



### Introduction to Text Mining

THE POWER TO KNOW<sub>®</sub> Tom Bohannon TAIR Conference February 2013

### **Objectives**

- Define text mining and identify text mining applications.
- Survey applications of text mining.
- Use an example to illustrate text mining concepts.
- Examine how text mining fits into modern data mining projects.

### What Is Text Mining ?

- Text mining is a process that employs a set of algorithms for converting unstructured text into structured data objects and the quantitative methods used to analyze these data objects.
- "SAS defines text mining as the process of investigating a large collection of free-form documents in order to discover and use the knowledge that exists in the collection as a whole." (SAS<sup>®</sup> Text Miner: Distilling Textual Data for Competitive Business Advantage)

### **Text Mining – Two General Goals**

- Pattern Discovery (Unsupervised Learning)
  - Identify naturally occurring groups (classification\*).
  - Derive convenient segments (clustering).
- Prediction (Supervised Learning)
  - Input variables are associated with values of a target variable.
  - Derive a model or set of rules that produces a predicted target value for a given set of inputs.
- \* Classification with a target variable is prediction.

### **Text Mining**

Text mining has the following characteristics:

- operates with respect to a corpus of documents
- employs a dictionary to identify relevant terms
- accommodates a variety of metrics to quantify the contents of a document within the corpus
- derives a structured vector\* of measurements for each document relative to the corpus
- employs analytical methods applied to the structured vector of measurements based on the goals of the analysis, for example, groups documents into segments
- \* Some text mining methods use a structured matrix.

### **Another View of Text Mining**



### **Application: Document Classification**



### **Document Categorization**

**Document categorization** 

Assign documents to pre-defined categories

Examples

- Process email into work, personal, junk
- Process documents from a newsgroup into "interesting", "not interesting", "spam and flames"
- Process transcripts of bugged phone calls into "relevant" and "irrelevant"

### **Application: Information Retrieval**



### Introduction

How can we retrieve information using a search engine?.

- We can represent the query and the documents as vectors (vector space model)
  - However to construct these vectors we should perform a preliminary document preparation.
- The documents are retrieved by finding the closest distance between the query and the document vector.

### **Application: Clustering**





### **Document Classification**

**Document classification** 

- Cluster documents based on similarity
- Examples
  - Group samples of writing in an attempt to determine author(s)
  - Look for "hot spots" in customer feedback
  - Find new trends in a document collection (outliers, hard to classify)

### **IR Applications Using Text Mining**

- Survey Analysis
- Analysis of Student Evaluations of Instructors
- Predictive Modeling

**Enrollment Models** 

**Retention Models** 

### **Predictive Modeling**



### **Obtaining the Prediction**

**Nominal Target** Binary/Categorical

#### **Example** Binary Response: Mail (Y/N)



### **Objectives**

- Explore the general concept of decision trees.
- Build a decision tree model.
- Examine the model results and interpret these results.

### **Fitted Decision Tree**



### **The Cultivation of Trees**

- Split Search
  - Which splits are to be considered?
- Splitting Criterion
  - Which split is best?
- Stopping Rule
  - When should the splitting stop?
- Pruning Rule
  - Should some branches be lopped off?

### **Benefits of Trees**

- Interpretability
  - tree-structured presentation
- Mixed Measurement Scales
  - nominal, ordinal, interval
- Regression trees
- Robustness
- Missing Values



### **Simple Prediction Illustration**

## Predict dot color for each $x_1$ and $x_2$ .

Training Data



### **Simple Prediction Illustration**

## Predict dot color for each $x_1$ and $x_2$ .

Training Data







24

....







# Calculate the *logworth* of every partition on input $x_1$ .





# Calculate the *logworth* of every partition on input $x_1$ .

0.52







1.0

0.9



0.52

## Select the partition with the maximum *logworth*.

0.9

1.0



max ogworth(x<sub>1</sub>) 0.95

Repeat for input  $x_2$ .













Create a partition rule from the best partition across all inputs.





Repeat the process in each subset.









max logworth(x<sub>1</sub>) 5.72













Create a second partition rule.





Create a second partition rule.





**Repeat to form a maximal tree.** 

### Example

- Two Year School on Texas & Mexico Border
- Strong in Mathematics and Sciences
- Weak in the Arts
- Half of the students are from newly emigrated families

### **Objective**

- Improve Graduation Rate
- Identify Students Most Likely Not to Graduate
- Collect Data and Build a Predictive Model
- Determine What Intervention is Approximate

### **Hypothetical Data**

- Sample of 1000 Students Entering Fall 2010
- Determine Which Students Had Left by Fall 2013
- Data Fields
- 1. Student ID
- 2. Age
- 3. Gender
- 4. Major
- 5. Population
- 6. School
- 7. Enrollment Statement
- 8. Target







### **Process Flow**



### **Decision Tree Model**



### **Fit Statistics**

#### Fit Statistics Model Selection based on Valid: Misclassification Rate (\_VMISC\_)

			Train:			Valid:
			Valid:	Average	Train:	Average
Selected	Model	Model	Misclassification	Squared	Misclassification	Squared
Model	Node	Description	Rate	Error	Rate	Error
Y	Tree2	With Text	0.10474	0.077413	0.09683	0.086645
	Tree	No Text	0.12469	0.069562	0.10017	0.089264

### **ROC Curve**





- Score Students Entering in Fall 2013 With Model
- Distribute Scoring Information to Approximate People
- Evaluate Model After Two Years