# Data Preparation for Predictive Modeling

Peggy Brinkmann, FCAS, MAAA
Actuary
Milliman, Inc.

April 1, 2014

**Milliman**

# Outline

- Why make a big deal about data prep?

- How to do a good data prep

- Case study

Milliman

# What is the big deal?

- You need good prep to meet actuarial Standards of Practice.

- You need good prep to get good results.

- It isn't as easy as you think.

Milliman

# ASOP 23 – Data Quality

3.3  Reliance on Data Supplied by Others

- "the accuracy and comprehensiveness of data supplied by others are the responsibility of those who supply the data"

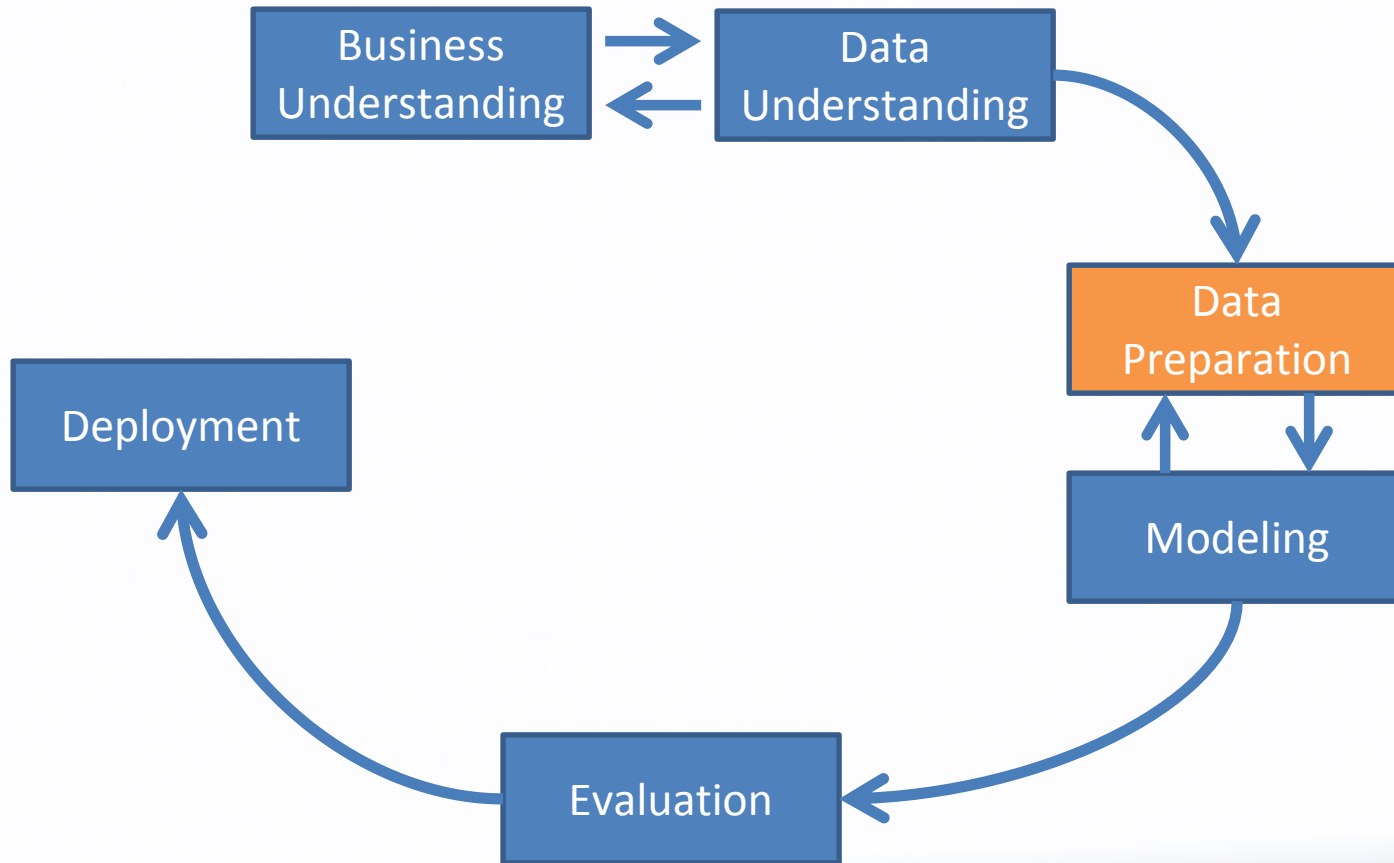**Milliman**

# ASOP 23 – Data Quality

3.5 Review of Data

- "the actuary should review the data for reasonableness and consistency, unless, in the actuary's professional judgment, such review is not necessary or not practical"

- Should consider the following
  - Data definitions
  - Questionable data values
  - Review of prior data

Milliman

# ASOP 23 – Data Quality

3.7 Use of Data

- Is the data sufficient for the analysis?

- Does it require enhancement?

- Are there material defects?

- If the data are inadequate, the actuary should obtain different data or decline the assignment

Milliman

# What can go wrong?

# Steps for good data prep

- Understand the business problem.

- Develop initial analysis plan.

- Review the raw data.

- Calculate the targets and predictors.

- Manually check calculations on examples.

- Review the prepped data.

- Document!

Milliman

# Understanding the business problem

- What is the problem area?

- What is the current solution?

- What are the business goals?

- How will the model(s) be used to achieve the goals?

- What are the implementation constraints?

- What is the desired timeline?

Milliman

# * CASE STUDY *

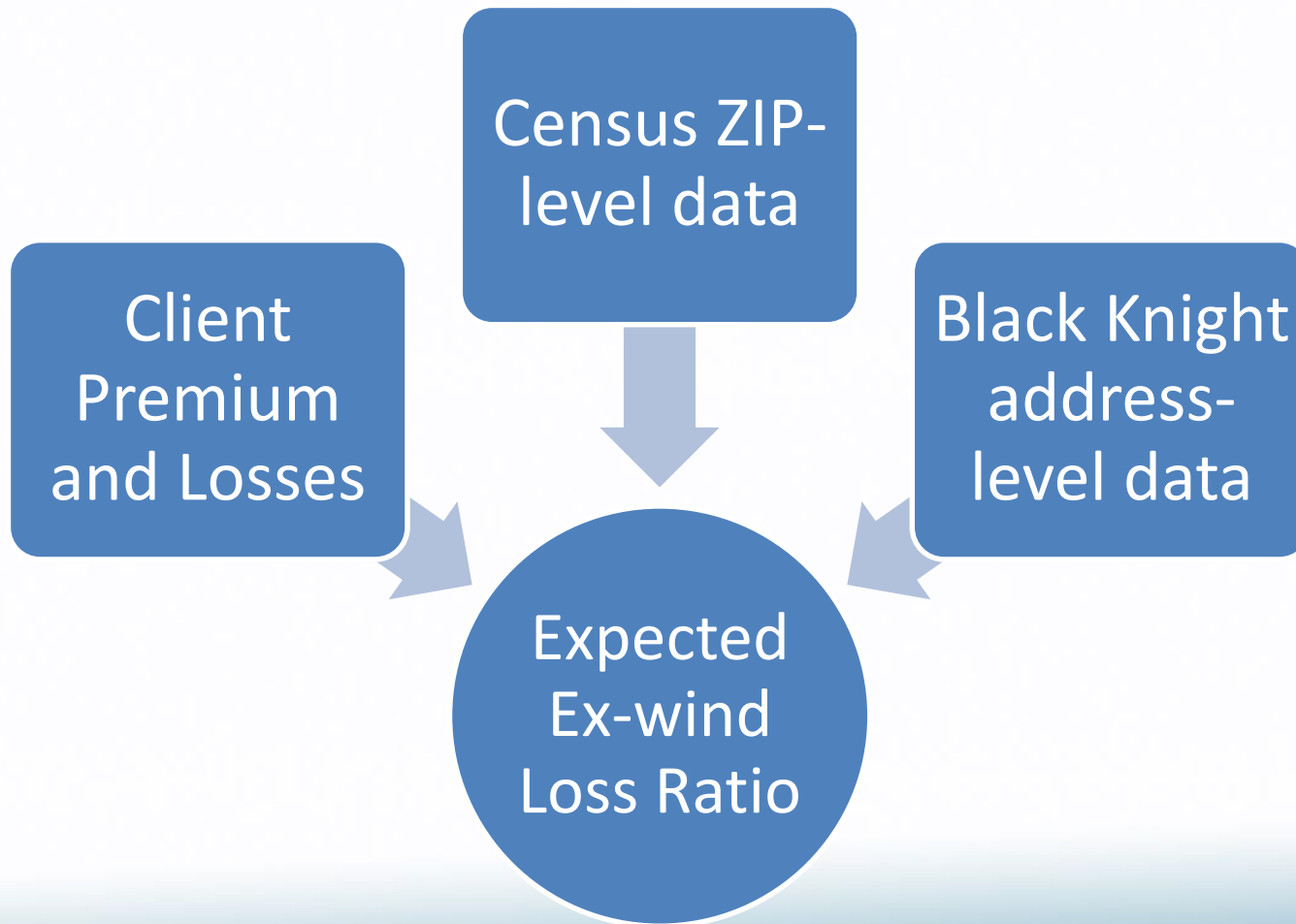| | |
|---|---|
| Problem area | Homeowners insurer in Florida wants to grow its HO3 book of business profitably. |
| Current solution | Appoint more agencies in profitable areas. |
| Business goals | Add HO3 policies in Broward, Duval, Lee, Orange, and Volusia counties. |
| Proposed approach | Build a model to estimate expected loss ratios given data available in a prospect database. Then score the prospect database, and provide list of most profitable prospects to current agencies. |
| Implementation constraints | Focus on higher-value homes, exclude current policyholders, do not use credit report data. |
| Timeline | ASAP! |

Milliman

# Develop initial analysis plan

- Is the data needed available?  Where is it stored, how to access?

- What is the sample size?

- What modeling technique(s) are going to be applied?  Using what software?

- How exactly are we defining the target variable(s)?  What adjustments will be needed?

Milliman

# * CASE STUDY *

- Need to match prospecting data to client loss ratios. Based on discussions with the client, we used Census data, as well as third party data from Black Knight, who can match onto client data by address

- Client has about 800,000 earned house years

- Use boosted scoring algorithm in EagleEye software

- Target variable will be ex-wind loss ratio at current rates. Data adjustments include
  - Extend exposures to get premium at CRL
  - Exclude wind losses

**Milliman**

# Data Overview



Census ZIP-level data

Client Premium and Losses

Black Knight address-level data

Expected Ex-wind Loss Ratio

Milliman

# Review the raw data

- Number of records received

- List of fields

- Sample of records

- Distributions of all variables – check for reasonable values

- Produce summary by year

**Milliman**

# * CASE STUDY *

QC Black Knight.txt        QC Policy Data.txt        QC Claim Data.txt

ASOP 23 questions:

- How is each data element defined?

- Are there questionable data elements?

- What enhancements/adjustments needed?

- Is the data sufficient (i.e. number of potential predictor variables, number of observations)?

- Is it consistent with data used for prior analysis?

**Milliman**

# Calculate targets and predictors

- One record per ???

- Definition of target variable

- List and describe of predictor variables

- Handling of missing values, bad data values

Milliman

# * CASE STUDY *

- One record per policy, per policy year

- Target = Ex-wind loss / Ex-wind earned premium

- Missing values – negative 1 for numeric, ~ for categorical

Milliman

# * CASE STUDY *

**Required Talon Attributes**

| Field Names | Type | Description | Notes/Calculations |
|---|---|---|---|
| POLICYNO | Character | Talon Required Field - Unique policy identifier | Combination of policynumber and effectivemonth and effectiveyear from the client file. |
| TOTEXPO | Numeric | Talon Required Field - Exposure based on the Policy | Calculated from effective, expiration, cancellation, and evaluation dates and assumes that 12 months of coverage per policy is 1 exposure. |
| PREMEA | Numeric | Talon Required Field - Premium and/or Loss associated with the Record. | earned premium calculated from onlevel premium and exposure as of 11/30/2012 |
| GROUP_ID | Numeric | Talon Required Field - Year associated with the Record such as Policy Year. | Year of effective date |
| SYS_EEA_PREMEA_Full | Numeric | Talon Required Field - PREMEA of record prior to earning to the data extraction date. | Written Premium provided on source file by policy. |
| SYS_PolicyEffectiveDate | Numeric | Talon Required Field - Policy Effective Date | Policy Effective Date put into Talon's required format |
| SYS_PolicyExpirationDate | Numeric | Talon Required Field - Policy Expiration Date | Policy Expiration Date put into Talon's required format |
| EffectiveDate | Numeric | Talon Required Field - Record Effective Date | Policy Effective Date put into Talon's required format |
| ExpirationDate | Numeric | Talon Required Field - Record Expiration Date | Policy Expiration Date put into Talon's required format |
| SYS_AgencyNo | Character | Talon Required Field - Agency Number | Agent Number |
| SYS_New_Business | Character | Talon Required Field - New Business Indicator. | Y if its effective year is the earliest for the policy, otherwise N |
| SYS_Renewed | Character | Talon Required Field - Policy Renewed Indicator. | Y if the policy renewed, N if the policy cancelled, blank if policy in force at evaluation date 11/30/2012 |
| SYS_POLICYNO | Character | Talon Required Field - Denotes the Policy number asssociated with each policy regardless of term. | Policynumber from the source file. Policynumber is not term specific and is the same across multiple terms. |
| Zipcode | Categorical | | Used to link in census data into Talon. |
| SYSRANDNUM | Numeric | Talon Required Field - System generated random number. Unique with each POLICYNO. For customer created interim data, simply include this field as a placeholder. | Blank field used as a placeholder |
| EEA_PolicyYear | Categorical | Talon Required Field - Year associated with the Record such as Policy Year. | Year of effective date |

| Field Names | Type | Description | Notes/Calculations |
|---|---|---|---|
| AOPPREMIUM_ONLEVEL | | | all other perils premium at current rate level |
| ONLEVELPREMIUM | | | |
| SINKHOLEPREMIUM_ONLEVEL | | | |
| WINDPREMIUM_ONLEVEL | | | |
| LPS_OutOfStateMail | | Flag if mailing address outside of FL | if MailState ^= "FL" then LPS_OutOfStateMail = 'Y'; |
| LPS_OwnerOccupied | | OwnerOccupied indicator | copy |
| LPS_LandUse | | LanduseDescription | Group into Residential vs non Residential |
| LPS_PurchaseDt | | LastArmsLength_RecordingDate | copy, if missing set to min(Loan1_startdt, Loan2_startdt, Loan3_startdt, Loan4_startdt); |
| LPS_PurchasePrice | | LastArmsLength_Price | copy |
| LPS_PctLandValue | | Estimated proportion of value in Lane | LPS_PctLandValue = round(MarketValueLand / MarketValueTotal * 100, 1); |
| LPS_YearBuilt | | YearBUilt | copy |
| LPS_LotSize | | LotUnit, LotSize | Restate acres to square feet |
| LPS_BuildingArea | | BuildingArea | copy |
| LPS_BuildingAreaInd | | BuildingAreaInd | copy |
| LPS_NoOfBuildings | | NoOfBuildings | copy |
| LPS_NoOfStories | | NoOfStories | copy |
| LPS_NoOfRooms | | NoOfRooms | copy |
| LPS_NoOfUnits | | NoOfUnits | copy |
| LPS_Bedrooms | | Bedrooms | copy |
| LPS_Baths | | Baths | copy |
| LPS_PartialBaths | | PartialBaths | copy |
| LPS_GarageType | | GarageType | copy |
| LPS_NoOfCars | | NoOfCars | copy |
| LPS_Pool | | Pool | copy |
| LPS_BuildingClass | | BuildingClass | copy |
| LPS_Style | | Style | copy |
| LPS_ConstructionType | | ConstructionType | copy |
| LPS_ExteriorWall | | ExteriorWall | copy |
| LPS_Foundation | | Foundation | copy |
| LPS_RoofCover | | RoofCover | copy |
| LPS_Heating | | Heating | copy |
| LPS_AirConditioning | | AirConditioning | copy |
| LPS_Elevator | | Elevator | copy |
| LPS_Fireplace | | Fireplace | copy |

Milliman

# Manual checks

- Select sample for checking – same policies from raw and prepped data
- Manually calculate target variables (e.g. earned exposures, earned premium, incurred losses, renewal) from raw data and compare to prepped
- Manually verify/calculate predictor variables from raw data

**Milliman**

# * CASE STUDY *

- Checks for onlevel premium calculation

- Checks for exposure and loss fields

- Checks for predictor variable fields

Milliman

# * CASE STUDY *

### Raw.xlsx



### Prepped.xlsx

# Review the prepped data

- Number of records received

- List of fields

- Sample of records

- Distributions of all variables – check for reasonable values

- Produce summary by year

Milliman

# * CASE STUDY *

QC Prepped Data.txt

# Documentation

- Raw data report

- Data layout

- Manual checks

- Prepped data report

- Reconciliation report (from raw to prepped)
  - Compare summaries, record counts
  - Document data adjustments (onlevel, trend, capping, etc.)
  - Document exclusions/filters (e.g. noncat loss only)

Milliman

# * CASE STUDY *

**Milliman**

650 California Street, 17th Floor
San Francisco, CA 94108-2702
USA

Tel   +1 415 403 1333
Fax   +1 415 403 1334

milliman.com

**MEMO**

**TO:** XXXXXXXXX
**FROM:** Peggy Brinkmann
**RE:** **Black Knight Data Reconciliation**

We have prepared this report to document the contents of the data files received from your client and the resulting analysis file that will be used for predictive modeling. Please review this with your clients to ensure that the data is correct before we begin the analyses.

**DATA FROM CLIENT**

Milliman received data from XXXXX for their Homeowners policies written from 2004 to November 30, 2012. The following data tables were transmitted to Milliman:

**Table 1 – Record counts**

| File Name | Date Transmitted | Record Count |
|---|---|---|
| TargetVariableFile.txt | 12/12/2012 | 1,355,137 |
| PredictiveVariableFile.txt | 12/12/2012 | 1,355,232 |
| ClaimFile.txt | 12/12/2012 | 2,255,417 |

**Milliman**

# Questions?

THANK YOU


peggy.brinkmann@milliman.com