

Data Mining Introduction

1. Nature of Data Sets
(Structured, unstructured, mixed)
2. Models vs. Patterns
3. Regression Example

Nature of Data Sets

- Structured Data
 - set of measurements from an environment or process
- Simple case
 - n objects with p measurements each: $n \times p$ matrix
 - p columns are called *variables, features, attributes or fields*

Structured Data and Data Types

US Census Bureau Data

Public Use Microdata Sample data sets (PUMS)

ID	Age	Sex	Marital Status	Education	Income
	Quantitative Continuous	Categorical		Categorical Ordinal	
248	54	Male	Married	High School grad	100000
	??	Missing data			Noisy data A guess?
249	??	Female	Married	HS grad	12000
250	29	Male	Married	Some College	23000
251	9	Male	Not Married	Child	0

PUMS Data has identifying information removed.

Available in 5% and 1% sample sizes. 1% sample has 2.7 million

Unstructured Data

1. Structured Data

- Well-defined tables, attributes (columns), tuples (rows)
- UC Irvine data set

2. Unstructured Data

- World wide web
 - Documents and hyperlinks
 - HTML docs represent tree structure with text and attributes embedded at nodes
 - XML pages use metadata descriptions
- Text Documents
 - Document viewed as sequence of words and punctuations
 - Mining Tasks
 - » Text categorization
 - » Clustering Similar Documents
 - » Finding documents that match a query
 - » Automatic Essay Scoring (AES)
 - Reuters collection is at <http://www.research.att.com/~lewis>

Representations of Text Documents

- Boolean Vector
 - Document is a vector where each element is a bit representing presence/absence of word
 - A set of documents
 - can be represented as matrix (d,w)
 - where document d and word w has value 1 or 0
(sparse matrix)
- Vector Space Representation
 - Each element has a value such as no. of occurrences or frequency
 - A set of documents represented as a document-term matrix

Vector Space Example

Document-Term Matrix

	t1	t2	t3	t4	t5	t6
D1	24	21	9	0	0	3
D2	32	10	5	0	3	0
D3	12	16	5	0	0	0
D4	6	7	2	0	0	0
D5	43	31	20	0	3	0
D6	2	0	0	18	7	16
D7	0	0	1	32	12	0
D8	3	0	0	22	4	2
D9	1	0	0	34	27	25
D10	6	0	0	17	4	23

t1	database
t2	SQL
t3	index
t4	regression
t5	likelihood
t6	linear

d_{ij} represents number of times that term appears in that document

Mixed Data: Structured & Unstructured

Medical Patient Data

- Blood Pressure at different times of day
- Image data (x-ray or MRI)
- Specialist's comments (text)
- Hierarchy of relationships between patients, doctors, hospitals

$N \times p$ data matrix is oversimplification of what occurs in practice

Transaction Data

List of store purchases: date, customer ID, list of items and prices

Web transaction log -sequence of triples: (user id, web page, time)

Can be transformed into binary-valued matrix

Individual

1	1		1	1	1				1
1					1				
1	1	1			1	1	1	1	
1							1		
1		1		1		1			1
1			1						1

Web Page Visited

Representations : Models and Patterns

- *Global Model*

- *Make a statement about any point in p-space*
 - *E.g., assign a point to a cluster*
 - *Even when some values are missing*
- *Simple model: $Y = aX + c$*
 - *Functional model is linear*
 - *Linear in variables rather than parameters*

- *Local Patterns*

- *Make a statement about restricted regions of space spanned by variables*
 - *E.g.1: if $X > thresh1$ then $Prob(Y > thresh2) = p$*
 - *E.g.2: certain classes of transactions do not show peaks and troughs (bank discovers dead people's open accounts)*

Data Mining Tasks

Different objectives

1. Exploratory Data Analysis

Visualizing Data (histograms)

2. Descriptive Modeling

Density Estimation, Cluster Analysis

3. Predictive Modeling

a) Classification (response variable is a class)

b) Regression (response variable is numeric)

4. Discovering Patterns and Rules

Association Rules

5. Retrieval by Content

Queries that are Textual(Google) or Images (QBIC)

Data Mining Algorithms

Four basic components in each algorithm*

1. Model or Pattern Structure

Determining underlying structure or functional form we seek from data

2. Score Function

Judging the quality of the fitted model

3. Optimization and Search Method

Searching over different model and pattern structures

4. Data Management Strategy

Handling data access efficiently

KDD Steps

Course
Focus:
Algorithms
for Modeling
the data

- Selecting Target Data
- Preprocessing the Data
- Transforming them
- Data Mining to extract Models or Patterns
- Interpreting/Assessing Discovered Structures

Process of Data Mining

1. Representation
2. Score function
3. Process to optimize score
4. Implementation:
data management,
efficiency

EXAMPLE of Model

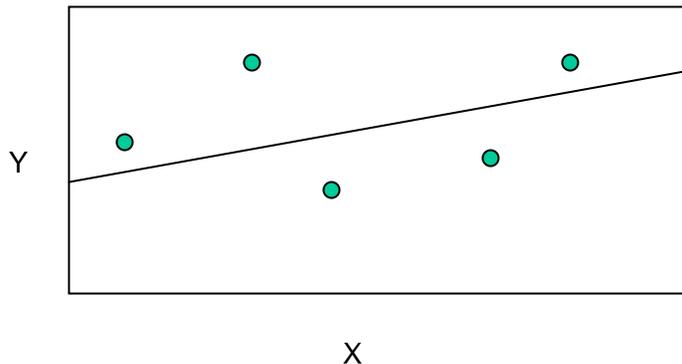
1. Regression:
 $Y = a + bX$
Predictor variable = X
(income)
Response variable = Y
(credit card spending)
2. Score: sum of squared errors

Linear Regression Process: Extracting a Linear Model

Linear regression with one variable

Data of the form (x_i, y_i) , $i = 1, \dots, n$ samples

Need to find a and b such that $Y = a + bX$



What is involved in calculating a and b
So that the line fits the points the best?

Data Representation

Y	X
1	3
8	9
11	11
4	5
3	2

Score: sum of squared errors

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - y'_i)^2 = \sum_{i=1}^n (y_i - a - bx_i)^2$$

Where y_i is the response value obtained from the model

We wish to minimize SSE

Minimizing SSE for Regression

Differentiating SSE with respect to a and b we have

$$\frac{\partial(SSE)}{\partial a} = -2 \sum_{i=1}^n (y_i - a - bx_i)$$

$$\frac{\partial(SSE)}{\partial b} = -2 \sum_{i=1}^n (y_i - a - bx_i)x_i$$

Setting partial derivatives equal to zero and rearranging terms we have equations

$$na + b \sum_{i=1}^n x_i = \sum_{i=1}^n y_i$$

$$a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i$$

Which we solve for a and b ,
the *regression coefficients*

Regression Coefficients

$$b = \frac{\sum_{i=1}^n (x_i - \text{mean}_x) \cdot (y_i - \text{mean}_y)}{\sum_{i=1}^n (x_i - \text{mean}_x)^2}$$

$$a = \text{mean}_y - b \cdot \text{mean}_x$$

To calculate a and b we need to find the *means* of the x and y values. Then we calculate b as a function of the x and y values and the means a as a function of the means and b

Application to data values

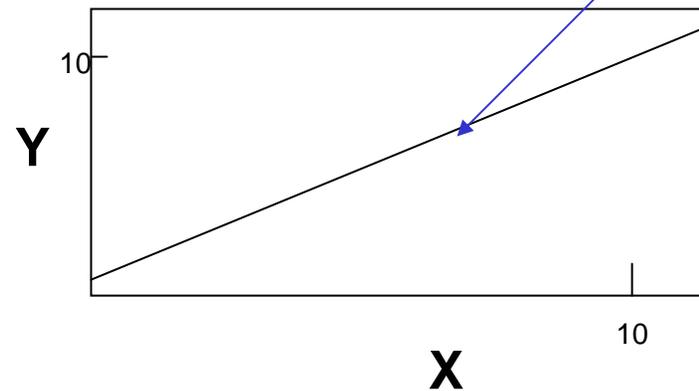
Y	X
1	3
8	9
11	11
4	5
3	2

$$\text{mean}_y = 5$$

$$\text{mean}_x = 6$$

$$a = 0.8, b = 1.04$$

Optimal regression line is
 $Y = 0.8 + 1.04X$



Linear Regression
For the data set

Multiple Regression

p predictor variables

$$y = a_0 + \sum_{i=1}^p a_j x_j$$

y	x_1	x_2	x_p
$y(1)$	$x_1(1)$			
$y(n)$	$x_1(n)$			

n objects

$X = n$ by $p+1$ matrix

Where a column of 1's are added to incorporate a_0 in model

$$y = Xa + e$$

y is a column vector, $a=(a_0, \dots, a_p)$
 e is a n by 1 vector containing residuals

Solution:

$$a = (X^t X)^{-1} \cdot X^t y$$

Implementation of Regression

$$y = a_0 + \sum_{i=1}^p a_j x_j$$

Solution:

$$a = (X^t X)^{-1} \cdot X^t y$$

Simple summaries of the data; sums, sums of squares and sums of products of X and Y are sufficient to compute estimates of a and b

Implies single pass through the data will yield estimates