Casualty Actuarial Society E-Forum, Fall 2016



The CAS E-Forum, Fall 2016

The Fall 2016 edition of the CAS *E-Forum* is a cooperative effort between the CAS *E-Forum* Committee and various other CAS committees, task forces, or working parties. This *E-Forum* contains the reports of the CAS Data & Technology and Bornhuetter-Ferguson Initial Expected Loss Ration Working Parties, and two independent research papers.

Data & Technology Working Party

Peter T. Bothwell, *Co-Chairperson* Mary Jo Kannon, *Co-Chairperson*

Benjamin Avanzi Joseph Marino Izzo Stephen A. Knobloch Raymond S. Nichols James L. Norris Ying Pan Dimitri Semenovich Tracy A. Spadola Linda M. Waite Dominique Howard Yarnell Cheri Widowski, *Staff Liaison*

Bornhuetter-Ferguson Initial Expected Loss Ratio Working Party

Nancy L. Arico Aaron Nicholas Hillebrandt Bertram A. Horowitz Lynne M. Bloom, Chairperson

Ziyi Jiao Douglas Robert Nation Michael J. Reynolds Xi Wu Karen Sonnet, *Staff Liaison*

CAS E-Forum, Fall 2016

Table of Contents

CAS Working Party Reports
CAS Data & Technology Working Party
Preface 1-4
Data Science and Analytics
Business Intelligence Technology and Tools: A Primer for Actuaries25-36
Data Quality Overview Actuarial Concepts in Data Quality
Databases
Bornhuetter-Ferguson Initial Expected Loss Ratio Working Party1-46
Independent Research
Escaping Hindsight: Case Reserve Development Using the Reserve Runoff Ratio
Joseph Boor, FCAS, PhD, CERA1-12
On Equality and Inequality in Stationary Populations David A. Swanson, Ph.D. and Lucky M. Tedrow, M.A1-16

E-Forum Committee

Dennis L. Lange, Chairperson Derek A. Jones, Chairperson-Elect Cara Blank Mei-Hsuan Chao Mark A. Florenz Mark M. Goldburd Karl Goring Donna Royston, Staff Liaison/Staff Editor Bryant Russell Shayan Sen Rial Simons Elizabeth A. Smith, Staff Liaison/Staff Editor John Sopkowicz Zongli Sun Betty-Jo Walke Qing Janet Wang Windrie Wong Yingjie Zhang

For information on submitting a paper to the *E-Forum*, visit <u>http://www.casact.org/pubs/forum/.</u>

PREFACE

The evolving definition of Advanced Analytics and the emergence of the Data Scientist

In its infancy, Actuarial Science operated at the leading edge of contemporary analytic capabilities and could be easily said to be employing "advanced analytics." Over the past 50 years, however, relentless data and technology breakthroughs have created modern analytic capabilities that far outstrip many of our traditional actuarial pricing and reserving methodologies. The role of "data scientist" has emerged as the holistic practitioner in advanced analytics. The Casualty Actuarial Society (CAS) has begun to address the need to update our methodologies with recent predictive modeling additions to the syllabus, but to function as data scientists, we still need additional data and technology capabilities as well.



Working effectively with Information Technology is key to advancing the goals of the Insurance Industry

Similarly, during this revolution of data and analytics capabilities, information technology (IT) departments and vendors have embraced "data & analytics," "big data," and "data science" as the new frontier for informed decision-making. In the P&C insurance industry, the CAS actuary is uniquely well-positioned to partner with IT to advance the potential of these disciplines to benefit the industry. In order to be a participant in the conversation, however, the actuary must have knowledge of the language, practices, tools and techniques of the technology supporting this revolution.

Data and Technology Narratives

Pursuant to three aims, namely:

- to introduce actuaries to concepts critical to the pursuit of data science;
- to encourage actuaries to take leadership/sponsorship roles in data governance;
- to familiarize actuaries with technical concepts important for working with IT professionals to evolve data-driven decision making in the insurance industry,

this collection of papers aims to address concepts that will inform the actuary on key terms and concepts underlying the data and technology disciplines. The ultimate goal of these papers is to identify the knowledge and skills actuaries must possess in order to participate in the changes brought about by rapidly evolving technology supporting data and analytics. The narratives are designed to provide brief descriptions of the key terms and concepts and then point to recommended publications that the reader should reference for a greater appreciation of the subject along with practical applications.

In order to apply structure and scope to the material, the key concepts were aggregated into four major categories: data science, business intelligence (BI), data quality, and databases. Although it is somewhat subjective what topics were assigned to which category, the categories align closely to the current usage of terms found within the P&C insurance industry.

The topics included in the major categories are outlined as follows:

"Data science" includes:

- A common definition of "data science"
- Other related sub-disciplines associated with the term "data science"
- Discussion on "big data"
- Mathematical modelling techniques
- Definitions of terms and recommended readings

"Business intelligence" includes:

- Business intelligence solutions supporting the actuarial process
- Description of current BI software tools and their application in an insurance company
- An actuary's role in the design and delivery of a BI project
- Definitions of terms and recommended readings

"Data quality" includes:

- A common definition of "data quality"
- Data management, governance and roles
- The use of metadata in various settings to control quality
- An actuary's role in the application of data quality best practices
- Definitions of terms and recommended readings

Preface

"Databases" includes:

- Comparison and contrast between a database and a data warehouse
- Actuarial considerations in the use of SQL and data tables
- High-level schematics and diagrams of data architectures
- Discussion on other structures
- Definitions of terms and recommended readings

The more you know, the more you know you don't know

Despite the different tone and structure of each paper, it is important to note that there is considerable overlap of perspective and terminology between the four topics. In fact, the interdependencies between these disciplines is what makes for compelling questions and unlimited opportunities for innovative solutions. For example, the reader should consider the following questions upon completing the readings:

- How do the databases that support data science differ from those that support actuarial process business intelligence deliverables?
- How does one build a business case for additional investment in data governance/management when competing with the forces that promote speed-to-market product development goals?
- What will the emergence of "self-service BI" mean for data warehousing strategy?
- How much of your actuarial analysis assumes your organization uses the same form of a reference data element (e.g., state code)? How will the lack of agreed-upon reference data format impact your database, data quality, business intelligence, and data science decisions?

With more formal education and research on data and technology topics, CAS actuaries will be better positioned to compete for data science roles and to partner with IT to use the combination of technology and analysis to develop innovative solutions to both long-standing and emerging challenges in the insurance industry and, likely, beyond. CAS Data and Technology Working Party Report

Data Science and Analytics

Tom Davenport published an article in the October, 2012 Harvard Business Review (HBR) titled "Data Scientist: The Sexiest Job of the 21st Century." He described the data scientist as "a hybrid of data hacker, analyst, communicator and trusted advisor." What happened? Wasn't Actuary the best job in America not long ago? So what is data science, what makes the data scientist different and why aren't actuaries ranked at the top anymore¹?

What is Data Science?

Let's start with a **definition** of data science. This is, unfortunately, not an easy task. There isn't an established professional or academic body to provide such a definition.

- Wikipedia defines it as: "an interdisciplinary field about processes and systems to extract knowledge or insights from data in various forms, either structured or unstructured, which is a continuation of some of the data analysis fields such as statistics, data mining, and predictive analytics, similar to Knowledge Discovery in Databases."
- TDWI asserts that data science "joins together contributions from several fields, including statistics, mathematics, operations research, computer science, data mining, machine learning (algorithms that can learn from data), software programming, and data visualization. It can cover the entire process of acquiring and cleaning data, methods for exploring the data and extracting value from it, and techniques for making insights actionable for humans and automated processes."
- Drew Conway provided² a popular view of the skills needed to be a data scientist using a Venn diagram, shown in the left panel of the figure below. For the purposes of this paper, we have created our own Venn diagram with labels that may be more familiar to the actuarial community, shown in the right panel of the following figure.

¹ In a recent careers survey, data scientist ranked first while actuary was at number ten: http://www.careercast.com/jobs-rated/jobs-rated-report-2016-ranking-200-jobs

² http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram



Conway posited that like Actuarial Science, Data Science is an empirical science. However, the difference between a traditional actuary and a data scientist is the addition of what Conway called "hacking skills," namely "being able to manipulate text files at the command-line, understanding vectorized operations, thinking algorithmically" More generally in this context, "hacking" can refer to data acquisition and transformation at scale together with coding expertise required to implement production ready prototypes of the mathematical models.

In popular use "hacking" carries pejorative connotations but its intent here is to indicate a certain degree of fluency in dealing with programmable systems³.

While "data science" has initially emerged as a label for analytics at web companies (Facebook and LinkedIn specifically), it is a reflection of deeper intellectual currents. A compelling account of the **history of data science** was recently given by a prominent academic statistician David Donoho [38], considering it in the context of the broader evolution of the practice of data analysis.

RELATION BETWEEN DATA SCIENCE AND OTHER ANALYTICS DISCIPLINES

It may at times seem difficult to differentiate between "data science" and more established fields of analytics. We believe that this is due to an increasing number of industries being affected by "digital transformation" – new business models facilitated by the ubiquitous availability of computing and telecommunication technologies. This "digital transformation" is to a large degree carried out by engineering / software-centric "web" companies following technology practices that have little overlap with traditional enterprise IT and adopting "data science" rather than traditional analytics.

³ See the definition of "hacker" in the Jargon File: http://www.catb.org/jargon/html/H/hacker.html

Due to the widening front of "digital transformation," the scope of "data science" is also expanding. We are already quite close to it becoming an umbrella term for what historically have been largely disparate areas of applied mathematical modelling in the commercial setting. Some of these are listed below.

Online advertising and website optimization: Online advertising has grown into a massive ecosystem over the last two decades providing critical revenue for the majority of online services. The nature of the medium is eminently accommodating of tracking and analytics, resulting in one of the more dramatic applications of "data science." Most sophisticated solutions (e.g. AdWords) are deployed by inventory providers and aggregators, such as Google and Facebook. Live A/B testing is also prevalent among online businesses – something which is still a rarity in traditional enterprise. This is at present the biggest area of employment for "data scientists."

Manufacturing quality control, statistical process control, lean manufacturing, Six Sigma – this is an area of analytics supporting manufacturing activities and has been progressively developed since at least the 1930s. Among the main objectives is monitoring and elimination of variability in manufacturing processes (e.g., part dimensions), ensuring that defect rates are thereby controlled.

Operations Research, industrial engineering, revenue management, mathematical optimization, management science. Operations research began as a scientific study of military operations (e.g., convoy composition, bomber interception protocols, and logistics) during the Second World War and the principles have been exported to many other industries in the following years, in particular manufacturing, travel and transportation. Main tools include mathematical optimization and stochastic processes.

Statistics really requires no introduction; perhaps its main focus of interest in applications has been analysis of government data, polls and surveys and support for design of experiments and evaluation of experimental results in life sciences and medicine.

Applied finance, financial engineering, algorithmic trading, HFT, portfolio management. There are close parallels between data science and quantitative finance in the 1980s and 1990s. This is not surprising, because in-market execution is a key part of any model driven trading strategy, placing a premium on "hacking skills." At present it is perhaps reasonable to view the majority of "data scientists" as "quants" of digital advertising.

Engineering control, control theory, signal processing. Successful engineering applications of control and information theories span from fly-by-wire systems to cellular networks and synthetic aperture radar. Much less ambitious in scope than AI, these systems work reliably and are by now absolutely ubiquitous.

Econometrics, mechanism design, causal inference (from observational data) – due to the difficulty and costs of real world experiments in economics, econometricians have developed tools and conceptual frameworks for causal inference with observational data [39]. Furthermore mechanism design and the study of auctions have had significant impact on the design of online marketplaces.

Business intelligence, database / warehouse design, dashboards – business intelligence is primarily an IT led activity to support descriptive and diagnostic analytics. Business Intelligence will be given its own discussion in another section of this paper.

Machine learning, natural language processing, computer vision, data mining – machine learning is a branch of computer science that initially focused on more tractable aspects of artificial intelligence, primarily by constructing models from example data using statistical methods rather than designing them by hand from general principles. Two large application areas are computer vision and natural language processing, including machine translation. By now the differences between theoretical machine learning and statistics communities are largely superficial, amounting to little more than preferences for different styles of analysis of statistical procedures. Machine learning research has also provided many of the tools used in analytics for online advertising and algorithmic trading. Data mining has originated from the databases research community and has also mostly converged with machine learning in terms of both objectives and methodologies. Notably, there is a significant community of machine learning researchers working at technology companies who self-identify as such rather than "data scientists."

DATA SCIENCE AND "BIG DATA"

Data science has come to be associated with so-called "big data" – in this section we argue that "big data" projects that some insurance companies have undertaken are only tangentially related to the success of data science and instead the key lessons that insurers can derive from the experience of web companies lie in an integrated approach to product management and design and the adoption of live market testing.

"Big data" projects

The focus of "big data" projects in insurance and consumer finance to date has largely been on data processing infrastructure – information from production systems (web servers, policy and claims management, finance systems etc.) is transferred in raw form into so called "data lakes" with the goal of subsequent "insight discovery."

In this sense much of the technology is a direct successor of the earlier generation of "business intelligence" (BI) or "data warehousing" solutions, with the key difference being the abandonment of

the fixed predetermined database schemas. Traditional BI architecture presupposed certain formats and relations to which all data was compiled, striving to present a "single source of truth" in one materialized data set.

Current big data tools (Hadoop, Spark etc.) replace this approach with computation; data views are not predesigned but are an output of a program run over the entire history of source system extracts. This approach is enabled by utilizing clusters of relatively cheap commodity server hardware⁴ and ideally ensures that no information is lost due to imposition of a schema and it is always possible to answer any (unanticipated) query addressable by historical data. This dramatically reduces both the upfront costs of data "ingestion" and transformation as well as making sure that the resulting system is potentially useful to many stakeholders in the organization, even those who have not been the key focus in its design. (For example, general purpose data warehouses developed internally by insurers often turn out to have limitations that still require actuarial teams to operate their own independent processes to meet pricing and valuation needs.)

This approach has the potential to dramatically simplify many of the reconciliation, reporting and model building activities, as all of the enterprise data can be collected on a single "computational substrate."

Another commonly cited benefit is the ability to construct a unified view of individual customer interactions with the company, records of which may be split across multiple systems. The data can then be used to both improve risk models and in some cases derive insights around other aspects of customer behavior. This is one area where diversified market participants, providing consumer services outside insurance, are at a clear advantage relative to traditional carriers.

Virtually all "big data" technologies originated from the need to support analytics at web companies. However, it is important to note that these are purely enabling technologies and are not essential for data science itself except in situations where associated data processing tasks cannot be accomplished by other means. It is perhaps these common origins that have created an association between "big data" and "data science."

Analytics solutions at web scale usually need to address challenges around the so called "four Vs" of data (nomenclature predominantly adopted by IT vendors):

- velocity: data is gathered at an increasing speed;
- variety: data is gathered in a large number of forms and ways;
- volume: exponentially increasing volumes of data are being gathered;
- veracity: it becomes increasingly difficult to guarantee the quality of the data;

⁴ Cost reductions of two orders of magnitude per terabyte relative to vendor BI solutions are sometimes claimed.

as well as satisfying certain operational properties:

- Automation: at scale, it becomes impossible to manually curate or even review models supporting operational decision. The entire pipeline needs to be fully automated. This is a challenging task if one wants results to be robust and credible.
- Speed of computation and algorithmic complexity of procedures involved come to the forefront with large volumes of data.
- Adaptability: special care needs to be taken in the design to allow adjustments to the analytics pipeline stemming from frequent changes to the front end systems.

It is important to remember that for most insurance companies with traditional product portfolios, issues relating to the scale of data are simply not present and data science solutions can be reasonably implemented on existing infrastructure, i.e., the link between "big data" and "data science" is very weak if it exists at all.

Limitations of "big data"

"Big data" technology is only a part of the solution to analytics-guided operational decision making - the standard operating practice of web companies. In this section we discuss another essential ingredient: live in -market testing.

Consider the typical online quoting process for a personal motor policy - the only interaction the customer has with the insurance company in this case consists of being presented a sequence of web forms. Who within the company is responsible for the overall customer experience? For some insurance companies the responsibilities may be separated as follows:

- A product team is responsible for policy options and associated wordings in the online world this translates into available check boxes and sliders on the quote screen.
- Pricing function is responsible for the actual quote amount displayed for a particular product configuration.
- Design, form layout and flow may be handled by a dedicated "channel" team.
- Banners or cross sell offers may be managed by the marketing function.
- Search engine campaigns directing traffic to the website are outsourced to a media agency.
- Underwriting may have input into what information is collected as well as business rules for generating referrals for manual processing.
- Finally, IT function would be responsible for the integration of the web front end and the "core" policy administration system.

While this structure is readily understood in historical context, it is often unclear who is ultimately accountable for the customer experience and any substantial change typically involves interdepartmental coordination which can further complicate or delay the process. In a modern web company, all of these responsibilities would be handled by a single "product" team, where "product" is not a particular policy wording but rather the software artifact that generates the customer experience with product options, wordings and prices all integral parts of the whole.

Traditional organizational structure is a major obstacle faced by established insurers seeking to adopt a "data science" approach to product management as it generally hinders rapid in-market testing of different variations of customer experience. "Big data" solutions are of little help in this environment as data captured from various systems will be by its nature observational - generated in the course of normal business operations - but only limited insight can be systemically obtained from observational data. For example, it is generally straightforward to estimate risk premium for a new cohort of business based on the history for a comparable book, but much more difficult to answer more pertinent questions around the impact of a proposed rate change on the expected business volume. The latter requires a model of demand elasticity, which is not identified⁵ without active intervention or external shocks (i.e. known changes in competitors' prices). The same applies to many business questions around the product offering and marketing strategies - few of them can be answered with any degree of credibility by analytics on historical data alone, ultimately requiring inmarket testing. We will revisit this in later sections.

Data quality considerations

Actuaries are quite familiar with data quality considerations when it comes to rate filing or reserving exercises and traditional data quality principles are discussed in detail in another section of this paper.

Different criteria, however, will apply when devising rules for operational decision making, e.g., choosing a particular version of an online quote form. While for a valuation missing data for 10% of policies would clearly be unacceptable, data missing (at random) for 10% of customers would have no significant effect on performance.

In the big data space, suitable determinations have to be made for each individual use case and it is at times necessary to significantly relax standards actuaries might be accustomed to. One issue specifically worth mentioning is the situation when the dataset used for analysis contains information not available at the time the decision needs to be made (e.g. due to a pre-processing step incorrectly incorporating knowledge of future transactions). This type of error has the potential to undetectably undermine the model validation protocol described in the following sections – underperformance only revealed once the model has been deployed in production.

OBJECTIVES OF DATA SCIENCE

At least one of the goals of data science is to bring rigor to optimizing operational decision making through integration of analytical and technological expertise. Additionally it seeks to incorporate rich new data sources such as text, audio, images and video into both analysis and decision making – this

⁵ http://en.wikipedia.org/wiki/Parameter_identification_problem

CAS Data and Technology Working Party Report

latter objective is made possible by advances both in the costs of hardware and progress made over the last two decades on the associated pattern recognition tasks (e.g. [41]).

Some questions around operational decisions can be answered effectively by constructing models (see next section) of observational data gathered in the course of normal business operation, sometimes referred to as "predictive analytics."

Situations where "predictive analytics" are directly applicable are not universal – insurance premium rating happens to be one such case, fraud detection in settings where notifications are reliably received from injured parties is another.

A great discussion of the limitations of "predictive analytics" approach in the context of evaluating effectiveness of advertising is given in [42]. The paper shows that the critically important questions of causality (in their case sales uplift from a particular advertising campaign) cannot be answered reliably from observational data alone without randomized intervention to fully remove confounding. Randomized interventions on the production systems are in many cases the only known way to reliably estimate and therefore optimize the effects of operational decisions. Sometimes this is called "prescriptive analytics" although the term is often also used in engineering applications where system dynamics can be reliably estimated from general theory and do not require ad hoc experimentation.

An example of "prescriptive analytics" in the insurance context would be so-called premium optimization. Demand based premium adjustments, however, is just one form of intervention and the exact same framework can be applied to evaluating the color and position in which the quoted premium is displayed vs. any loading applied to the amount itself.

Indeed it is in the design, execution and analysis of live market tests of this type that technical expertise of a data scientist is often crucial, quoting R. A. Fisher:

"To consult the statistician after an experiment is finished is often merely to ask him to conduct a post mortem examination. He can perhaps say what the experiment died of."

In practice this means that some formal metric needs to be defined that can be estimated in a relatively short period of time - it could be conversions, click through rates, retention, net promoter scoring and so on – as well as the size of the test or an appropriate stopping rule⁶.

Design of large scale sequential experiments and analysis of resulting data is an active area of research in the machine learning community, with [43] offering the most accessible introduction to date. Pervasive testing is likely to prove the key analytics lesson to be adopted from the consumer tech

⁶ https://en.wikipedia.org/wiki/Sequential_analysis

Data Science and Analytics

companies - Google, for instance, runs hundreds⁷ of parallel experiments on its search product alone, as does Microsoft⁸.

Business analytics is also sometimes said to follow a maturity model. While there are many sources with slight variations around the same theme, we have included the Gartner model in this paper.



Here we also see the progression from "predictive" to "prescriptive" analytics. While actuarial analysis is generally "predictive" there is considerable room for advancement when it comes to model validation, and live testing in insurance remains exceedingly rare.

METHODOLOGY: SOME SPECIFIC TECHNIQUES

Mathematical modeling techniques

There is dramatic variation in mathematical tools used in different areas of "analytics." It is not infrequently observed that different groups of practitioners develop solutions that are operationally very similar, while diverging significantly in ideology and mathematical apparatus. Despite significant overlaps it can still be useful to consider the major approaches as they form major components of respective intellectual traditions.

Summary tables are pervasive in business reporting, and while this aspect is usually ignored, one must make implicit assumptions about the underlying data generating process in order to make inferences from such information.

⁷ http://research.google.com/pubs/pub36500.html

⁸ http://www.exp-platform.com/Documents/2014\%20experimentersRulesOfThumb.pdf

Basic probability and statistics are quite familiar to actuaries, perhaps the biggest gap being hypothesis testing, should experimental methodology become more widely adopted in the insurance industry. In particular, the general confusion between Fisherian, Neyman-Pearson and Bayesian points of view⁹ in introductory texts make it difficult to acquire fluency in the correct application of standard methods.

Parametric conditional models – these include traditional regression tools, like generalized linear models, quantile regression etc. Multiple useful extensions have been developed, including regularization, random effects and additive models, all of which are closely linked to actuarial credibility. It is also sometimes possible to "predict" more complex objects than just a single dependent variable, such as part of speech labelling for an entire sentence (e.g. so called conditional random fields).

Dynamical systems – dynamical system models seek to incorporate the temporal aspects of the phenomenon or control process under consideration. These are particularly effective when the evolution of the system can be somewhat reliably predicted (e.g. on the basis of physical laws). Time series models are without exception special cases of dynamical systems.

Mathematical optimization – many inference problems, from testing to regression and beyond rely on solving optimization problems (e.g. maximum likelihood). Mathematical optimization studies both properties of such problems as well as computational procedures that can be used to find or approximate solutions.

"Model free control" – in many situations it is not possible to construct a reasonable model of the system to be optimized from general principles; this includes the majority of applications of data science in analysis and control of live experiments. Such settings necessitate joint estimation and control. For example, in online advertising, before a click through rate for a new ad can be estimated it has to be displayed a certain number of times in different contexts. Investigation of methods for doing this efficiently while simultaneously optimizing for an overall objective, such as revenue, is a central problem in "reinforcement learning," a subfield of machine learning.

Bayesian modelling – it is possible to consider most of the above settings from the Bayesian point of view. Notoriously difficult computationally, algorithmic advances (geometric MCMC, variational methods) and the availability of open source software make this approach tractable for a growing range of practical problems. Some aspects of Bayesian analysis are known to actuaries as credibility theory.

⁹ R. Christensen, Testing Fisher, Neyman Pearson and Bayes, <u>http://www.stat.ualberta.ca/~wiens/stat665/TAS%20-</u> %20testing.pdf

"Non-parametric" conditional models – gradient boosting machines, support vector machines, much of "deep learning" or neural networks – these types of methods are most commonly associated with machine learning or data science. They extend the usual regression models by introducing non-linear dependence of output on the input variables while still maintaining the ability to control overall model complexity.

With the proliferation of specific analytical methods and implementations, it is important for actuaries to be able to place specific methods into a theoretical framework to evaluate their relative merits and specific applicability. For example it would be valuable to understand connections between actuarial credibility and penalized regression and ensemble methods developed in machine learning and computational statistics literatures [46].

There exists a multitude of such frameworks at various levels of abstraction. One particularly useful viewpoint is that of optimization [19] – it is generally very difficult to understand whether two statistical procedures are related, especially if they are presented in the form of algorithms. Understanding what objective function is minimized or maximized by a given procedure allows us to readily appreciate similarities between methods. As an example, it turns out that the maximum likelihood estimator for logistic regression is almost identical to the optimization problem solved by "support vector machines" popular in machine learning. Many examples of optimization models in premium rating are given in [44].

Algorithmic thinking

Increasing volumes of data brings to the forefront computational issues around mathematical modelling and data processing. Beyond certain problem sizes, algorithms with second or higher degree polynomial complexity simply stop working (i.e., they do not terminate in any reasonable time). Actuaries must be aware of this possibility and some common workarounds where they exist.

Finally, we should point out that all popular environments for cluster computing (Hadoop, Spark, etc.) impose significant limitations on the user in terms of how the computation needs to be structured relative to using a single computer. Understanding these limitations and how they can impact common tasks require both familiarity with distributed system architectures as well as the underlying algorithms.

Visualization and exploratory data analysis

Sanity checks on the available data and trying to understand how recorded observations relate to the generating process are the core activity in "data science" and indeed among actuaries. John Tuckey has referred to this "Exploratory Data Analysis" in his influential book [45]. This type of investigation is made particularly important when working with heterogeneous data originating from rapidly evolving systems.

Advances in theory [47] and computer software (e.g. ggplot2) have made advanced visualization [48] readily available.

Model validation

Data driven model validation has emerged as one of the central themes in "data science." These methods have seen relatively limited use in ratemaking to date due to the labor-intensive nature of model construction (model validation requires fitting multiple models to different subsets of data).

Basic model validation involves dividing the available modelling data set into three disjoint subsets:

- 1. Training this is a (random) subset of the data used to construct successive iterations of the model.
- 2. Test this (random) subset of the data is not used in model fitting but only to evaluate model performance. If a model iteration performs well (relative to all other model iterations as well as some known baseline) on the test set according to some formal metric, such as AUC or RMSE, that iteration is declared the winner. It is usual to consider dozens if not hundreds of iterations in a course of a modelling project¹⁰.
- 3. Validation this set of the data is withheld for a final validation of the model that passes the testing process. Often most reliable validation is, in fact, not a random subset of the data at all, but an "out of time" dataset that more accurately approximates live deployment.

This approach can then be naturally extended to multiple participants and has enabled steady progress in a number of applications of machine learning, particularly computer vision and natural language processing, with the more general framework as follows:

- 1. A dataset is made available (perhaps publicly) containing for each observation a value to be predicted (these can be numeric, categorical, or more complex structures altogether).
- 2. An objective function which the prediction rule or model is to optimize is communicated to the participants.
- 3. A referee who is able to evaluate models on a separate dataset whose objective values are not visible to the participants and report back the scores.

The goal of the participants is the construction of model which minimizes deviations from the objective values as reported by the referee. Beyond academic research, this is the mechanism that is used by Kaggle, a company that provides "crowd sourced" modelling solutions to companies willing to share their data publicly.

In the view of the authors, the key to success in applying "hard" data science to business problems is the creation of appropriate evaluation frameworks that can rigorously evaluate the quality of decision rules – sometimes historical observational data alone is sufficient (e.g., for models of claim costs) and sometimes live market testing may be required.

¹⁰ To get more accurate estimate of out of sample performance when limited data is available, it is common to repeat the process over multiple training/test splits, e.g. so called "cross-validation."

DATA SCIENCE RESOURCES

In what follows we list some particularly noteworthy graduate and undergraduate courses that could help develop a broad fundamental understanding of computing and mathematical modelling. These could be argued to be core "data science" skills for addressing future business problems, with increasing number of processes and low-level operational decisions subject to automation. In compiling these resources we have intentionally stayed away from "flavor of the month" or introductory offerings, focusing instead on fundamentals.

Analytics at web companies

To get an impression of what the future of insurance analytics might look like, it is worthwhile to review some of the courses offered by people with experience implementing analytics solutions for the leading web companies. Examples include CS281B "Scalable Machine Learning" at UC Berkeley [25] by Alex Smola (formerly of Yahoo) and "Big Data, Large Scale Machine Learning" at NYU [26] by Yan LeCunn (currently at Facebook). In particular the first course offers an interesting insight into the importance of understanding systems, numerical methods and statistics to develop analytics solutions at web scale.

Prerequisites for this material include linear algebra, basic probability and statistics and, ideally, convex optimization and an introduction to machine learning, as discussed next.

Mathematical background and numerical computing

Numerical linear algebra is the most essential tool in applied mathematics. The majority of computational procedures for solving mathematical models ultimately reduce to iteratively solving systems of linear equations.

An excellent introductory treatment of linear algebra is given by Gilbert Strang in MIT 18.06 [2]. The material is further developed in MIT 18.085 [3] and 18.086 [4], demonstrating a very broad range of applications across engineering subfields. The observation that the differential operator can be discretized as e.g. a tri-diagonal matrix (the so called "finite differences" method) is the key connection between linear algebra, traditional calculus (in the form of integral and differential equations) and computing.

Another take on the material is given in Stanford EE263 taught by Stephen Boyd - in addition to basic linear algebra, the course gives a highly intuitive exposition to least squares regression, regularization, singular value decomposition and linear dynamical systems (which can be viewed as a generalization of a wide class of time-series models in the CAS syllabus). The material above should provide sufficient background to appreciate some of the technology behind modern robotics platforms, such as those formerly developed at Boston Dynamics, now part of Google (MIT 6.832

Underactuated Robotics [12]).

Finally, the Fourier transform is one of the most famous special cases of a linear operation– an intuitive introduction to the subject and its multitude of applications, including the Central Limit Theorem, is given in Stanford EE261 [17].

Optimization

Beyond differential equations, one of the main applications of linear algebra is in mathematical optimization or "mathematical programming." Optimization based models are pervasive in analytics, whether it be maximum likelihood estimation, "empirical risk minimization," Neyman-Pearson hypothesis testing, optimal control, Markowitz portfolio theory or option pricing.

Prof. Stephen Boyd's course EE364A Convex Optimization [19,20] not only gives a solid grounding in the theory but also considers many of the above-mentioned examples. Convex optimization is widely seen as the foundation of modern statistics, machine learning and signal processing. Familiarity with theory and algorithms will enable the practitioner to identify and implement solutions to a very wide range of problems across industries.

There is also an interesting connection between mathematical optimization and classical algorithms studied in undergraduate computer science courses (e.g. [7]) - many of the problems such as sorting, shortest path, max flow, etc. turn out to be special cases of linear programming (itself a special case of convex optimization).

The follow up course EE364B [21] provides more detailed background on scalable and distributed optimization as well as the clearest introduction to the General Equilibrium theory of microeconomics you are likely to find. The background for these courses is limited to linear algebra [2,18] and basics of multivariable calculus (gradient, Hessian) [1].

Probability, statistics, machine learning, information theory.

There are few unequivocally great introductory probability and statistics courses publicly available, at least at the moment. MIT 6.041 [9] is a useful probability refresher. A worthwhile follow up is MIT 6.262 [10] "Discrete Stochastic Processes."

When it comes to statistics, or at least a take on the topic that is more attuned to analytics applications, Stanford Statistical Learning [21] is a solid introduction from the authors of the well-known book. A closely related subject area is machine learning, with the introductory course by Andrew Ng [23] and a much more in depth treatment by Alex Smola [24]. So called "deep networks" are a recent "hot" topic in machine learning, providing state of the art performance for many recognition tasks. This material is covered in [27].

Data Science and Analytics

Information theory provides perhaps one of the most successful and widely used applications of probability. There are also important connections to statistics and machine learning (as efficient compression requires effective conditional probability estimation). MIT 6.450 "Principles of Digital Communications I" [11] is an excellent course by the pioneer of digital communications Rob Gallager, who invented one of the most effective known coding schemes and was a founding engineer at Qualcomm where he designed the first 9600 baud modern. Information theory is an essential foundation of all digital information processing technology.

Another excellent discussion of information theory is given in the course taught by David MacKay at Cambridge [36], bringing together topics from coding theory, statistics and machine learning.

Convex optimization provides a very helpful background for the courses in this section even if it is not explicitly alluded to.

Programming

There exists a very wide range of high quality introductory programming courses. Perhaps the Stanford sequence deserves a particular mention [15,16]. Alternatives include the introductory courses at MIT [5,8].

MIT 6.001 [6] (now superseded) is the most celebrated introductory programming course of all, with the textbook "Structure and Interpretation of Computer Programs" used in dozens of top universities. While Scheme, the language that it uses for teaching programming concepts, has for long time been considered less than practical, over the recent years there has been a dramatic resurgence of popularity of the related body of ideas called "functional programming," underpinning many of the latest "big data" technologies.

Beyond the introductory courses, "Programming Paradigms" [17] gives a useful overview of design choices behind a variety of programming languages and [31] offered on Coursera by the University of Washington, provides a more advanced grounding in the functional programming paradigm.

An introduction to Scala, an increasingly popular compatible replacement of Java, is available from its creator on Coursera [32].

No such list would be complete without an algorithms class [7]. Conceptual links with optimization or "mathematical programming" offer a connection back to the material in the earlier sections.

Finance, economics and social science

While the exact relation between actuarial pricing and financial economics is not clearly set out in the actuarial curriculum, it has been understood in the academic literature for some time as the so called "incomplete markets" setting. An introductory discussion of the modern theory of finance (CAPM, option pricing, etc.) from this more advanced point of view is given in John Cochrane's (University of Chicago) class "Asset Pricing" on Coursera [30].

A useful generalization of the concept of an optimization problem (see e.g. Stanford EE364A [19,20]) is offered by game theory. Instead of considering a "central planning" problem where all the decisions are taken by a single agent, game theory looks at situations where there are multiple self-interested parties involved. Coursera classes [28,29] provide an introduction to a range of topics, including auctions and mechanism design. Applications of game theoretic methods to the study of social insurance, optimal taxation and related ideas are given in the Harvard course "Public Economics" [34].

Problems addressed by "business analytics" are not dissimilar to those found in the social sciences, especially when it comes to identifying what is sometimes called "actionable insights" - a social scientist may instead talk about "policy targets." While causal attribution is oftentimes not necessary, it is important to be aware of limitations of analyses carried out purely on observational data. One example in social science where large-scale experiments have been possible is "development economics." The MIT course 14.73 [35] offers an in depth discussion of considerations that go into designing a convincing experimental study. A broad introduction to the design of quantitative methods that are directly applicable to the question being studied is given in Gary King's excellent methodology course at Harvard [33].



REFERENCES

- [1.] MIT 18.02 Multivariate Calculus http://ocw.mit.edu/courses/mathematics/18-02sc-multivariable-calculus-fall-2010/
- [2.] MIT 18.06 Linear Algebra http://ocw.mit.edu/courses/mathematics/18-06-linear-algebra-spring-2010/
- [3.] MIT 18.085 Computational Science and Engineering I http://ocw.mit.edu/courses/mathematics/18-085-computational-science-and-engineering-i-fall-2008/
- [4.] MIT 18.086 Mathematical Methods for Engineers II http://ocw.mit.edu/courses/mathematics/18-086-mathematical-methods-for-engineers-ii-spring-2006/
- [5.] MIT 6.00 Introduction to Computer Science and Programming <u>http://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-00sc-introduction-to-computer-science-and-programming-spring-2011/</u>
- [6.] MIT 6.001 Structure and Interpretation of Computer Programs <u>http://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-001-structure-and-interpretation-of-computer-programs-spring-2005/</u>
- [7.] MIT 6.06 Introduction to Algorithms <u>http://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-006-introduction-to-algorithms-fall-2011/</u>

CAS Data and Technology Working Party Report

- [8.] MIT 6.01 Introduction to Electrical Engineering and Computer Science <u>http://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-01sc-introduction-to-electrical-engineering-and-computer-science-i-spring-2011/</u>
- [9.] MIT 6.041 Probabilistic Systems Analysis and Applied Probability <u>http://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-041sc-probabilistic-systems-analysis-and-applied-probability-fall-2013/</u>
- [10.] MIT 6.262 Discrete Stochastic Processes <u>http://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-262-discrete-stochastic-processes-spring-2011/</u>
- [11.] MIT 6.450 Principles of Digital Communications I <u>http://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-450-principles-of-digital-communications-i-fall-2006/</u>
- [12.] MIT 6.832 Underactuated Robotics <u>http://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-832-underactuated-robotics-spring-2009/</u>
- [13.] UNSW COMP1917 Higher Computing http://www.youtube.com/playlist?list=PL6B940F08B9773B9F
- [14.] Stanford CS106A Programming Methodology http://see.stanford.edu/see/courseinfo.aspx?coll=824a47e1-135f-4508-a5aa-866adcae1111
- [15.]Stanford CS106B Programming Abstractions http://see.stanford.edu/see/courseinfo.aspx?coll=11f4f422-5670-4b4c-889c-008262e09e4e
- [16.] Stanford CS107 Programming Paradigms http://see.stanford.edu/see/courseinfo.aspx?coll=2d712634-2bf1-4b55-9a3a-ca9d470755ee
- [17.] Stanford EE261 Fourier Transform and its Applications <u>http://see.stanford.edu/see/courseinfo.aspx?coll=84d174c2-d74f-493d-92ae-c3f45c0ee091</u>
- [18.]Stanford EE263 Introduction to Linear Dynamical Systems <u>http://see.stanford.edu/see/courseinfo.aspx?coll=17005383-19c6-49ed-9497-2ba8bfcfe5f6</u>
- [19.] Stanford EE364A Convex Optimization http://see.stanford.edu/see/courseinfo.aspx?coll=2db7ced4-39d1-4fdb-90e8-364129597c87
- [20.] Stanford CVX101 Convex Optimization https://class.stanford.edu/courses/Engineering/CVX101/Winter2014/about
- [21.] Stanford Statistical Learning https://class.stanford.edu/courses/HumanitiesScience/StatLearning/Winter2014/about
- [22.] Stanford EE364B Convex Optimization II http://see.stanford.edu/see/courseinfo.aspx?coll=523bbab2-dcc1-4b5a-b78f-4c9dc8c7cf7a
- [23.] Stanford CS229 Machine Learning http://see.stanford.edu/see/courseinfo.aspx?coll=348ca38a-3a6d-4052-937d-cb017338d7b1
- [24.] CMU 10-701 Introduction to Machine Learning http://alex.smola.org/teaching/cmu2013-10-701/
- [25.] UC Berkeley CS281B Scalable Machine Learning Slides - <u>http://alex.smola.org/teaching/berkeley2012/</u> Videos - <u>http://www.youtube.com/playlist?list=PLOxR6w3fIHWzljtDh7jKSx_cuSxEOCayP</u>
- [26.]NYU Big Data, Large Scale Machine Learning http://cilvr.cs.nyu.edu/doku.php?id=courses:bigdata:start
- [27.]NYU Deep Learning http://cilvr.cs.nyu.edu/doku.php?id=courses:deeplearning:start

- [28.] Coursera Stanford/UBC Game Theory https://www.coursera.org/course/gametheory
- [29.] Coursera Stanford/UBC Game Theory II: Advanced Applications <u>https://www.coursera.org/course/gametheory2</u>
- [30.] Coursera University of Chicago Asset Pricing http://www.coursera.org/course/assetpricing
- [31.] Coursera University of Washington Programming Languages https://www.coursera.org/course/proglang
- [32.] Coursera EPFL -Principles of Functional Programming in Scala <u>https://www.coursera.org/course/progfun</u>
- [33.] Harvard Gov 2001 Quantitative Research Methodology http://projects.iq.harvard.edu/gov2001/home
- [34.] Harvard Econ 2450a http://obs.rc.fas.harvard.edu/chetty/public_lecs.html
- [35.] MIT 14.73 The Challenge of World Poverty http://ocw.mit.edu/courses/economics/14-73-the-challenge-of-world-poverty-spring-2011/
- [36.] Cambridge Information Theory, Pattern Recognition and Neural Networks http://www.inference.phy.cam.ac.uk/itprnn/Videos.shtml
- [37.] https://www.soa.org/library/research/transactions-of-society-of-actuaries/1990-95/1995/january/tsa95v475.pdf
- [38.]D. Donoho, 50 Years of Data Science http://courses.csail.mit.edu/18.337/2015/docs/50YearsDataScience.pdf
- [39.] J. Agrist, J. Pischke, Mostly Harmless Econometrics: an empiricist's guide, PUP, 2009
- [40.] Gill Press, Forbes Contributor, A Very Short History of Data Science, 5/28/2013
- [41.] Stanford CS231n Convolutional Neural Networks for Visual Recognition https://www.youtube.com/watch?v=NfnWJUyUJYU&feature=youtu.be
- [42.]B. Gordon et al., A Comparison of Approaches to Advertising Measurement: Evidence from Big Field Experiments at Facebook

http://www.kellogg.northwestern.edu/faculty/gordon_b/files/kellogg_fb_whitepaper.pdf

- [43.] Agrawal et al., Multi-World Testing: A System for Experimentation, Learning and Decision-Making <u>http://research.microsoft.com/en-US/projects/mwt/mwt-intro.pdf</u>
- [44.]D. Semenovich, Applications of Convex Optimization in premium rating, CAS E-Forum 2013 <u>https://www.casact.org/pubs/forum/13spforum/Semenovich.pdf</u>
- [45.] J. Tukey, Exploratory Data Analysis, 1977
- [46.] H. Miller and P. Mulquiney. Credibility, penalized regression and boosting; let's call the whole thing off, 2011 <u>https://www.casact.org/education/infocus/2011/handouts/AM2-Fry.pdf</u>
- [47.] L. Wilkinson, The Grammar of Graphics, 1999
- [48.] E. Tufte, The Visual Display of Quantitative Information, 2001

CAS Data and Technology Working Party Report

Business Intelligence Technology and Tools: A Primer for Actuaries

BUSINESS INTELLIGENCE DEFINED

While the term "Business Intelligence" is widely used, it doesn't have a single, widely accepted definition. However, several authoritative sources have each published definitions close enough conceptually to provide a good starting point for any discussion of the subject.

"Business Intelligence" has been defined in the following ways:

"...a set of concepts and methodologies to improve decision making in business through the use of facts and fact-based systems"¹

"BI is neither a product nor a system. It is an architecture and a collection of integrated operational, as well as decision-support, applications and databases that provide the business community easy access to business data."²

"Business intelligence encompasses data warehousing, business analytic tools and content knowledge management." 3

"...the ability to transform data into useable, actionable information for business purposes. BI requires:

- Collections of quality data and metadata important to the business
- The application of analytic tools, techniques and processes
- The knowledge and skills to use business analysis to identify/create business information
- The organizational skills and motivation to develop a BI program and apply the results back to the business"⁴

"... an umbrella term that combines architectures, tools, databases, analytical tools, applications, and methodologies that allows actuaries: (1) to have interactive access (sometimes in real time) to data, (2) to manipulate data, and (3) to conduct appropriate analyses. The process of BI is based on the transformation of data to information, then to decisions, and finally to actions."5

INTRODUCTION

The purpose of this narrative is to introduce our actuarial audience to business intelligence terms and tools. The technologies involved are not all new, but are evolving as insurers use more and more computations and analytics to drive their business.

¹ Howard Dresner, Gartner Group

² Larissa T. Moss and Shaku Arte, Business Intelligence Roadmap, Pearson Edition, 2003

³ David Loshin, Business Intelligence: The Savvy Manager's Guide, Addison Wesley, 2003

⁴ TDWI Course titled, "Business Intelligence Fundamentals..."

⁵ Sharda, Ramesh; Delen, Dursun; Turban, Efraim; Aronson, Janine; Liang, Ting Peng (2014-01-14). Business Intelligence and Analytics: Systems for Decision Support (Page 14). Pearson Education. Kindle Edition

Business Intelligence Technology and Tools: A Primer for Actuaries

Casualty actuaries have been concerned with data used in decision making since the beginning of the Casualty Actuarial Society (originally the Casualty Actuarial and Statistical Society of America) or CAS. Actuaries used computers as personal support tools to complete their work, but today business intelligence systems are coupled with data mining and automated processes to greatly expand an insurer's understanding of market trends.

When looking to improve their processes with increased use of the business intelligence landscape, the actuary needs to have a good working relationship with the IT⁶ department who often maintains the data and tools used for the actuarial work. The nature of the relationship relies on the actuary having a clear idea of the business goal. Together they will determine what data and tools the actuaries will provide and what data and tools the IT department will provide as inputs to the actuarial work. The relationship and how they communicate together often determines if these tools will be successful.

Source	Author	Publisher
Business Intelligence Roadmap	Larissa T. Moss and Shaku Arte	Pearson
Business Intelligence: The Savvy Manager's Guide	David Loshin	Addison Wesley
Business Intelligence Fundamentals	TDWI Course	TDWI
Business Intelligence and Analytics: Systems for Decision Support	Sharda, Ramesh; Delen, Dursun; Turban, Efraim; Aronson, Janine; Liang, Ting Peng	Pearson Education
Decision support and Business Intelligence Systems	Turban, Sharda & Delen	Prentice Hall
Understanding Actuarial Management: the Actuarial Control Cycle, 2 nd Ed.	Bellis, Lyon, Klugman and Shepard. Ed.	Institutes of Actuaries of Australia and the Society of Actuaries

Sources for further information

⁶ This paper takes the perspective that the actuary is supported by an IT department within his or her organization. For those actuaries who do not have that benefit, "IT" can be also thought of as an external software or hardware vendor.

Source	Author	Publisher
Understanding Actuarial Practice	Stuart Klugman, Ed	Society of Actuaries
Deloitte – "How to Build a Successful BI Strategy"	Prashant Pant	Deloitte
Applied Insurance Analytics	Patricia Saporito	Pearson

BUSINESS INTELLIGENCE SOLUTIONS IN THE CONTEXT OF ACTUARIAL PROCESSES

Business Intelligence solutions are subject to data governance processes which provide for "an enterprise-wide data governance body, a policy, a set of processes, standards, controls, and an execution plan for managing the data."⁷ While data governance processes control the flow and allocation of data as a resource, it is up to the actuary to determine how these resources are turned into actuarial work products. In order to ensure a smooth process of data-to-information and to determine the proper use of actuarial information, the actuary must have an active working relationship with the process, people and technology of their Business Intelligence landscape. This includes familiarity with the BI development/delivery lifecycle as well as the various capabilities of the larger organization's accepted BI toolset. Embracing with academic curiosity this decidedly non-actuarial discipline and developing these relationships can often be the critical factors in the success or failure of an actuarial department.

Actuarial resources, products and workflow must be managed like any other business activity. This activity can be described as the Actuarial process. This process can be broken into distinct components.

- Define the Problem
- Design the Solution
- Monitor the Results

Each of these components involves large and small tasks. Large tasks, like measuring and reporting annual statement loss reserves; and small tasks, like answering a regulator's inquiry about a rate filing, all require historical information. The source of required data can be internal to the corporation, such as that arising from policyholders and claimants, or external, such as economic and political data.

⁷ Deloitte - "How to Build a Successful BI Strategy"

Business Intelligence Architects design and direct the flow of data into the Actuarial Process

BI solution architects and managers are charged with managing the corporate data as a resource. They are responsible for delivering:

- Faster computations
- Improved communications and collaboration
- Increased productivity of supported groups
- Improved data management
- Better management of actuarial data reservoirs
- Improved quality of decisions
- More agile support
- Increased cognitive limits
- Using the web
- Anytime, anywhere support

Optimally, the IT architects and the actuaries work together to decide on the best tools and data for supporting the actuarial process.

Sources for further information

Source	Author	Publisher
Decision support and Business Intelligence Systems	Turban, Sharda & Delen	Prentice Hall
Understanding Actuarial Management: the Actuarial Control Cycle, 2 nd Ed.	Bellis, Lyon, Klugman and Shepard. Ed.	Institutes of Actuaries of Australia and the Society of Actuaries
Understanding Actuarial Practice	Stuart Klugman, Ed	Society of Actuaries
Loss Models: From Data to Decisions	Klugman, Panjer & Willmott	Wiley & SOA
Introduction to Scientific Computation and Programming	Daniel Kaplan	Thompson
Intelligence and Other Computational Techniques in Insurance: Theory and Applications	Shapiro & Jain, ed.	World Scientific

THE TECHNOLOGY OF ACTUARIAL PROCESSES AND BUSINESS INTELLIGENCE SOLUTIONS

Actuaries, underwriters, claims examiners and all others within an insurance organization are part of a complex data-driven system aimed at making accurate assessments to advance the goals of their management. Technology is applied at both divisional and enterprise levels to achieve these goals. Each organization balances this technology distribution between division and enterprise differently depending on the goal targeted and/or their appetite for standardization.

The Technology of Actuarial Processes

As professionals, actuaries work through an actuarial process that constantly looks at loss experience from the past, adjusts that experience to current conditions and then reports the findings for actionable decisions. Problems are defined, solutions are designed and reported, decisions and actions are monitored, all in a cycle of activity. The actuary leverages multiple technologies for this monitoring, adjusting, and analysing past experience.

Actuarial Technology Characteristics and Capabilities are frequently characterized by the following:

- They can be standalone.
- They are used for multiple levels of management.
- They are adaptable and flexible.
- They are interactive and easy to use.
- Actuaries control the process.

Actuarial Systems

Certain actuarial processes can rely heavily on technology to maintain standardization. Reserving and capital allocation processes are typical examples where heavy investment in the build of a structured Actuarial System is not uncommon. Best practices of Actuarial Systems suggest they include a database management subsystem, a model management subsystem, and a user interface subsystem.

The database management subsystem consists of a decision support database with a DBMS, a Data dictionary, and a Query facility. The data in the decision support database (primarily premium, exposure and loss data) are non-volatile, cleaned, in a standard format and not used in a transaction processing environment. IT manages the directory, data quality, query facility, data integration, data scalability and data security.

The second component of a structured Actuarial System is a model management subsystem which contains strategic planning models, capital allocation models, pricing and reserving modes and predictive models.

The final component of structured Actuarial System is a user interface system which provides for clear communication of the results of analysis. Of growing importance in user interface systems is data visualization. Data visualization capabilities are a growing expectation for both structured Actuarial Systems as well as smaller actuarial deployments of technology.

Modeling Languages – Many languages can be used in Actuarial Processes. Earlier languages were FORTRAN, Basic, and APL. SQL has replaced many of the processes previously driven by these languages. Today, statistical languages like R and SAS are becoming more common.

Despite the breadth of modeling languages available, it is not uncommon for many actuarial process to be supported exclusively through Microsoft Office applications like Access and Excel spreadsheets including VBA and Excel Add-ins like those from Palisades Corporation.

The Technology of Business Intelligence Solutions

As discussed, actuarial processes may or may not depend on the technology commonly associated with Business Intelligence solutions. With the desire to include more data into the actuarial process, it has become important to consider the use of more sophisticated and scalable Business Intelligence solutions into actuarial processes. In this section, specific tools are mentioned that have been known to create successes within financial institutions like insurance companies. It is important, however, to recognize that new tools are introduced to the market every year which may or may not improve on the capabilities of existing products.

Business Intelligence & Analytics Software Tools

This term refers to the software that supports the business processes, methodologies, and metrics, used by the insurance enterprise to measure, monitor, interpret and forecast business performance. From a cost perspective it is desired to have a single software platform support these processes. In practice, however, the need to balance the varied metrics, timing and detailed appetite of different business support personnel generally requires that different software tools be implemented throughout the enterprise. IBM Cognos and SAP Business Objects are seen as leaders in the space of standard reporting delivered across an enterprise. Their highly structured design, however, creates challenges to meeting new expectations in the marketplace. Frequently they add capabilities through acquisition of innovative competitors. Despite the perceived lack of flexibility of reporting, these giant technology companies benefit from their software deployments by early adopters of the enterprise BI movement in the early 2000's. They are likely the mainstay of many financial organizations where consistency and a "single source of truth" are critical to forming and delivering on expectations. As the desire for "data discovery" continues to widen to support functions that seek

trend and correlation insight over precise measurement, tools that specialize in "visualization" have moved to the forefront. Tableau and Qlik, once the newcomers are now becoming established providers looking to maintain market share against steady upstart competition. Microsoft Excel, easily considered the current tool of choice of financial analysts and actuaries, continues to evolve to compete with these mid to large scale BI tools. Depending on the controls required for the business process, it may still be the optimal choice.

Data Warehouses

Data warehouses can but frequently do not have an operational role within an organization. That is, traditionally the role of a data warehouse is to provide control and efficient storage for timely reporting of financial values recognized widely throughout the organization. Data warehousing tools are designed to store and refresh large volumes of structured data (discrete values with a limited number of bytes). With the onset of "big data," the concept of the data warehouse has become less *en vogue*. Despite its current departure from popularity, data warehouses (in one form or another) are the foundational data structures of many organizations. This is particularly true for financial organizations that require stringent internal and regulatory controls on financial information, both historical and prospective. Database tools such as SQL Server, Oracle, IBM DB2 are strongly associated with traditional data warehouses.

Data Mining Applications

These are technology that use statistical, mathematical and artificial intelligence techniques to extract and identify useful information and patterns obtained from large sets of data.

Any tool that can parse text (which includes Excel and SQL), can technically be considered a data mining tool. That notwithstanding, an unstructured dataset can easily exceed a terabyte which generally calls for a tool that can comfortably accommodate such volume. Both SAS and Hadoop are seen as successful tools for data mining projects. A look at Gartner's 2015 *Magic Quadrant for Business Intelligence and Analytics Platforms* drives home the rapid expansion of BI tools that can manage large databases although there are few tools that appear to be a "one-size-fits-all" solution. In fact, many tools in this space are configured to work with other tools recognizing the varied goals of data mining projects.

Sources for further information

Source	Author	Publisher
Decision support and Business Intelligence Systems	Turban, Sharda & Delen	Prentice Hall
Excel 2010 Data Analysis and Business Modeling	Wayne Winston	Microsoft Press
VBA for Modelers: Developing Decision Support Systems with Microsoft Office Excel, 4 th Ed.	s. Christian Albright	South-Western Centage Learning
@Risk: Advanced Risk Analysis for Spread Sheets		Palisade Corporation
Computational Actuarial Science with R	Arthur Charpentier, ed.	CRC Press
Session 57 L: Business Intelligence for Actuaries	Rigby & Levine	Society of Actuaries
Practical Management Science, 4 th ed.	Winston & Albright	South-Western Centage Learning
Magic Quadrant for Business Intelligence and Analytics Platforms	Gartner	Gartner
Wikipedia: "Database" "Business intelligence tools"	Various	NA
Making Successful Presentations: A Self- Teaching Guide	Terry Smith	Wiley Press
Practical Data Science with R	Zumel & Mount	Manning
Applied Insurance Analytics	Patricia Saporito	Pearson
THE ROLE OF THE ACTUARY IN BUSINESS INTELLIGENCE PROJECTS

Business Intelligence Project Challenges

Despite the millions of dollars spent on BI projects, many BI projects fail or, at the very least, fail to meet their expected potential. There is no shortage of publications lamenting this observation. Even in the financial industry (which includes insurance), where BI projects are most attempted, success stories are limited. A post-mortem of a project's diminished success frequently boils down to two critical gaps:

- Sufficient and sustainable support at the sponsorship level
- Concrete recognition and consensus on the return on investment

Of course, these two causes are frequently related. Sometimes BI projects are embraced with only a fuzzy understanding of the potential or, sadly, because "everyone else is doing it." Many are justified by the assumption that current processes are so fractured and inefficient that a BI project can only improve the organization.

BI projects are generally costly (in the millions) and lengthy, frequently scheduled over years to sufficiently spread the costs. Unless the BI project's business objectives are clear and core to the long term success of the organization, the risk is great that the originally envisioned expectations of the project will fall victim to changing short term priorities or unexpected forces.

Even when sufficiently clear goals and an expected ROI are embraced at the onset, technological advances unveiled during the course of a project can easily sway an impatient sponsor (the CFO or CEO, for example) into changing directions that appear to be faster and less costly. Similarly, a change in the "c-suite" can derail a BI project mid-stream as the new chief officer may want to recast the future BI landscape to their liking and not to that of their predecessor's.

Given this bleak track record, you might wonder why organizations continue to pursue large BI projects instead of operating locally with the power of spreadsheet technology. The answer remains the same as it always has been: the decision-making confidence gained by centralized and integrated data exceeds the pain and cost of a BI project. In fact, as the need to integrate data from external sources and stored descriptive (non-financial) data proliferates, BI projects are more critical than ever if an organization is to remain competitive.

Envisioning improvements to the Actuarial Process within a Business Intelligence Project Scope

Similar to the standards employed in developing an actuarial opinion, creating and delivering a BI report to be used by an organization has an accepted set of appropriate steps or "best practices."

Without these best practices, the delivered report is likely to not meet the expectations of the users in all of the familiar disappointing ways: the report will not be available as scheduled, it will not address all of the expected needs of the user (aka the actuary), nor will it be as flexible to use as desired. Given the investment of time and efforts of the multiple individuals involved (often far more than the creation of an actuarial opinion), it only makes sense that actuaries embrace these best practices as they would uphold the standards of their own work product.

The BI Development/Delivery Life Cycle can be arranged in varying ways within an organization. Some methodologies require formal, detailed meetings and documents with sign off responsibilities for each invested party. Others are less formal and may rely on multiple iterations of prototypes in order to build out the final product. Even when the planning and management is driven more by the business than IT, it is important to align resources, timelines and the efforts of others not directly involved in the project in order to make the best use of resources. It would do well for the actuary to inquire and adopt the expected approach and terminology used by those planning the project.

Despite the variations of approach (waterfall or agile⁸), BI solution development has a commonly understood set of required activities (although the terms used to describe the activities may vary), all of which are critical in order to achieve the goals of the product. It is unlikely that an actuary will be 100% dedicated to a BI project targeted to meet actuarial needs, but the actuary has a defined role in each of these activities.

Planning

"Planning" can be thought of as a constant drive to innovate and/or improve processes. Throughout their career, actuaries should seek out opportunities to brainstorm with co-workers in IT and other areas (finance, operations) on system or reporting improvements that would benefit all parties. That way, when there is an appetite to invest in technology, there is a better chance that the actuary's needs are known across the various support functions. If, on the off chance, the ensuing BI project targets the actuarial department's needs directly, that actuary is likely to have influence on the ultimate design.

Analysis

The use of the term "analysis" in this context generally refers to "Business Requirements Analysis" which is the gathering and documenting of the BI needs and expectations of the targeted users. The actuary will be well served to take the term quite literally from an actuarial science perspective and

⁸ In the phases of a traditional waterfall development arc, you move to the next phase only when the previous one is complete.

However, [in an agile development model] instead of tackling all the steps for all of your product features at once, you break the project into *iterations* (smaller segments of the overall project), called *sprints*. Mark C. Layton, *Agile Project Management for Dummies*

craft out any statistics to support business impacts that would assist in prioritizing needs, either initially or later during the project.

It can be a challenge to sort through the needs and wishes of everyone to determine core deliverables that would enable the users to experience desired results. One recommended approach is to ask "why?" repeatedly until solutions that meet the common needs of all engaged users are uncovered. It is not uncommon for any technical analyst to balk at this process. It can seem patronizing and a waste of time to discuss with non-actuaries what an actuary does and to what end. It is important during this process for the actuary to remember the values of methodical exploration, numerical or verbal, and to recognize that their organization is investing significant resources to the project and resistance to an accepted methodology is the true waste of time.

Design

This includes the finalization of the technical design and specifications. Although it may be discussed during Analysis activities, during Design the technology (software package and infrastructure) will be determined. As discussed earlier, it would be optimal for the actuary to inquire extensively on the software options available to the project. If IT is leading the project, they might assume that, like other support areas, an actuary would have little interest in the varying options of one software package or another. An actuarial department, however, is likely to rely heavily on Excel and custom queries in the existing process and should provide considerable input into the required capabilities of the new solution.

Development

During Development it is tempting for an actuary to "get back to their real job" and wait until they are beckoned again by the project leads. It is during the development process, however, that the hard decisions and compromises are made. Those closer to the project during Development are likely to have the most influence on the initial output. Additionally, knowing where the initial project design had to be "tweaked" due to unforeseen development, will provide the actuary insight later on during the testing phase. It is widely accepted that modifications during this phase are considerably less effort (aka cheaper) than changes made near the end of the project.

Test

Unfortunately, Testing is frequently considered the "last call" for any changes to the BI deliverable. In fact, testing should be performed throughout the project to constantly validate the requirements. As such, the actuarial area should strive to develop test cases that will not only challenge the speed of performance but should also seek to test the unusual but valid request. Some test failures will be more critical than others. It is important to classify the degree of failure as some fixes are likely to be postponed to later releases. The decomposition of these failures is likely to lead to the discovery of a poorly understood business requirement or software capability not captured earlier in the project.

Implementation

Assuming testing was thorough, this process should be relatively painless for the actuary. If, however, the testing process was abbreviated, it is likely to be the most painful. Even if implementation goes smoothly, the actuary will likely need to advocate for the new solution and demonstrate to others its contributions to the Actuarial Control Process.

Source	Author	Publisher
Business-Driven Business Intelligence and Analytics: Achieving Value through Collaborative Business/IT Leadership	David Stodder	TDWI
Seven Strategies for Creating High-Performance BI Teams	Wayne Eckerson	TDWI
Wikipedia "Systems development life cycle"	Various	NA
Agile Project Management for Dummies	Mark C. Layton	Wiley

Sources for further information

CONCLUSION

Providing management data and information for supporting insurance financial systems is the main task for professional actuaries. The data, the models, the communications and professional advice are all resources that actuaries use to support decisions.

For all actuaries currently working, from the actuarial student to those in the final years of their professional life, one constant is and has been the rapidly evolving technology that supports the collection and analysis of data and the proliferation of data sources available. Actuaries in all practice areas and certainly CAS actuaries should stay vigilant to the opportunities and problems brought about by this rapid evolution in technology.

Data Quality Overview Actuarial Concepts in Data Quality

DATA QUALITY PRINCIPLES

In its simplest terms, Data Quality can be defined as data "fit for its intended use." In other words, Data Quality is measured in terms of how well it fits the data consumers' expectations. The categories in which Data Quality is evaluated include:

- Validity is the information captured in correct formats, with codes or values that are appropriate for the business? For example: Is Zip Code 10019 is valid for the state of New York?
- Accuracy does the information captured truly reflect the business information? Continuing with the above example, although valid, the data is not accurate if it is for a risk in Upstate NY while the intended zip code was meant to describe New York City.
- Completeness is used to measure the breadth of the data. Is all of the data that is supposed to be in the file or analysis included? What may have been excluded or duplicated?
- Timeliness is associated with the 'freshness' or time-lag of the data. If we need to support near real time customer service calls, data that is a month old may not meet the quality expectations of the consumer.
- Reasonability refers to the consistency or materiality of the data, given the business conditions. For example, a significant shift in the distribution or profile of a company's book of business may be reasonable, if the company has entered into a new territory or market.
- Data Lineage is a newer category of Data Quality. It includes ability to transparently trace the data path from creation to reporting, including data transformations. This path provides information about the reliability of the data.

Managing Data Quality

Data Quality is typically managed through the development and monitoring of metrics. These metrics must be measurable and should be quantifiable within a discrete range. Note, however, that while there are many things that can be measured, not all translate into useful metrics, implying the need for business relevance. Therefore, every data quality metric should demonstrate how meeting its acceptability threshold correlates with business expectations. The above data quality dimensions should frame the business requirements for data quality. Quantifying how quality is measured along the identified dimension provides hard evidence of data quality levels. The determination of whether the quality of data meets business expectations can be based on specified acceptability thresholds; if the score is equal to or exceeds the acceptability threshold, the quality of the data meets business expectations. Any measurable characteristic of information that is suitable as a metric should reflect some controllable aspect of the business. In other words, the assessment of the data duality metric's value within an undesirable range should trigger some action to improve the data being measured.

If the score is below the acceptability threshold, the appropriate data steward must be notified, and some action must be taken. These data quality metrics can be organized by the data quality dimensions noted above. For example:

- Completeness: Metric to monitor whether total dollars on claim records (e.g., a loss run) balances to a total on a control report;
- Timeliness: Is the claims upload delivery in the agreed upon time range agreed upon among stakeholders completed?
- Validity: What percentage of zip codes in the data are actual valid US zip codes?
- Integrity: What percentage of policyholder records contain a missing or null field?

Quantifiable metrics enable an organization to measure data quality performance improvement over time. Tracking helps in monitoring activities within the scope of data quality service level agreements and demonstrates the effectiveness of improvement activities. Once an information process is presumed to be stable, tracking enables the institution of statistical control processes to ensure predictability with respect to continuous data quality.¹

Sources for further information

For more information on the overall principles of data quality, please see:

Source	Author	Publisher	Link (if applicable)
Data Quality - The Field Guide	T. Redman, Ph.D.	Digital Press	
Data Quality Assessment	A. Maydanchik	Technics Publications	
Risk Management and Insurance Operations - CAS Course 1 For Preparation For Exam CA1 - Assignment 16 - Actuarial Data Management	Casualty Actuarial Society	The Institutes	

Outside of the basic principles of data quality, the concepts which many actuaries would benefit being aware of fall under an umbrella of two broad topics, joined by their aims to promote a higher standard of information quality within the organization. As discussed in detail in the following sections, these are:

- Data Governance Concepts
- Data Documentation Concepts

¹ IDMA Tools for Managing Data Effectively One Day Class 2014 Insurance Data Management Association

DATA GOVERNANCE CONCEPTS

Quick Glossary of Important Terms

	Actuaries Should Know:
Term	Definition
	Data governance is "the exercise of authority, control, and shared decision making (planning, monitoring, and enforcement) over the management of data assets." ²
Data Governance	Similarly:
Governance	"Formalized behavior associated with data. Includes execution and enforcement of authority over the management of data and data-related assets/processes." ³
Data Stewardship	"Formalized accountability over the definition, production, and use of data and data-related assets/processes." ⁴
Data Governance Committee	A data governance council or committee (DGC) is a cross-functional group with members from both IT and the organization's operational side. Members of the DGC generally include the Chief Information Officer (CIO), Chief Data Officer (CDO), the Data Management (DM) leader, and a business executive who acts as Chief Data Steward. It is not uncommon for this group to include executives representing other functions, such as actuarial, underwriting, and claims. The DGC makes high level, strategic decisions about data governance as an integrated function within the organization.

² [1] The DAMA Guide to the Data Management Body of Knowledge (DAMA-DMBOK) First Edition, Mark Mosley Editor, Technics Publications LLC, New Jersey, copyright 2009 DAMA International, p., p. 37

³ Robert S. Seiner, TDAN.com (The Data Administration Newsletter)

⁴ Robert S. Seiner, TDAN.com (The Data Administration Newsletter)

Data Governance: Overview

As described above, the primary goal of any corporate data governance initiative is to manage data for the purpose of delivering accurate, valid, timely, and complete data which can be used to inform decisions across the company.

Data governance, therefore, may encompass elements extending beyond simply data. For example, there are the classic people, policy, process, and technology dimensions to data governance. Each must be individually defined with goals in order to successfully support an organization's data governance strategy. Let's examine the elements of each pillar:

- Process Data governance processes which provide for "an enterprise-wide data governance body, a policy, a set of processes, standards, controls, and an execution plan for managing the data."⁵
- People Clearly defining roles and responsibilities across the data managers or influencers in the organization. These may include:
 - Business Analysts Those who actively utilize the data. Those who utilize the data are
 often the best to provide feedback on the data required for robust analysis. Actuaries,
 for example, may be considered business analysts.
 - IT Architects Those who design and direct the flow of data within the organization.
- Technology As will be discussed in more detail in a subsequent section of this paper, this involves deploying the best suited tools and data infrastructure needed to support the objectives of the organization. For example, if real-time data is necessary, data architects must design technology appropriately suited to deliver the information to the analysts in real-time.

Under the umbrella of these pillars, "Data governance focuses on the delivery of trustworthy, secure information to support informed business decisions, efficient business processes, and optimal stakeholder interactions. It is therefore not an end in itself, but merely the means: data governance supports your most critical business objectives."⁶

While the paper is written from the perspective of the need for data quality, there are business considerations which must be considered in this pursuit. Data quality initiatives involve significant commitments of time, resources, IT architecture, and capital. While actuaries may strive for the absolute best data quality, this must be balanced with an analysis of the marginal returns of greater and greater quality initiatives. For example, does spending \$1m for a marginal increase in data quality warrant the business benefits? Actuaries must be prepared to discuss with management the cost implications for the data quality initiatives discussed in this paper.

⁵ Deloitte – "How to Build a Successful BI Strategy"

⁶ https://www.informatica.com/content/dam/informatica-com/global/amer/us/collateral/white-paper/metadatamanagement-data-governance_white-paper_2163.pdf

Data Governance Committees [People and Process]

Overview

"The Data Governance Council's (DGC) primary duty is to ensure responsibility, accountability and sustainability of data practices. The framework for effective data governance planning contemplates the personnel, technology and policies and procedures necessary to ensure the preservation, availability, security, confidentiality and usability of the company's data. Furthermore, a DGC encourages strategic thinking and the creation of opportunities surrounding the appropriate use of data within the organization."⁷

In simpler terms, the DGC manages all projects related to and impacting data within the firm. The advantages of establishing a DGC are numerous and include:

- Enhanced Deployment of Resources. Various stakeholder groups may have similar data request needs which are completed or attempted in parallel tracks, resulting in overlap or inefficient deployment of resources. The result may be a patch work of data flows rather than a logical, consolidated flow. A DGC coordinates IT efforts from a focal point of authority.
- **Singular Management of Data.** Data has typically been collected and managed at the business unit or stakeholder level based on individual needs. With increasing needs to integrate or exchange data, organizations need coordinated management to effectively:
 - o Migrate data from legacy platforms to current, more advanced solutions;
 - Integrate various systems which may speak a different language (e.g., essentially contain the same information such as dates but in different formats);
 - o Determine the accepted definitions of the data for reporting and analysis needs;
 - Report data consistently and ensure its fit for use.

A DGC which centrally manages the data assets of an organization ultimately improves the quality of data and information across the organization. Via consistent naming standards, definitions, formal metrics and calculations, there is an improved understanding of data. This facilitates better communication and understanding of the data, which in term aids in the ability to share or re-use data.

• **Representation of Stakeholders**. The DGC has cross-functional representation and works to understand the needs of corporate stakeholders from a data perspective. Understanding the needs of stakeholders is a key component in creating a synergistic data strategy.

Membership and Actuarial Roles in Data Governance

As insurance companies begin to establish formal Data Governance Committees, representatives

⁷ http://www.insidecounsel.com/2014/01/27/inside-establishing-a-data-governance-committee-as

from both the business and technical data stakeholder groups are often included:

- Senior Executives (e.g., Chief Information Officer, Chief Data Officer, and/or Chief Technology Officer)
- Business Stakeholders from:
 - o Financial Reporting;
 - o Underwriting;
 - o Claims;
 - o Actuarial
- Technical stakeholders:
 - o Data architects
 - o Business analysts
 - o Project managers

Actuaries usually have a seat or two at the DGC table. Many organizations choose an actuary to be the Chief Data Steward if there is not a formal chief data officer.

Roles and Responsibilities

The roles and responsibilities of a DGC include, but are not limited to:

- Clearly establishing senior authority over data streams which cross organizational boundaries;
- Evaluating all internal data-related projects for coherence with the overall corporate data strategy, architecture, and overlapping work-streams to reduce inefficiencies, redundancies, and conflicting data streams;
- Designing the data related controls and processes for which data travels throughout the organization;
- Monitoring the compliance of data processes with controls mandated by various internal and external authorities (e.g., regulators, auditors, Sarbanes-Oxley);
- Responding to data process or compliance issues by prioritizing resources, approving remediation or strategic plans, and approving the data architecture utilized to support the data processes; and
- Conduct annual audits of strategic data processes.

Tools Available for Data Governance [Technology]

Data Governance Councils rely on the following tools or artifacts:

- **Policies and Procedures**: These are the 'rules' the organization has established for how data and information is processed and analyzed. They may include both internal and external Standards and Guidelines for how data and information is to be coded, processed or exchanged.
- Enterprise Data Models: These are the 'blueprints' for how data is organized in the various databases and systems. These are used to understand how the data is related to other data and processes and may be a source for data quality rules.
- **Collaboration Tools:** Most DGCs use a variety of collaboration tools to capture discussions, histories and revised policies and procedures. These aid in the recordkeeping and documentation of important decisions and actions.

	A . 1	D 1 1 1	Link (if
Source	Author	Publisher	applicable)
Stewardship Approach to Data Governance	Robert S. Seiner	The Data Administrators Newsletter (TDAN)	http://tdan.com/the-data- stewardship-approach-to- data-governance-chapter- 1/5037
Deloitte – "How to Build a Successful BI Strategy"	Prashant Pant	Deloitte	http://www.loria.fr/~ssidh om/UE909R/1_BI_strateg y.pdf
Establishing a Data Governance Committee as part of 2014 strategic priorities	David Katz	Inside Counsel	http://www.insidecounsel.c om/2014/01/27/inside- establishing-a-data- governance-committee-as
Defining Organizational Structures	Gwen Thomas	The Data Governance Institute	http://www.datagovernanc e.com/defining- organizational-structures/
Risk Management and Insurance Operations - CAS Course 1 For Preparation For Exam CA1 - Assignment 16 - Actuarial Data Management	Casualty Actuarial Society	The Institutes	

Sources for further information

DATA DOCUMENTATION CONCEPTS

Introduction

There are varying levels of data management and documentation, all of which an actuary can play an integral role in for an organization. The following glossary of terms are important in the subsequent documentation discussion. The interdependence of documentation and governance are concepts that will be explored further in this section.

Term	Definition
Big Data	Big Data is high-volume, high-velocity and high-variety information assets that
	demand cost-effective, innovative forms of information processing for enhanced
	insight and decision making. ⁸
Master Data	Master Data represents the business objects which are commonly agreed to and
	shared across the organization. Customer and / or Product IDs are examples of
	Master Data. ⁹
Master Data	Master Data Management (MDM) is a comprehensive method of enabling an
Management	enterprise to link all of its critical data to one file, called a master file that provides
	a common point of reference. When properly done, MDM streamlines data
	sharing among personnel and departments. In addition, MDM can facilitate
	computing in multiple system architectures, platforms, and applications. ¹⁰
Metadata	Metadata are business and technical information about an organization's data.
	They help put data in context, reveal their meaning and make them accessible. ¹¹
	Metadata is structured information that describes, explains, locates, or otherwise
	makes it easier to retrieve, use, or manage an information resource. Metadata is
	often called data about data or information about information. ¹²
	Metadata summarizes information about data for the purpose of making that data
	easy to find and work with.

Quick Glossary of Important Terms Actuaries Should Know

⁸ www.gartner.com/it-glossary/**big-data**

⁹ Wikipedia

¹⁰ http://searchdatamanagement.techtarget.com/definition/master-data-management

¹¹ IDMA Curriculum Rewrite Task Force - Course 201 Assignment 6

¹² http://www.techrepublic.com/blog/it-security/is-metadata-collected-by-the-government-a-threat-to-your-privacy/

Term	Definition					
Metadata	A Metadata Depository is a database and software used to capture, manage and					
Repository	access metadata. It is where an organization collects, integrates, standardizes,					
	consolidates, organizes, controls, and stores its metadata, and makes them available					
	for shared general use. ¹³					
	A Metadata Repository is a special type of database containing information about					
	another database, e.g., how the data in the other database was collected,					
	transformed, and formatted, how frequently it is updated, and generally anything					
	that can be useful to analysts that need to query data from that database. ¹⁴					
Data Flow	A mapping of data flows from source systems (e.g., Policy Admin System,					
Chart	Enterprise Data Warehouse, etc.) to intermediate data stores (if any) and finally to					
	end user.					
Stakeholder	A comprehensive review of the data requirements of the enterprise stakeholders,					
Data	both for analysis and reporting purposes. This often comprises the first step in					
Analysis	formulating an overall enterprise data strategy.					
Data	"A data dictionary is a tool for displaying metadata to business and technical					
Dictionary	personnel. A data dictionary is important for expediting the transfer of knowledge					
	regarding the meaning of data values stored in the data fields." ¹⁵					

 ¹³ DMA Curriculum Rewrite Task Force - Course 201 Assignment 5
 ¹⁴ http://www.casact.org/pubs/forum/05wforum/05wf274.pdf
 ¹⁵ https://www.casact.org/pubs/forum/05wforum/05wf274.pdf

Data Documentation and the relationship to Data Governance

The key to superior data governance is the processes supporting the management and documentation of data. As will be discussed, there are varying levels of data management and documentation, all of which an actuary can play an integral role in for an organization.

"Maintenance of adequate documentation describing the data can help avoid problems associated with relying exclusively on people's memories of what is contained in the data. As actuaries we can help persuade our business and data management partners that system documentation is vital to the actuarial work product." Documentation is the cornerstone for well performing data governance – in that sense, data documentation and data governance are not mutually exclusive. Data won't govern itself.

Stakeholder Analysis

One of the first items of documentation a Data Governance Committee may seek to help formulate strategy, priorities, and objectives of the data governance function is a stakeholder analysis.

Stakeholders across different departments demand information presented in a way which aligns with their operational objectives – understanding what these stakeholder operational and reporting objectives are help define the data an organization needs to collect, store, and process.

As a simplistic example, personal auto coverage is often differentiated by actuaries into property damage and bodily injury components due to the differing development characteristics of these pieces. However, the profit center managers for personal auto do not view these components in isolation - they understand composite results and rates, and thus demand combined metrics. Furthermore, actuaries may look at the business pieces on a countrywide basis, while profit centers require a state by state breakdown of the results. As such, it's necessary for a data organization to understand what its stakeholders require from a data perspective – in this case, the actuaries would require losses and premiums on a state by state basis, split by property damage and bodily injury, for both analysis and allocation purposes.

Data and Process Flow Diagrams

Data and process flow diagrams are used to document the lifecycle that data goes through in the organization (again, one may hear the term data lineage when discussing the data lifecycle). It visually represents the various systems (input and outputs) that are involved the creation, consumption and transformation of data. It helps to ensure that all systems and processes are accounted for and is used to manage data lineage and data impacts.

Sources for further information

Source	Author	Publisher	Link (if applicable)
The Data Governance Institute			http://www.datagovernan ce.com/
The Data Administration Newsletter		Robert S. Seiner	http://tdan.com/
Risk Management and Insurance Operations - CAS Course 1 For Preparation For Exam CA1 - Assignment 16 - Actuarial Data Management	Casualty Actuarial Society	The Institutes	

Metadata: Overview and Technical Concepts

Metadata is one of the more topical and least understood data documentation concepts.

Metadata provides the context and descriptions of the data (the type, what it means, where it is located, how it used, etc.)¹⁶ Metadata is important from a data lineage perspective – the ability to trace the data through its various stages and transformations is a key to ensuring data quality. In the broadest sense,

Metadata can be considered the documentation of the contents of a database. In addition to the information about the data itself, metadata contains information about business rules and data processing.¹⁷

As defined in the CAS 2007 Winter Forum and in reference to Data Quality: the Accuracy Dimension by Jack Olson,

Metadata is a term used by data management professionals for information about the data such as definitions, a description of permissible values and business relationships that define the data in a database. <u>Comprehensive metadata is a prerequisite for good information quality</u>.¹⁸

To that extent, Metadata can be classified into subtypes.

Types of Metadata

There are three types of metadata the actuary should be aware of:

- 1. Technical (also known as Structural)
- 2. Business (also known as Descriptive)

¹⁶ http://dataqualitypro.com/data-quality-pro-blog/data-quality-through-metadata-strategy-anne-marie-smith)

¹⁷ https://www.casact.org/pubs/forum/08wforum/actuarialIQ.pdf

¹⁸ https://www.casact.org/pubs/forum/07wforum/07w279.pdf

3. Operational (also known as Administrative)

Technical Metadata

Technical metadata assists in understanding the format and definition of the data collected. For example, if one were sending a letter, the information contained on the envelope can be considered the data. The "data describing data" about that envelope may include a name and address and specify the format that "data about the data" is in.¹⁹ This may include about the addressee:

- Surname coded as 20 digits, alphabetic
- U.S. State coded as two digit alphabetic code defined by US Postal Service
- Zip Code coded as five digit numeric as defined by the US Postal Service

Technical metadata helps interpret the raw data – for example, we would know that NJ is a state abbreviation for New Jersey through structural metadata. Structural metadata is thus a key component in data quality – it gives context to the data (e.g., we know that NJ is a reference to the US State of New Jersey, how it should be coded – in this case with a 2 letter abbreviation, and its meaning).

Technical metadata includes the source table information so a user can know exactly where the information is sourced from. This is important in using metadata to generate data flow charts or maps.

Business Metadata

Business (or descriptive) metadata assists in describing a resource for purposes such as discovery and identification.²⁰ This type of metadata provides context around the data, including but not limited to the data field's name, definition of contents, related data, as well as the applicable business rules.

Continuing with our letter mailing example from above, business metadata may log characteristics of the transaction around sending the letter. This includes the addressee, sender, post-mark date, post offices handling, date of delivery, etc.²¹

Operational Metadata

Operational (or administrative) metadata provides administrative "data about the data", including date of last update, date of last access, user who last modified, movement from source to target, availability and usage. Operational metadata may also include a description of the types of data control or quality checks performed on the data, and where in the data processes these occur – this is a metric that allows for audit trails providing proof of compliance for data related controls.²²

¹⁹ http://www.riskandinsurance.com/the-whodunit-of-big-data/

²⁰ http://www.niso.org/publications/press/UnderstandingMetadata.pdf

²¹ http://www.riskandinsurance.com/the-whodunit-of-big-data/

²² https://www.casact.org/pubs/forum/08wforum/actuarialIQ.pdf

Metadata as a Key to Successful Data Governance

Key to successful data governance is the management of metadata – the frame of reference giving data its context and meaning. Effectively governed metadata provides a view into the flow of data, the ability to perform an impact analysis, a common business vocabulary and accountability for its terms and definitions, and finally an audit trail for compliance.

High Level Metadata Example: Metadata in the News

Metadata's association with "data about data" is best seen through an example which is rather topical, albeit non-insurance related. In 2013, President Obama addressed the current data collection practices of the NSA with a reference to metadata:

What the intelligence community is doing is looking at phone numbers and durations of calls. They are not looking at people's names, and they're not looking at content. But by sifting through this so-called **metadata**, they may identify potential leads with respect to folks who might engage in terrorism.²³

Continuing with this example, suppose you make a phone call to a friend: The phone call's conversation (e.g., what was said) is raw data. Data without context may have little to no value – for example, a relatively mundane conversation may not generate any actionable information. However, the data about this data (i.e. the data about the conversation) may be classified as metadata and create the context needed to generate actionable information. For example, this metadata may include the date and time you called somebody, the duration of the phone call, the phone numbers involved, or the location of the participants.²⁴

Using this information, as President Obama noted in his press conference, is how the NSA has generated links to those involved in terrorist activities.

High Level Metadata Example: Metadata in Health Care

An example of the power of metadata in cross-referencing data sources can be seen in the health care industry where vast amounts of patient information is collected, often by different users or systems. Metadata is a key component in tying this information together – "allowing resources to be found by relevant criteria, identifying resources, bringing similar resources together, distinguishing dissimilar resources, and giving location information."²⁵

This information may be tied together via a metadata repository which consolidates the metadata from various source systems, and from there integrate with the system through which end users query.

 ²³ http://blogs.wsj.com/washwire/2013/06/07/transcript-what-obama-said-on-nsa-controversy/
 ²⁴ http://www.theguardian.com/technology/interactive/2013/jun/12/what-is-metadata-nsa-

surveillance#meta=1111111

²⁵ http://www.niso.org/publications/press/UnderstandingMetadata.pdf

In health care, it can be used to link patient information across sources. For instance, "If the physician prescribes the patient aspirin for a chronic headache, metadata could be used to retrieve other patient information, alerting the physician that the patient currently takes a blood thinner." ²⁶

Metadata Promotes Data Quality

Data quality intuitively can be measured by how well insurers can cross-reference, analyse, interpret and capitalize on the vast amounts of data collected. As seen above, robust metadata is a powerful tool in creating connections between data (and explicitly enhancing its quality and usability). Metadata is key to organizing the data collected, reducing confusion, and enhancing the usability and cross referencing ability of the data.

"Without structural metadata, both descriptive metadata and, ultimately, the data content of the transaction, have no context."²⁷

"Good metadata management can lead to good data quality since having and relying on the metadata can identify poor data / incorrect data / missing data. Also, having good metadata shows an understanding of data management and shows that the organization is committed to good data – hence an improvement in data quality almost always follows." ²⁸

In the simplest sense, defining exactly how data is to be recorded within a database (format, character types, permissible values, etc.) is crucial to reducing the amount of time needed to scrub the data.

Metadata in the P&C Insurance Context

Let's expand the examples of metadata to an insurance context. One traditional compilation of metadata which we often take part in are the ISO and NCCI statistical plans. Statistical Plans are developed with a goal of providing a data base of homogeneous experience for comparable policies that fulfils both a regulatory need and a business need to correctly price the insurance product.

- For regulatory purposes, the statistical plans collect historical insurance company experience by state, by class, and by coverage. The minimum requirements for the regulatory needs are included in the National Association of Insurance Commissioners (NAIC) Statistical Handbook of Data Available to Insurance Regulators.
- For the business purpose of pricing the insurance product, Statistical Plans go beyond the regulatory mandated data elements (or standard data elements) and collect both additional detail within the standard data elements and additional or new data elements to perform research and development to better refine the rating or classification of insurance policies and to provide advisory prospective loss costs. By aggregating the data together from many insurers,

²⁶ http://library.ahima.org/xpedio/groups/public/documents/ahima/bok1_049357.hcsp?dDocName=bok1_049357

²⁷ http://www.riskandinsurance.com/the-whodunit-of-big-data/

 $^{^{28}\} http://dataqualitypro.com/data-quality-pro-blog/data-quality-through-metadata-strategy-anne-marie-smith$

the resulting ISO data base provides a larger, more credible data base than any one insurer can do alone.²⁹

The Statistical Plans are essentially the rules that described how the data is to be captured, including the format and codes or values that to be used. These descriptions and instructions help to form the metadata, definitions and business (and data quality) rules for the data.

Furthermore, metadata is crucial for insurers to capitalize on the advantages of big data. Robust metadata can be used to make connections between disparate data sources for use in analytics.

Metadata in Predictive Analytics

"Metadata within data infrastructures enables us to locate and combine data, and to analyze its lifecycle and history. Consider, for instance, the addition of weather, geographical and social media data to the daily sales figures for a retail chain. It is easy to conceive that correlations with peaks and troughs in sales could be elicited: perhaps with good weather, word-of-mouth trends or road accessibility. With sufficient data, some of these events might even be found to be predictive of sales."³⁰

Metadata to Detect Insurance Fraud

According to the Insurance Information Institutes, Property / Casualty insurance fraud amounts to about \$32 Billion a year. From quote, policy issuance and even claims reporting, more and more insurance transactions are conducted online and through various mobile applications. Each internetenabled device, which can include computers, tablets and cell phones, has metadata associated with it. The metadata can include which email accounts are associated with it, what the IP address is, etc. This metadata can be used to identify whether this electronic device is connected to an account (email) or other device with a known history of fraud. Technology is now available that can use this metadata, check for anomalies, attributes and activity levels of the device to determine a "Reputation" of the device. This device-based intelligence, gathered through metadata, is helping insurers identify fraud at the point of first contact.

Metadata in Property Risk Modeling

Without accurate address data, you can't property risk model without making assumptions about the risk characteristics of a particular property. Those assumptions usually take the form of using the average or modal value for a particular characteristic at the ZIP code level.

The Actuary's Role in Metadata

Actuaries are key individuals in developing the enterprise metadata and helping to promulgate its

²⁹ CAS STUDY NOTE: IS0 STATISTICAL PLANS, by Virginia R. Prevosto, FCAS, MAAA

³⁰ http://www.forbes.com/sites/edddumbill/2013/12/31/big-data-variety-means-that-metadata-matters/

usage and acceptance within the organization. Perhaps most importantly, actuaries play a role in the definition of the business metadata (e.g., providing clear definitions to the data to avoid confusion and misunderstanding across functional users or business units).

As an example, the definition of "loss" may significantly vary among business users depending on the context the term is used in. For a data field called "loss," it is important to exactly define what comprises these values. For instance, if we're talking in the context of Workers' Compensation, loss may include indemnity only loss, medical only loss, allocated loss adjustment expense, unallocated loss adjustment expense, etc. For a user in finance computing a "loss ratio" using the loss field as the numerator, potential confusion and erroneous indications may result without a clear context of what defines the loss field. Metadata is a key governance solution to avoid this type of confusion and provide business users a clear context for the data utilized in analysis or decision making.

Actuarial IQ has a well-defined starting list of questions an actuary can ask when helping IT create the enterprise metadata. The following is borrowed directly from the paper:³¹

- Are all the data elements listed?
- Has the source of each data element been provided?
- Is there a special value that is used to indicate missing data?
- Are there any transformations being applied to the data? (Note: data clean up such as filling in missing values should be considered data transformation).
- Have the contents and use of each data element been properly described?
- Have all the categorical values of each data element been properly described?
- In the case of numeric data, has the range of possible values for each data element been provided?
- Has the valuation date of all data been provided?
- Has a schedule of planned updates to the data been provided?
- Has the business process changed during the experience period?
- Have any of the data definitions changed during the experience period?

A good place to start is with our own actuarial work product. In many instances, we may produce or maintain databases that underlie our analyses. How well documented are these systems? How well understood are the sources that feed the actuarial systems? Once the actuarial systems are understood, one can start to drill back into the source systems. Along the way, missing metadata can be identified. The benefits and costs of producing the metadata can be weighed and ownership could be assigned.

³¹ https://www.casact.org/pubs/forum/08wforum/actuarialIQ.pdf

Sources for further information

Publishe				
Source	Author	r	Link (if applicable)	
Data Quality: The Field Guide	Thomas Redman	Digital Press, 1st Edition, 2011		
Actuarial I.Q. (Information Quality)	CAS Data Management Educational Materials Working Party	CAS	https://www.casact.org/pubs/forum/08wforum/actuarialIQ.pd f	
Metadata Management for Holistic Data Governance	Informatica Whitepaper	Informatica Whitepaper	https://www.informatica.com/content/dam/informatica- com/global/amer/us/collateral/white-paper/metadata- management-data-governance_white-paper_2163.pdf	
Understanding Metadata	NISO			
Survey of Data Management and Data Quality Texts	CAS Data Management Educational Materials Working Party	CAS	https://www.casact.org/pubs/forum/07wforum/07w279.pdf	
Actuarial Data Management In A High- Volume Transactional Processing Environment	Joseph Strube and Bryant Russell, Ph.D., ACAS, MAAA	CAS	https://www.casact.org/pubs/forum/05wforum/05wf274.pdf	
Risk Management and Insurance Operations - CAS Course 1 For Preparation For Exam CA1 - Assignment 16 - Actuarial Data Management	Casualty Actuarial Society	The Institutes		
Actuarial Standard of Practice 23 - Data Quality	Actuarial Standards Board (ASB)	ASB	http://www.actuarialstandardsboard.org/wp- content/uploads/2014/02/asop023_141.pdf	

CAS Data and Technology Working Party Report

Databases

Why an Actuary Needs to Know about Databases

Insurance organizations are increasingly becoming data driven. While most insurance companies have had Chief Actuaries for some time, the advent of the Chief Data Officer is a recent phenomenon, and increasingly common. In some instances, the Chief Actuary and Chief Data Officer are one-and-the- same. In other instances, Actuaries are leading the Data Management and Governance processes. Actuaries have traditionally been the prime users of data for analytical purposes at insurers, and as primary users of data, Actuaries should be at the forefront of the new data-driven culture. As this data culture grows, the discussions about data will become increasingly technical. To be a valuable and influential participant, it will be important that the actuary is reasonably fluent in data terminology, and knowledgeable about how data is stored. Databases come in many forms that vary by intended use; this paper will provide an overview of some of the more common forms.

APPLICATION DATABASES VS DATA WAREHOUSES



What's the difference between an application database (ADB) and a data warehouse (DW)? They are both databases, but they serve very different purposes.

As an actuary, you may rely on reports that will be as of yesterday, even though you know the data is already in "the system" today. Why can't your report include the data you know is there?

One reason is likely to be that your application data is in the ADB, but your report references the DW which may only be updated nightly. So why bother to have a DW when the data already exists in the ADB?

ADBs are designed to store transactional data efficiently and allow IT to efficiently add new transactions and update existing transaction data, while DWs are designed to serve as a source for reports and data analysis. Unfortunately, the underlying structures that optimize these functions are significantly different. For example, the following table illustrates a subset of data you might download in a report on current policyholders:

Insured ID	Insured Name	Premium	State Name	City Name
1	Aaron	90,000	Illinois	Chicago
2	Brian	30,000	Illinois	Chicago
3	Chris	30,000	Illinois	Chicago
4	David	40,000	Illinois	Springfield
5	Eddie	80,000	Wisconsin	Madison
6	Frank	10,000	Wisconsin	Madison
7	Gary	20,000	Wisconsin	Milwaukee
8	Henry	20,000	New York	New York
9	Isaac	40,000	New York	New York
10	John	30,000	New York	Albany
11	Kevin	90,000	New York	Albany

This data is easy to understand and once in Excel it is ready to serve as the source for a PivotTable or aggregating functions like SUMIFS(), as you might want to see the Premium at a state or city level. The format of this data is more likely to appear in DWs.

Data Table Normalization

The data above would likely come from an ADB, where it would probably be formatted as follows.

Insured ID	Insured Name	Premium	City ID
1	Aaron	90,000	1
2	Brian	30,000	1
3	Chris	30,000	1
4	David	40,000	2
5	Eddie	80,000	3
6	Frank	10,000	3
7	Gary	20,000	4
8	Henry	20,000	5
9	Isaac	40,000	5
10	John	30,000	6
11	Kevin	90,000	6

Databases

State ID	State Name	City ID	State ID	City Name
1	Illinois	1	1	Chicago
2	Wisconsin	2	1	Springfield
3	New York	3	2	Madison
		4	2	Milwaukee
		5	3	New York
		6	3	Albany

The data has been "normalized" for use in an ADB. When data is normalized, it is reorganized so that it is as parsimonious as possible. As you can see, the data above was reorganized from the original insured table and has been split into an insured table, a state table, and a city table, and ID fields have been included. Since the City ID in the insured table maps the City Name, and the State ID in the city table maps the State Name, all the information of the original insured table is preserved in the normalized data tables.

But why would normalized data tables be of use to IT in an ADB? One reason is to note that the addition of new data only requires a City ID instead of a State Name and City Name, since the City ID encapsulates both. Since new policyholders are added to the ADB more frequently than cities or states, less work is required to upload the same level of information.

Another reason to normalize data is to easily update entries across multiple tables using "primary keys," which have unique entries for every record in a table. For instance, let's say an insured's name was entered incorrectly and it must be updated in every table in the ADB. If the data tables are normalized, the insured's name can be corrected in the insured table, and any other table that references the insured's name does so via the Insured ID, so the update effectively takes place across all data tables. Otherwise, the insured's name would have to be tracked down across all the various data tables, which could be very time consuming—and imagine the problem when two insureds happen to have the same name! By using IDs as primary keys, IT is able to uniquely identify relevant data and update it efficiently.

In addition to making it easy for IT to make update data, tables with unique entries can be used to ensure data is entered consistently into the system. In order to prevent users from manually entering "New York City" instead of "New York," IT can reference the city table to validate the data before it goes into the system.

So if there's a good reason to normalize data tables in the ADB, why bother to de-normalize it in the DW? Why doesn't I'T simply provide actuaries with reports directly from the ADB instead of reformatting the data in the DW?



Why bother with a Data Warehouse?

One reason is that a DW can get data from several ADBs associated with different business functions, like policy administration and claims handling. Furthermore, single business functions can rely on several ADBs that have evolved separately to handle specific lines of business or regulatory requirements. Extracting data from these ADBs with varying data formats and transferring it to a DW dramatically simplifies the development of reports for end users.

Structured Query Language (SQL)

Even if a DW represents data from a single ADB, the process of de-normalizing data makes it much easier for IT to develop reports. For instance, let's say you wanted to get the insured data presented in the first, de-normalized table. Since data in relational databases¹ is most often manipulated using Structured Query Language (SQL), the SQL statement written to query the table, "Denormalized," would look as follows:

SELECT [Insured Name], [Premium], [State Name], [City Name]

FROM Denormalized

Even if you're not familiar with SQL, the statement above is pretty straightforward. Now let's look at what the statement would look like to get the same data from three normalized tables, "Insureds," "States," and "Cities."

SELECT i. [Insured Name], i. [Premium], s. [State Name], c. [City Name]

FROM Insureds i

¹ Relational databases (like Microsoft Access and SQL Server) are based on tables ("relations") that are linked by their primary keys. While there are databases that make use of different formats, you can usually assume that when IT talks about a database, it will be a relational database unless they specifically qualify it.

LEFT JOIN Cities c ON i.[City ID] = c.[City ID] LEFT JOIN State s ON c.[State ID] = s.[State ID]

By referencing normalized data tables, the query must explicitly state how the tables are related through JOIN statements, and also identify the tables that contain the various field names. And while the statement above uses LEFT JOIN statements, there are also RIGHT JOIN, INNER JOIN, OUTER JOIN, and FULL JOIN statements that might be appropriate depending on the nature of the data. You could think of the transformation of data to denormalized tables as a way of baking in the joins so that downstream queries don't need to deal with them. As the SQL statements increase in complexity, the value of referencing denormalized tables becomes clear².

Extract, Transform, and Load

The process of migrating data from ADBs to a DW is often referred to as ETL, which stands for Extract, Transform, and Load.

Extracting data from multiple sources can be challenging when those sources are not relational databases, and therefore require methods beyond traditional SQL statements. The extraction process can include a validation step that halts the process unless the data conforms to certain standards, preventing complications further downstream.

Transforming data can involve a number of adjustments, including the aggregation of data. Since the reports generated from the DW may not require the level of detail that exists in an ADB, performance could be enhanced by aggregating data during the transfer. For instance, a DW used by actuaries to evaluate experience by territory might only need data aggregate by ZIP code, county, or state, rather than the exact address of each insured. While street names and addresses would be lost in the course of aggregating the data, reports referencing this data in the DW would run faster. Another alternative would be for IT to preserve the detail for the data in the DW and implement an index, which keeps track of groups of records in a table, enabling more efficient retrieval of specific records. If indexes were implemented in an ADB, they would need to be maintained and updated with the addition of each new record, and if records are added more often than they are read, the cost in computing resources would probably outweigh the benefit. Other manipulations that may occur when moving data from an ADB to a DW include

- Adding accounting periods
- Calculating unearned premium reserves and earned premium
- Assigning loss development months for the creation of loss triangles
- Calculating reinsurance premium and loss from direct and assumed business
- Mapping premium and loss to lines of business defined by various regulatory regimes
- Mapping data to general ledger codes
- General scrubbing of data

The final step, Loading, is (hopefully) executed in such a way as to leave an audit trail, so that any data that looks odd can be traced back to the source ADB and verified. The loading process must also determine which of the preexisting data to overwrite, update, or leave alone. This step can be challenging, for what if an underwriter moves from California to New York, and a report aggregates

² For a useful reference on SQL statements, visit: <u>http://www.w3schools.com/sql/default.asp</u>

premium written by underwriter location? Should the underwriter's premium be re-classified as New York premium, so that the current reports won't be consistent with previous reports? One solution to this Slow Changing Dimension (SCD) issue is to duplicate the underwriter record and add fields that indicate the start date and end date for when that record is valid. This way the current reports are accurate but users can still generate reports as of earlier dates that tie to the original versions.

Levels of Normalization

The data formats discussed so far have been categorized as "normalized" or "denormalized," which is actually a convenient simplification. Data tables can go through increasing stages of normalization, with almost all of them in a state of at least "first normal form," or "1NF," and most being at "third normal form" or "3NF."³ Meanwhile, denormalized tables used in DWs are referred to as "dimension tables" that contain various combinations of categories used to evaluate data held in "fact tables." DWs with more denormalized dimension tables can be described as having "star schemas" whereas less denormalized DWs with less denormalized dimensions tables can be described as having "snowflake schemas." Whatever the structure of your ADBs and DWs, the goal is to take data from an ADB that's optimized to be regularly updated with new data, and move it to a DW that's optimized to provide periodic reports. If you're interested in learning some of the more technical aspects of understanding and interacting with databases and data warehouses, the Microsoft Virtual Academy has tutorial videos intended to assist students preparing to take exams necessary for certification in SQL Server. You can access these videos here: https://mva.microsoft.com/colleges/mcsa-sql.

Development of a Data Warehouse / Business Information System (DW/BI)

For a DW / BI system to be successful, the business community must accept it. For them to accept it, information needs to be presented consistently, be easily accessible, timely, secure, authoritative and trustworthy.⁴ A transaction schema focused on business processes and related measureable events at the most granular level allows for maximum flexibility when data is extracted for analysis. For insurance companies, core business processes could include policy issuance, premium collection and claim processing. A claim processing measurable event could include a claim payment or reserve change.

Facilitating a good system requires a team effort and should include Data Governance, IT, actuarial and other heavy users of analytic data. Actuaries need to be somewhat familiar with IT terminology to be effective members of the design team.

Data Architecture Terminology

Design features described in previous sections are again described here with a few common data warehouse architecture terms. 5

• **Star Schemas** – refers to the architecture of a dimensional model implemented in a relational database management system. It includes a fact table at the core surrounded by multiple dimension tables joined by keys in a star-like formation.

³ There are higher order normal forms, but 3NF is usually sufficient for database administration purposes.

⁴ [2] Kimball R, Ross M. 2013. The Data Warehouse Toolkit. Third Edition. Indianapolis (IN): John Wiley & Sons. Chapter 1, Goals of Data Warehousing and Business Intelligence; p. 3-4.

⁵ [2] Kimball R, Ross M. 2013. The Data Warehouse Toolkit. Third Edition. Indianapolis (IN): John Wiley & Sons. Chapter 1, Dimensional Modeling Introduction; p. 7-18.

Databases



• Snowflake Schema - represents a dimensional model which is also composed of a central fact table and a set of constituent dimension tables which are further normalized into subdimension tables. Because snowflake schemas are normalized, they are easier to maintain, but harder to query.



• Online Analytical Processing (OLAP) cubes – presentation area containing aggregations and precalculated summary tables. Data is loaded from the tables after additional processing.



- Fact Table stores the performance measurements resulting from an organization's business process events. Most useful facts are numeric, additive and continuously valued. A non-additive fact would be rate per unit of exposure. Semi-additive facts like claim reserve balances cannot be summed across the time dimension. These tables consume the most storage. Potential facts might include premium in a policy transaction Fact Table; claim dollars in a claim transaction Fact Table.
- **Grain** level of detail in each Fact Table row; three categories include:

- Transaction like reserve changes
- **Periodic snapshot** like triangles
- Accumulating snapshot like claim reserve balances
- **Dimension Table** integral companion to a fact table; contains the textual context associated with a business process measurement event; often have many columns or attributes. Source of query constraints, groupings and report labels. Good practice is to minimize the use of codes in dimensional table for the sake of consistency and clarity across business processes. Potential Dimensions applicable might include policy effective date, policyholder, coverage, covered item, claimant, date of loss.
- **Keys** join fact with dimension table. Generated when the fact table is created as sequential integers.
 - **Primary Key –** a key in a relational database that is unique for each record in a table
 - Foreign Key a field in one table that uniquely identifies a row in another table
 - Natural Key a key that uses its naturally occurring value as its unique identifier (e.g. Telephone Number)
 - **Surrogate Key** a key that has to be created to uniquely identify a row in a table (e.g. Policy Number)

Steps in Dimensional Design Process for Insurance Company⁶

- Select the business process (e.g. ratemaking)
- Declare the grain the lowest level of detail that will be stored (e.g. coverage)
- Identify the facts the quantitative data that will be measured (e.g. premium)
- Identify the dimensions the attributes of the facts in a dimensional database (e.g. state)

One of the collaboration tools commonly used in a Data Warehouse development process is a bus matrix.⁷ The matrix is simply a table with the core business processes as rows and core dimensions as columns. It is a useful as a communication and documentation tool for DW / BI team participants. Below is an example:⁸

	Date	Policyholder	Covered Item	Coverage	Employee	Policy	Claim	Claimant
Policy Transactions	Х	Х	Х	Х	Х	Х		
Premium Snapshot	Х	Х	Х	Х	Х	Х		
Claim Transactions	Х	Х	Х	Х	Х	Х	Х	Х

⁶ Kimball R, Ross M. 2013. The Data Warehouse Toolkit. Third Edition. Indianapolis (IN): John Wiley & Sons.

⁷ Kimball R, Ross M. 2013. The Data Warehouse Toolkit. Third Edition. Indianapolis (IN): John Wiley & Sons. Chapter 4, Enterprise Data Warehouse Bus Architecture; p. 123-130.

⁸ Kimball R, Ross M. 2013. The Data Warehouse Toolkit. Third Edition. Indianapolis (IN): John Wiley & Sons. p. 389.

Databases

Conformed fact and dimension tables keep naming conventions and defined calculations consistent across all departments. DW/BI system should:

- Deliver data that is understandable to the business users.
- Deliver fast query performance.

OTHER TYPES OF DATABASES

Database types are dictated by their intended purpose. So far we have focused on application databases and Data Warehouses. Operational applications handle data input one transaction at a time; data warehouse business intelligence (DW/BI) systems maintain historical data for analysis by the business community. Other database types include:

Flat and Wide (F&W)

Flat and wide (or F&W) is a common data type given to actuaries from IT, usually in an Excel spreadsheet. It is denormalized and typically relatively small.

Operational Data Store (ODS)

An operational data store (or "ODS") is a database designed to integrate data from multiple sources for additional operations on the data. Unlike a master data store, the data is not passed back to operational systems. It may be passed for further operations and to the data warehouse for reporting.

Because the data originate from multiple sources, the integration often involves cleaning, resolving redundancy and checking against business rules for integrity. An ODS is usually designed to contain low-level or atomic (indivisible) data (such as transactions and prices) with limited history that is captured "real time" or "near real time" as opposed to the much greater volumes of data stored in the data warehouse generally on a less-frequent basis.

Columnar Databases

A column-oriented DBMS (or columnar database) is a database management system (DBMS) that stores data tables as columns rather than as rows. Practical use of a column store versus a row store differs little in the relational DBMS world. Both columnar and row databases use traditional database languages like SQL to load data and perform queries. Both row and columnar databases can become the backbone in a system to serve data for common ETL and data visualization tools. However, by storing data in columns rather than rows, the database can more precisely access the data it needs to answer a query rather than scanning and discarding unwanted data in rows. Query performance is often increased⁹ as a result, particularly in very large data sets.

Another benefit of columnar storage is compression efficiency.¹⁰ It is well known that a row of similar data, dates for example, can be compressed more efficiently than disparate data across rows. It's for this reason, columnar databases are well-known for minimizing storage and reducing the amount of I/O needed to read data and answer a query. Columnar databases most often are paired with Massively Parallel Processing (MPP) capability to allow for it to share the analytical workload across a cluster. They may also leverage Hadoop MPP capability for this purpose.

⁹ Ventana; et al. (2011). "Ins and Outs of Columnar Databases."

¹⁰ Ventana; et al. (2011). "Ins and Outs of Columnar Databases."

NoSQL Databases

A NoSQL (originally referring to "non SQL" or "non-relational")¹¹ database provides a mechanism for storage and retrieval of data which is modelled in means other than the tabular relations used in relational databases. Such databases have existed since the late 1960s, but did not obtain the "NoSQL" moniker until a surge of popularity in the early twenty-first century,¹² triggered by the needs of Web 2.0 companies such as Facebook¹³, Google¹⁴ and Amazon.com¹⁵. NoSQL databases are increasingly used in big data and real-time web applications.¹⁶ NoSQL systems are also sometimes called "Not only SQL"¹⁷ to emphasize that they may support SQL-like query languages.¹⁸

Graph Databases

In computing, a graph database is a database that uses graph structures for semantic queries with nodes, edges and properties to represent and store data.

Most graph databases are NoSQL in nature and store their data in a key-value store or documentoriented database. In general terms, they can be considered to be key-value databases with the additional relationship concept added. Relationships allow the values in the store to be related to each other in a free form way, as opposed to traditional relational databases where the relationships are defined within the data itself. These relationships allow complex hierarchies to be quickly traversed, addressing one of the more common performance problems found in traditional keyvalue stores. Most graph databases also add the concept of tags or properties, which are essentially relationships lacking a pointer to another document.

¹¹ NoSQL DEFINITION: Next Generation Databases mostly addressing some of the points: being non-relational, distributed, open-source and horizontally scalable. See: <u>http://nosql-database.org/.</u>

¹² Leavitt, Neal (2010). "<u>Will NoSQL Databases Live Up to Their Promise?</u>" (PDF). IEEE Computer.

¹³ Mohan, C. (2013). <u>History Repeats Itself: Sensible and NonsenSQL Aspects of the NoSQL Hoopla (PDF)</u>. Proc. 16th Int'l Conf. on Extending Database Technology.

¹⁴ "Dynamo Clones and Big Tables" http://www.eventbrite.com/e/nosql-meetup-tickets-341739151.

¹⁵ Garling, Caleb (2012). "<u>Amazon helped start the "NoSQL" movement</u>." Wired Magazine.

¹⁶ "<u>RDBMS dominate the database market, but NoSQL systems are catching up</u>". DB-Engines.com. 21 Nov 2013. Retrieved 24 Nov 2013.

¹⁷ "NoSQL (Not Only SQL)". NoSQL database, also called Not Only SQL

¹⁸ Fowler, Martin. "<u>NosqlDefinition</u>." Many advocates of NoSQL say that it does not mean a "no" to SQL, rather it means Not Only SQL.

Databases

Unstructured Databases

To this point we have discussed the databases most familiar to actuaries. In the era of "Big Data", unstructured data is starting to play a substantial role in the world of data and analytics. Unstructured data includes any data whose structure is not compatible with the data warehousing structures discussed above. The incompatibility may be due to:

- Volume: databases too lard for data warehousing (e.g. weather data)
- Variety: data with formats incompatible with traditional warehousing (e.g. images)
- Velocity: data generated and delivered too frequently to be organized into a warehouse (e.g. telematics)

XML and geospatial data are often called "semi-structured" data as they contain their own inherent structure, but as their structure requires transformation before it can be combined with fully structured data.

	Organization	
Source Description	Name	Web Link
Introduction to Databases, a	Platform by	https://www.coursera.org/course/db
free online course that covers	Coursera, class	
database design and the use of	taught by Jennifer	
database management systems	Widom of Stanford	
for applications.	University	
SQL Tutorial that teaches you	W3Schools.com	http://www.w3schools.com/sql/default.a
how to use SQL to access and		sp
manipulate data in: MySQL,		
SQL Server, Access, Oracle,		
Sybase, ADB2, and other		
database systems.		
SQL Server Certification and	Microsoft Virtual	https://mva.microsoft.com/colleges/mcs
Training videos for	Academy	<u>a-sql</u>
professionals working toward		
earning a Microsoft Certified		
Solutions Analyst (MCSA):		
SQL Server certification		
Kimball R, Ross M. 2013. The	John Wiley & Sons.	http://www.kimballgroup.com/
Data Warehouse Toolkit.		
Third Edition. Indianapolis		
(IN)		

	Organization	
Source Description	Name	Web Link
Wikipedia	Various	 <u>https://en.wikipedia.org/wiki/Op</u><u>erational_data_store</u> <u>https://en.wikipedia.org/wiki/Col</u><u>umn-oriented_DBMS</u> <u>https://en.wikipedia.org/wiki/Gra</u><u>ph_database</u> <u>https://en.wikipedia.org/wiki/No</u><u>SQL</u> <u>https://en.wikipedia.org/wiki/Uns</u><u>tructured_data</u>

Bornhuetter-Ferguson Initial Expected Loss Ratio Working Party Paper

This paper is the culmination of effort of the working party over the span of several years. Listed below are the working party members that were part of the group during the survey that was created and conducted, during the initial presentation at the Casualty Loss Reserve Seminar, during the subsequent presentation at the CAS Annual Meeting or during the writing of the this paper. Special thanks to Lynne Bloom, Kelly Moore and Chandu Patel for being instrumental in bringing this paper to completion.

Nancy Arico Lynne Bloom Aaron Hillebrandt Bertram Horowitz Dennis Lange Kelly Moore Douglas Nation Chandu Patel

INTRODUCTION

"The Actuary and IBNR" was published in 1972 by Ronald Bornhuetter and Ronald Ferguson [1]. The methodology from this paper has exploded into a veritably universal methodology used by actuaries and commonly referred to as the "Born Ferg" or "BF" method. The technique and its application are included in the syllabus for the CAS actuarial exams and the use of the technique is pervasive in both the reserving and pricing worlds.

The method involves the selection of an "Initial Expected Loss Ratio" or "IELR" for which the selection criteria varies greatly and a great degree of latitude is permitted to the practitioner for "actuarial judgment." Given the widespread use of this method and its impact on financial reporting, the Bornhuetter Ferguson Initial Expected Loss Ratio Working Party set out to glean an understanding of general industry practices surrounding the selection of the IELR used in this method.

A survey was conducted across the CAS membership and the results of that survey are presented in this paper.

Along with the survey, the paper also explores several alternative methods to selecting the initial expected loss ratios, their relative strengths and weaknesses and their relative predictive value when applied to historical data. Carried reserves versus the outcome of several alternative methods for selecting the IELR are also explored to determine the effectiveness of industry practices.

The Basic BF Method

The Bornhuetter-Ferguson ("BF") expected loss projection method based on reported loss data relies on the assumption that remaining unreported losses are a function of the total expected losses rather than a function of currently reported losses. The expected losses used in this analysis are generally based on a review of previous accident year ultimate loss ratios and the company's business plan. The expected losses are multiplied by the unreported percentage to produce expected unreported losses at a point in time. The unreported percentage is calculated as one minus the reciprocal of the selected cumulative reported loss development factor ("LDF") for the segment under review. Finally, the expected unreported losses are added to the current reported losses to produce the estimated ultimate losses.

The calculations underlying the Bornhuetter-Ferguson expected loss projection method based on paid loss data are similar to the reported Bornhuetter-Ferguson calculations with the exception that paid losses and unpaid percentages replace reported losses and unreported percentages.

Alternative Choices for Initial Expected Loss Ratios

A critical assumption within the framework of the Bornhuetter-Ferguson method is the Initial Expected Loss Ratio ("IELR"). The IELR can be determined using several methods. The most frequently used methods for determining IELR for long-tailed lines are as follows:

- Pricing Loss Ratio
- Prior Analysis Ultimate Loss Ratios
- Industry Aggregates
- Cape Cod
- Prior Accident Years' projected loss ratios
- Prior Accident Years' loss ratios adjusted for rate changes and trends
- Judgment

METHODS

Pricing or Plan Loss Ratio

This method uses a pricing target loss ratio from the pricing actuary or a plan loss ratio from the company's financial plan as the IELR.

A refinement to this method is to adjust the target or plan loss ratio for the difference between
Bornhuetter-Ferguson Initial Expected Loss Ratio Working Party Paper

actual and target or planned pricing. Companies often have price monitoring systems that monitor actual price level compared to target price level. Actuaries can also track actual rate changes implemented compared to planned rate changes. For example, if the plan loss ratio is 60% and included a planned earned price change of 5%, but the company actually achieved an earned price change of 3%, then the IELR would be calculated as:

Other adjustments to the plan assumptions could be reflected as well. For example, if the actuary has an updated estimate of loss trend compared to the loss trend assumptions used in the plan, the IELR could be adjusted accordingly. If an operational or regulatory change is implemented that wasn't anticipated in the plan, then the expected impact of that change could be reflected in the IELR as an adjustment to the plan loss ratio.

Below are some advantages and disadvantages of this method. Similarly, advantages and disadvantages will be listed for each method in subsequent sections.

Advantages

- It is straightforward.
- It will be generally understood and accepted by management and staff in other departments.
- It includes information from multiple departments.

Disadvantages

• Pricing targets and plan loss ratios can be aspirational and therefore may not reflect the true expected loss ratio;

• Plan loss ratios are often derived by subtracting expense and profit provisions from a target combined ratio. The target combined ratio often reflects optimistic estimates for the impact of rate/pricing changes and underwriting actions. If the rate and underwriting effects do not materialize, the plan ratio can be significantly understated.

Rate Indication Adjusted for Rate Changes and Trends

Another way to use a pricing loss ratio as the basis for the IELR is to start with the indicated loss ratio from a rate indication/pricing study and adjust for rate changes and loss trend from the prospective proposed effective period to the appropriate accident year. An example of this method is shown in Exhibit 1 of Appendix A. In this example, in the latest rate indication the pricing actuary has projected the loss ratio for policies effective from 7/1/2017 through 6/30/2018, and we start with that loss ratio to estimate the IELR for accident year 2016. First, we adjust for the net loss and

Bornhuetter-Ferguson Initial Expected Loss Ratio Working Party Paper

premium trend from the projection period in the rate indication back to accident year 2016. Next, we adjust for any rate changes reflected in the rate indication that haven't been fully earned in accident year 2016. In the example, there were two such rate changes, one effective on 7/1/2015 and one effective on 7/1/2016.

Other adjustments to the projected loss ratio from the rate indication could be reflected as necessary, for example, the impact of operational, regulatory, and/or underwriting changes.

Advantages

- It leverages the work already done by the pricing actuary.
- It reflects the expected impact of trend and rate changes.

• Indicated loss ratios from rate indications have generally already been smoothed for large losses and catastrophes.

Disadvantages

• Rate indications are often done at a lower level of detail than reserve analyses, for example by state or business unit, so this method may require aggregation before use.

Prior Analysis Ultimate Loss Ratio

Another method that can be used is to select the ultimate loss ratios from the prior reserve analysis as the IELRs in the current reserve analysis. For example, if the company does semiannual reserve reviews, we would use the ultimate loss ratio for accident year 2015 from the 6/30/2016 reserve review as the IELR for accident year 2015 in the 12/31/2016 reserve review.

Advantages

- It is straightforward.
- It leverages work already done to arrive at a best estimate of the ultimate loss ratios.

Disadvantages

• This method will increase the responsiveness of the BF method to the extent that the ultimates from the prior reserve review reflect the actual loss emergence, which may be a disadvantage in some cases, for example when an accident year has experienced unusually high or low large losses.

Industry Aggregates

The IELR may be based on industry aggregate loss ratios. Sources for industry results include the following:

• A.M. Best

- NCCI
- SNL
- ISO
- Internal benchmarks

This approach may be especially appropriate when a company is writing a new type of business and doesn't have the historical data necessary to use many of the other methods, or when a company has a small book of business and doesn't have credible historical data.

Advantages

• It reflects the whole industry, so results are based on a credible volume of data.

• Industry results reflect the aggregate impacts of price changes, loss trend, and the underwriting cycle.

Disadvantages

• There is a lag in receiving industry results and the selected IELR is usually based on dated information;

• It doesn't reflect factors specific to the company's book of business that can impact the loss ratio, such as pricing, underwriting, and mix of business.

Prior Accident Years

Another method is to select an IELR based on the loss ratios for prior accident years for the same book of business. Averages of the loss ratios from several years can be used to smooth out or exclude abnormal variations in the results.

Advantages

- It is straightforward.
- It is easy to explain.

Disadvantages

• It doesn't reflect changes in pricing, loss trend, and underwriting that can impact the loss ratio.

Prior Accident Years Adjusted for Rate Changes and Trends

In this method, the IELR is based on estimates for prior accident years adjusted for rate changes and loss trends. Examples of this method are shown in Appendix A, Exhibit 4, which uses on-level earned premiums and loss ratios, and Appendix A, Exhibit 5, which uses exposures and pure premiums. In both examples, we start with ultimate losses from the prior reserve review for accident years 2007 through 2015 and use them to estimate the accident year 2016 IELR.

This method is similar to the rate indication adjusted for rate changes and trends method in that both start with an indicated loss ratio and adjust for rate changes and trends to get the IELR for the accident year in question. However, this method starts with the estimated loss ratios for prior accident years and adjusts forward to the appropriate accident year, while the rate indication method starts with an indicated loss ratio for a prospective proposed effective period from a rate indication and adjusts back to the appropriate accident year.

In Exhibit 4, we calculate the ultimate loss ratios from the prior reserve review by dividing the ultimate losses from the prior review by the earned premiums. Then, we adjust each of the loss ratios for accident years 2007 through 2015 to the accident year 2016 level. The on-level premium factors are calculated based on the rate change history and the loss trend factors are calculated based on selected annual loss trends. We apply the on-level and loss trend adjustments to the loss ratios for accident years 2007 through 2015 to arrive at various estimates of the accident year 2016 IELR. Then, we calculate various averages of the indicated IELRs and make a selection.

Exhibit 5 shows a similar calculation except using pure premiums instead of loss ratios. We calculate the ultimate pure premiums from the prior reserve review by dividing the ultimate losses from the prior review by the earned exposures. We then apply the pure premium trend adjustments to the pure premiums for accident years 2007 through 2015 to arrive at various estimates of the accident year 2016 pure premiums. Next, we calculate various averages of the indicated pure premiums and select an expected pure premium for accident year 2016. Finally, we convert the selected accident year 2016 expected pure premium to an expected loss ratio.

Advantages

• It reflects the expected impact of trend and rate changes.

• By using several accident years and taking averages, random variation in loss results should be smoothed.

Disadvantages

• It requires either rate change or exposure information, which in practice is sometimes not available.

Cape Cod

The Cape Cod or Stanard-Buhlmann method (Stanard [2]) calculates the expected loss ratio based

Bornhuetter-Ferguson Initial Expected Loss Ratio Working Party Paper

on the reported loss experience for all accident years. First, reported losses are trended and earned premiums are adjusted for rate changes such that they are at an equivalent point of evaluation. Then, the "used-up" or "reported" on-level earned premiums are calculated as the on-level earned premiums times the percent of losses expected to be reported, which is equivalent to the on-level earned premiums divided by the cumulative loss development factor. The IELR for each accident year is calculated as the weighted average of the IELR for each year using the "reported" on-level earned premiums as weights.

Gluck [3] introduced a decay factor to the Cape Cod method in order to give more weight to those accident years that are closer in time to the accident year whose IELR is being estimated. This refinement recognizes that the results for more remote accident years are less relevant to estimating an IELR for a given accident year since trend and on-level estimates are not perfect and there may have been changes in the book of business over time due to mix or underwriting changes. The decay factor is between zero and one, with lower factors being more appropriate for books of business with more stable experience and higher factors being more appropriate for books of business with more volatile experience. A decay factor of one results in the original Cape Cod method.

Examples of this method are shown in Appendix A, Exhibit 2, which uses on-level earned premiums and loss ratios, and Appendix A, Exhibit 3, which uses exposures and pure premiums.

This method can also be done ignoring both rate changes and trend in losses under the assumption that pricing changes are reflective of loss changes. Later in this paper, we used the method both ways to demonstrate the impact with industry data.

Advantages

- It uses all the available reported loss experience to develop the IELR.
- It can reflect the expected impact of trend and rate changes.
- It is very responsive to experience.

Disadvantages

• It may require a complete history of either rate change or exposure information, which in practice is sometimes not available.

• Each accident year is treated as similar experience if decay factors are not used; decay factors are difficult to program in practice.

Judgment

The actuary could use judgment to select the IELR, incorporating knowledge of the book of

Bornhuetter-Ferguson Initial Expected Loss Ratio Working Party Paper

business including underwriting and pricing, information on industry and company results in similar types of business and awareness of the underwriting cycle. Similar to industry aggregates, this approach may be especially appropriate for new or small books of business lacking the credible historical data necessary to use many of the other methods.

For books of business with credible experience, judgment is used in the application of the other methods described above, for example, in deciding which methods to use, what adjustments are appropriate and what selections to make when faced with varying indicated IELRs from different methods or various averages.

Advantages

- It doesn't require any specific data.
- It allows the actuary to apply knowledge gained through experience.

Disadvantages

• It may be more difficult to document and support the selection.

THEORETICAL ROBUSTNESS OF THE BF METHOD

The Bornhuetter-Ferguson method is most useful as an alternative to other models for immature accident years. For these immature years, the amounts reported or paid may be small and unstable and therefore not predictive of future development. Therefore, future development is assumed to follow an expected pattern that is supported by more stable historical data or by emerging trends. This method is also useful when changing reporting patterns or payment patterns distort historical development of losses and for lines of business with volatile reporting and payment patterns. For example, it is effective for lines of business such as aviation where a dominant large loss may distort current paid and reported loss experience and render it unusable for the reported and paid loss development methods. It can also be very useful for lines of business with significantly long reporting periods. For example, in high excess casualty occurrence lines of business, paid and reported loss activity may be zero for decades and losses may manifest themselves many years after the policy has been issued. In this instance, the reported and paid loss development methods cannot be applied in any meaningful manner.

SURVEY RESULTS

Major Observations / Conclusions from the BF IELR Survey

A complete list of questions posed and responses received is included in Appendix B. Discussed below are the highlights from the survey.

Extent of Use

Not surprisingly, the BF methodology is used extensively within the industry; over 75% of survey respondents use the BF methodology for all lines of business analyzed. The BF methodology is used to analyze loss as well as ALAE/DCC (often in combination with loss); it is not commonly used to analyze ULAE/AAO. Although the methodology is used extensively, there was a fair amount of negative feedback regarding the misuse of the methodology, particularly in the selection of the IELR:

"In the vein of coming up with a best estimate using all available information, the rationale for using some initial expected loss ratio in the analysis despite information that suggests that initial expected loss ratio was either too high or too low is a flawed approach."

"I do see abuse and unsupported BF selections frequently on the low side as a reviewer."

"Although my decisions are independent, I feel pressure from management, and I can't imagine an actuary working for a client that doesn't."

The testing in the next section of this paper is geared toward addressing possible industry biases.

Choice of Method

For long-tailed lines of business, the most prevalent method for determining the IELR is prior accident year loss ratios adjusted for rate changes and loss trend. Second is the ultimate loss ratio from the prior analysis. Cape Cod is the third most popular, but less than 10% of respondents use it. Within the reinsurance industry, the pricing/plan loss ratio is the most popular, consistent with long-tailed lines of business. For short-tailed lines of business, prior accident year loss ratios adjusted for rate changes and loss trend is most popular.

Other Common Practices

- 1) It is very common to use the BF method to estimate loss ratios for the most recent accident year; for older accident years, the use of BF drops off rapidly.
- 2) There is wide degree of variation in beliefs about whether the IELR should be changed, if new data indicates that a change is necessary based on either higher or lower actual loss experience. A few responders believe that once picked, the IELR should not be changed; on the other hand, the majority of responders believe that it is necessary to change the IELR once the new experience indicates that a change is necessary and over 60% of responders change the selection annually.

- Although most respondents considered their selection of IELR to be independent, for a significant amount of respondents, management plays a role in reviewing/guiding the actuary in the process of selecting the IELR.
- 4) A majority of the respondents do not place minimum boundaries in the selection of the IELR; for example, they keep the IELR the same even if paid or reported loss ratios exceed the previously selected IELR.
- 5) A majority of the respondents used an internal peer review process and/or actual versus expected analysis to test the reasonability of the selected IELR.

INDUSTRY TESTING

The survey gave us a snapshot of what respondents were doing in practice, but we also wanted to understand the methods in the context of financial statements and real world data. There is a great degree of cynicism surrounding the use of the BF methods and the partially judgmental selection of IELR. Using actual reported loss data, we sought to glean what the general industry practice was and how well it was working.

Both the efficacy and accuracy of the method itself are important aspects of our study. Therefore, we tested the industry use of the BF method with the following questions:

- 1. How do actual carried reserves compare to the various forms of BF method?
- 2. How well do the various BF methods compare to hindsight reserves?

To answer these questions, we used Schedule P Data, an industry rate change index for commercial lines (CIAB) and industry claim cost inflation trends (Towers Watson). With the available data we were able to test three forms of the BF method:

- 1. Prior evaluation (using past carried)
- 2. Cape Cod (used with and without inflation and rate change information)
- 3. Trended rate-adjusted loss ratio (using Schedule P carried)

For commercial lines, we had aggregate rate change information and tested Workers Compensation, General Liability (Claims Made and Occurrence), Medical Malpractice (Claims Made and Occurrence) and Commercial Auto. We also tested personal lines (Private Passenger Auto and Homeowners) and Commercial Multiple Peril, but we did not have aggregate rate change information.

Observations on Carried Reserves

Commercial Lines

Commercial lines that we studied included:

Commercial Auto

Workers Compensation

General Liability (Both Occurrence and Claims Made)

Medical Malpractice (Both Occurrence and Claims Made)

Note, we are not including CMP in this aggregation due to the lack of available rate change information.

Findings for All Commercial Lines - Current Accident Year

Actual carried and indicated net loss ratios for the industry for accident year 2012 evaluated as of December 31, 2012 are as follows:



The graph displays the paid and incurred LDF methods, where LDF's are selected based on weighted averages along with an industry tail factor. These LDF's are also used to project the Cape Cod, the trend / rate change adjusted Cape Cod methods and the trended loss ratio BF method on a paid and incurred basis. For accident year 2012, for Commercial business in aggregate, the industry

booked loss ratio is at a level commensurate with trended paid and incurred loss ratio BF methods, higher than the LDF or Cape Cod methods but lower than the trend adjusted Cape Cod methods. In this case, the prior ultimate method would not be applicable since we are evaluating the current accident year.

Conclusion: For accident year (2012), the industry aggregate commercial lines booked net loss ratio most closely matches indications from the trended loss ratio BF method.

Findings for Trended Loss Ratio BF Method - All Accident Years



The graph allows us to focus on the trended loss and DCC ratio BF method for all accident years as of December 31, 2012. Although current carried reserves are close to these methods for the latest accident year and are slightly lower than this method for prior accident years, we can see that initial carried reserves in the 2003 through 2006 period are much higher than this method. It is clear that initial carried amounts reflected more pessimism about loss ratios at the time and may be indicative of the market cycle change and the higher loss ratio experience around the year 2000.

Conclusion: Initial carried reserves in hindsight appear to reflect the market cycle more than the BF indications.



Findings for Cape Cod Paid Adjusted Method - All Accident Years (Loss Plus DCC Ratio)

If we look at just the Paid Cape Cod method adjusted for rate changes and trend across accident years, we can see that the current carried is more optimistic for older accident years and more pessimistic for more recent accident years. This demonstrates that carried reserves lean more toward methods that view accident years separately rather than gravitating toward a long-term average, which is what the Cape Cod method does.

Conclusion: The industry may look at accident year results in more isolation than Cape Cod methods would imply.



Findings for All Methods - All Accident Years

Based on the graph above (methods are done using most recent data), one can see that other than the trended adjusted Cape Cod Methods, the methods and carried reserves are very close until the latest accident year. It is also clear that the initial carried amounts do not move with the methods, implying a tendency not to deviate from a set level of reserves or a tendency not to react to market conditions.

Conclusions: In selecting the initial carried amounts, reactions to market cycles appear to play a more prominent role than actuarial indications.

Findings for DCC to Loss Ratio - Current Accident Year

The previous graphs all consider the projection of losses and DCC together. We also examined the BF method in the context of separate loss and DCC analyses and using the projection of ultimate losses as an exposure base for DCC.



DCC Method results are slightly more erratic than loss plus DCC methods. For the 2012 accident year, actual carried reserves fell close to an incurred Cape Cod method. This makes sense since the ratio of DCC to loss might be expected to be more stable than losses during market changes.

Conclusion: DCC ratio may not be as greatly impacted by the market cycle.



Findings for Gross Loss and DCC Ratio - Current Accident Year

In general, the pattern of methods is very similar to what we observed on a net basis. One difference is that carried reserves are closer to the higher end of the methods on a gross basis, perhaps suggesting that carried reserves are swayed more by balance sheet considerations than exact methodology.

Conclusion: Ceded reserves may not be as impacted by preconceived reserve level expectations.

Findings for Individual Lines of Business - Commercial Lines

The findings for individual lines of business were very similar to the all lines indications. A notable exception is the medium-tailed Commercial Auto line of business, where the market cycle effects are less pronounced. Shown below are the results for Commercial Auto.



Personal Lines

For Personal Lines, we did not have the benefit of a rate change index, but we did have trend information for the two personal lines studied, Personal Auto and Homeowners, as well as CMP.





The graph is very similar to the graph above for commercial lines, with carried reserves falling in line with most methods but below Cape Cod adjusted methods. It is interesting that using trend in losses but not rate change for these lines would seem to overstate indications versus that of other methods in the same manner it does for commercial lines, where rate change was incorporated. This would suggest that the available rate change information does not account for all the changes in loss ratio due to market cycle effects.

Conclusion: Rate changes do not necessarily compensate for loss trends during a market cycle; this would suggest that changes in terms and conditions play a significant role in the final outcome of the results.



Findings for Trended Loss Ratio BF Method - All Accident Years

This graph allows us to focus on the trended loss and DCC ratio BF method for all accident years as of December 31, 2012. Similar to findings for shorter-tailed lines in the Commercial segment, there is much less variation in carried reserves over time and less deviation of carried reserves from a specific method. However, the cycle effect on initial carried reserves is still present, even though to a lesser degree.

Conclusion: The choice of IELR even in the latest accident has very little significance on shorter-tailed lines.





Based on the graph above, one can see that other than the trended adjusted Cape Cod Methods, the methods and carried reserves are very close until the latest accident year. It is also clear that the initial carried amounts do not move with the methods, implying a tendency not to deviate from a set level of reserves or a delay in reacting to market swings. This effect is minimized in these shorter-tailed lines.

Conclusions: The impact of market cycles is present in Personal Lines as well; however the impact is less pronounced than commercial lines.

Findings for DCC to Loss Ratio - Current Accident Year

The previous graphs all consider the projection of losses and DCC together. We also examined the BF method in the context of separate loss and DCC analyses and using the projection of ultimate losses as an exposure base for DCC.



DCC Method results are slightly more erratic than loss plus DCC methods. For the 2012 accident year, actual carried reserves fell close to an incurred Cape Cod method. This makes sense since the ratio of DCC to loss might be expected to be more stable than losses during market changes. These findings are very similar to the findings for Commercial lines.

Conclusion: Consistent with Commercial Lines, DCC ratio is not as impacted by the market cycle.



Findings for Gross Loss and DCC Ratio - Current Accident Year

In general, the pattern of methods is very similar to what we observed on a net basis. One difference is that carried reserves are closer to the higher end of the methods on a gross basis, perhaps suggesting that carried reserves are swayed more by balance sheet considerations than exact methodology. Once again, we see personal lines results are similar to Commercial lines results.

Conclusion: Ceded reserves may not be as impacted by preconceived reserve level expectations.

Findings for Individual Lines of Business - Personal Lines

The findings for individual lines of business were very similar to the all lines indications. The findings for Commercial Lines with regards to shorter-tailed lines have a more profound effect. The Homeowners line demonstrates convergence of carried reserves to methods very quickly.

Shown below are the results for Homeowners.



Conclusion: As expected, choice of initial expected loss ratio has very little effect on the Homeowners line of business.

A full set of graphs is available in Appendix C.

Hindsight Testing

Commercial Lines

Findings for Commercial Lines - 2003 Accident Year

The following graph represents results of methods as they would appear at the end of 2003. 2003 was coming off of a very severe soft market, where overall loss ratios peaked in 1999 at over 100%. In 2003, the industry picked the carried loss ratio ignoring Cape Cod information from the prior years. Although Cape Cod methods, adjusted for large rate increases, came closer to the hindsight 2003 ratio (as carried in 2012), they still overstated the loss ratio. Even trending from 2002 overstated the actual loss ratio achieved during 2003. Here "Oldest Estimate" is the earliest carried amount we have data for.



Conclusion: During a year not affected by the soft market cycle but following the soft market, adjusted (on-leveled) methods for Cape Cod tend to overstate ultimate losses.



Findings for Commercial Lines - All Accident Years Adjusted Cape Cod

The effects of the market cycle and reserving practices appear clearly on this graph. We can see that initial carried amounts were furthest from method indications and 2012 carried levels during the high point of the cycle. The adjusted Paid and Incurred Cape Cods (performed with information as of year-end 2003) did a good job of matching the losses coming off the soft market.

Conclusion: During accident years in a soft market cycle, adjusted (on-leveled) Cape Cod methods do a good job of predicting ultimate losses.



Findings for Commercial Lines – All Accident Years All Methods

As mentioned above it becomes apparent that the soft market renders the trended and adjusted methods most useful whereas in years following the soft market, these methods will overstate losses.

Conclusion: Knowledge of the market cycle is critical to establishing an appropriate IELR; in many instances, knowledge of the market cycle is more important than the variety of methods used.



Findings for Commercial Lines - DCC to Loss Ratio 2003

Final DCC to loss ratio was higher than initially carried. Either Paid or Incurred Cape Cod methods would have provided a better estimate than actual booked reserves. The source of the low carried amounts is unclear, but it is possible that the uncertainty in loss amounts makes the DCC prediction less predictable.

Findings for Commercial Lines - DCC to Loss Ratio All Accident Years



Looking at all accident years for commercial lines, DCC booked ratios were deficient in the years following the soft cycle. Cape Cod methods were more accurate. DCC booked ratios seemed to go down when booked loss ratios went up.

Conclusion: Carriers held lower DCC reserves than necessary following a soft market cycle, but Cape Cod methods would have predicted DCC more accurately.

Findings for Commercial Lines - Gross Loss and DCC Ratio 2003 Accident Year



In the case of Gross losses, the carried loss ratios were overstated. Similar to the more recent loss ratios above, gross carried loss ratios are less affected by market considerations and are more commensurate with the results of Cape Cod methods.



Findings for Commercial Lines - Gross Loss and DCC Ratio All Accident Years

Similarly to above, booked gross losses were closer to final estimates and adjusted Cape Cod methods. This has a serious implication, in that the methodology is adequate but insurers chose to book lower net reserves in a soft market cycle.

Personal Lines

For Personal Lines, we did not have the benefit of a rate change index, but we did have trend information for the two personal lines studied, Personal Auto and Homeowners, as well as CMP.

Findings for All Lines - Current Accident Year



For personal lines, all methods and carried reserves overstated losses. This is more similar to gross reserves on commercial lines.





As seen above, the shorter-tail personal lines demonstrate a convergence of methods as well as carried reserves.

Findings for All Methods - All Accident Years

			\sim										
1994	1995	1996	1997	1998	1999	2000	2001	2002	2003				
76.26%	72.24%	74.65%	66.68%	71.26%	73.53%	78.71%	78.30%	69.09%	64.25%				
76.26%	72.24%	74.65%	66.68%	71.26%	73.53%	78.73%	78.28%	68.94%	63.03%				
		[73.53%	78.75%	78.73%	72.04%	68.78%				
					73.53%	78.77%	78.43%	70.22%	64.57%				
76.26%	72.14%	74.52%	66.84%	71.37%	73.55%	78.34%	78.50%	69.81%	66.27%				
76.26%	72.24%	74.65%	66.67%	71.26%	73.53%	78.74%	78.27%	68.85%	62.18%				
76.22%	72.26%	74.63%	66.75%	71.19%	73.35%	78.18%	77.68%	69.89%	67.37%				
76.24%	72.24%	74.64%	66.69%	71.20%	73.43%	78.54%	78.05%	69.39%	64.48%				
					73.16%	78.11%	77.68%	70.76%	70.57%				
					73.34%	78.47%	77.99%	69.57%	65.52%				
	1994 76.26% 76.26% 76.26% 76.26% 76.22% 76.24%	1994 1995 76.26% 72.24% 76.26% 72.24% 76.26% 72.14% 76.26% 72.24% 76.26% 72.24% 76.26% 72.24% 76.26% 72.24% 76.26% 72.24% 76.26% 72.24%	1994 1995 1996 76.26% 72.24% 74.65% 76.26% 72.24% 74.65% 76.26% 72.14% 74.52% 76.26% 72.24% 74.65% 76.26% 72.24% 74.65% 76.26% 72.24% 74.65% 76.26% 72.24% 74.65% 76.22% 72.26% 74.63% 76.24% 72.24% 74.64%	1994 1995 1996 1997 76.26% 72.24% 74.65% 66.68% 76.26% 72.24% 74.65% 66.68% 76.26% 72.24% 74.65% 66.68% 76.26% 72.24% 74.65% 66.84% 76.26% 72.24% 74.65% 66.7% 76.26% 72.24% 74.65% 66.67% 76.22% 72.26% 74.63% 66.75% 76.24% 72.24% 74.64% 66.69%	1994 1995 1996 1997 1998 '6.26% 72.24% 74.65% 66.68% 71.26% '6.26% 72.24% 74.65% 66.68% 71.26% '6.26% 72.14% 74.52% 66.84% 71.37% '6.26% 72.24% 74.65% 66.67% 71.37% '6.26% 72.24% 74.65% 66.67% 71.26% '6.26% 72.24% 74.63% 66.67% 71.26% '6.22% 72.24% 74.64% 66.69% 71.20% '6.24% 72.24% 74.64% 66.69% 71.20%	1994 1995 1996 1997 1998 1999 '6.26% 72.24% 74.65% 66.68% 71.26% 73.53% '6.26% 72.24% 74.65% 66.68% 71.26% 73.53% '6.26% 72.24% 74.65% 66.68% 71.26% 73.53% '6.26% 72.14% 74.52% 66.84% 71.37% 73.53% '6.26% 72.14% 74.65% 66.67% 71.26% 73.53% '6.26% 72.24% 74.65% 66.67% 71.36% 73.53% '6.26% 72.24% 74.65% 66.67% 71.26% 73.53% '6.22% 72.26% 74.63% 66.75% 71.19% 73.35% '6.24% 72.24% 74.64% 66.69% 71.20% 73.43% '6.24% 72.24% 74.64% 66.69% 71.20% 73.43% '6.24% '73.34% '73.34% '73.34%	1994 1995 1996 1997 1998 1999 2000 76.26% 72.24% 74.65% 66.68% 71.26% 73.53% 78.71% 76.26% 72.24% 74.65% 66.68% 71.26% 73.53% 78.73% 76.26% 72.24% 74.65% 66.68% 71.26% 73.53% 78.73% 73.53% 78.75% 73.53% 78.75% 73.53% 78.75% 73.53% 74.65% 66.68% 71.37% 73.53% 78.75% 76.26% 72.14% 74.52% 66.84% 71.37% 73.55% 78.34% 76.26% 72.24% 74.65% 66.67% 71.26% 73.53% 78.74% 76.26% 72.24% 74.63% 66.75% 71.19% 73.35% 78.18% 76.24% 72.24% 74.64% 66.69% 71.20% 73.43% 78.54% 72.24% 74.64% 66.69% 71.20% 73.43% 78.14% 73.16% 78.11% 73.3	1994 1995 1996 1997 1998 1999 2000 2001 76.26% 72.24% 74.65% 66.68% 71.26% 73.53% 78.71% 78.30% 76.26% 72.24% 74.65% 66.68% 71.26% 73.53% 78.73% 78.28% 76.26% 72.24% 74.65% 66.68% 71.26% 73.53% 78.73% 78.28% 73.53% 78.75% 78.73% 78.73% 78.73% 78.28% 73.53% 78.75% 78.73% 78.73% 78.73% 78.73% 74.65% 66.68% 71.37% 73.55% 78.74% 78.50% 76.26% 72.24% 74.65% 66.67% 71.26% 73.53% 78.74% 78.27% 76.26% 72.24% 74.65% 66.67% 71.26% 73.53% 78.18% 77.68% 76.22% 72.26% 74.63% 66.75% 71.19% 73.35% 78.18% 77.68% 76.24% 72.24% 74.64% 6	1994 1995 1996 1997 1998 1999 2000 2001 2002 26.26% 72.24% 74.65% 66.68% 71.26% 73.53% 78.71% 78.30% 69.09% 26.26% 72.24% 74.65% 66.68% 71.26% 73.53% 78.73% 78.28% 68.94% 26.26% 72.24% 74.65% 66.68% 71.26% 73.53% 78.73% 78.28% 68.94% 2 73.53% 78.75% 78.73% 72.04% 2 73.53% 78.77% 78.43% 70.22% 26.26% 72.14% 74.52% 66.84% 71.37% 73.53% 78.74% 78.50% 69.81% 26.26% 72.24% 74.65% 66.67% 71.26% 73.53% 78.74% 78.27% 68.85% 26.26% 72.24% 74.63% 66.75% 71.19% 73.35% 78.18% 77.68% 69.89% 26.26% 72.24% 74.64% 66.69% 71.20% 73.43% </td				

For short-tailed personal lines, only the most recent two years show any material variation in methods. In hindsight, the Incurred Cape Cod methods seemed to be the closest for both years.



Findings for Personal Lines - DCC to Loss Ratio 2003

In this case the initial carried for DCC overstated the latest carried. This is the opposite effect observed on the commercial lines. If anything, this underscores the unpredictability of DCC after a soft market cycle. Unlike the commercial lines, none of the methods would have predicted the right level of DCC.

SUMMARY OF OBSERVATIONS AND CONCLUSIONS

In selecting the initial carried amounts, reactions to market cycles appear to play a more prominent role than actuarial indications. During accident years in a soft market cycle, adjusted (on-leveled) Cape Cod methods often do a more accurate job of predicting ultimate losses, but could easily overestimate losses during the period following a soft market. In this case, accurate rate changes (and changes in terms and conditions) may not be available to properly adjust the method. Overall it appears the industry selects accident year loss ratios more uniquely than the Cape Cod results, which would weight in a more long-term average.

The influence of market cycle in deflating net reserves during a soft market is not seen on the gross side, which suggests carriers approach net and gross reserves differently during the market cycle. In addition, DCC reserves tend to be deflated during a soft cycle, despite the fact that actuarial indications such as Cape Cod are not distorted by the cycle.

Most of the observations above impact long-tailed commercial lines and not surprisingly, have a lesser effect on short-tailed or personal lines.

Knowledge of the market cycle is critical to establishing an appropriate IELR; in many instances, knowledge of the market cycle is as important as the appropriateness of the methods used to select the IELR. Based on hindsight testing, it is apparent that methods that reflect rate changes, loss trends and give appropriate weights to the on-level loss ratios (the best example being Cape Cod) tend to perform better than methods that do not. However, it is evident from the survey results that the use of the Cape Cod is not prevalent within the industry. Although the use of appropriate models can play a role in improving the accuracy of the booked reserves, changing business conditions and business considerations are also factors that have an important impact.

REFERENCES

- [1] Bornhuetter, Ronald L. and Ronald E. Ferguson, "The Actuary and IBNR," *PCAS*, 1972, Vol. LIX, 181-195.
- [2] Stanard, James N., "A Simulation Test of Prediction Errors of Loss Reserve Estimation Techniques," *PCAS*, 1985, Vol. LXXII, 124-148.
- [3] Gluck, Spencer M., "Balancing Development and Trend in Loss Reserve Analysis," *PCAS*, 1997, Vol. LXXXIV, 482-532.

Appendix A – Method Examples

Adjusted Rate Indication MethodExhibit 1Estimating Accident Year 2016 Initial Expected Loss Ratio at 12/31/2016Exhibit 1

(1)	Projection Period	Policies Effective from 7/1/2017 to 6/30/2018
(2)	Indicated Ultimate Loss Ratio for Projection Period	65.3%
(3)	Net Annual Loss/Premium Trend	3.0%
(4)	Average Earned Date for Projection Period	6/30/2018
(5)	Midpoint of Accident Year 2016	6/30/2016
(6)	Number of Years of Trend	2.0
(7)	Detrend Factor	0.943
(8)	2015 Rate Change	2.0%
(9)	Effective Date of 2015 Rate Change	7/1/2015
(10)	Portion of 2015 Rate Change Not Earned in 2016	12.4%
(11)	Unearned 2015 Rate Change Adjustment	1.002
(12)	2016 Rate Change	2.0%
(13)	Effective Date of 2016 Rate Change	7/1/2016
(14)	Portion of 2016 Rate Change Not Earned in 2016	87.4%
(15)	Unearned 2016 Rate Change Adjustment	1.017
(16)	Selected IELR	62.8%

Notes:

(1), (2), (3), (8), (9), (12) and (13) from rate indication (6) = ((4) - (5)) / 365(7) = $(1 / (1 + (3)))^{6}$ (10) = $(((9) + 365 - 12/31/2015)/365)^{2}/2$ (11) = $1 + (8) \times (10)$ (14) = $1 - ((12/31/2016 - (13))/365)^{2}/2$ (15) = $1 + (12) \times (14)$ (16) = (2) x (7) x (11) x (15)

Cape Cod Method Using On-Level Earned Premiums (\$000's) Estimating Accident Year 2016 Initial Expected Loss Ratio at 12/31/2016

			Cumulative	On-Level	On-Level	Annual	Cumulative	Loss	Trended		"Reported"	Trended		
Accident	Reported	Earned	Rate	Premium	Earned	Loss	Loss Trend	Trend	Reported	Percent	On-Level	Developed	Decay	
Year	Losses	Premium	Index	Factor	Premium.	Trend	Index	Factor	Losses	Reported	Premium	Loss Ratio	Weight	<u>Weight</u>
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)	(13)	(14)	(15)
2007	68,000	120,000	1.004	1.275	152,981		1.000	1.409	95,819	98.0%	149,982	63.9%	0.075	11,261
2008	69,476	123,152	1.030	1.242	152,981	5.0%	1.050	1.342	93,236	97.1%	148,497	62.8%	0.100	14,866
2009	71,765	126,846	1.061	1.206	152,981	5.0%	1.103	1.278	91,723	95.2%	145,585	63.0%	0.133	19,433
2010	75,217	130,652	1.093	1.171	152,981	5.0%	1.158	1.217	91,557	93.3%	142,730	64.1%	0.178	25,403
2011	73,397	134,571	1.126	1.137	152,981	5.0%	1.216	1.159	85,088	88.9%	135,934	62.6%	0.237	32,258
2012	70,124	139,994	1.159	1.104	154,511	3.0%	1.252	1.126	78,925	82.3%	127,123	62.1%	0.316	40,223
2013	65,882	145,636	1.194	1.072	156,056	3.0%	1.290	1.093	71,991	73.5%	114,638	62.8%	0.422	48,363
2014	56,643	152,814	1.228	1.042	159,177	3.0%	1.328	1.061	60,092	60.2%	95,845	62.7%	0.563	53,913
2015	41,603	156,056	1.255	1.020	159,177	3.0%	1.368	1.030	42,851	42.3%	67,259	63.7%	0.750	50,445
2016	27,981	159,177	1.280	1.000	159,177	3.0%	1.409	1.000	27,981	28.2%	44,84 0	62.4%	1.000	44,84 0
Total	620,087	1,388,899			1,553,004				739,263		1,172,431			341,004
(16) Selecte	d Decay Fac	ctor	0.75											

(17) Selected IELR 62.9%

Notes:

(2), (3), (4), (7) and (11) from company data
(5) = ((4) for Accident Year 2016) / (4)
(6) = (3) x (5)
(8) cumulative index based on (7)
(9) = ((8) for Accident Year 2016) / (8)
(10) = (2) x (9)
(12) = (6) x (11)
(13) = (10) / (12)
(14) = (16)^{(2016 - (1))}
(15) = (12) x (14)
(16) judgmentally selected
(17) weighted average of (13) using (15) as weights
Cape Cod Method Using Earned Exposures (\$000's) Estimating Accident Year 2016 Initial Expected Loss Ratio at 12/31/2016

Cumulative Trended "Reported" Trended Annual Accident Reported Earned Loss Loss Trend Loss Trend Reported Percent Earned Developed Decay Trend Index Factor Losses Reported Exposures Pure Prem. Weight Weight Year Losses Exposures (1) (2)(3) (4) (5) (6)(7)(8)(9)(10)(11)(12)2007 68,000 100,000 1.000 1.409 95,819 98.0% 98,039 977 0.075 7,361 2008 69,476 100,000 5.0% 1.050 1.342 93,236 97.1% 97,069 961 0.100 9,718 2009 71,765 100,000 5.0% 1.103 1.278 91,723 95.2% 95,165 964 0.133 12,703 2010 75,217 100,000 5.0% 1.158 1.217 91,557 93.3% 93,299 981 0.17816,605 2011 73,397 100,000 5.0% 1.216 1.159 85,088 88.9% 88,856 958 0.237 21,086 2012 70,124 101,000 3.0% 1.252 1.126 78,925 82.3% 83,097 950 0.316 26,292 65,882 2013 102,010 3.0% 1.290 1.093 71,991 73.5% 74,936 961 0.422 31,614 2014 56,643 104,050 1.061 60,092 60.2% 62,651 959 0.563 35,241 3.0% 1.328 2015 41,603 1.030 43,966 975 0.750 32,974 104,050 3.0% 1.368 42,851 42.3% 27,981 104,050 3.0% 1.409 1.000 27,981 28.2% 29,311 1.000 29,311 2016 955 Total 620,087 1,015,161 739,263 766,389 222,906 (13) Selected Decay Factor 0.75 (14) Selected Expected Pure Premium 962 (15) Accident Year 2016 Earned Premium 159,177 (16) Selected IELR 62.9%

Notes: (2), (3), (4), (8) and (15) from company data (5) cumulative index based on (4) (6) = ((5) for Accident Year 2016) / (5) (7) = (2) x (6) (9) = (3) x (8) (10) = (7) / (9) (11) = (13)^(2016 - (1)) (12) = (9) x (11) (13) judgmentally selected (14) weighted average of (10) using (12) as weights (16) = (14) x ((3) for 2016) / (15)

Prior Accident Year Loss Ratios Trended and Rate Adjusted (\$000's) Estimating Accident Year 2016 Initial Expected Loss Ratio at 12/31/2016

Exhibit 4

			Estimated	Estimated	Cumulative	On-Level	Annual	Cumulative		Estimated
Accident	Earned	Earned	Ultimate	Ultimate	Rate	Premium	Loss	Loss Trend I	Loss Trend	Expected
Year	<u>Exposures</u>	Premium	Loss	Loss Ratio	Index	Factor	Trend	Index	Factor	Loss Ratio
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)
2007	100,000	120,000	69,360	57.8%	1.004	1.275		1.000	1.409	63.9%
2008	100,000	123,152	71,574	58.1%	1.030	1.242	5.0%	1.050	1.342	62.8%
2009	100,000	126,846	75,411	59.5%	1.061	1.206	5.0%	1.103	1.278	63.0%
2010	100,000	130,652	80,619	61.7%	1.093	1.171	5.0%	1.158	1.217	64.1%
2011	100,000	134,571	82,602	61.4%	1.126	1.137	5.0%	1.216	1.159	62.6%
2012	101,000	139,994	85,231	60.9%	1.159	1.104	3.0%	1.252	1.126	62.1%
2013	102,010	145,636	89,686	61.6%	1.194	1.072	3.0%	1.290	1.093	62.8%
2014	104,050	152,814	94,071	61.6%	1.228	1.042	3.0%	1.328	1.061	62.7%
2015	104,050	156,056	98,458	63.1%	1.255	1.020	3.0%	1.368	1.030	63.7%
2016	104,050	159,177	99,331	62.4%	1.280	1.000	3.0%	1.409	1.000	

Total 1,015,161 1,388,899 846,343

(12) Average Estimated Expected Loss Ratios	
Average All Years	63.1%
Average Latest 7 Years	63.0%
Average Latest 5 Years	62.8%
Average Latest 3 Years	63.1%
(13) Selected IELR	63.1%

Notes:

(2), (3), (6) and (8) from company data
(4) from prior reserve review valued at 6/30/2016
(5) = (4) / (3)
(7) = ((6) for Accident Year 2016) / (6)
(9) cumulative index based on (8)
(10) = ((9) for Accident Year 2016) / (9)
(11) = (5) x (10) / (7)
(12) simple averages of (11)

(13) selected based on (11) and (12)

Prior Accident Year Pure Premiums Trended and Rate Adjusted (\$000's)Exhibit 5Estimating Accident Year 2016 Initial Expected Loss Ratio at 12/31/2016Exhibit 5

			Estimated	Estimated	Annual	Cumulative		Estimated
Accident	Earned	Earned	Ultimate	Ultimate	Loss	Loss Trend	Loss Trend	l Expected
Year	<u>Exposures</u>	Premium	Loss	ure Premiur	Trend	Index	Factor	ure Premiur
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
2007	100,000	120,000	69,360	694		1.000	1.409	977
2008	100,000	123,152	71,574	716	5.0%	1.050	1.342	961
2009	100,000	126,846	75,411	754	5.0%	1.103	1.278	964
2010	100,000	130,652	80,619	806	5.0%	1.158	1.217	981
2011	100,000	134,571	82,602	826	5.0%	1.216	1.159	958
2012	101,000	139,994	85,231	844	3.0%	1.252	1.126	950
2013	102,010	145,636	89,686	879	3.0%	1.290	1.093	961
2014	104,050	152,814	94,071	904	3.0%	1.328	1.061	959
2015	104,050	156,056	98,458	946	3.0%	1.368	1.030	975
2016	104,050	159,177	99,331	955	3.0%	1.409	1.000	
Total	1,015,161	1,388,899	846,343					
(10) Avera	ge Estimated	d Expected	Pure Premi	lums				0 / F
Av	erage All Ye	ars						965
Av	erage Latest	7 Years						964
Av	erage Latest	5 Years						960
Av	erage Latest	3 Years						965
(11) Selecte	ed Expected	Pure Prem	ium					965
(12) Selecte	ed IELR							63.1%
NT /								
Notes: (2) (2) and	(() (_					
$(2), (3)$ and (4) from π	(0) from cc	mpany data	ad at $\frac{20}{20}$	2016				
(4) from p (5) - (4) /	(2)	review value	ed at 0/ 50/	2010				
(3) - (4) / (7)	(4)	and on (f)						
(7) cumula	ive maex D?	$\frac{1}{2} \sum_{n=1}^{\infty} \frac{1}{2} \sum_{n=1}^{\infty} \frac{1}$	(7)					
(0) - ((/) I (0) - (5)	(9)	1 cai 2010)	/ (/)					
(2) - (3) X	(0)	(0)						
(10) simple	: averages 01	. (2)						

(11) selected based on (9) and (10)

(12) = (11) x ((2) for Accident Year 2016) / ((3) for Accident Year 2016)

Bornhuetter-Ferguson Initial Expected Loss Ratio Working Party Paper

Appendix B – Survey Results

Choice of Method - Long-Tailed Lines

43.6%
27.6%
9.6%
9.3%
6.1%
2.3%
1.5%

Choice of Method - Short-Tailed Lines

Prior Accident Years Adjusted for Rate Changes and Trends	34.3%
Prior Analysis Ultimate Loss Ratios	31.8%
Pricing Loss Ratio	11.4%
Cape Cod	8.6%
Prior Accident Years	8.0%
Judgment	3.4%
Industry Aggregates	2.5%

Choice of Method – Additional Considerations

In addition to the above, the actuary may also consider the following in selecting the IELR: (select all that apply)

Maturity of accident year	78.0%
Homogeneity of portfolio	48.3%
Credibility of development factors	46.6%
Size of Book	45.9%
Size of development factors	33.8%

BF Used to Develop (select all that apply)

ALAE/DCC	81.0%
Claim Counts	51.4%
Salvage and Subrogation	31.2%
ULAE/AAO	6.5%

How is DCC Treated

Analyze Loss and Expense combined	30.7%
Assume an Expense/Ultimate Loss Ratio that varies by year	23.2%
Don't use BF on expenses	22.2%
Assume a fixed percent to losses/premium for all years as IE	15.0%
Assume an Expense/Premium IE that varies by year	8.2%
Use a claim count method to determine ultimate expenses	0.7%

For Current AY, BF is

Always used	49.6%
Sometimes used	40.5%
Rarely used	7.1%
Not used	2.8%

For Other than Current AY, BF is

Sometimes used	76.1%
Always used	14.1%
Rarely used	9.3%
Not used	0.6%

How Often is IELR Reselected?

Annually	61.1%
Quarterly	31.4%
Every 2 - 3 years	2.9%
Every 3 -5 years	2.3%
Never	2.3%

Restrictions on IELR?

No boundaries put in place	62.1%
Higher than reported losses	26.8%
Higher than paid losses (excluding high salvage situations)	11.1%

Use of Cape Cod with

Don't use Cape Cod	65.0%
Loss trend	29.3%
Rate changes	25.9%
A decay factor	18.2%
Rate Changes considered with	
A price monitor	63.6%
Not considered	20.2%
Planned changes	16.2%
Sources of Industry LR Benchmarks	
Not considered	50.6%
AM Best	13.5%
Internal benchmarks	13.1%
NCCI	9 3%
	2.570
SNL	9.0%
SNL ISO	9.0% 4.5%
SNL ISO Management Influence	9.0% 4.5%
SNL ISO Management Influence My decisions are completely independent	9.0% 4.5% 50.7%
SNL ISO Management Influence My decisions are completely independent Management points out factors that I consider in my analysis	9.0% 4.5% 50.7% 42.2%

I feel pressure from management

1.4%

Reasonability Checks of IELR? (select all that apply)

Internal Peer Review	82.6%
Comparison of expected losses to actual emerged losses to date	65.9%
Hindsight Tests of accuracy of methodology	36.8%
External Peer Review	32.1%
Audit controls under SOX	14.7%
Audit controls under Model Audit Rule	6.8%

Escaping Hindsight: Case Reserve Development Using the Reserve Runoff Ratio

By Joseph Boor, FCAS, PhD, CERA

Abstract: The common calculation used in developing case reserves are based on "hindsight" from a separate development test, thus they are based on data that already reflects judgment. A method is presented for estimating development factors for case reserves that strictly uses data within the standard loss development triangles, primarily the paid loss to case reserve disposed or "runoff" ratios. This method is thus, a truly independent view of the case development factors.

Keywords: case reserve development, hindsight

1. INTRODUCTION

Developing case reserves of older years instead of using chain-ladder or other common methods to estimate the ultimate losses for those years has received at least one laudatory review. A 2009 paper by Jing, Lebens and Lowe suggests that case reserve development is often the best method of the alternatives for some maturities. It does have a weakness, though. Developing case reserves often uses "hindsight" methods that begin with the same chain-ladder, etc. methods and then compute the case reserve development factors that would have been needed in the years above the diagonal if the estimated ultimate is an accurate estimate. The case reserve factor for the current diagonal is estimated from the results. Aside from the development beyond the last diagonal being driven by a potentially misjudged ultimate loss¹, ultimate losses developed using this hindsight case reserve development process may be prone to match the beginning chain-ladder ultimate loss. So using the hindsight case development method could result in either repeating a misjudgment or (maybe also) a misleading confidence in the results.

1.1 The Benefits

That being said, developing case reserves on mature years has high potential to estimate the loss on mature years. They potentially reflect whether a large number of claims remain open, or whether few claims, or only small claims, are open at present. So, it is advisable to have a case reserve development process that is not based on an initial ultimate loss estimate.

¹ In the context of this paper, the word "loss" is used to represent whatever data is being developed, whether that is "loss", "loss and defense and cost containment", or some similar type of data.

2. THE METHOD

This paper presents an alternative method that strictly uses information inside the triangles, and does not involve any external judgment. Rather it is a slight extension of the "runoff ratio" presented in Sherman 2006. That in turn stems from the "paid loss to reserve disposed of ratio" used in Sherman 1984, Boor 2006, and the Report of the CAS Tail Factor Working Party (Herman, et al 2013).

That process is fairly simple. It will be illustrated by an example that follows the process from start to finish. First, each incremental paid loss by development cell is computed by subtracting the adjacent-to-the-left paid loss from the paid loss in each corresponding cell. Secondly, the reserve disposed of uses a negative process. The value in each cell of the case reserve triangle is instead subtracted from the value to the left. So, the incremental paid loss value is the actual costs in the cell. Dividing by the case reserve produces a measure of the actual cost (in that cell) of disposing of a dollar of case reserve. That is the core calculation behind this set of correction factors for case reserves.

However, one factor from that analysis will not make a proper correction factor for the case reserves. The changes between the key processes of the claims department must be considered. As said in Boor 2006 and repeated in the 2014 Report of the Tail Factor Working Party

"It is important to consider the primary activity within each development stage.

When using multiple periods to estimate a tail factor, it is relatively important that the periods reflect the same general type of claims department activity as that which takes place in the tail. For example, in the early 12 to 24 month stage of workers compensation, the primary development activity is the initial reporting of claims and the settlement and closure of small claims. The primary factors influencing development are how quickly the claims are reported and entered into the system, and the average reserves (assuming the claims department initially just sets a 'formula reserve', or a fixed reserve amount for each claim of a given type such as medical or lost time) used when claims are first reported.

In the 24 to 36-48 month period, claims department activity is focused on ascertaining the true value of long-term claims and settling claims. After 48-60 months most of the activity centers on long-term claims. So, the 12-24 link ratio has relatively little relevance for the tail, as the driver behind the link ratio is reporting and the size of initial formula reserves rather than the handling of long-term cases. Similarly, if the last credible link ratio in the triangle is

the 24 to 36 or 36 to 48 link ratio, that triangle may be a poor predictor of the required tail factor."

Of course, the exact maturities at which the stages change may not match a particular reserving situation, but the progression through the stages likely will be an issue².

So, in summary, the key concerns dictate using the paid loss/reserve disposed of and being able to target the activity in a stage are key. As one might surmise, the first step is to develop a triangle of paid/disposed ratios.

Thankfully, such a triangle can be computed from the standard paid and reported loss reserving triangles. For example, given the following sample paid and reported loss data triangles,

Accident Year	12	24	36	48	60	72	84	96	108	120
										-
1999	2,065	4,759	8,883	11,832	13,005	13,290	13,502	13,508	13,510	13,510
2000	1,915	6,662	13,952	17,899	19,406	19,796	20,066	20,068	20,140	
2001	3,976	12,534	21,164	26,134	29,416	32,098	32,942	33,074		
2002	3,906	11,115	18,526	30,371	41,207	44,158	48,138			
2003	7,619	21,043	41,439	58,151	72,731	79,336				
2004	10,376	19,406	39,902	58,127	69,684					
2005	9,662	23,869	32,016	38,311						
2006	9,225	18,106	24,546							
2007	3,062	8,751								
2008	2,278									

Table 1: Cumulative Paid Loss

² It would seem, though, that for claims-made products the first stage might either not occur or have a short duration.

Accident										
Year	12	24	36	48	60	72	84	96	108	120
1999	5,605	8,126	10,710	12,586	13,378	13,495	13,508	13,516	13,516	13,516
2000	7,074	11,431	16,681	19,466	20,178	20,223	20,132	20,082	20,140	
2001	9,913	18,493	25,986	29,141	31,815	32,906	33,154	33,300		
2002	9,979	17,277	23,366	38,199	45,036	47,100	48,804			
2003	22,625	34,198	58,881	70,094	79,626	82,086				
2004	23,770	37,119	54,780	68,054	73,579					
2005	20,019	32,326	40,188	39,814						
2006	20,176	28,624	29,439							
2007	9,080	13,335								
2008	6,011									

Table 2: Cumulative Reported Loss

one may readily compute the incremental paid loss by subtracting values in adjacent columns in the cumulative paid loss data (Table 1).

Accident										
Year	12	24	36	48	60	72	84	96	108	120
1999	2,065	2,694	4,124	2,950	1,172	285	212	6	2	0
2000	1,915	4,747	7,290	3,948	1,507	390	270	2	72	
2001	3,976	8,558	8,630	4,970	3,282	2,682	844	132		
2002	3,906	7,209	7,411	11,845	10,836	2,951	3,980			
2003	7,619	13,424	20,396	16,712	14,580	6,605				
2004	10,376	9,030	20,496	18,225	11,557					
2005	9,662	14,207	8,148	6,295						
2006	9,225	8,881	6,440							
2007	3,062	5,690								
2008	2,278									

Table 3: Incremental Paid =Costs of Disposing of Case (Table 1 Value – Value in Previous Table 1 column)

Next, the case reserves disposed of in each cell must be computed. The first step, of course, is to compute the case reserves.

Accident										
Year	12	24	36	48	60	72	84	96	108	120
1999	3,540	3,367	1,827	754	374	205	6	8	6	6
2000	5,159	4,769	2,730	1,566	771	427	66	14	0	
2001	5,937	5,959	4,822	3,007	2,399	808	212	226		
2002	6,073	6,162	4,840	7,827	3,829	2,942	666			
2003	15,006	13,156	17,442	11,943	6,896	2,750				
2004	13,394	17,713	14,878	9,927	3,895					
2005	10,357	8,458	8,172	1,503						
2006	10,950	10,518	4,893							
2007	6,019	4,584								
2008	3,733									

Table 4: Case Reserves (Table 2 – Table 1)

Then, the reserve disposed of is computed using the additive inverse of the process used to compute the incremental paid loss. In other words, instead of subtracting the value in the previous column from the value in the current3 column, one would subtract the value in the current column from the value in the previous column. That is logical since case reserves tend to decrease after some point in the triangle whereas paid loss would increase. Thus, one would compute the case disposed of using the outline bin Table 5.

Table 5: Case Reserves Disposed of

Accident										
Year	12	24	36	48	60	72	84	96	108	120
1999	-3,540	173	1,540	1,074	380	169	199	-2	2	0
2000	-5,159	390	2,039	1,163	795	345	361	52	14	
2001	-5,937	-22	1,137	1,815	608	1,591	596	-14		
2002	-6,073	-89	1,322	-2,987	3,998	887	2,276			
2003	-15,006	1,851	-4,286	5,498	5,047	4,145				
2004	-13,394	-4,319	2,836	4,951	6,032					
2005	-10,357	1,899	286	6,668						
2006	-10,950	433	5,625							
2007	-6,019	1,435								
2008	-3,733									

³ The is perhaps an unusual phrase to some readers. The "value in the current column" would be the value in the cell with the same maturity and accident (or report for some coverages) year as the cell being computed.

Once the reserves are computed, it is easy to compute the ratio of paid loss to case reserves disposed of (the "runoff ratio").

Accident											
Year	12	24	36	48	60	72	84	96	108	120	
1999	-0.5833	15.5812	2.6777	2.7471	3.0870	1.6918	1.0646	-3.0000	1.0000	1.0000	
2000	-0.3712	12.1603	3.5747	3.3942	1.8955	1.1301	0.7487	0.0385	5.1429		
2001	-0.6697	-385.8499	7.5895	2.7382	5.4000	1.6855	1.4168	-9.4286			
2002	-0.6432	-80.5960	5.6059	-3.9656	2.7101	3.3284	1.7484				
2003	-0.5077	7.2531	-4.7588	3.0394	2.8886	1.5934					
2004	-0.7746	-2.0908	7.2282	3.6810	1.9159						
2005	-0.9329	7.4801	28.4870	0.9440							
2006	-0.8425	20.5332	1.1449								
2007	-0.5086	3.9647									
2008	-0.6100										
	12	24	36	48	60	72	84	96	108	120	Tail
Averages:											
Column ^{\$\$} Woightod	0.6746	12 5299	7 8006	2 5719	2 5464	1 8004	1 5/61	2 8880	4 6250	NI/A	
	-0.0740	42.3200	7.0990	5.5710	2.3404	1.0094	1.5401	3.0009	4.0250	IN/A	
3 Col. Centered \$\$											
Weighted				4.1898	2.8637	2.2295	1.7313	1.5845	4.1154		
5 Col. Centered \$\$											
Weighted					3.7254	2.7655	2.2331	1.7356			
All-Time Unweighted	-0.6444	-44.6182	6.4436	1.7969	2.9828	1.8858	1.2446	-4.1300	3.0714	1.0000	
Selected Values	0.0000	0.0000	7.0000	3.0000	3.0000	2.2500	2.5000	3.0000	4.0000	4.0000	4.0000
Notes: Early factors	set at zero as t	hey clearly do	not involve u	ipward deve	elopment in	existing cla	ims-they ap	ppear to ofte	en show		1 1
disposed of,	in those period	stead of decrea	ises. The 12 a	and 24 paid,	/ disposed f	atios were s	et at zero s	ince case res	serves are ci	early being	built, not
1							-				
Selections ge	enerally relied he	eavily on the 3	column aver	rage, then th	ne 5 column	i, although s	some crede	nce was give	en to consis	stency with	the single

Table 6: "Runoff Ratio"----Paid Loss to Case Reserve Disposed of (Table 3 Value/Table 5 Value)

Where no 5 column or 3 column averages existed, the nearest ones were considered.

column dollar weighted, especially at earlier maturities

Also, note, efforts were made to round for consistency and to show consistent patterns of increase and decrease.

This gives the core information of the process --- how much it costs to eliminate a dollar of case reserve. These may be called runoff ratios, since they represent the true value of closing or "running off" each dollar of case reserve. Note that comments on how the runoff ratios were selected are included. However, as one may see, each of the selected runoff ratios only covers the activity during a twelve month period of development. Since a case reserve will pay out over multiple development periods, it is necessary to use a weighted average of the appropriate set of runoff ratios. There is more than one way to compute the weights. One could analyze the average decrease in case reserves (the case at the end of the period divided by the case at the beginning of the period) through each twelve month stage of development, and use that decrease pattern to determine the weights for the various value factors. That process is shown in Appendix A. When one is performing a reserve review, it is likely that reported loss and paid loss development patterns, and corresponding patterns of the percentages of loss reported and paid, have already been estimated. Using those, one may back into the case reserves at each stage as a percentage of ultimate loss at the various twelve month stages. That allows for a calculation of the expected decay in case reserves.

The table below begins with paid and reported loss development factors arising from mechanical selection of all-time weighted average link ratios and Sherman-Boor tail analysis (using the runoff ratio expected for the tail in Table 6). Then, case reserves development factors are computed by weighting the runoff ratios using the case decay rates. Of course, since the disposition will take in subsequent periods, one must use the data in subsequent columns, not the column in question. Note that all that is required to both compute the runoff ratios and the weights are the paid and reported loss triangles. This example does use paid and reported development patterns as a start for computing the decay rates. But one could just as readily compute decay rates by dividing adjacent values in Table 4 (see Appendix A). Therefore, this reserve runoff approach to case development is not heavily affected by the other development tests.

	12	24	36	48	60	72	84	96	108	120	Tail
Selected Runoff Ratios (A)	0.0000	0.0000	7.0000	3.0000	3.0000	2.2500	2.5000	3.0000	4.0000	4.0000	4.0000
% Incrd to-Date per ILDFs	63.21%	72.70%	82.14%	91.08%	97.62%	99.01%	99.92%	99.86%	100.00%	99.40%	100.00%
% Paid to-Date per PLDFs	12.05%	31.74%	55.69%	75.95%	90.06%	95.39%	98.87%	99.02%	99.21%	99.21%	100.00%
Case @Date (B)	51.16%	40.96%	26.45%	15.13%	7.55%	3.62%	1.04%	0.83%	0.79%	0.20%	0.00%
Decay in Period (C)	80.07%	64.57%	57.22%	49.92%	47.87%	28.90%	79.85%	95.11%	25.00%	0.00%	100.00%
=(B next)/(B)											
Cumulative Case											
Development Factors (D)	3.4701	4.3341	2.8713	2.7750	2.5493	2.8752	3.7985	4.0000	4.0000	4.0000	
(next A)*(1.0-(C))+(D next)*(C)											

 Table 7: Calculation of Case Reserve Development Factors

At this stage, one has factors that could plausibly be used to develop case reserves at 36 months maturity and could more plausibly be used to develop case reserves of 48 or more months of maturity. On that basis, this process could be used to convert case reserves into estimated loss liabilities, at least for some years. However, in practice most actuaries first develop ultimate loss, and then develop the loss reserve/liability indication from the ultimate losses.

Therefore, it is necessary to show how the ultimate loss may be estimated using case development. It should be clear that to estimate the ultimate loss of a given accident or report year, one need only develop the case reserves, then add in the paid loss to date. Table 8 shows the calculations.

(1)	(2)	(3)	(4)	(5)	(6)	(7)
	(Table 1))	(Table 4)	(Table 7)	(3)*(4)	(2)+(5)	
						Are Claims
Begin of	Paid to	Case	Case	Estimate of	Estimate of	Developed
Accident	Date @	Reserves @	Reserve	Ultimate	Total	Enough to
Year	12/31/08	12/31/08	LDF	Reserve	Ultimate	Be Usable?
1999	13,510	6	4.0000	24	13,534	Yes
2000	20,140	0	4.0000	0	20,140	Yes
2001	33,074	226	4.0000	904	33,978	Yes
2002	48,138	666	3.7985	2,530	50,668	Yes
2003	79,336	2,750	2.8752	7,908	87,244	Yes
2004	69,684	3,895	2.5493	9,929	79,612	Yes
2005	38,311	1,503	2.7750	4,172	42,483	Yes
2006	24,546	4,893	2.8713	14,049	38,595	Yes
2007	8,751	4,584	n/a	n/a	n/a	No
2008	2,278	3,733	n/a	n/a	n/a	No
	337,768	22,257		39,516	366,255	

Table 8: Final Case Reserve Development and Ultimate Loss

3. CONCLUSIONS

A process of case reserve development without relying on ultimate loss estimates from other methods is presented above. Hopefully, this method will achieve wide adoption and improve the quality of ultimate loss estimates, especially for the years near the tail.

Appendix A– Computing the Case Reserve Decay (and Consequent Case Reserve Development Factors) from the Case Reserve Triangle

If one desires to more completely isolate the case reserve development results from the paid and incurred loss development tests, one need only compute the decay in the case reserves from the triangle of case reserves. Then, one may weight the runoff ratio in a column with the decay in the case reserve over that period, and assign the remaining weight to the (composite) case reserve development factor for the next maturity. The process proceeds as follows:

Begin of Accident										
Year	12	24	36	48	60	72	84	96	108	
1999	0.9512	0.5426	0.4124	0.4961	0.5488	0.0292	1.3333	0.7500	1.0000	
2000	0.9243	0.5724	0.5739	0.4924	0.5531	0.1547	0.2121	0.0000		
2001	1.0037	0.8092	0.6235	0.7978	0.3367	0.2625	1.0660			
2002	1.0147	0.7855	1.6171	0.4892	0.7684	0.2264				
2003	0.8767	1.3258	0.6848	0.5774	0.3988					
2004	1.3224	0.8399	0.6672	0.3923						
2005	0.8166	0.9662	0.1840							
2006	0.9605	0.4652								
2007	0.7616									
	12	24	36	48	60	72	84	96	108	Tail
Averages										
All-Time \$Weighted	0.9771	0.8502	0.6677	0.5186	0.4998	0.1332	0.8732	0.2727	1.0000	
3 Year \$Weighted	0.8621	0.7616	0.5773	0.4923	0.4953	0.2260	0.8732	0.2727	1.0000	
5 year \$Weighted	0.9767	0.8968	0.6821	0.5191	0.4998	0.2168	0.8732	0.2727	1.0000	
All-Time Unweighted	0.9591	0.7884	0.6804	0.5409	0.5212	0.1682	0.8705	0.3750	1.0000	
Selected (A)	0.8621	0.8502	0.6677	0.5186	0.4998	0.2000	0.2000	0.2000	0.2000	1.0000
Incremental Runoff Ratio	(B)									
	0.0000	0.0000	7.0000	3.0000	3.0000	2.2500	2.5000	3.0000	4.0000	4.0000
Case Reserve Development	nt Factor (C)									
= (B)*[1-(A)]+(A)(B next)	t)		4.2131	2.8258	2.6641	2.3280	2.6400	3.2000	4.0000	4.0000

Table A: Case Decay Next Maturity Case(Table 4 Value)/Current Case(Table 4 Value)

Note that because slightly different decay rates are used here, the case reserve development factors differ slightly from those in Table 7.

4. REFERENCES

- [1] Boor, Joseph A., 'Estimating tail development factors: what to do when the triangle runs out', Casualty Actuarial Society Foru,m Casualty Actuarial Society, Arlington, Virginia 2006: Winter, pp. 345-390
- [2] Herman, Steven C., et al., 'The estimation of loss development tail factors: a summary report', Casualty Actuarial Society Forum, Casualty Actuarial Society, Arlington, Virginia 2013: Fall, Vol. 1, pp. 31-111
- [3] Jing, Y. Lebens, J., and Lowe, S., 'Claim reserving: performance testing and the control cycle', Variance, Casualty Actuarial Society, Arlington, Virginia 2009: Vol. 03, Issue 02, pp. 161-193
- [4] Sherman, R., 'Extrapolating, smoothing, and interpolating development factors', Proceedings of the Casualty Actuarial Society, Casualty Actuarial Society, Arlington, Virginia 1984: Vol. LXXI pp. 122-155
- [5] Sherman, R., 'Techniques for projecting claims costs', Business Insurance, Chicago, Illinois, Virginia April 6, 2006

Biography of the Author

Joseph Boor is an actuary at the Office of Insurance Regulation in Florida. He has a Baccalaureate degree in Mathematics from Southern Illinois University at Carbondale, and Master's and Doctoral degrees in Financial Mathematics from Florida State University. He is a Fellow of the CAS and is a Chartered Risk Analyst. Over a long and varied career he has had roles as diverse as regulator, Chief Actuary, consultant, and regional actuary. He also contributed significantly to the CAS literature on topics such credibility procedures, tail factors, interpolation, and the commercial market cycle.

By David A. Swanson and Lucky M. Tedrow

Abstract. Although it is an analytic construct important in its own right, a stationary population is an integral component of a life table. Using this perspective, we discuss well-known and not-so-well known equalities that are found a stationary population as well as a set of inequalities. There are two parts to the set of inequalities we discuss. The first (theorem 1) is that at any given age x, the sum of mean years lived and mean years remaining exceeds life expectancy at birth when x is greater than zero and less than the maximum lifespan (When x = zero or x = maximum lifespan, then the sum of mean years lived and mean years remaining is equal to life expectancy at birth). The second inequality (theorem 2) is a generalization of the first, namely that for the entire population, the sum of mean years lived and mean years remaining exceeds life expectancy at birth. It may be that the inequality we identify as Theorem 1 is common knowledge in some circles. However, we have found no formal description of it and believe that Theorem 1 represents a contribution to the literature. Similarly, it may be the case that one would expect that Theorem 2 would hold, given Theorem 1, but we also have not found a formal description of this in the literature and believe that it also represents a contribution. Finally, we note we have not found any discussion of an equality we found embedded in Theorem 1 (when age = 0 and when age = ω , then $\lambda_x + e_x = e_0$ and believe that the identification of this equality represents a contribution. We provide illustrations of the two inequalities and discuss them as well as selected equalities.

Keywords. Carey's Equality Theorem, Two Inequality Theorems, Mean years lived, mean years remaining, life expectancy at birth, sum or mean years lived and mean years remaining, mean age at death, variance in age at death

1. INTRODUCTION

Although many of them are apparent and some that are not so apparent have been described, equalities represent a defining characteristic of stationary populations (Kintner 2004). In addition to the obvious equalities such as the crude birth rate and crude death rate, research has revealed that: (1) mean years lived is equal to mean years remaining; and (2) the distribution of age composition is equal to the distribution of remaining lifetimes(Carey et al. 2008; Rao and Carey 2014, Vaupel 2009). To these equalities, the following can be added: (1) mean age is equal to mean years lived (Rao and Carey 2014); and (2) mean age is equal to mean years remaining (Kim and Aron 1989).

As we show in this paper, mean age can be expressed as a function of total years lived by the stationary population and its life expectancy at birth, which implies that for a given stationary population, its mean age can be expressed as a function of its crude birth rate as well as its crude death rate. In turn, because mean age is equivalent to mean years lived and mean years remaining, it also can be expressed as a function of total years lived and, respectively, life expectancy at birth, the crude birth rate and the crude death rate.

To these equalities, we add a set of inequalities by demonstrating: (1) that at any given age x, the sum of mean years lived and mean years remaining exceeds life expectancy at birth in a given stationary population, where $0 < x < \omega$ (maximum lifespan); and (2) that for a stationary population as a whole, the sum of mean years lived and mean years remaining exceeds life expectancy at birth. We discuss this set of inequalities and provide an empirical illustration of them.

Before proceeding, it is worth noting that while a stationary population is an analytic construct in its own part, it is an integral component of a life table [1]. As such, the equalities and inequalities we identify and discuss apply to life tables and their construction. As our main findings, we offer: (1) Theorem 1 and provide a proof for it that shows that for a given age x, the sum of mean years lived (λ_x) and mean years remaining (e_x) exceeds life expectancy at birth where $0 < x < \omega$; (2) Theorem 2 as a generalization of Theorem 1 to all ages and provide a proof for it; and (3) an equality we found embedded in Theorem 1, namely that when age = 0 or when age $= \omega$, then $\lambda_x + e_x = e_0$

1.1 Equalities in a Stationary Population

Let the size of a stationary population be T_o

where

 $T_0 = ke_0$

and

 \mathbf{k} = radix of the life table (i.e., \mathbf{k} = 100,000) = I_0

 e_0 = life expectancy at birth (Mean years remaining at birth)

Extending the notation used by Vaupel (2009), the age distribution of a stationary population of size T_o can be described by: (1) the probability density function c(a), the distribution of years lived; (2) the probability density function $\lambda(a)$; and (3) the distribution of years remaining be described by the probability density function r(a). Note that by definition, $c(a) = \lambda(a)$. Using this notation, we can define the total number of years lived by individuals currently alive in the stationary population (T_a) and the total number of years remaining to them (T_c) , respectively, as:

(1)
$$T_{\lambda} = \int_{n}^{\omega} \alpha c(\alpha) = T_{0}\mu_{\lambda}$$

(2) $T_{r} = \int_{0}^{\omega} \alpha r(\alpha) = T_{0}\mu_{n}$

Because, as we noted earlier, $c(\alpha) = \lambda(\alpha)$,

then
$$T_{c} = \int_{0}^{\omega} \alpha c(\alpha) = T_{\lambda} = \int_{0}^{\omega} \alpha \lambda(\alpha)$$

Kim and Aron (1989) provide a proof that mean age in a stationary population is equal to mean expected years remaining. Because Vaupel (2009) demonstrated that that the mean number of years lived in a stationary population is equal to the mean expected years remaining, we can see that the three means are equivalent, using the notation just described:

(3)
$$\mu_c = \mu_r = \mu_\lambda$$

where

 $\mu_{c} = \text{mean age} = \int_{0}^{\omega} \alpha c(\alpha) d\alpha$ $\mu_{r} = \text{mean years remaining} = \int_{0}^{\omega} \alpha r(\alpha) d\alpha$ and

 μ_{λ} = mean years lived = $\int_{0}^{\omega} \alpha \lambda(\alpha) d\alpha$ Because $T_{0} = ke_{0}$, then it follows that

$$(4) T_c/T_0 = \mu_c$$

Because $\mu_c = \mu_r = \mu_\lambda$, then it follows that

$$(5) \quad \boldsymbol{T_c}/\boldsymbol{T_0} = \boldsymbol{\mu_r} = \boldsymbol{\mu_{\lambda}}$$

And because $T_0 = ke_0$, μ_c can be expressed as

(6)
$$\mu_c = T_c/ke_0$$

then it follows that

(7)
$$T_c = \mu_c k e_0$$

and

(8)
$$T_c/k = \mu_c e_0$$

In verbal terms, equation (8) states that when divided by the radix of the life table, k, the total number of years lived by those alive in the stationary population, T_c , is equal to the product of the mean age of the stationary population, μ_c , and its life expectancy at birth, e_a . When divided by the radix of the life table, the total number of years lived by those alive in the stationary population also is equal to: (1) the product of the mean number of years lived by those alive in the stationary population, μ_a , and life expectancy at birth, e_a , and (2) the product of the mean number of years remaining to those alive in the stationary population, μ_a , and life expectancy at birth, e_a and life expectancy at birth, e_a .

Further,

(9)
$$e_0 = T_c / k \mu_c$$

and because $1/e_0 = b = d$

where

b = the crude birth rate in the stationary population (k/T_0)

d = the crude death rate in the stationary population (k/T_0)

then it follows that the relationship, $\mu_c = T_c/ke_0$ can be expressed as

(10)
$$\mu_c = (T_c b)/k$$

In verbal terms, equation (9) states that when divided by the radix of the life table, k, the product of the total number of years lived by those alive in the stationary population, T_c , and the population's crude birth rate, b, is equal to the mean age of the individuals currently alive in the stationary population. This equality is the product of the force of fertility and the total years lived by those alive. Because b = d, the equality can also be viewed as the product of the force of mortality and the total years lived by those alive. These equalities should not be surprising because for a population to be stationary, the force of increments is equal to the force of decrements. Similarly, it should not be surprising that specific values of mean years lived, μ_{λ} , and mean years remaining, μ_r , also result from the specific equality of the force of increments and the force of decrements acting in concert with the total years lived in a given stationary population.

1.2 A Set of Inequalities

Theorem 1

when $0 < x < \omega$, then $\lambda_x + e_x > e_0$

Definition

 $\lambda_{\mathbf{x}} = (T_{\theta} - T_{\mathbf{x}})/I_{\theta}$ = mean years lived to age x

and

 $e_x = T_x/I_x$ = mean years remaining at age x

Corollary

when x = 0 then $\lambda_x + e_x = e_0$ since

 $(T_{\theta} - T_{\theta})/I_{\theta} + T_{\theta}/I_{\theta} = 0 + e_{\theta} = e_{\theta}$

and when $x = \omega$ then $\lambda_x + e_x = e_0$ since

 $(T_{0} - T_{\omega})/l_{0} + T_{x}/l_{x} = (T_{0} - T_{\omega})/l_{0} + T_{\omega}/l_{\omega} = (T_{0} - 0)/l_{0} + 0 = e_{0} + 0 = e_{0}$

Proof

Let
$$\lambda_{x} = (T_{0} - T_{x})/l_{0} = (e_{0}l_{0} - T_{x})/l_{0} = e_{0} - T_{x}/l_{0}$$

then $\lambda_x + e_x = e_0 - T_x/l_0 + T_x/l_x$

and except when x = 0, so that $T_x/I_0 = T_0/I_0 = e_0$

and when $T_x/l_x = T_0/l_0$ so that $e_0 - T_0/l_0 + T_0/l_0 = 0 + e_0 = e_0$

and except when $x = \omega$, so that $T_x/l_0 = T_{\omega}/l_0$

and when $T_x/I_x = T_\omega/I_\omega$, so that $e_0 - T_\omega/I_0 + T_\omega/I_\omega = e_0 - \theta/I_0 + \theta/0 = e_0 - \theta + \theta = e_0$

then $T_x/I_0 < T_x/I_x$ because $I_0 > I_x$ when x >0

Thus, $\lambda_x + e_x > e_0$ because

 $e_0 - T_x/l_0 + T_x/l_x > e_0$

Theorem 2

 $\mu_{\lambda} + \mu_r > e_0$

Proof

Because $\mu_c = \mu_r = \mu_\lambda$

then it follows that $\mu_{\lambda} + \mu_{c} = 2\mu_{c} = 2\mu_{x} = 2\mu_{\lambda}$

Because
$$e_0 = Tc/k\mu_c$$

then it follows that

 $e_0/2 = T_c/k2\mu_c$

and since $e_0/2 < e_0$

then

 $(\mu_{\lambda} + \mu_{r}) > e_{0}$

Once we have T_c and μ_c , both of which are easily obtained when $c(\alpha)$ is determined, we can determine life expectancy at birth by dividing total years in the stationary population by the product of k (remember $k = I_0$) and the mean age of the population. Because of the equalities shown earlier, e_0 also can be determined when either $r(\alpha)$ or $\lambda(\alpha)$ is found. And, of course, once e_0 is obtained, b and d can be determined, as can T_0 .

It is useful to note here that Pressat (1972: 479-480) examined the relationship between mean age of a stationary population and life expectancy at birth and found (in the notation we use):

(11)
$$\mu_c = \frac{1}{2}(e_0 + (\sigma^2/e_0))$$

where

 μ_c = mean age of the stationary population

 e_0 = life expectancy at birth

and

 σ^2 = variance in age at death

Pressat's identification of equation (11) was independently re-discovered by Morales (1989) and identified as a re-discovery by Preston (1991).

Equation (11) is particularly useful here because it provides the basis for an interpretation of the inequality given in Theorem 2, namely that $\mu_{\lambda} + \mu_r > e_0$ First, recall that as shown earlier, the mean age of the stationary population is equal to mean years lived and to mean years remaining: $\mu_c = \mu_r = \mu_{\lambda}$ and, therefore $= 2\mu_c = 2\mu_r = 2\mu_{\lambda}$. Thus, if we multiply μ_c by 2, then equation (11) can be restated as

(12)
$$2\mu_c = 2(\frac{1}{2}(e_0 + (\sigma^2/e_0))) = e_0 + (\sigma^2/e_0)$$

Because $2\mu_c =$ mean years lived (μ_{λ}) plus mean years remaining (μ_r) and because $2\mu_c = e_0 + (\sigma^2/e_0)$, we can see that the sum of mean years lived and mean years remaining is equal to the sum of life expectancy at birth and the ratio of variance in age at death to life expectancy at birth: $\mu_{\lambda} + \mu_r = e_0 + (\sigma^2/e_0)$. Further, where $\sigma^2 > 0$, then it follows that $\mu_{\lambda} + \mu_r > e_0$ and where $\sigma^2 = 0$, then $\mu_{\lambda} + \mu_r = e_0$.

Because we also know that life expectancy at birth is equivalent to mean age at death, we also can state equation (12) as:

(13)
$$2\boldsymbol{\mu_c} = \boldsymbol{\mu_d} + (\boldsymbol{\sigma^2}/\boldsymbol{\mu_d})$$

where

 μ_d = mean age at death and μ_c and σ^2 are defined as before.

Because $2\mu_c = \mu_{\lambda} + \mu_r$ we can re-express (13) as:

(14)
$$\mu_{\lambda} + \mu_{r} = \mu_{d} + (\sigma^{2}/\mu_{d})$$

where

Casualty Actuarial Society E-Forum, Fall 2016

all of the terms are as previously defined.

Thus, the sum of mean years lived and mean years remaining is equal to mean age at death plus the ratio of the variance in age at death to mean age at death. Further, where $\sigma^2 > 0$, then it follows that $\mu_{\lambda} + \mu_r > \mu_d$ and where $\sigma^2 = 0$, then $\mu_{\lambda} + \mu_r = \mu_d$.

Equation 12 provides a shortcut method for calculating the variance in e_0 (and its equivalent, mean age at death):

(15)
$$\sigma^2 = [\boldsymbol{e}_{\boldsymbol{\partial}^*}(\boldsymbol{\mu}_{\boldsymbol{\lambda}} + \boldsymbol{\mu}_{\boldsymbol{r}})] - \boldsymbol{e}_{\boldsymbol{\partial}^2}$$

This approach to calculating is simpler to implement than others (Hakkert 1987, Hill 1993, Wrycza 2014) (e.g., one can simply multiply mean age (μ_c) by 2 and substitute this in the right had side of equation [15] in place of $\mu_{\lambda} + \mu_r$). This approach also provides a meaningful estimate of σ^2 that among other desirable characteristics includes mortality at all ages (see Wryzca 2014 for a discussion of this issue), which has a range of applications (see, e.g., Schindler et al. 2012). Appendix Table 1 provides a set of such estimates using the information found in Table 1.

1.2.1 Illustration of Theorem 1

Using a 1990 USA Life Table (both sexes combined) from the Human Mortality Database (2009) as an illustration of a stationary population, we examine λ_x , e_x , and $\lambda_x + e_x$ by age, where $\omega = 110.5$ (which we set as the maximum life span; nobody lives beyond this age). Our examination is displayed by Figure 1, which provides a scatterplot of the relationship between age (x axis) and $\lambda_x + e_x$, the sum of mean years lived and mean years remaining (y axis). Life expectancy at birth for this population is 75.40 years. As shown in Figure 1, when age (x) = 0, $\lambda_x + e_x = e_0$ and when age (x) = 110.5, $\lambda_x + e_x = e_0$ The scatterplot shows that $\lambda_x + e_x$ rises non-monotonically from 75.40 years (e_0) when age = zero, reaches a maximum of 79.82 years at age 78.5, remains at this maximum to age 79.5, then monotonically declines back to 75.40 (e_0), at the maximum possible age, 110.5. As it increases, the curve is steepest from age 45 to age 79 and the decline from age 79 is steep all the way to age 110.5.



On Equality and Inequality in Stationary Populations

1.2.2 Illustration of Theorem 2

In order to empirically illustrate the inequality provided by Theorem 2 and the relationship linking it to variance in age at death (see equations (11) through (14)), we selected a (non-random) sample of complete USA life tables for years ending in zero and five from the Human Mortality Database (2009), which has an online collection of these life tables annually from 1933 to 2013. Table 1 provides these 16 empirical examples of this inequality, $\mu_{\lambda} + \mu_{r} > e_{0}$.

TABLE 1. DIFFERENCE BETWEEN THE SUM OF MEAN YEARS LIVED & MEAN YEARS REMAINING AND LIFE EXPECTANCY AT BIRTH: SELECTED USA LIFE TABLES FOR BOTH SEXES COMBINED, 1935 TO 2010 (N=16)							
YEAR	E ₀ (1)	MEAN YRS LIVED (2)	MEAN YRS REMAINING (3)	TOTAL MEAN YRS LIVED & REMAINING (4)	DIFFERENCE: (4) - (1)		
1935	60.89	35.47	35.47	70.94	10.05		
1940	63.23	35.86	35.86	71.72	8.49		
1945	65.58	36.55	36.55	73.10	7.52		
1950	68.07	37.12	37.12	74.24	6.17		
1955	69.56	37.62	37.62	75.24	5.68		
1960	69.83	37.66	37.66	75.32	5.49		
1965	70.24	37.81	37.81	75.62	5.38		
1970	70.74	38.00	38.00	76.00	5.26		
1975	72.54	38.67	38.67	77.34	4.80		
1980	73.74	39.09	39.09	78.18	4.44		
1985	74.67	39.39	39.39	78.78	4.11		
1990	75.40	39.75	39.75	79.50	4.10		
1995	75.89	39.90	39.90	79.80	3.91		
2000	76.86	40.20	40.20	80.40	3.54		
2005	77.63	40.60	40.60	81.20	3.57		
2010	78.85	41.14	41.14	82.28	3.43		

Source of data discussed in text. Calculations by authors.

As can be seen in Table 1, the difference between $\mu_{\lambda} + \mu_{r_{5}}$ on the one hand, and e_{0} on the other, declines (although not monotonically) as e_{0} increases from 1935 to 2010. The mean difference over all 16 observations is 5.37 years, with a standard deviation of 1.90. Because of Theorem 2 we know that the difference will remain positive from the re-expressed form of equation (12), namely, $\mu_{\lambda} + \mu_{r} = e_{0} + (\sigma^{2}/e_{0})$. The trend in the sample confirms that the relationship is curvilinear as expected from this same re-expressed equation. To empirically illustrate this, we constructed scatter plots of different equations and variable transformations that seemed promising using the NCSS package, version 8 (2016) and found that a quadratic model of the following form fit well: (difference 2) = A + B*(ln(e_{0})) + C*(ln(e_{0}))2, where A = 25498.4. B = -11685.8 and C = 1339.6, with R² = .9965. This model was estimated in 21 iterations with a random seed of 2695. A scatterplot of the relationship between difference and e_{0} along with the fitted model's trend line is shown in Figure 2.

In verbal terms, the explanation for the empirical illustration of the relationship found in Figure 2 and specified in the non-linear equation given by $\mu_{\lambda} + \mu_{r} = e_{0} + (\sigma^{2}/e_{0})$, is that the sum of mean years lived (μ_{λ}) and mean years remaining (μ_{r}) is equal to the mean age at death (μ_{d}) plus the ratio of the variance in age at death to mean age at death (σ^{2}/μ_{d}). Recalling that mean age at death is equal to life expectancy at birth (e_{0}), we can see that if the variance in age at death remained relatively constant (or, relatively speaking, did not increase as much as life expectancy) from 1935 to 2010 while life expectancy increased, then the difference, μ_{λ} + $\mu_{r} - e_{0}$, would decrease during the same period, which is what is shown in Figure 2. To some extent, the trend found in Figure 2 likely reflects this because other than the initial effect of the baby boom (1946-64), the US population aged between 1935 and 2010 and holding all else constant, one would expect that variance in age at death would not increase as a population ages because deaths become more concentrated in the older population, which, in turn, would be reflected in life tables constructed from such a population.





2. RESULTS AND DISCUSSION

Using Carey's equality Theorem (Carey et al. 2008, Rao and Carey 2014, Müller et al. 2004) and a 2005 life table for the United States, Vaupel (2009) estimates that more than 48 percent are 41 years or older, which implies that nearly half of the life table population will be alive in 2050, assuming that the 2005 life table holds to 2009. Using the same US life table and corresponding stationary population, we find that on average the population lived 40.60 years and will live another 40.60 years on average. If we assume that the 2005 life table applied to 2009 as did Vaupel, then on average the members will live to almost 2050, which is in agreement with Vaupel's estimate. Even without such an assumption, it is the case that on average the 2005 population lived 40.6 years and will, on average, live an another 40.6 years, or 81.3 years in total, which is 3.67 years more than their life expectancy at birth of 77.63 years. While the actual differences may vary, the proof shown earlier for Theorem 2 shows that mean years lived + mean years remaining is greater than life expectancy at birth ($\mu_{\lambda} + \mu_r >$ eo). If we apply this line of reasoning to the actual 2010 US life table, we find that on average the 2010 population lived 41.14 years and will, on average, live another 41.14 years, or 82.28 years in total, which is 3.43 years longer than this population's life expectancy at birth of 78.85. Notice that as shown in Figure 2, that this difference is less than the difference found for the 2005 life table, which is consistent with the model shown in Figure 2 and discussed at the end of the preceding section.

Vaupel (2009) notes that in regard to work by Müller et al. (2004) and Müller et al. (2007) on wildlife population dynamics, Carey's equality Theorem could be used to estimate population age structure. In regard to this application, we add that if a representative age structure is obtained for a stationary population (or one that can be made stationary with adjustments suggested by Müller et al. (2004) and Müller et al. (2007), through Vaupel's suggestion or from another method, such as a sample, then its mean age, mean years lived, and mean years remaining can be determined as can its life expectancy at birth, its crude birth rate and its crude death rate. If a representative age structure is obtained from a random sample then interval estimates of these parameters can be constructed for the stationary population in question.

In the form of λ_x and e_x , Carey's Equality Theorem also manifests itself in the data displayed as Figure 3, although somewhat imperfectly because the data are discrete rather than continuous.¹As can be seen in Figure 3, the plotted values of λ_x by age are nearly a mirror

image of the plotted values of e_x by age. The two curves cross at 39.75 years, which is the average number of years lived for this population and, also, the average number of years remaining.



Theorem 1 shows that for a given age x, the sum of mean years lived (λ_x) and mean years remaining (e_x) exceeds life expectancy at birth where $0 < x < \omega$. Theorem 2 generalizes Theorem 1 to all ages. As shown in equations (12) through (14) and the discussion directly related to these equations, we have an explanation for the inequality demonstrated in theorem 2, which is linked to the variance in age at death. For example, if variance in age at death is held constant and life expectancy (mean age at death) increases then the inequality described by theorem 2 decreases; if variance in age at death increases and life expectancy is held constant then the inequality described by theorem 2 increases.

The explanation provided for the inequality described by theorem 2 can be extended to theorem 1 by looking at the variance in age at death up to and including a given age. For example, if we are interested in the inequality found at age x, we will find that if variance in age at death up to and including age x is held constant and life expectancy (mean age at death)

increases, then the inequality described by theorem 1 decreases; if variance in age of death up to and including age x increases and life expectancy is held constant then the inequality described by theorem 1 increases.

One implication of these two related theorems is that the average longevity of all of the "living" members of a given stationary population exceeds the average number of years lived expected at birth. From a different perspective, Pressat (1972: 480) recognizes this inequality by stating that "the mean age of a stationary population is greater than half of the expectation of life." He follows this with an important observation, namely that this inequality is due to variation in individual lengths of life. This variation is why the sum of mean years lived and mean years remaining exceeds life expectancy at birth. This inequality suggests that when a life table is used for planning the future, it is worthwhile to keep in mind that life expectancy at birth understates average longevity for the "living" members of the life table population relative to the non-linear relationship found in the ratio of variance in age at death to life expectancy at birth.² As such, when this ratio is elevated then it may be preferable to use the sum of mean years lived and mean years remaining instead of life expectancy at birth in some applications. For a similar reason, this also suggests that at a given age, it may be preferable to use the sum of mean years lived to that age and mean years remaining at that age instead of simply using life expectancy at the age in question.³ Although it does not directly take into account the inequalities we have demonstrated here, work by others such as Canudas-Romo and Zarulli (2016) and Canudas-Romo and Engelman (2016) recognizes similar implications involving years lived and years remaining.

Acknowledgments

We are grateful to Robert Schoen, Anatoliy Yashin and Jakub Bijak for comments on earlier versions of this paper and to comments from participants in the 51st ARC breakout session on mortality modeling held on July 30th, 2016 and the reviewers of this submission to *CAS E-Forum*.

3. ENDNOTES

- 1. Villavicencio and Riffe (2016) provide a complete and formal proof of Carey's equality in a discrete-time framework.
- 2. In addition to Pressat (1972), Morales (1989), and Preston (1991), among others, Canudas-Romo and Engelman (2016) have examined the sum of mean years lived and mean years remaining. However, none of these authors describes the inequalities demonstrated here in the forms of theorems 1 and 2.
- 3. The ratio, σ^2/e_0 is equivalent to the coefficient of variation, as is σ^2/μ_d . As such, when making comparison across stationary populations in regard to variation in e_0 or μ_d , it is more appropriate to use these measures, respectively, instead of σ^2 . Following the observations of Pressat (1972: 480), it is worthwhile to note here that when any subject is examined from the perspective of "longevity," the inequalities we have identified

will be found where there is variation in individual longevity. Among many others, these subjects include, for example, duration of first marriage (Schoen 1975), length of working life (Yusuf, Martins, and Swanson 2014: 222-224), length of the second birth interval (Swanson 1985, 1986), length of product reliability (Ebeling 2010), age and length of time to product substitution (Martins, Yusuf, and Swanson 2012: 169-189), duration of disability (Office of the Chief Actuary 2002), and the longevity of species other than humans (Carey and Judge 2000).

4. REFERENCES

- Canudas-Romo, V., and V. Zarnulli. (2016). Am I halfway? Years Lived = Expected Life. pp. 33 50 in R. Schoen (Ed.) *Dynamic Demographic Analysis*. Springer. Dordrecht, The Netherlands.
- [2] Canudas-Romo, V., and M. Engelman. (2016). Maximum life expectancies: Revisiting the best practice trends. *Genus* 65 (1): 59-79.
- [3] Carey, J. R. and D. Judge. (2000). Longevity Records: Life Spans of Mammals, Birds, Reptiles, Amphibians and Fishes. Odense University Press. Odense, Denmark.
- [4] Carey, J.R., R., Papadopoulos, H-G Müller, B. Katsoyannos, B., N. Kouloussis, J-L Wang, K. Wachter, W. Yu, and P. Liedo. (2008). Age structure and extraordinary life span in wild medfly populations. *Aging Cell* 7: 426-437.
- [5] Ebeling, C. (2010). An Introduction to Reliability and Maintainability Engineering, 2nd edition. Waveland Press. Long Grove, II.
- [6] Hakkert, R. (1987). Lifetable transformations and inequality measures: some noteworthy formal relations. *Demography* 23: 615-622.
- [7] Hill, G. (1993). The entropy of the survival curve: an alternative measure. *Canadian Studies in Population* 20: 43-57.
- [8] Human Mortality Database. (2009) University of California, Berkeley, and Max Planck Institute for Demographic Research. <u>www.mortality.org</u>.
- [9] Kim, Y. and J. Aron. (1989). On the equality of average age and average expectation of remaining life in a stationary population. SIAM Review 31 (1): 110-113.
- [10] Kintner, H. (2004). The Life Table. pp. 301-340 in J. Siegel and D. Swanson (Eds.) The Methods and Materials of Demography, 2nd Edition. San Diego, CA: Elsevier Academic Press.
- [11] Martins, J., F. Yusuf, and D. Swanson. (2012). Consumer Demographics and Behaviour. Springer. Dordrecht, The Netherlands.
- [12] Morales, V. (1989). Mean age and life expectancy at birth in stationary populations: Research Note. Social Biology 36 (1-2): 114
- [13] Müller, H-G., J-l Wang, J. Carey, E. Caswell-Chen, C. Chen, N. Papadoupoulos and F. Yao. (2004). Demographic window to aging in the wild: Constructing life tables and estimating survival functions from marked individuals of unknown age. *Aging Cell* 3(3): 125–131.
- [14] Müller, H-G., J-L Wang, W. Yu, A. Delaigle, and J. Carey. (2007) Survival in the wild via residual demography. *Theoretical Population Biology* 72 (4): 513–522.
- [15] NCSS (2016). "General Description of NCSS 8 and its system requirements." http://www.ncss.com/download/ncss/updates/ncss-8/
- [16] Office of the Chief Actuary. (2002). Canada Pension Plan Experience Study of Disability Beneficiaries. Actuarial Study no. 1. Office of the Superintendent of Financial Institutions, Government of Canada. Ottawa, Ontario, Canada. <u>http://www.osfi-bsif.gc.ca/eng/docs/cpp_disability_paper.pdf</u>.
- [17] Pressat, R. (1972). Demographic Analysis: Methods, Results, Applications. Aldine-Atherton: Chicago.
- [18] Preston, S. (1991). Mean age and life expectancy at birth in stationary populations: Comment. Social Biology 38 (1-2): 154
- [19] Rao, A., and J. Carey. (2014). Generalization of Carey's equality and a Theorem on stationary population." *Journal of Mathematical Biology*. DOI 10.1007/s00285-014-0831-6. http://entomology.ucdavis.edu/files/203430.pdf
- [20] Schindler, S., S. Tuljapurkar, J. Gaillard, and T. Coulson (2012). Linking the population growth rate and the age-at-death distribution. *Theoretical Population Biology* 82(4): 244-252
- [21] Schoen, R. (1975). California divorce rates by age at first marriage and duration of first marriage. *Journal of Marriage and Family* 37 (3): 548-555.

- [22] Swanson, D. (1985). The Timing of Fertility: Mathematical Models and Socio-demographic Aspects of the Second Birth Interval. Unpublished Ph.D. dissertation, University of Hawai'i. Honolulu, HI.
- [23] Swanson D. (1986). Timing the second birth: Fecundability models for selected race and age groups in Hawai'i. Janasamkhya 4 (December):82-113.
- [24] Vaupel, J. (2009). Life lived and left: Carey's equality. (2009). Demographic Research 20: 7-10.
- [25] Villavicencio, F., and T. Riffe. (2016). Symmetries between life lived and left in finite stationary populations. *Demographic Research* 35: 381-398.
- [26] Wrycza, T. (2014). Variance in age at death equals average squared remaining life expectancy at death. Demographic Research 30 (50): 1405-1412
- [27] Yusuf, F., J. Martins, and D. Swanson (2014). *Methods of Demographic Analysis*. Springer. Dordrecht, The Netherlands.

Biographies

David A. Swanson is professor of sociology at the University of California Riverside, where he teaches courses on demography and statistics. His Ph.D. is from the University of Hawai'i. He is a member of the Academic Central Program of CAS. He can be reached at <u>dswanson@ucr.edu</u>

Lucky M. Tedrow is Director of the Demographic Research Laboratory at Western Washington University and holds an M.A. from the same institution. He also teaches courses on computer usage and demography. His email is Lucky.Tedrow@wwu.edu.

APPENDIX TABLE 1. ESTIMATE OF VARIANCE (σ^2) IN e_0 (MEAN AGE AT DEATH) : VARIANCE = $(e_0^* (MEAN YEARS LIVED + MEAN YEARS)$							
REMAINING)) - e_0^2							
	TOTAL MEAN YRS						
	LIVED &	VARIANCE (σ^2) IN					
e ₀	REMAINING	e ₀ (MEAN AGE AT	STANDARD				
(1)	(4)	DEATH)	DEVIATION ((σ)				
60.89	70.94	611.94	24.74				
63.23	71.72	536.82	23.17				
65.58	73.1	493.16	22.21				
68.07	74.24	419.99	20.49				
69.56	75.24	395.10	19.88				
69.83	75.32	383.37	19.58				
70.24	75.62	377.89	19.44				
70.74	76	372.09	19.29				
72.54	77.34	348.19	18.66				
73.74	78.18	327.41	18.09				
74.67	78.78	306.89	17.52				
75.40	79.5	309.14	17.58				
75.89	79.8	296.73	17.23				
76.86	80.4	272.08	16.49				
77.63	81.2	277.14	16.65				
78.85	82.28	270.46	16.45				