# C-14 FINDING THE RIGHT SYNERGY FROM GLMS AND MACHINE LEARNING

CAS Annual Meeting
November 7-10

---

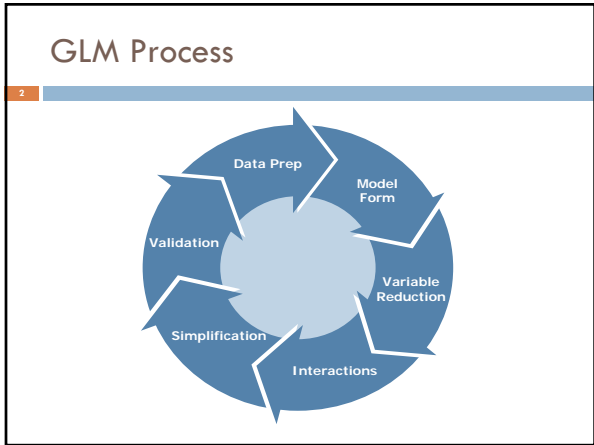## GLM Process

Data Prep

Model Form

Variable Reduction

Interactions

Simplification

Validation

---

## GLM Process

Opportunities for incorporating machine learning

Data Prep

Model Form

Variable Reduction

Interactions

Simplification

Validation

## GLM best practices
### Gathering data

- Know the data
  - Analyze one-way cuts (exposures, frequencies, severities, loss ratios, etc.)
  - Monitor changing distributions over time
  - Identify outliers and determine whether capping or removing extreme values might be appropriate

- Clean data
  - Validate between sources
  - Consider impact of nulls
  - Account for changing level definitions (i.e. territory boundary redefinition)

---

## Incorporating machine learning
### Mining the data

> Incorporate data mining to supplement data knowledge and identify adjustments.

- Often requires little up-front data prep
- Provides valuable insight into your data

---

## Incorporating machine learning
### Using decision trees

| Advantages | Disadvantages |
|---|---|
| • Easy to understand<br>• Relatively quick to run<br>• Makes no prior assumptions about the data<br>• Able to process both numeric and categorical data | • Tree structure is unstable (not in terms of prediction)<br>• Can be complex<br>• Prediction is not smooth/continuous |

## Incorporating machine learning
### Identifying important variables

| Variable | Importance |
|---|---|
| Age | 100 |
| Limits | 98 |
| Prior Accidents | 93 |
| Tier | 72 |
| Vehicle Symbol | 26 |
| Prior Convictions | 22 |
| Mileage | 21 |
| Territory | 19 |
| Model Year | 18 |
| Gender | 16 |
| … | |

□ Data mining output usually includes variable importance, useful for:
- ▫ Communication
- ▫ Determining where to start GLM
- ▫ Gaining insight into new variables

Data Prep

---

## GLM best practices
### Selecting the model form

□ Goal of GLM is to identify the signal and remove the noise.

Signal          Noise

$$y = h(\text{Linear Combination of Variables}) + \text{Error}$$

Model Form

---

## GLM best practices
### Selecting the model form

□ Error term:
- ▫ Reflects the noise of the process.
- ▫ Typically within the exponential family.

| Distribution | Common Uses |
|---|---|
| Poisson | Frequency |
| Gamma | Severity (left skewed) |
| Inverse Gaussian | Severity (right skewed) |
| Tweedie | Pure Premium |
| Binomial | Response |

□ The link function determines how variables relate to one another.

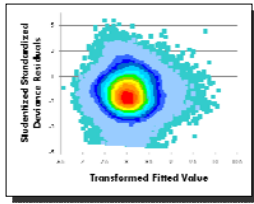| Link Function | Common Uses |
|---|---|
| Log | Multiplicative algorithm |
| Identity | Additive algorithm |
| Logit | Retention/Close Rate Studies |

## GLM best practices
### Validating the error term

- Residuals should:
  - Be symmetric around 0
  - Have constant variance across fitted values
  - Be pre-grouped for frequency ("crunched")

Model Form
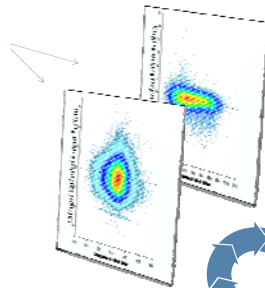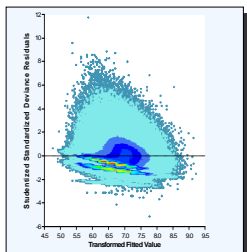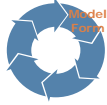
## GLM best practices
### Example: Modeling heterogeneous risks

- Bi-modal distribution can be corrected by modeling perils separately

Model Form

## GLM best practices
### Building the GLM

A balanced model is both refined and robust.



Under-fit          Over-fit

**Most robust model:**
- Mean only
- No explanatory power
- No noise

**Most explanatory model:**
- Restating history
- Full explanatory power
- Full of noise

Variable Reduction

## Slide 13

# GLM best practices
## Example: Deciding to include/exclude variables

| Variable | Chi-squared % | Deviance/AIC/BIC | Business judgment | Confidence intervals | Standard errors | Time consistency | ... | Decision |
|---|---|---|---|---|---|---|---|---|
| Model Year | | | | | | | | |

| Statistic | Impact |
|---|---|
| Chi-squared % | 0.0% |
| Deviance | -995 |
| AIC | -1026 |
| BIC | -1019 |

Variable Reduction

## Slide 14

# GLM best practices
## Example: Deciding to include/exclude variables

| Variable | Chi-squared % | Deviance/AIC/BIC | Business judgment | Confidence intervals | Standard errors | Time consistency | ... | Decision |
|---|---|---|---|---|---|---|---|---|
| Model Year | | | | | | | | |

| Statistic | Impact |
|---|---|
| Chi-squared % | 0.0% |
| Deviance | -995 |
| | 026 |
| | 019 |

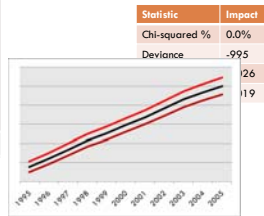Variable Reduction

## Slide 15

# GLM best practices
## Example: Deciding to include/exclude variables

| Variable | Chi-squared % | Deviance/AIC/BIC | Business judgment | Confidence intervals | Standard errors | Time consistency | ... | Decision |
|---|---|---|---|---|---|---|---|---|
| Model Year | | | | | | | | |

| Statistic | Impact |
|---|---|
| Chi-squared % | 0.0% |
| Deviance | -995 |
| | 026 |
| | 019 |

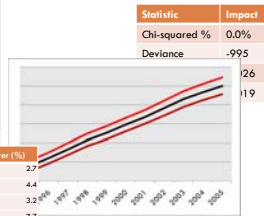| | Standard Error (%) |
|---|---|
| 1995 | 2.7 |
| 1996 | 4.4 |
| 1997 | 3.2 |
| 1998 | 7.7 |
| 1999 | 12.0 |
| 2001 | 13.5 |
| 2002 | 6.1 |
| 2003 | 4.6 |
| 2004 | 4.0 |
| 2005 | 2.8 |

Variable Reduction

# GLM best practices
## Example: Deciding to include/exclude variables

16

| Variable | Chi-squared % | Deviance/AIC/BIC | Business judgment | Confidence intervals | Standard errors | Time consistency | ... | Decision |
|---|---|---|---|---|---|---|---|---|
| Model Year | | | | | | | | |

| Statistic | Impact |
|---|---|
| Chi-squared % | 0.0% |
| Deviance | -995 |
| | 26 |
| | 19 |

Standard Error (%)

2.7
4.4
3.2
7.7
12.0
13.5
6.1
4.6
4.0
2.8

Variable Reduction

---

# GLM best practices
## Example: Inclusion/exclusion decisions

17

| Variable | Chi-squared % | Deviance/AIC/BIC | Business judgment | Confidence intervals | Standard errors | Time consistency | ... | Decision |
|---|---|---|---|---|---|---|---|---|
| Model Year | | | | | | | | |
| Vehicle Size | | | | | | | | |
| Vehicle Class | | | | | | | | |
| Horsepower | | | | | | | | |
| ... | | | | | | | | |

Variable Reduction

---

# Incorporating machine learning
## Bulk variable reduction

18

Machine learning can be particularly helpful in reducing a long list of potential variables.

- Examples:
  - Principal Components Analysis (PCA)
  - k-Means clustering
  - Decision trees
  - Pairwise correlations
  - Forward Stepwise Regression
  - …

Variable Reduction

## Incorporating Machine Learning
### Example: Bulk reduction of new variables

| | |
|---|---|
| **Initial Variables** | • Hundreds of new variables |
| **Need** | • Bulk reduction of variables that may or may not be predictive |
| **Response** | • Frequency<br>• Severity |
| **Primary Method** | • **Pairwise correlations** to identify variables that pass some minimum R^2 threshold |
| **Secondary Method** | • **Forward Stepwise Regression** |
| **Output** | • <100 variables selected for further study in GLM |

Variable Reduction

---

## Incorporating Machine Learning
### Example: Reducing highly correlated variables

| | |
|---|---|
| **Initial Variables** | • Thousands of geo-demographic variables |
| **Need** | • Bulk reduction of highly correlated variables |
| **Primary Method** | • **Principal Components Analysis** to explain the majority of the variation with a few variables |
| **Secondary Method** | • **k-Means Clustering** to pick the 'best' representative variable of the components |
| **Output** | • Few hundred unique variables that capture the majority of the variation |

Variable Reduction

---

## GLM best practices
### Identifying interactions

- Develop a list of potential interactions:
  - Brainstorm with business partners
  - Use filed rating manuals to investigate what the competition is doing
  - Study most predictive variables, especially with a wide range of predicted values
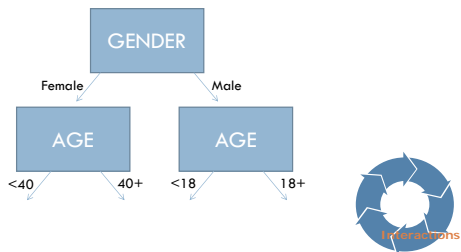- Guess and check!

Interactions

## Incorporating machine learning
### Example: Identifying interactions with trees

Decision trees can be used to identify potential interactions.

GENDER

Female | Male

AGE | AGE

<40 | 40+ | <18 | 18+

Interactions

---

## GLM best practices
### Simplifying variables and interactions

- A good GLM describes the signal with as few parameters as possible.
- Reduce the number of parameters by fitting curves to continuous variables and logically grouping categorical variables.
- Diagnostics help validate simplification decisions.

Simplification

---

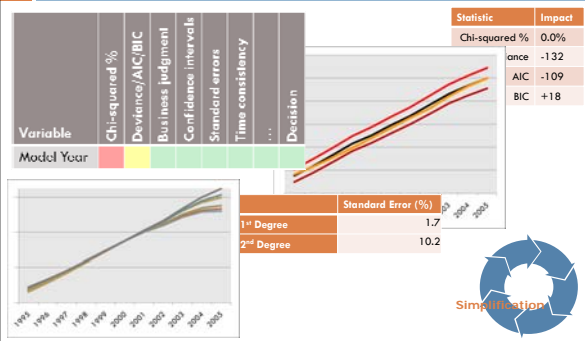## GLM best practices
### Example: Simplifying variables

| Variable | Chi-squared % | Deviance/AIC/BIC | Business Judgment | Confidence intervals | Standard errors | Time consistency | ... | Decision |
|---|---|---|---|---|---|---|---|---|
| Model Year | | | | | | | | |

| Statistic | Impact |
|---|---|
| Chi-squared % | 0.0% |
| ance | -132 |
| AIC | -109 |
| BIC | +18 |

| | Standard Error (%) |
|---|---|
| 1st Degree | 1.7 |
| 2nd Degree | 10.2 |

1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005

Simplification

## Incorporating machine learning
### Variable Simplification

> Machine learning can help determine how best to simplify the GLM.

- Examples:
  - Identify potential binning of categorical variables
  - Test whether groupings or curves are more appropriate for continuous variables
  - How to handle nulls

Simplification

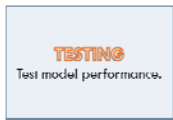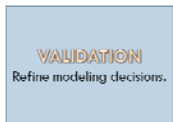## GLM best practices
### Validating with hold-out samples

- Split data for modeling and validation.

TRAINING
Build initial models.

VALIDATION
Refine modeling decisions.

TESTING
Test model performance.

Validation

## GLM best practices
### Starting with a solid model

TRAINING
Build initial models.

- Employ GLM best practice techniques.
- Be prepared to make multiple iterations through your model.

Validation

## GLM best practices
### Validating modeling decisions

**28**



VALIDATION
Refine modeling decisions.

Training —Validation

- Fit the **model structure** developed to a new set of data.
- Compare predicted values between validation and training.
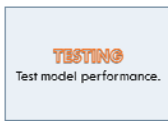- Refine and validate model decisions.

Validation

---

## GLM best practices
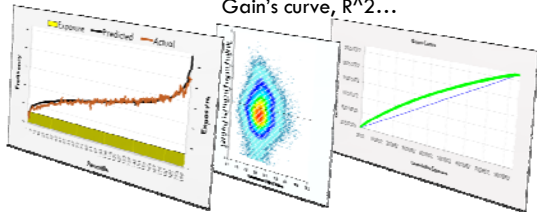### Validating model performance

**29**

TESTING
Test model performance.

- Test **model prediction** on a new set of data.
- Consider an out-of-time hold-out.
- Validate the model through a comparison with actuals, residuals, Gain's curve, R^2…



---

## Incorporating machine learning
### Validating against neural net

**30**

- Compare predictions from a well-built GLM to a neural network.

Lift        Transparency

- The neural net should offer some additional lift at the cost of transparency.
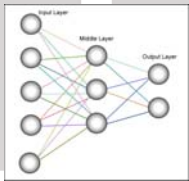- If lift is significant, consider further complication of GLMs.

Validation

## Incorporating machine learning
### Using neural networks

| Advantages | Disadvantages |
|---|---|
| • Very flexible | • Can easily be over-fit |
| • Will pick up important interactions even if not specified | • Black box, difficult to decipher useful info. |
| • Can model complex non-linear relations | |
| • Should always give equal or better model fit than a GLM | |

Validation

## Incorporating machine learning
### Example: Validating GLM vs. neural net

- Used neural net to validate performance of a less complex GLM.
- The comparison resulted in the following conclusions:
  - The less complicated model held up well.
  - The slight improvement in predictive power achieved using a neural net was not enough to justify further complicating the model.

| Validation Statistics | Additional lift using neural net |
|---|---|
| Log Likelihood | 2.7% |
| Classification rate | 1.3% |

Validation

## Summary

- •Decision Trees to investigate and adjust data

**Data Prep**

**Model Form**

Various method of variable reduction:
•Principal Components Analysis
•K-Means Clustering
•Decision Trees
•…

•Neural nets for validation

**Validation**

**GLM PROCESS**

**Variable Reduction**

**Simplification**

**Interactions**

•Data mining for potential ways to simplify

•Decision Trees to identify interactions