ChoicePoint

# The Rest of the Story: Applications and Practical Considerations of GLM & Predictive Modeling

## Fun with Factors

CAS Ratemaking Seminar 2005

New Orleans Marriott

# Raw Variable Data

## Types of Variables

- Continuous
  - Numerical Values or Statistics (Age, Miles, etc)
- Categorical
  - Discrete Classes or Groups (Symbols, Gender, etc)
- Ordinal
  - Ranks or Scores
- Spatial/Temporal
  - Boundary or Point Spatial Data (Lat/Lon, Zip, etc)

# Why create factors?

- **Challenges**
  - Low credibility in distinct classifications, either univariately or multivariately
  - Sparse values across range of continuous or ordinal values
  - Unmanageable number of distinct classifications

- **Solutions**
  - Split or combine variable(s) to create new groupings which are better suited to the GLM modelling process

**ChoicePoint**

# Creating Predictive Factors

- Methods
  - Ad Hoc
    - EDA to manually combine factor levels
  - Decision Trees/Partitioning
    - Tree based methods to split or combine variables
    - Use either leaf as final classification, or splits and combines from tree
  - Cluster Analysis
    - Combine level(s) of factors(s) to form new groupings
  - Genetic Algorithms
    - Genetic rules select "best" combination of splits or combines
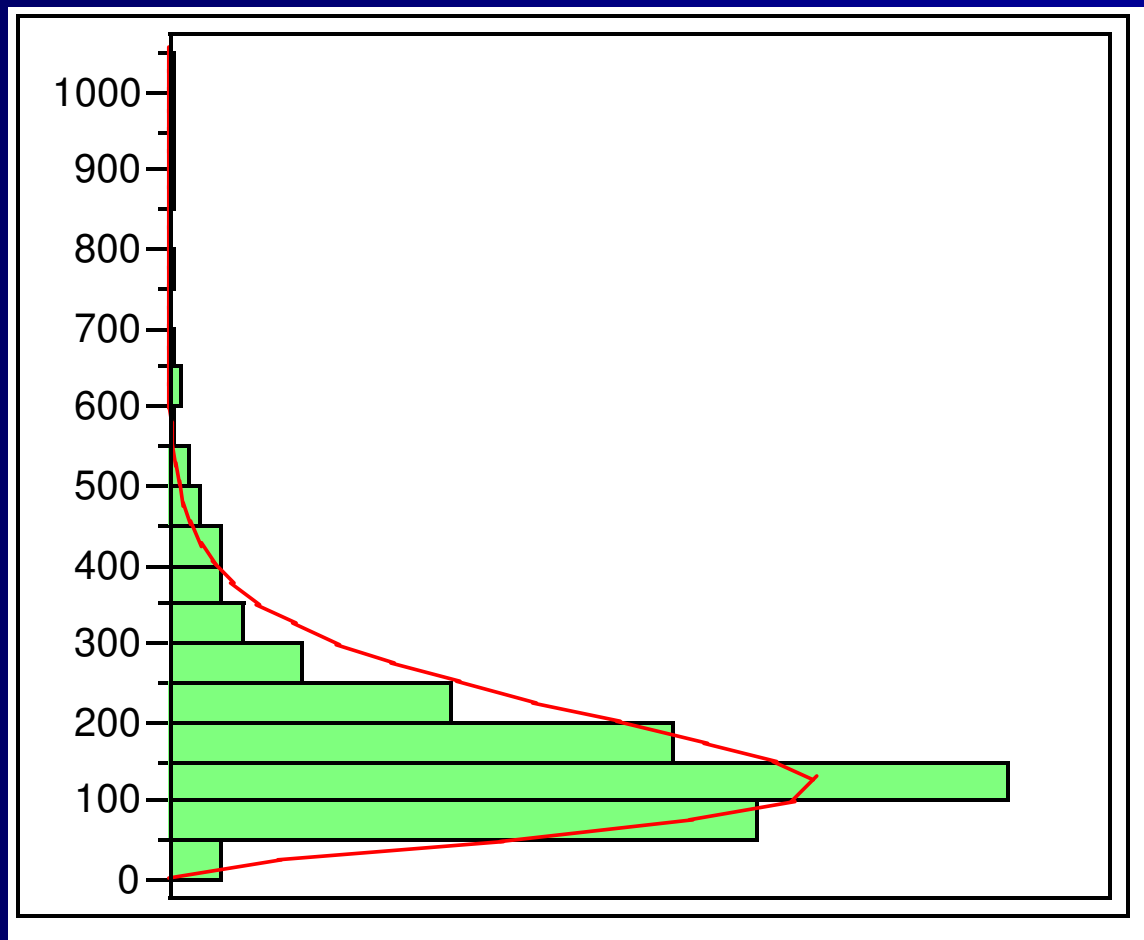
# Examples

- **Fundamentals**
  - Any type of dependent variable (loss, frequency, severity, response) can be used
  - Can accommodate measures of credibility appropriate to method
- **Data**
  - Personal Auto Loss Data
  - Modeling Losses
  - Multiple potential predictive variables
  - Low multivariate credibility and predictive power
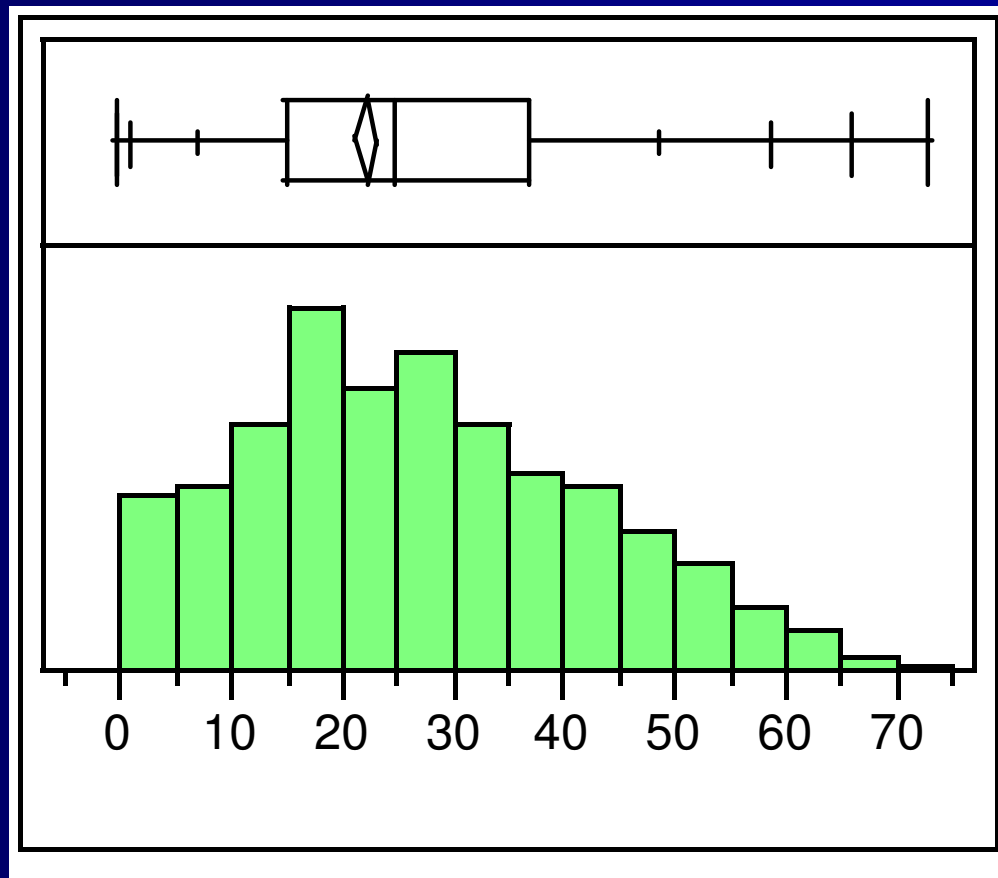
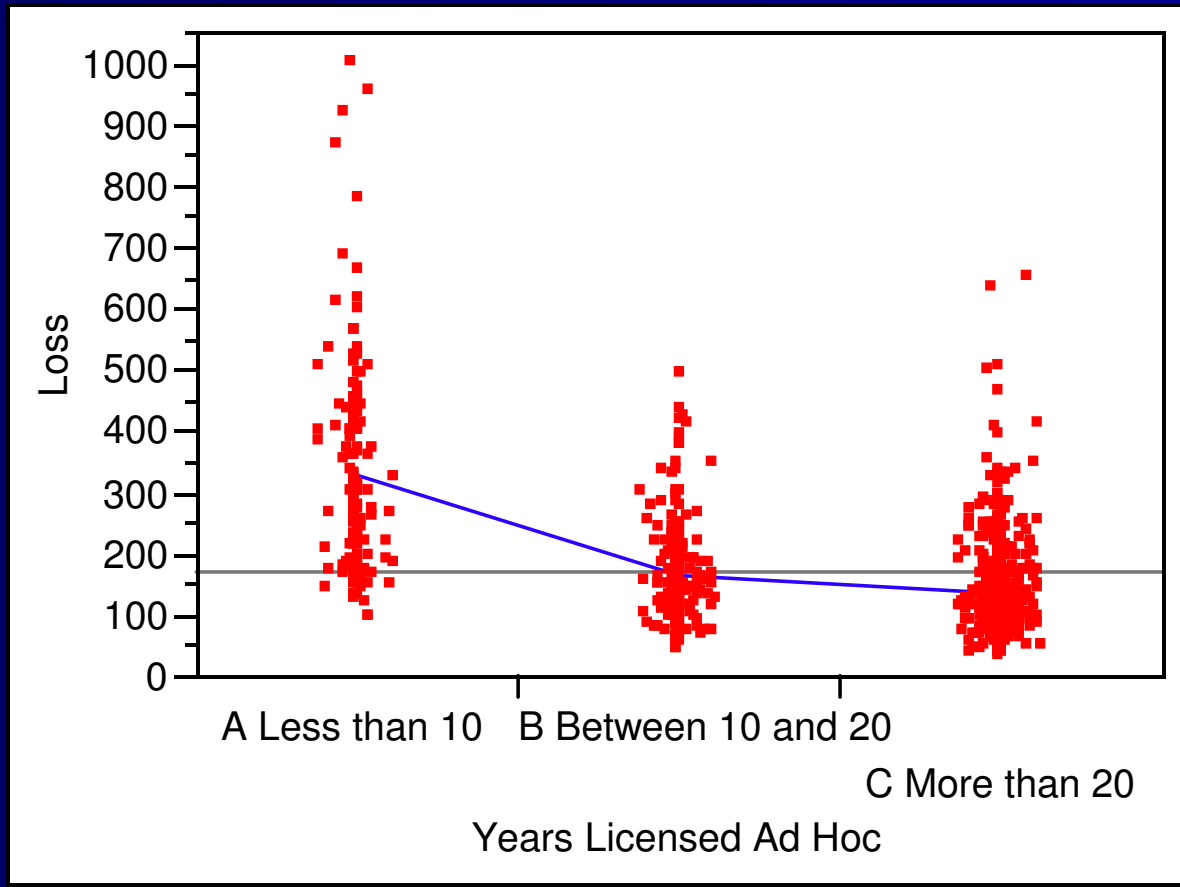# Loss Distribution

- Gamma distributed Auto losses
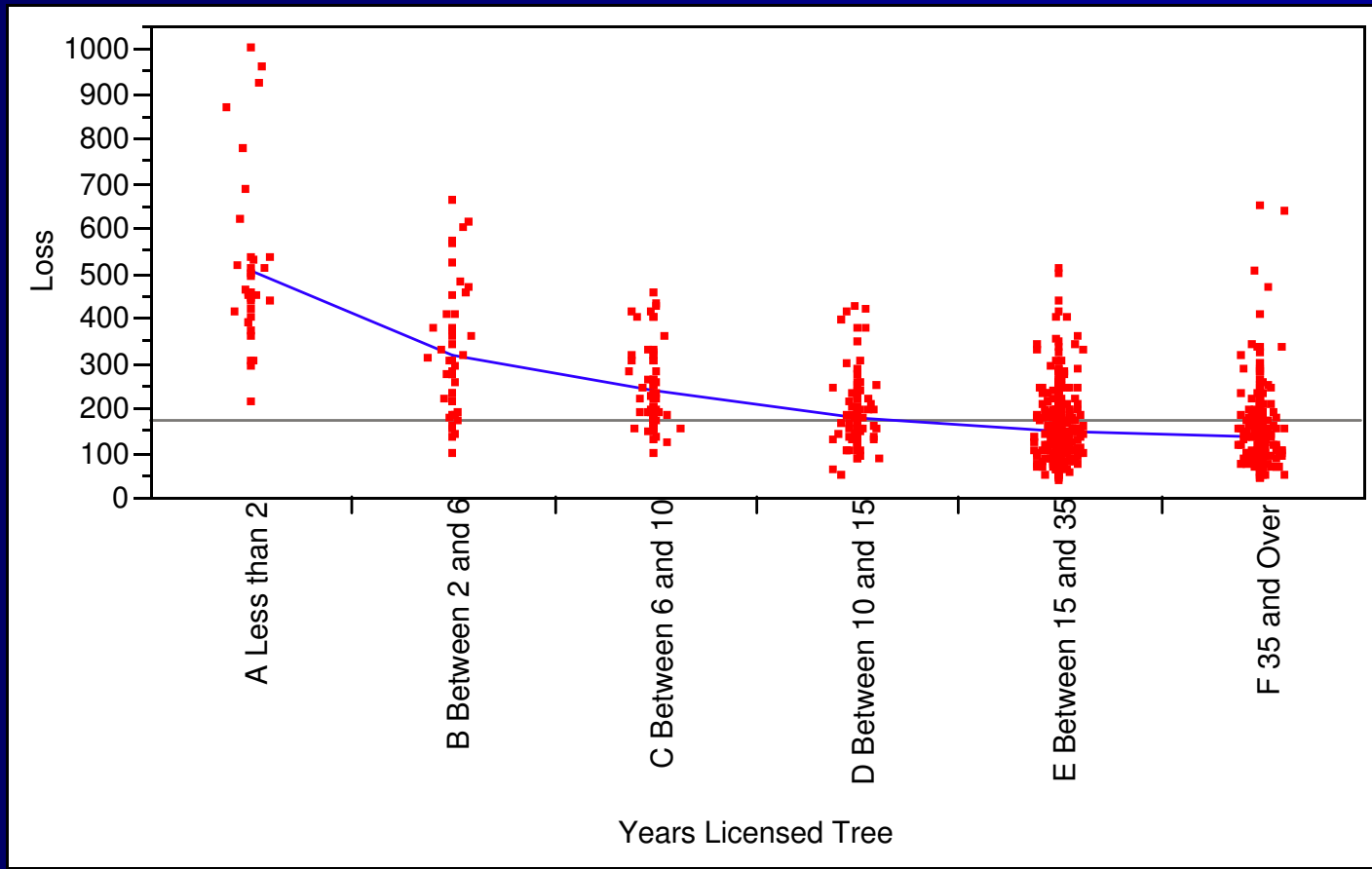
# Continuous Variable: Years Licensed

Min=0, Max=73, Mean=22.4, Median=25

# Ad Hoc Categorization

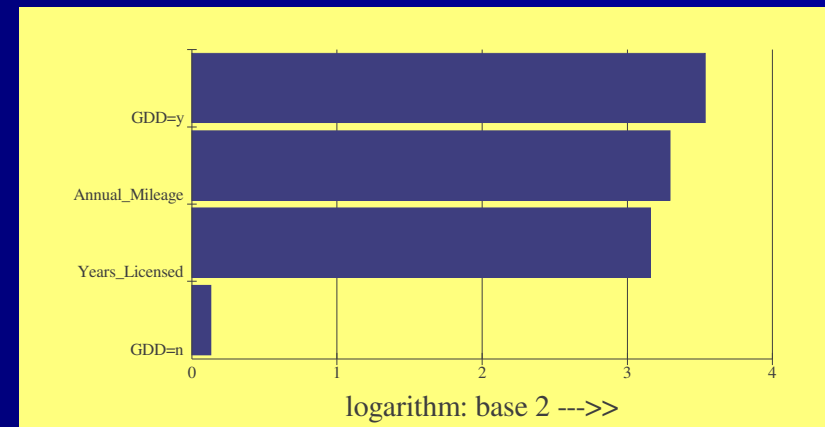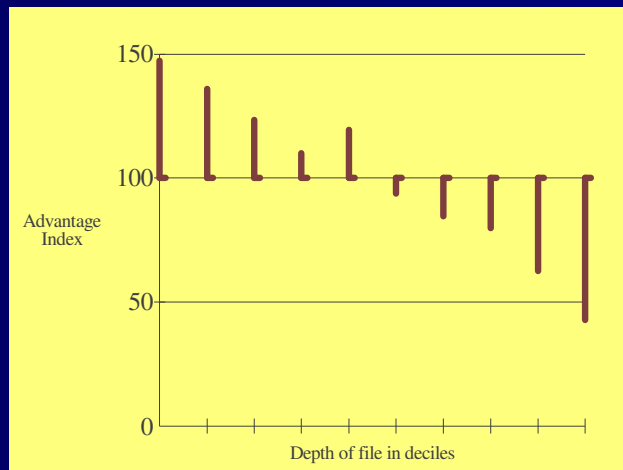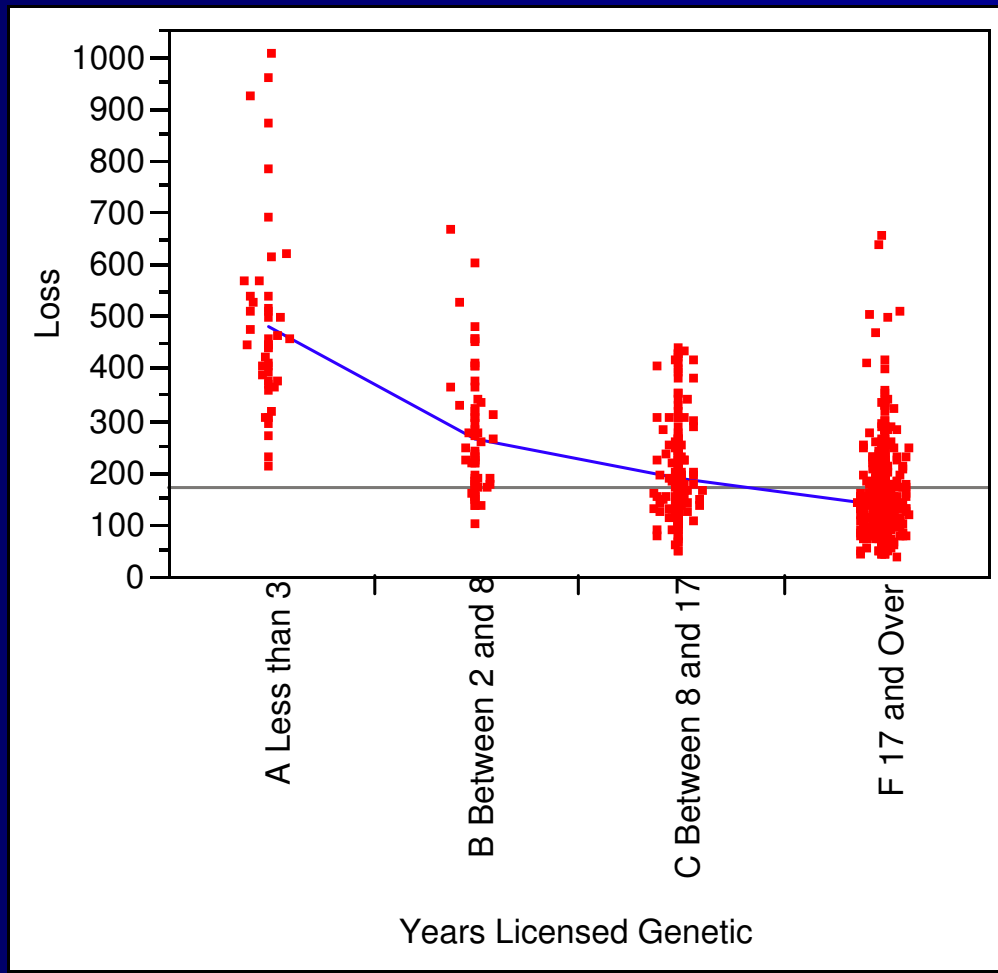# Partition – Multivariate with GDD & Miles



9

# Partition Categorization

# Genetic – Multivariate with GDD & Miles



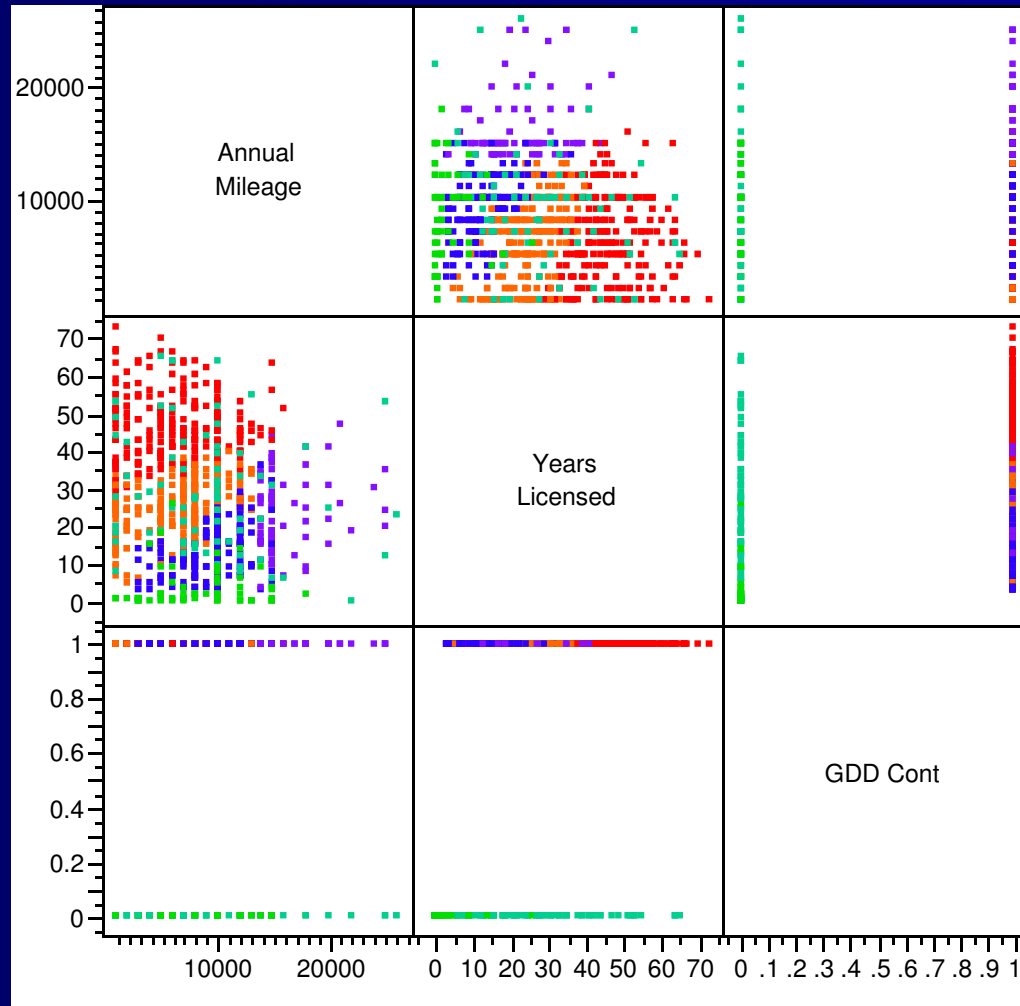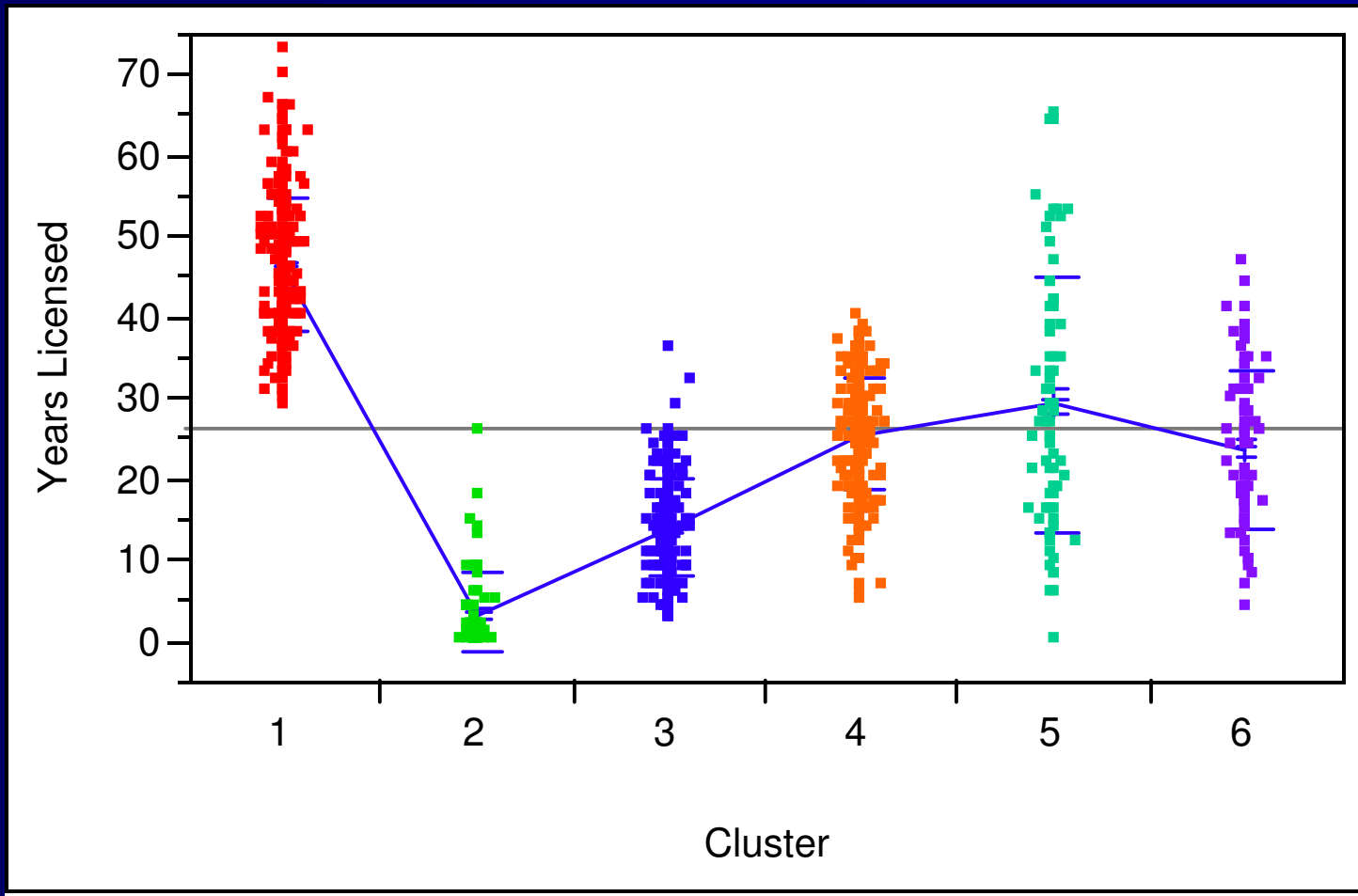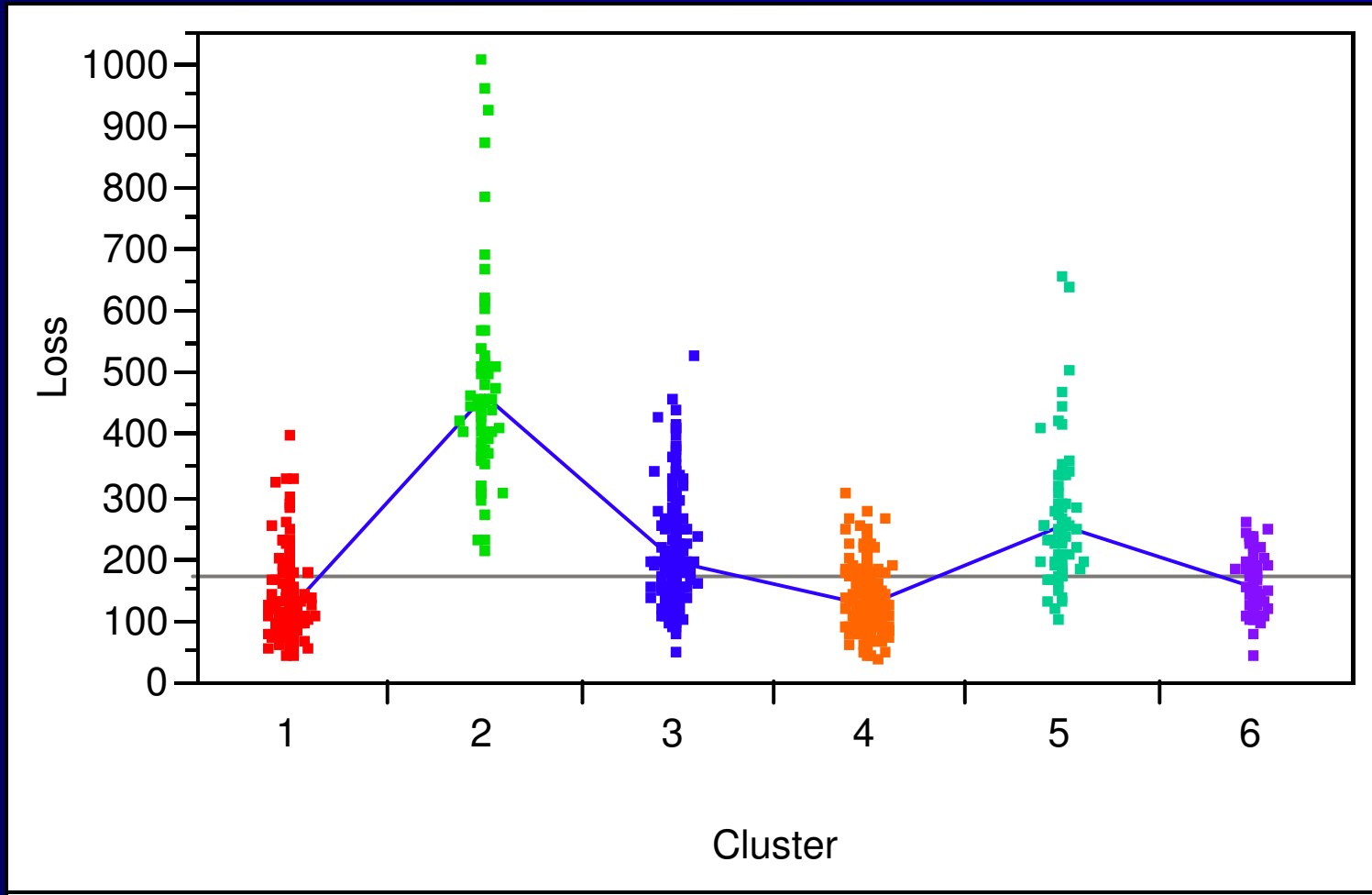|  | Predicted Profit Cum. Gain | | Avrg Profit Min score | | Cum. Avrg Max score | |
|---|---|---|---|---|---|---|
| Top | 1 | 0.01 | 0.01 | 141 | 1.006 | 1.067 |
| 2nd | 1 | 0.01 | 0.01 | 128 | 1.003 | 1.006 |
| 3rd | 1 | 0.01 | 0.01 | 114 | 1.001 | 1.003 |
| Bottom | 1 | 0.00 | 0.01 | 100 | 0.000 | 1.001 |

# Genetic Categorization

# Cluster – Multivariate with GDD & Miles
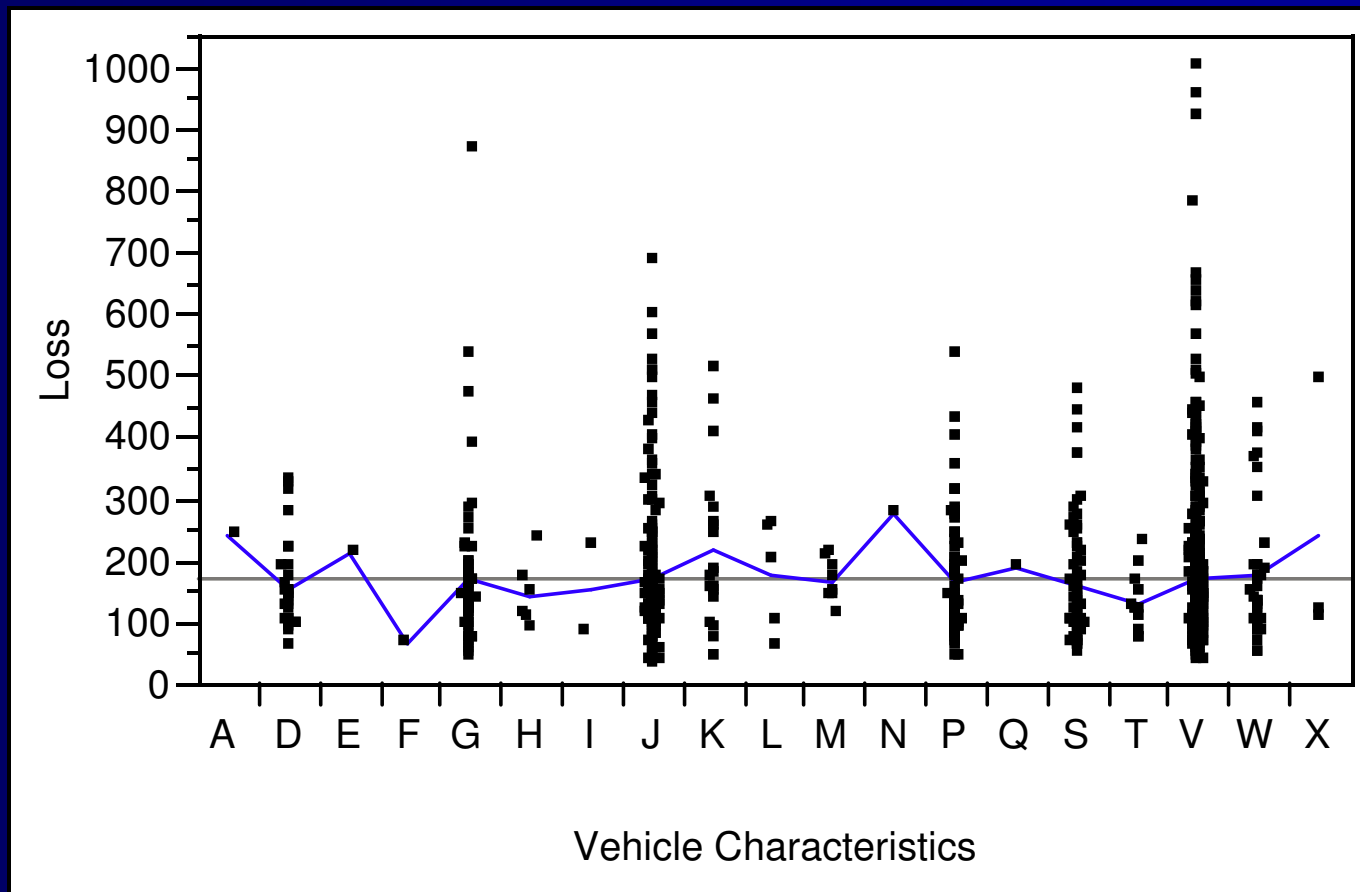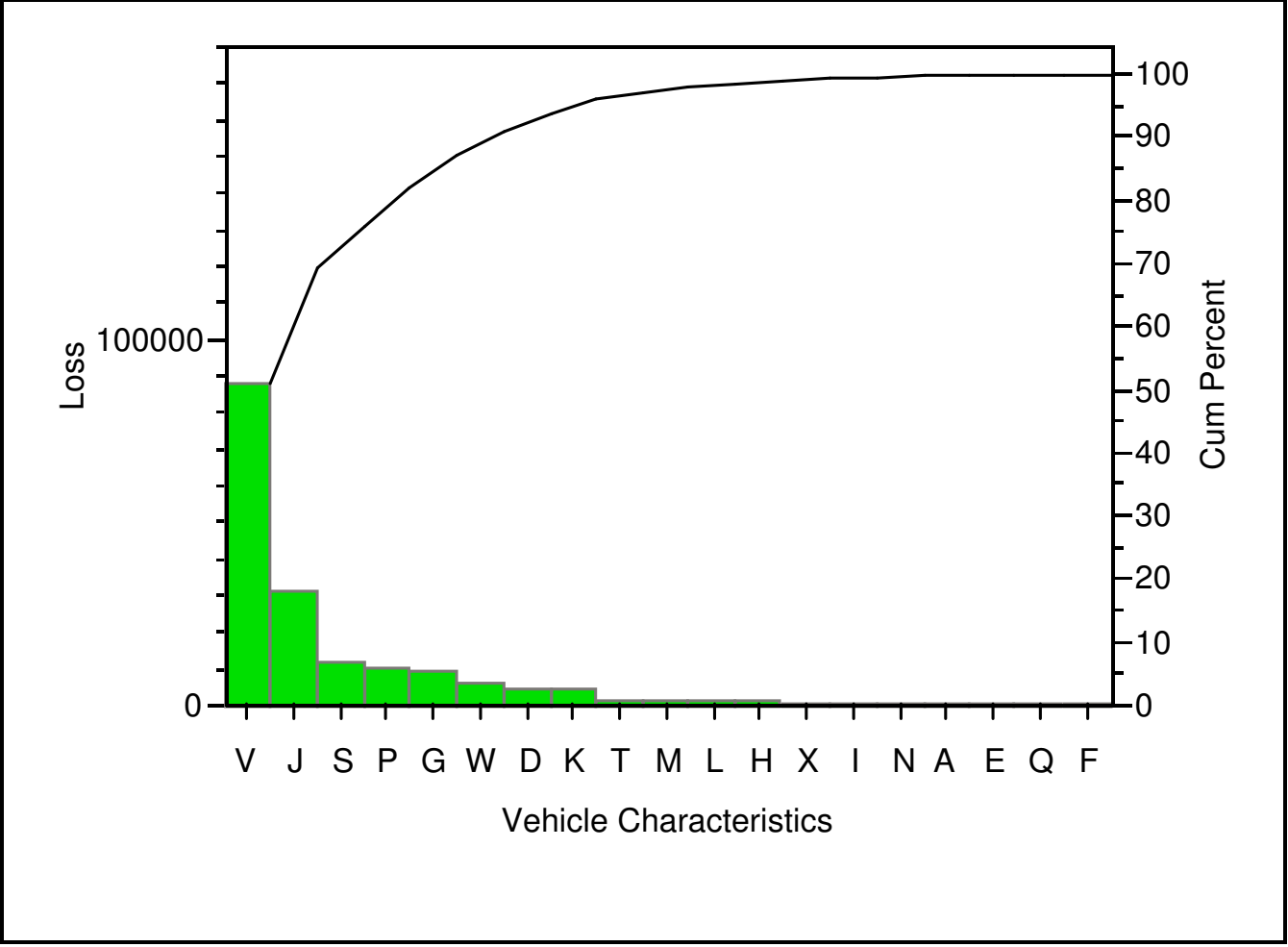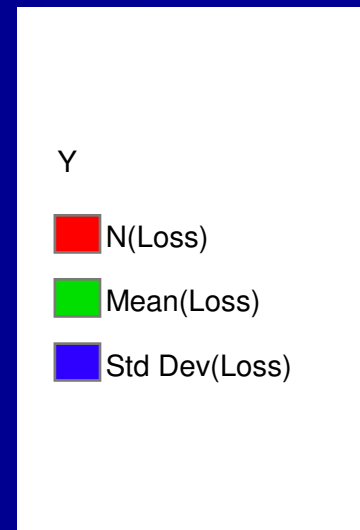
# Cluster View
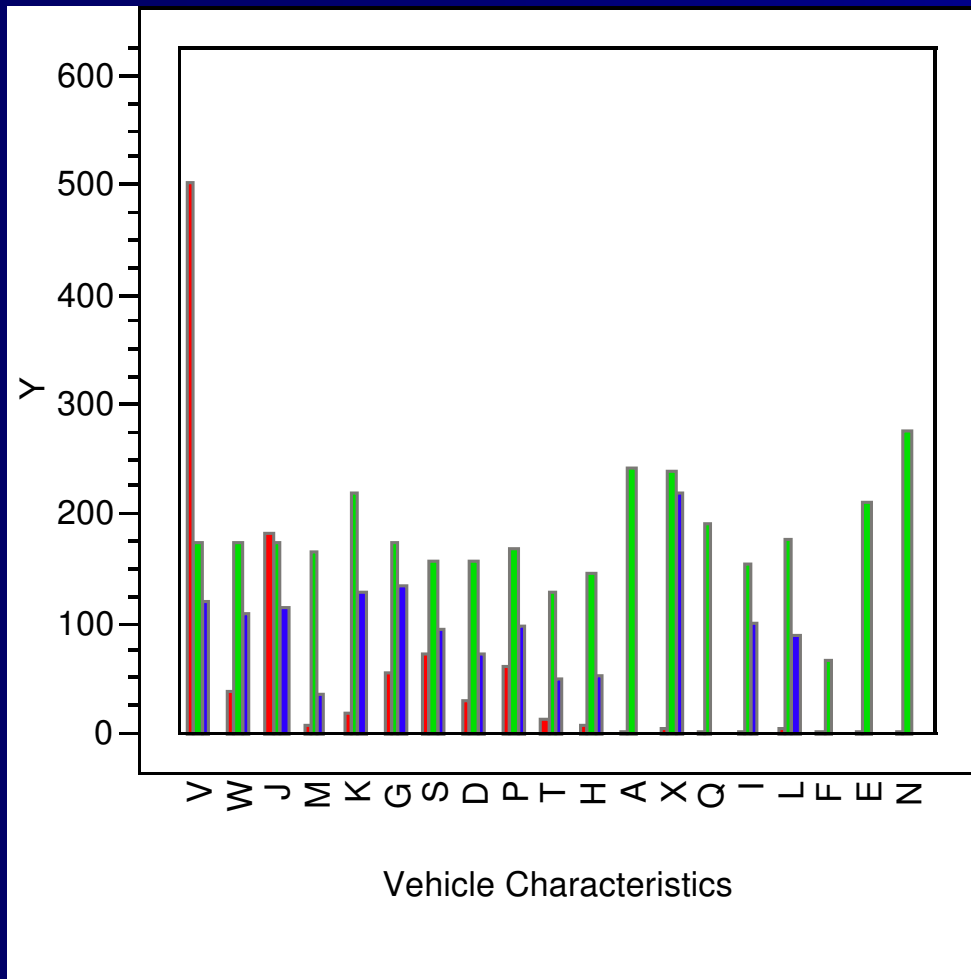
- Wards

14

# Cluster Categorization

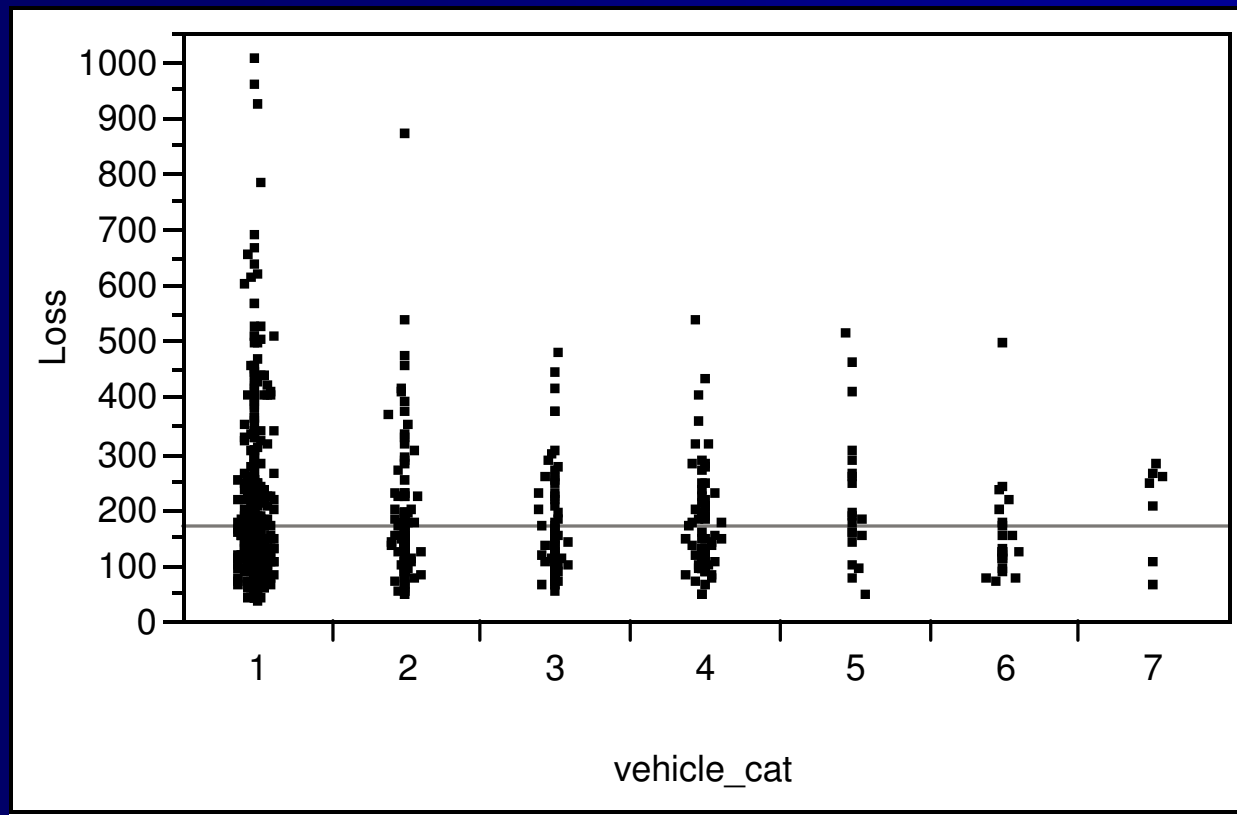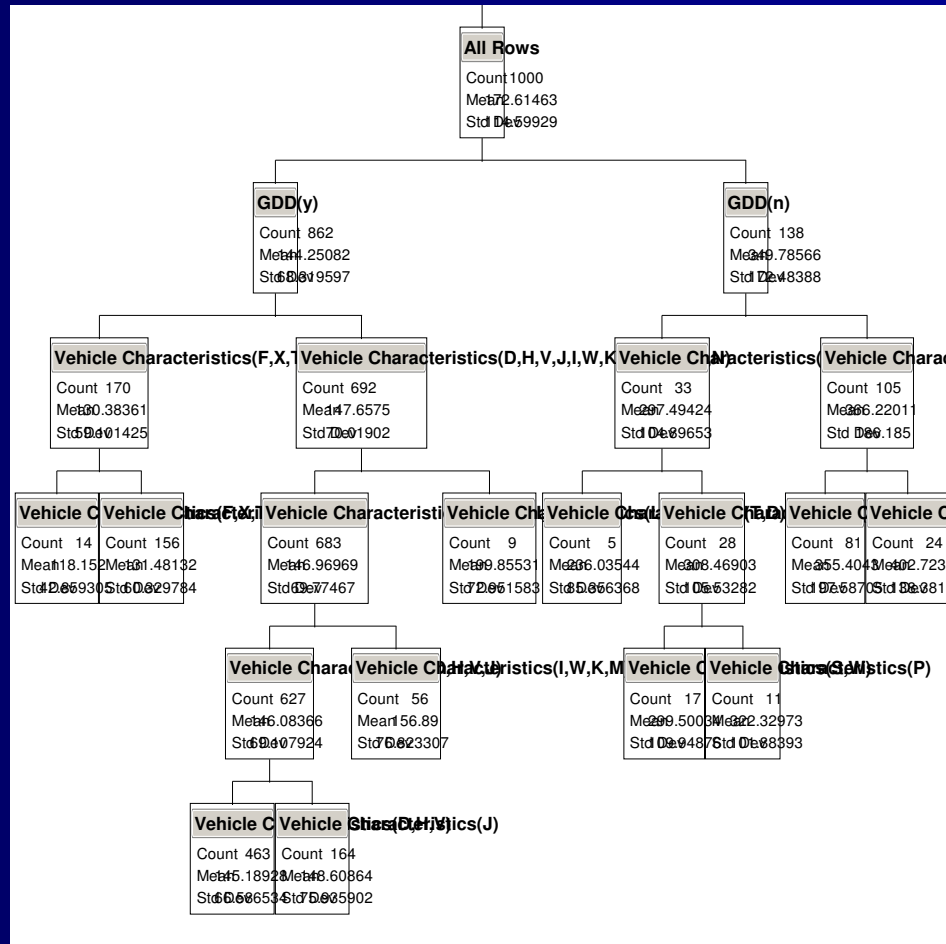# Categorical Variable: Vehicle Characteristics
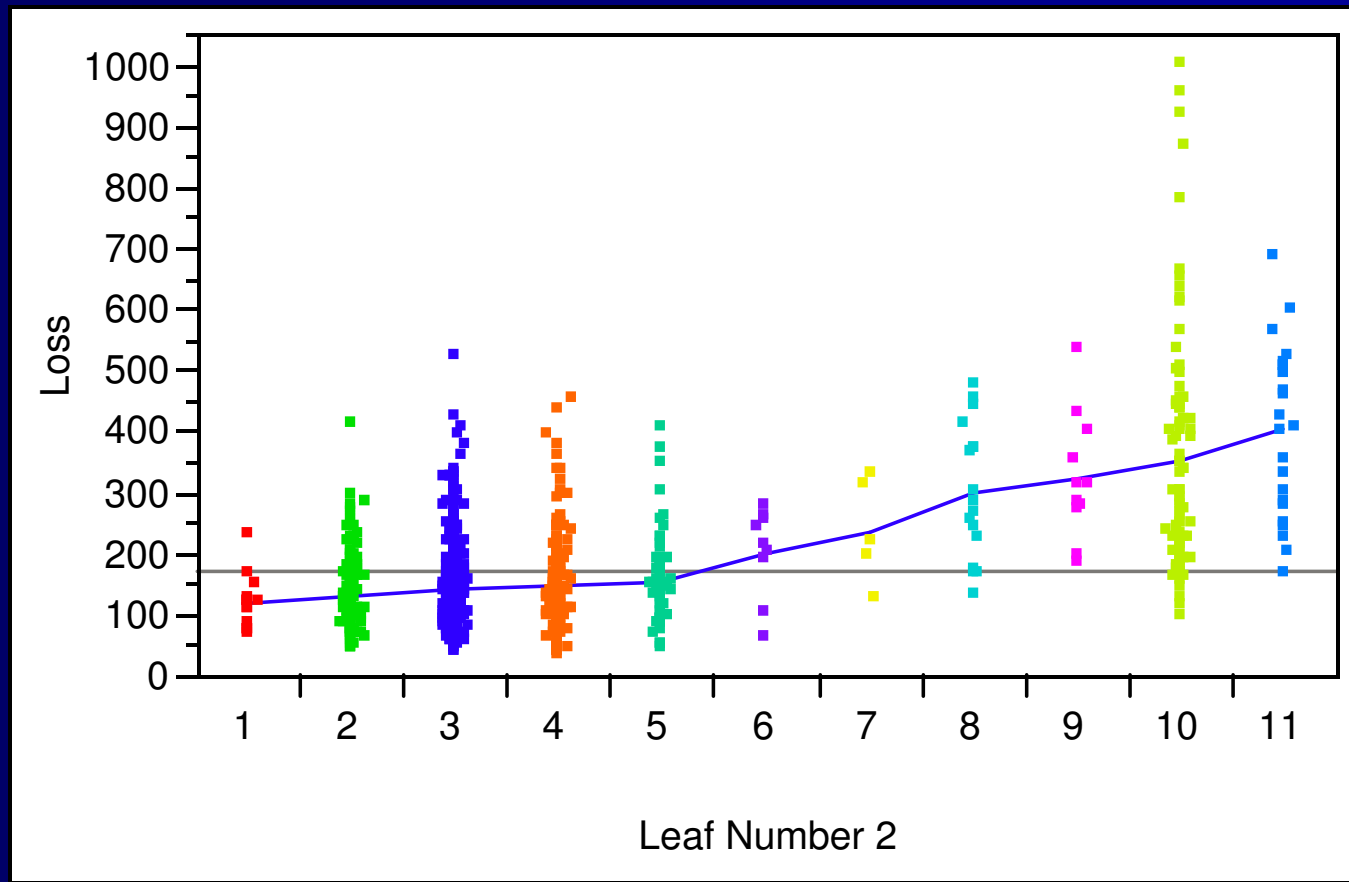
# Ad Hoc

# Empirical View

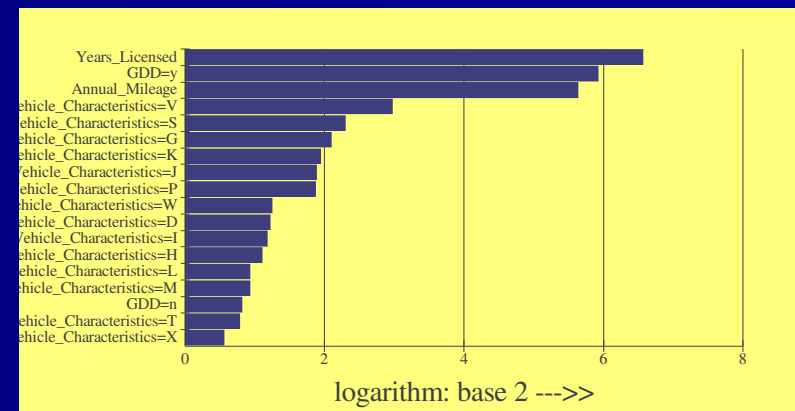# Ad Hoc Classification

# Partition – Multivariate with GDD

# Partition Classification

# Genetic – Multivariate

- **GDD, Annual Miles, Years Licensed**





| | Predicted Profit Cum. Gain | | Avrg Profit Min score | | Cum. Avrg Max score | |
|---|---|---|---|---|---|---|
| Top | 1 | 0.01 | 0.01 | 140 | 0.0098 | 0.1258 |
| 2nd | 1 | 0.01 | 0.01 | 126 | 0.0052 | 0.0098 |
| 3rd | 1 | 0.01 | 0.01 | 114 | 0.0023 | 0.0052 |
| Bottom | 1 | 0.00 | 0.01 | 100 | 0.0000 | 0.0023 |

# Genetic Classification

# Methods for Other Factor Types

- **Spatio-Temporal**
  - Correlated in one or two dimensions
  - Spatial smoothing
  - Geostatistical methods (Variogram modeling, kriging)
- **Ordinal**
  - Specialized Methods
  - Marketing Research Preference Studies

# Final Notes

- Model continuous variables if sufficient range and credibility
- Check at least bivariate with best predictor prior to modeling
- Best modeling practices and full model evaluation

# Software Used

- **SAS/JMP (Graphics, Partition, Cluster)**
  - http://www.jmp.com
- **ANGOSS Knowledgeseeker (Partition)**
  - http://www.angoss.com
- **SAS/STAT (Data manipulation, GLM, Cluster)**
  - http://www.sas.com
- **Minetech GMax (Genetic Selection)**
  - http://www.minetech.com

# Contact Information

Kate Phinney

Insurity

170 Huyshope Avenue

Hartford, CT 06106

860-616-7413

Kate.Phinney@Insurity.com

# Thank you!