**TOWERS PERRIN**

**TILLINGHAST**

# Techniques for Dimension Reduction – Variable Selection with Clustering

## CAS Special Interest Seminar on Predictive Modeling

**Robert Sanche**

**October 5, 2006**

# Contents

- **Predictive Variables**

  - Economics of Data Storage

  - Sources of Data

  - Redundancy of Variables

- **Dimension Reduction**

  - Goals of Predictive Modeling

  - Model Generalization

  - Clustering Analysis for Dimensions Reduction

- **Variable Clustering**

  - Description of Variable Clustering

  - Selection of the Cluster Representative

  - Example of Variable Clustering

- **Conclusion and Benefits of Variable Clustering**

# Economics of Data Storage

"In 1956, IBM sold its first magnetic disk system, RAMAC (Random Access Method of Accounting and Control). It used 50 24-inch metal disks, with 100 tracks per side. It could store 5 megabytes of data and cost $10,000 per megabyte. (As of 2005, disk storage costs less than $1 per gigabyte)."
http://en.wikipedia.org/wiki/History_of_computing_hardware

- 1 gigabyte = 130 numeric characteristics

  — for 1 million policies

  — for $1.00

# Sources of Data

- New data sources

    - Data warehousing (coverage and claims)

    - External sources

        — Geo-demographics

        — Meteorological

        — Policyholder, household, business owner, company information or agent

    - Other

# External Data (Census)

# Census (Geo-demographics)

- Population

  - Average household size

  - Median household size

  - Population density

  - Proportion of household with more than 4

  - Etc.

# Meteorological (Environmental Canada)

# Meteorological (Temperature)

| Temperature | Days with Minimum Temperature |
|---|---|
| Daily Average (°C) | > 0 °C |
| Standard Deviation | <= 2 °C |
| Daily Maximum (°C) | <= 0 °C |
| Daily Minimum (°C) | < -2 °C |
| | < -10 °C |
| **Degree Days** | < -20 °C |
| Above 24 °C | < - 30 °C |
| Above 18 °C | |
| Above 15 °C | **Days with Maximum Temperature** |
| Above 10 °C | <= 0 °C |
| Above 5 °C | > 0 °C |
| Above 0 °C | > 10 °C |
| Below 0 °C | > 20 °C |
| Below 5 °C | > 30 °C |
| Below 10 °C | > 35 °C |
| Below 15 °C | |
| Below 18 °C | |

# Meteorological (Precipitation)

| Precipitation | Days with Rainfall |
|---|---|
| Rainfall (mm) | >= 0.2 mm |
| Snowfall (cm) | >= 5 mm |
| Precipitation (mm) | >= 10 mm |
| Average Snow Depth (cm) | >= 25 mm |
| Median Snow Depth (cm) | |
| Snow Depth at Month-end (cm) | **Days With Snowfall** |
| | >= 0.2 cm |
| **Days with Precipitation** | >= 5 cm |
| >= 0.2 mm | >= 10 cm |
| >= 5 mm | >= 25 cm |
| >= 10 mm | |
| >= 25 mm | **Days with Snow Depth** |
| | >= 1 cm |
| | >= 5 cm |
| | >= 10 |
| | >= 20 |

# Redundancy of Variables

- External sources of data are highly redundant

- Note that the data is almost exclusively numeric

    - This fact is primordial in order to use variable clustering

# Goals of Predictive Modeling

- **Predictive model**

  - $Y = \alpha_1 X_1 + \ldots + \alpha_n X_n + \beta$

    — **n is universe of all available predictors**

- Goal of predictive modeling

  - Obtain coefficients for $\alpha$'s and $\beta$

- Additional goal

  - Predictive of future results

  - Model generalizes well over time

# Model Generalization

- As the number of variables increases and the model complexity increases, the potential of <u>over-fitting</u> the input data increases

- Dimensions reduction

  - Clustering (K-Means)

    — Rows

  - variable clustering

    — Columns

    — Alternatives (Factor, PCA, One-way)

# Clustering Analysis for Dimensions Reduction

- "**Cluster Analysis** is a set of methods for constructing a sensible and informative classification of an initially unclassified set of data, using the variable values observed on each individual"
  B.S. Everitt , *The Cambridge Dictionary of Statistics*, 1998

- Divide set of data (variables) into groups of similar characteristics

- Unsupervised learning technique

- Useful only when there is **redundancy** in the data

# Description of Variable Clustering

- **Variable clustering** divides a set of <u>numeric </u>variables into clusters.

- A large set of variables can be replaced by a single member (cluster representative).

- Reduce the number of variables

    - More difficult to identify irrelevant variables than redundant variables

- $Y = \alpha_1 X_1 + \ldots + \alpha_m X_m + \beta$

    - **where m<n**

# Selection of the Cluster Representative

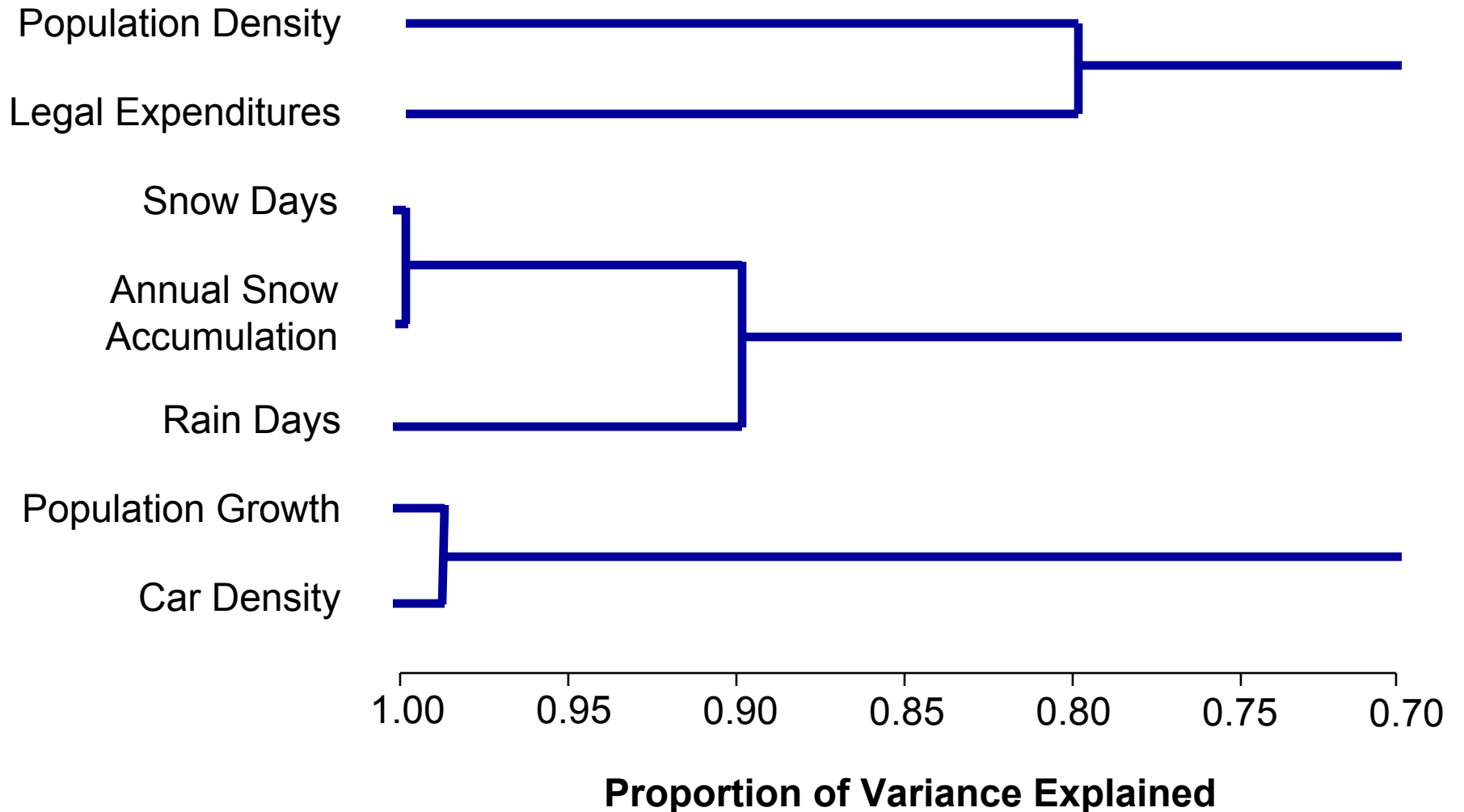$$1\text{-}R^2{}_{ratio} = (\ 1\text{-}R^2{}_{own}\ )/(\ 1\text{-}R^2{}_{nearest}\ )$$

- Intuitively, we want the cluster representative to be as closely correlated to its own cluster ($R^2{}_{own} \rightarrow 1$) and as uncorrelated to the nearest cluster ($R^2{}_{nearest} \rightarrow 0$).

- Therefore, the optimal representative of a cluster is a variable where $1\text{-}R^2$ ratio tends to zero

# Example of Variable Clustering

| 3 CLUSTERS | | R-SQUARED WITH | | 1-R² Ratio |
|---|---|---|---|---|
| Cluster | Variable | Own Cluster | Next Closest | |
| Cluster 1 | Rain Days | 0.5995 | 0.0426 | 0.4183 |
| | Snow Days | 0.8976 | 0.0317 | 0.1095 |
| | Annual Snow | 0.8940 | 0.0314 | 0.1095 |
| Cluster 2 | Population Density | 0.9804 | 0.0228 | 0.0201 |
| | Car Density | 0.9804 | 0.0113 | 0.0199 |
| Cluster 3 | Population Growth | 0.6459 | 0.0911 | 0.3896 |
| | Legal Expenditures | 0.6459 | 0.0013 | 0.3546 |

# Conclusion and Benefits of Variable Clustering

- Variable clustering reduces the amount of variables available for predictive modeling (GLM, etc.)

- The predictive modeling process using variable clustering

  - Produces a model that generalize well over time

  - Increases interpretability of the results

  - Reduces time spend on variables selection