



CS 584

Data Mining

Professor Jessica Lin

Introduction

1/19/16



Basics

Instructor: Dr. Jessica Lin

Contact Info:

Email: [jessica \[at\] cs \[dot\] gmu \[dot\] edu](mailto:jessica[at]cs[dot]gmu[dot]edu)

Homepage: <http://www.cs.gmu.edu/~jessica>

Office: Engineering Building Room 4419

Phone: (703)993-4693

Office Hours: Tuesday 2-4pm

Class Meeting: Tuesday 4:30-7:10pm
Innovation Hall 206

Pre-requisites: (strictly enforced) C or better in CS 310 and STAT 344

TA: Monjura Afrin Rumi ([mrumi \[at\] gmu \[dot\] edu](mailto:mrumi[at]gmu[dot]edu))

TA Office Hours: Monday 7:30-8:30pm, Thursday 3:30-4:30pm

Outline

- Course syllabus
- Introduction to Data Mining



Administration Trivia

- Class webpage:

http://www.cs.gmu.edu/~jessica/cs584_s16.html

- In most cases, I will put the slides online the night before the lecture.
- You are 100% responsible for any announcements and updates on the class webpage and Piazza (more on this later), so visit the page(s) frequently.

Textbook

- Required:
 - Introduction to Data Mining by Tan, Steinbach and Kumar
- Recommended:
 - Data Mining and Analysis by Mohammed Zaki
 - online pdf version:
<http://www.dataminingbook.info/pmwiki.php/Main/BookDownload>

Grading

- Assignments: 20%
 - Project: 30%
 - Midterm Exam: 20%
 - Final Exam: 30%
-
- There will be one midterm and one final exam.
 - Final exam is comprehensive.
 - Both exams are closed-book, closed-notes.
 - Exams must be taken at the scheduled time and place, unless prior arrangement has been made with the instructor.
 - Missed exams cannot be made up.

Grading Scale

- We will use the following grading scale:

A: 89-100

B: 78-88

C: 67-77

D: 60-66

F: < 60

- Pluses and Minuses may be used as well.

Honor Code System

- GMU honor Code
<http://oai.gmu.edu/the-mason-honor-code-2/>
- Make sure you read the full honor code within the link.
- In addition, the CS Department has specific honor code policies for programming projects, etc.:
<http://cs.gmu.edu/wiki/pmwiki.php/HonorCode/CSHonorCodePolicies>
- For this class
 - You may work in a team of 2 for the project.
 - Homework: individual effort
 - Exams: individual effort, closed books/notes

Tools we will use for the class

- Weka (more on this later in the semester)
- Piazza
 - A free online class Q&A platform. I will send out an invitation to sign up once the the class roster is finalized.
- Blackboard



Piazza


- Think before posting the question (e.g. is the answer in the book or the lecture slides?). This is a discussion forum, not Siri. You are encouraged to answer each other's questions but do not give out answers for assignments.
- **You are responsible for all discussions, announcements and updates posted on the discussion forum**, so make sure your notification setting is set to “real time” – this ensures that you receive notifications in the timely manner.

Homework & Project Submission

- <https://mymasonportal.gmu.edu/webapps/portal/frameset.jsp>
- Login with your GMU student account
- Click on the Courses tab on the upper right hand corner
- Choose CS584
- Use Blackboard for:
 - Electronic submission of assignments and project
 - Checking grades
 - Getting course materials such as homework solutions
- All homework and project, unless otherwise specified, are due at 4:30pm (start time of the class) on the due date.
- No late submissions will be accepted.
- Please bring a hardcopy to class.



Email Policy

- Please email me from your official GMU email. If you must email me from another account, you must state your full name and your official GMU email address.
 - Please put [CS584] in the beginning of the subject line
- 

Class Schedule

- See class website for the tentative class schedule:

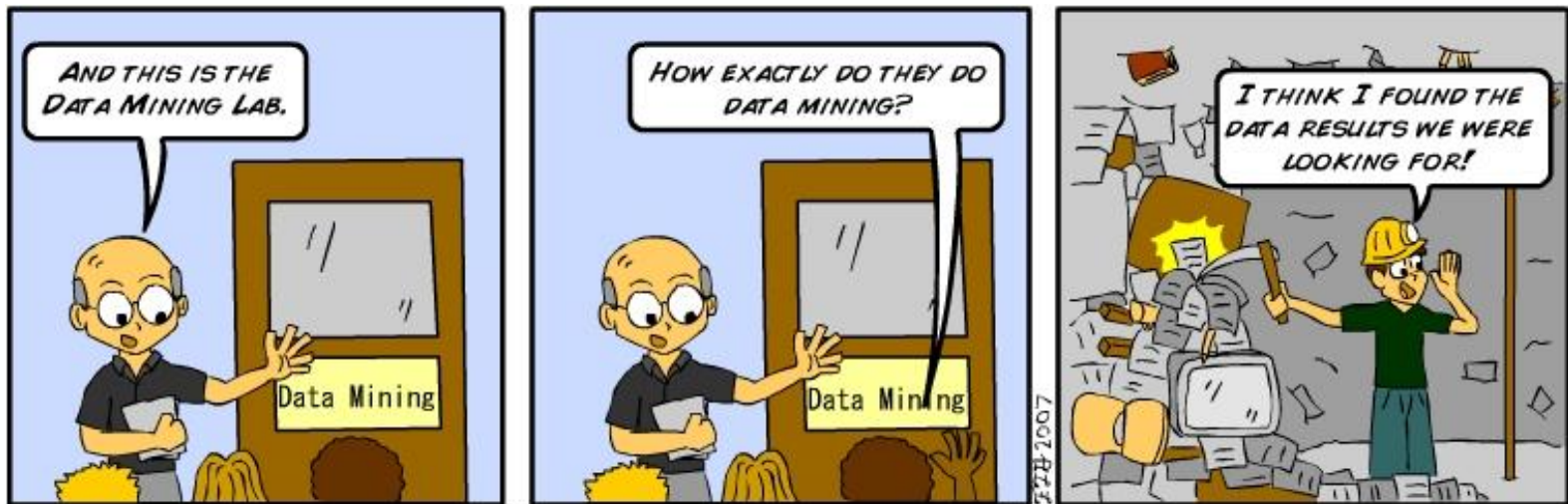
https://cs.gmu.edu/~jessica/cs584_s16.html





Any Questions?

What do you think of data mining?

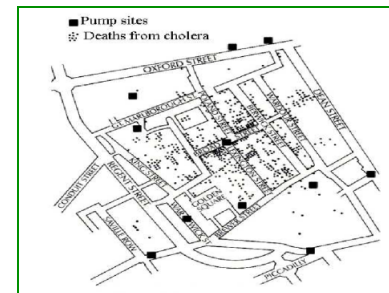


Large-scale Data is Everywhere!

- Tremendous data growth
- New mantra
 - Gather whatever data you can whenever and wherever possible.
- Expectations
 - Gathered data will have value either for the purpose collected or for a purpose not envisioned.



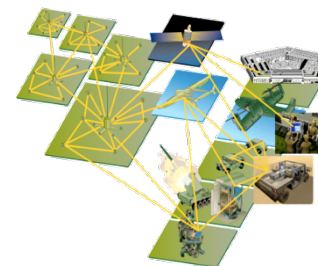
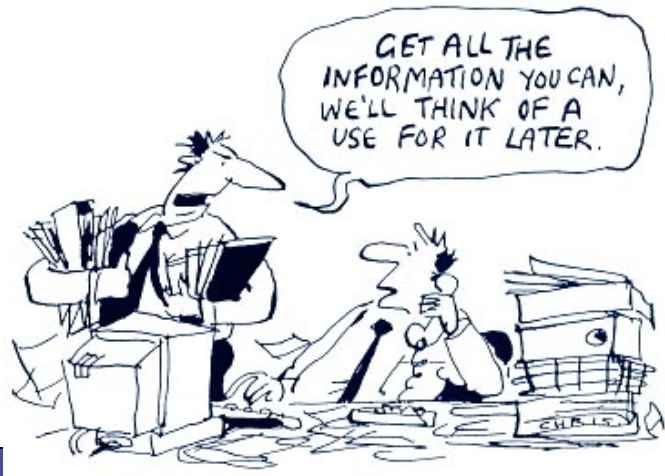
Homeland Security



Geo-spatial data



Business Data



Sensor Networks



Social Media & the Internet

An example business problem

- TelCo, a major telecommunications firm, wants to investigate its problem with customer attrition, or “churn”
- Lets consider this for now as a marketing problem only



How would you go about targeting some customers with a special offer, prior to contract expiration? Think about what data should be available for your use.

From data & business to strategy



HURRICANE FRANCES was on its way, barreling across the Caribbean, threatening a direct hit on Florida's Atlantic coast. Residents made for higher ground, but far away, in Bentonville, Ark., executives at Wal-Mart Stores decided that the situation offered a great opportunity for one of their newest data-driven weapons, something that the company calls predictive technology.

A week ahead of the storm's landfall, Linda M. Dillman, Wal-Mart's chief information officer, pressed her staff to come up with forecasts based on what had happened when Hurricane Charley struck several weeks earlier. Backed by the trillions of bytes' worth of shopper history that is stored in Wal-Mart's data warehouse, she felt that the company could "start predicting what's going to happen, instead of waiting for it to happen," as she put it.

What is Data Mining?



Copyright © 1999 United Feature Syndicate, Inc.

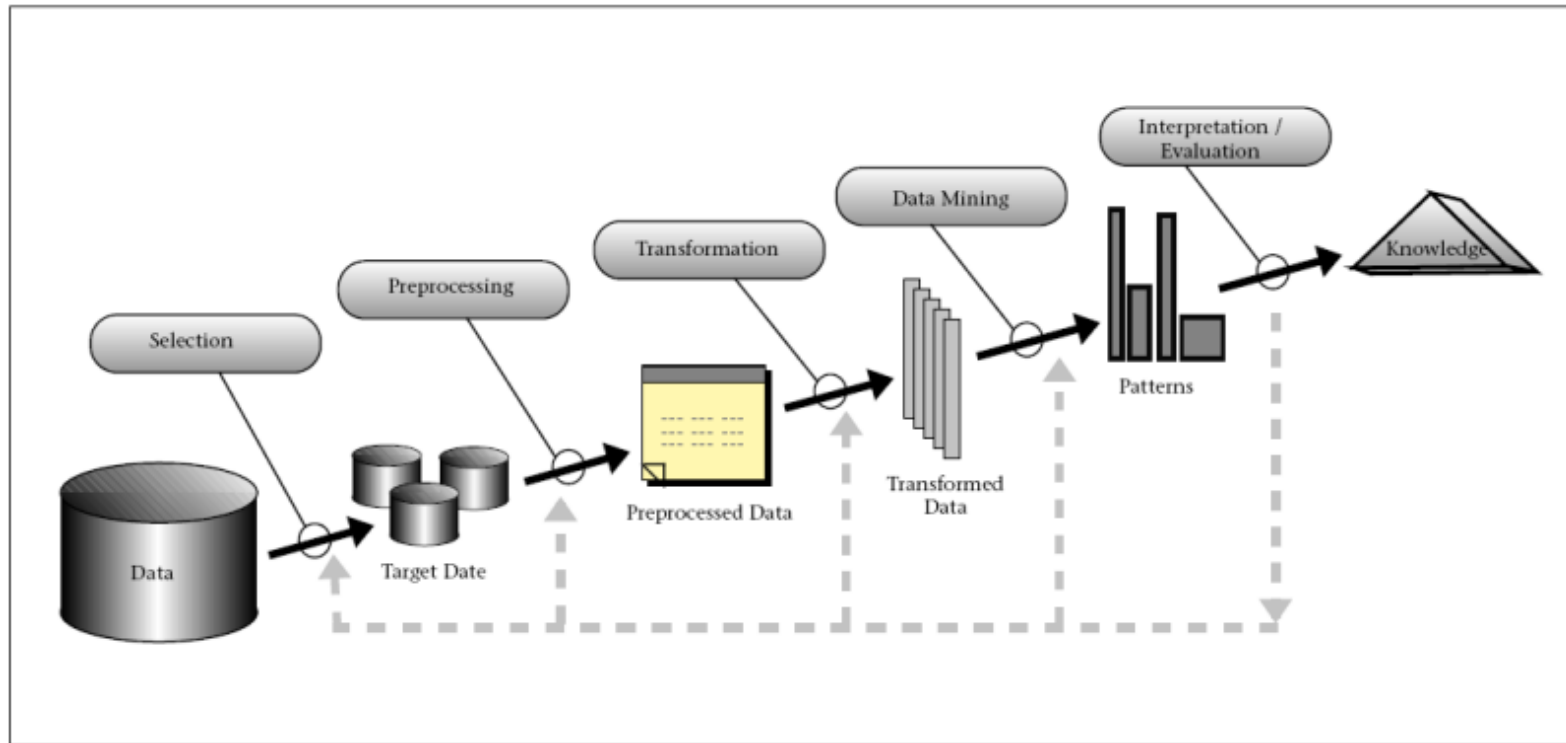
What really *is* data mining?

- *A process for using information technology to extract useful (non-trivial, hopefully actionable) knowledge from large bodies of data*

KDD Process

- KDD: (Knowledge Discovery in Databases)

CONVERTING RAW DATA TO USEFUL INFORMATION.



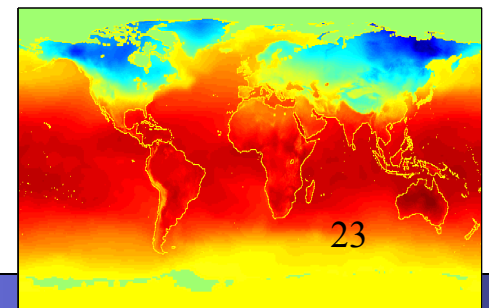
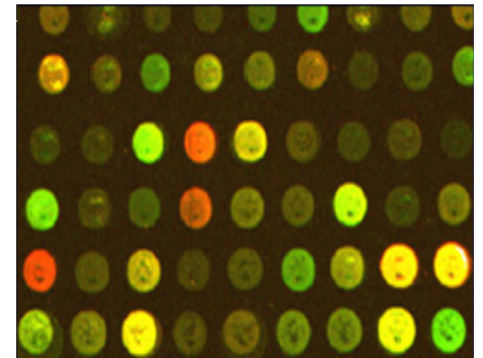
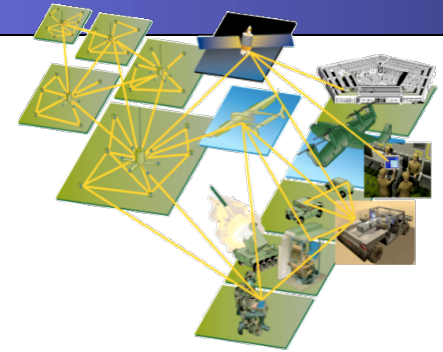
Why Data Mining? Commercial Viewpoint

- Lots of data is being collected and warehoused
 - Web data (as of 2014/15, most likely outdated by now)
 - Google processes 20 PB/day
 - Facebook has 955M active users
 - Twitter has more than 400M tweets/day
 - purchases at department/grocery stores, e-commerce
 - Amazon has 42 TB of data
 - Bank/Credit Card transactions
- Computers have become cheaper and more powerful
- Competitive Pressure is Strong
 - Provide better, customized services for an edge (e.g. in Customer Relationship Management)



Why Data Mining? Scientific Viewpoint

- Data collected and stored at enormous speeds
 - remote sensors on a satellite
 - NASA EOSDIS archives over 1-petabytes of earth science data / year
 - telescopes scanning the skies
 - Sky survey data
 - LSST (Large Synoptic Survey Telescope) project: 20 PB science data & 100 PB image archive
 - High-throughput biological data
 - scientific simulations
 - terabytes of data generated in a few hours
- Data mining helps scientists
 - in automated analysis of massive datasets
 - in hypothesis formation



What is (not) Data Mining?

- What is not Data Mining?
 - Look up phone number in phone directory
 - Query a Web search engine for information about “Amazon”
- What is Data Mining?
 - Certain names are more prevalent in certain US locations (O’ Brien, O’ Rurke, O’ Reilly... in Boston area)
 - Group together similar documents returned by search engine according to their context (e.g. Amazon rainforest, Amazon.com)

Some Data Mining Examples

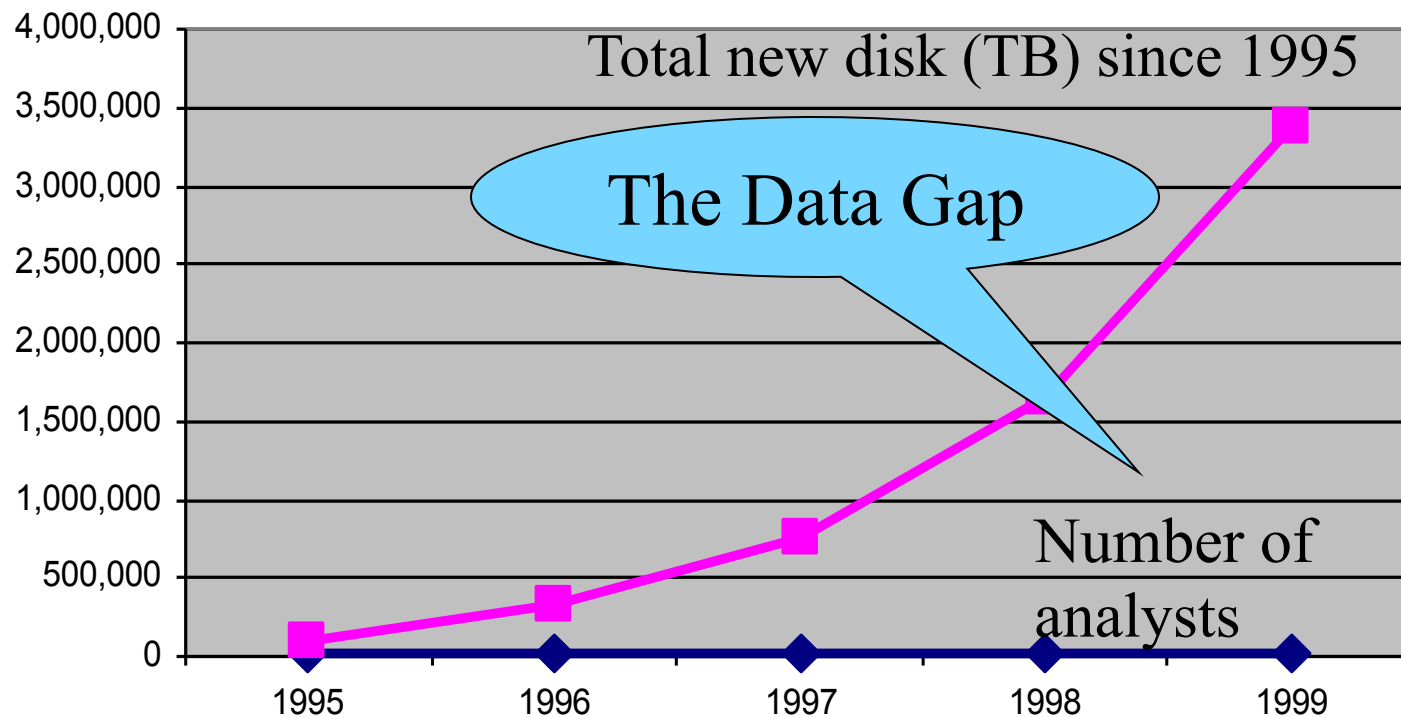
- Amazon.com, Google, Netflix
 - Personal Recommendations
 - Profile-based advertisements
- Spam Filters/Priority Inbox
 - 94B spam messages sent daily; cost society \$20B annually
- Scientific Discovery
 - Finding potential causes of cancer
 - Prediction of weather and natural disasters
- Security
 - Fraud detection, Network Traffic monitoring
- Politics
 - Obama campaign: how data mining helped Obama win the election

Some of my Data Mining Projects

- Intel manufacturing data
 - Compare “good” chips and “bad” chips
 - What are the potential causes of “bad” chips?
- ICU medical alarms
 - “Alarm fatigue”
 - How to reduce false alarms
- Trajectory data mining
 - Find repeated routes
 - Anomaly detection, route recommendation, etc.
- Power load prediction

Mining Large Data Sets - Motivation

- There is often information “hidden” in the data that is not readily evident
- Human analysts may take weeks to discover useful information
- Much of the data is never analyzed at all



From: R. Grossman, C. Kamath, V. Kumar, “Data Mining for Scientific and Engineering Applications”

Data Mining Tasks

- Prediction Methods
 - Use some variables to predict unknown or future values of other variables.
- Description Methods
 - Find human-interpretable patterns that describe the data. Make good inferences from the data.

Applications of Data Mining

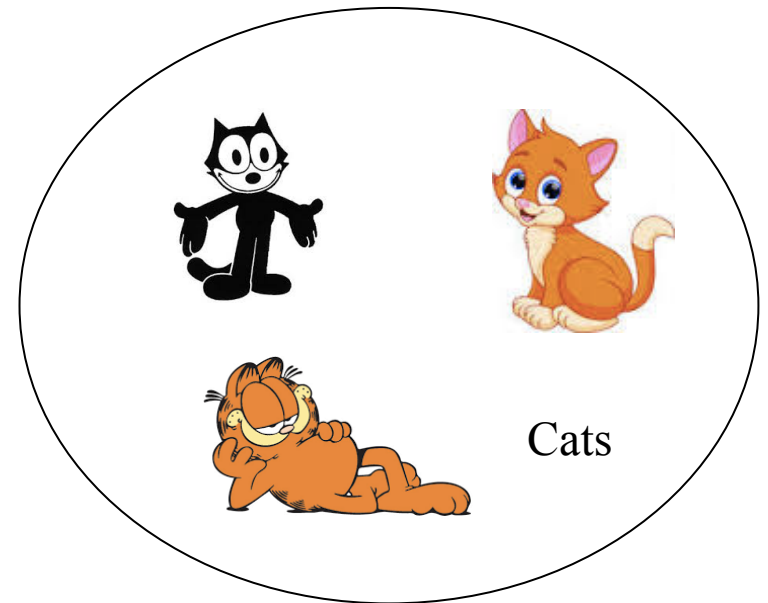
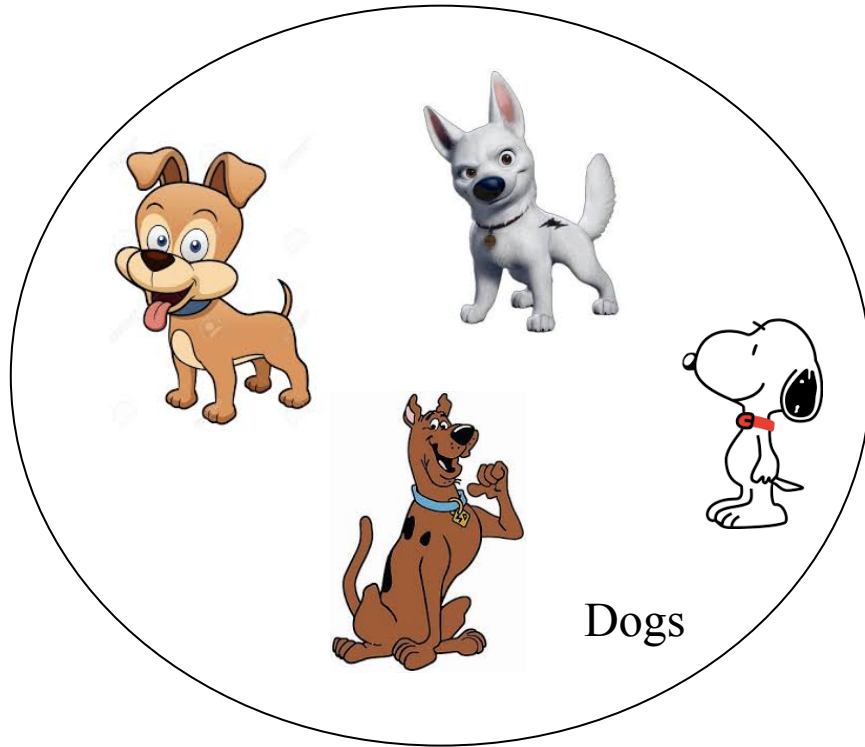
- **Prediction** based on past history
 - Predict if a credit card applicant poses a good credit risk, based on some attributes (income, job type, age, ..) and past history
 - Predict if a customer is likely to switch brand loyalty (churn)
 - Predict if a customer is likely to respond to “junk mail”
 - Predict if a pattern of phone calling card usage is likely to be fraudulent
- Example of prediction mechanisms:
 - **Classification** - Given a training set consisting of items belonging to different classes, and a new item whose class is unknown, predict which class it belongs to
 - **Regression** – Predicting the value of one variable from one or more variables.

Applications of Data Mining (Cont.)

- **Descriptive Patterns**

- **Associations** – Find items that are often bought by the same customers. If a new customer buys one such item, suggest that he buys the others too.
 - Associations may also be used as a first step in detecting **causation**
 - E.g. association between exposure to chemical X and cancer, or new medicine and cardiac problems
- **Clusters** – Finding natural grouping in data
 - E.g. Do my customers form natural groups?
 - E.g. typhoid cases were clustered in an area surrounding a contaminated well

Classification Example

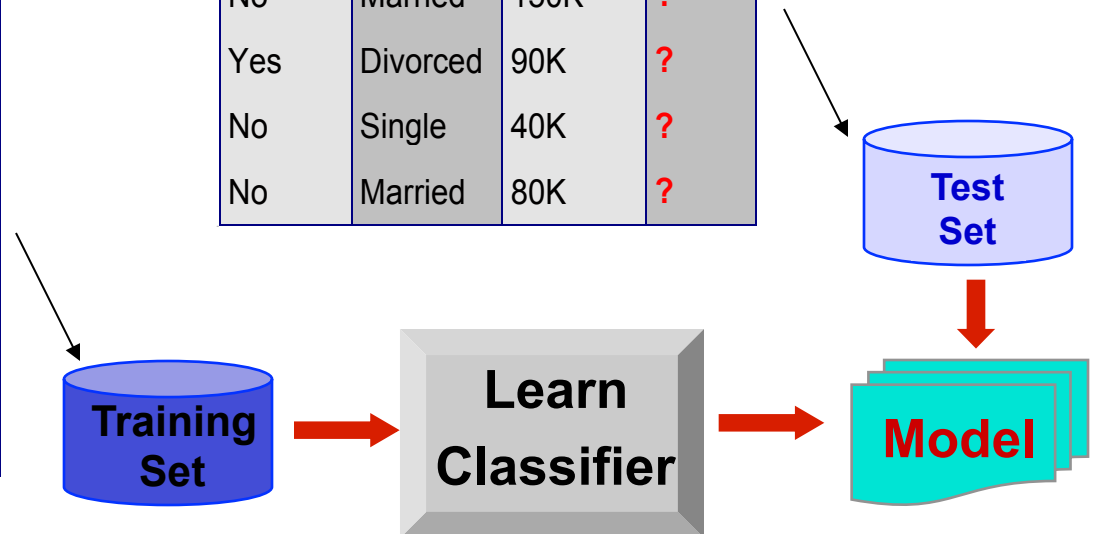


Classification Example

categorical
categorical
continuous
class

<i>Tid</i>	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Refund	Marital Status	Taxable Income	Cheat
No	Single	75K	?
Yes	Married	50K	?
No	Married	150K	?
Yes	Divorced	90K	?
No	Single	40K	?
No	Married	80K	?



Classification: Definition

- Given a collection of records (*training set*)
 - Each record contains a set of *attributes*, one of the attributes is the *class*.
- Find a *model* for class attribute as a function of the values of other attributes.
- Goal: previously unseen records should be assigned a class as accurately as possible.
 - A *test set* is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.
- Also known as *Supervised* learning

Classification: Direct Marketing

- Direct Marketing
 - Goal: Reduce cost of mailing by *targeting* a set of consumers likely to buy a new cell-phone product.
 - Approach:
 - Use the data for a similar product introduced before.
 - We know which customers decided to buy and which decided otherwise. This *{buy, don't buy}* decision forms the *class attribute*.
 - Collect various demographic, lifestyle, and company-interaction related information about all such customers.
 - Type of business, where they live, how much they earn, etc.
 - Use this information as input attributes to learn a classifier model.

Classification: Fraud Detection

- Fraud Detection
 - Goal: Detect fraudulent cases in credit card transactions.
 - Approach:

Classification: Fraud Detection

- Fraud Detection
 - Goal: Detect fraudulent cases in credit card transactions.
 - Approach:
 - Use credit card transactions and the information on its account-holder as attributes.
 - When does a customer buy, what does he buy, how often he pays on time, etc
 - Label past transactions as fraud or fair transactions. This forms the class attribute.
 - Learn a model for the class of the transactions.
 - Use this model to detect fraud by observing credit card transactions on an account.

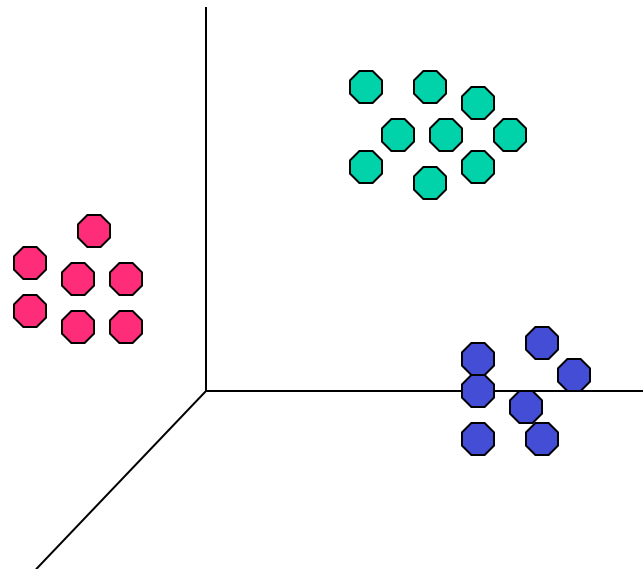
Classification: Your Turn

- Churn prediction
 - Goal: Predict which customers will terminate contracts soon after they expire
 - Approach:
- Can you think of 3 more applications for classification?

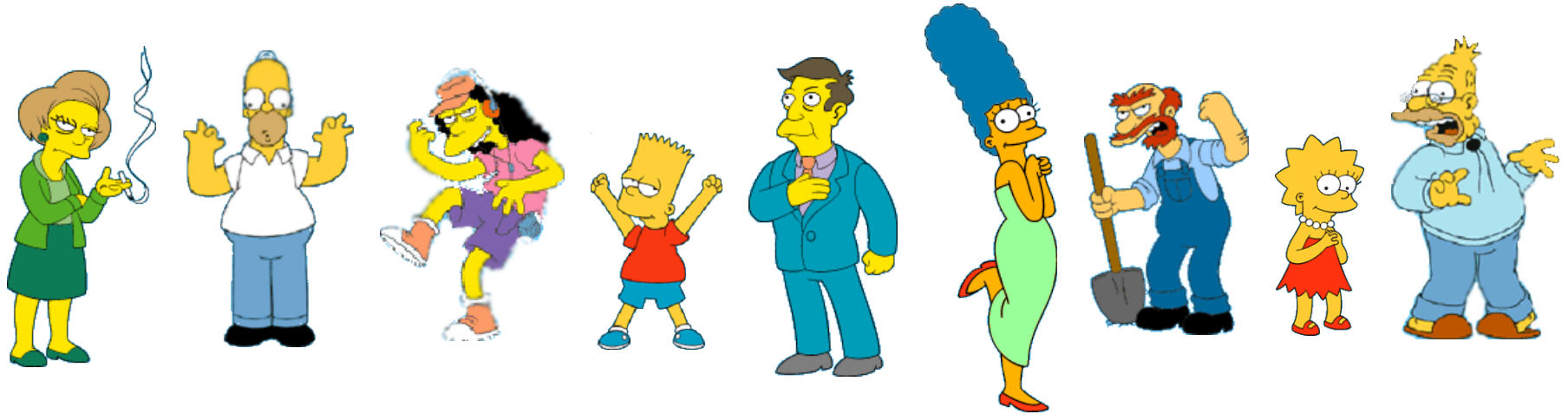
Clustering

- Given a set of data points, each having a set of attributes, and a similarity measure among them, find clusters such that
 - Data points in one cluster are more similar to one another.
 - Data points in separate clusters are less similar to one another.

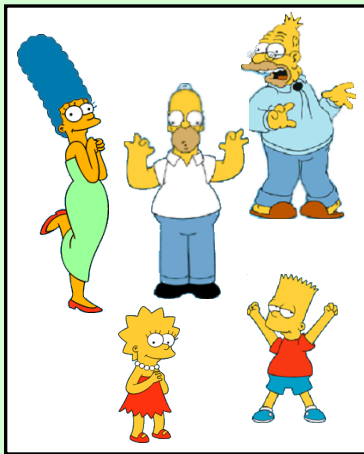
Illustrating Clustering



What is a natural grouping among these objects?



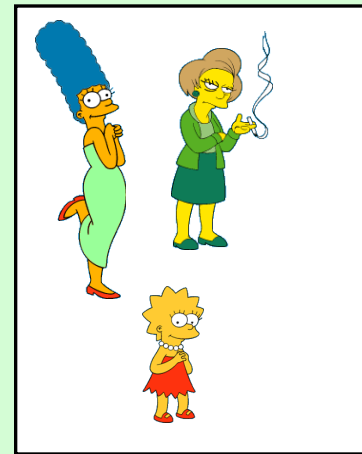
Clustering is subjective



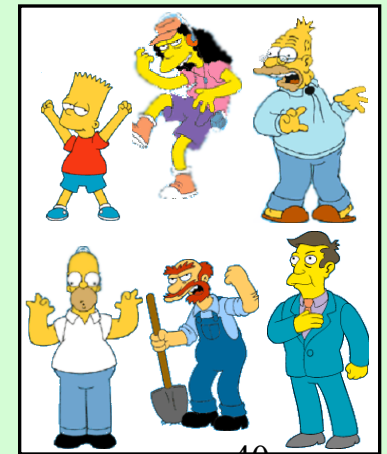
Simpson's Family



School Employees



Females



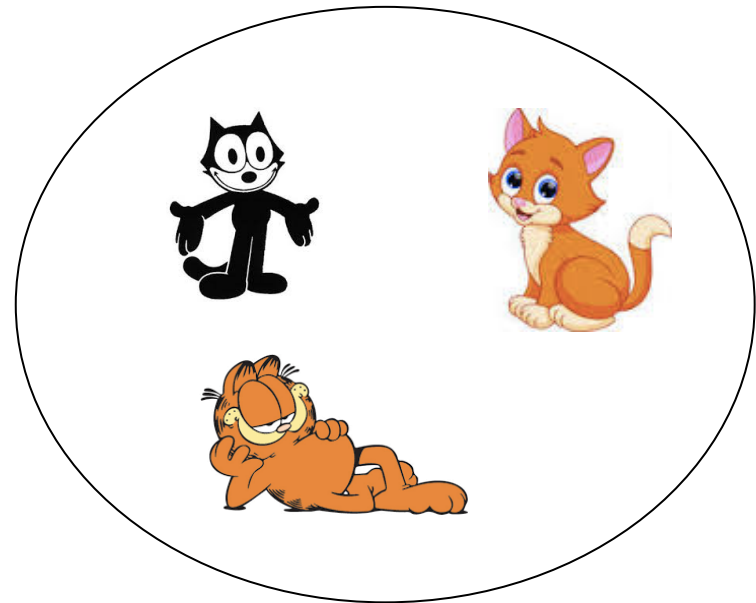
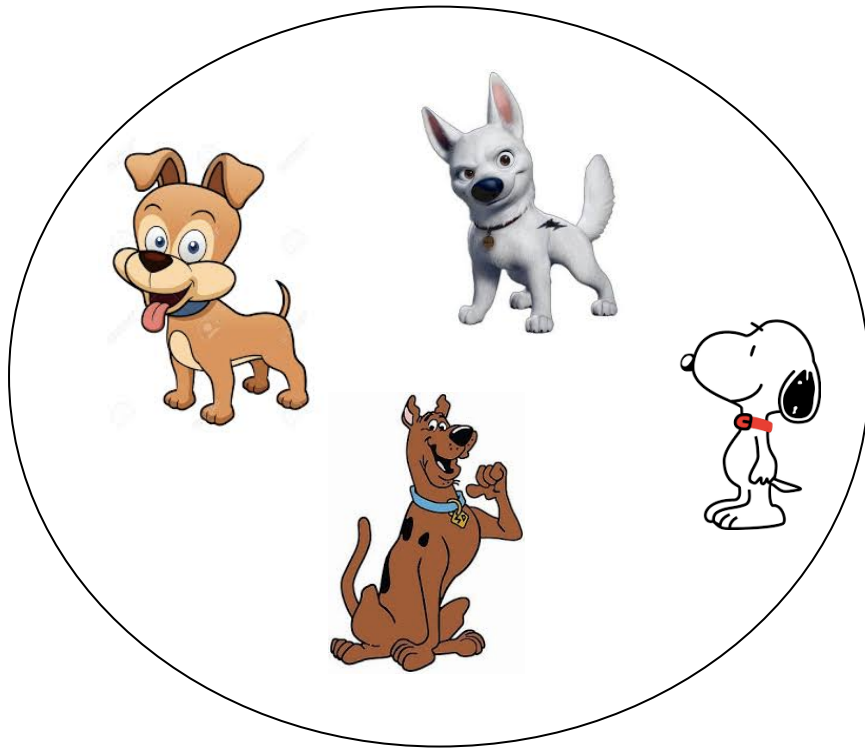
Males

Think point ?

- Differences between classification and clustering?

Think point ?

- Differences between classification and clustering?



Association Rule Discovery: Definition

- Given a set of records each of which contain some number of items from a given collection;
 - Produce dependency rules which will predict occurrence of an item based on occurrences of other items.

<i>TID</i>	<i>Items</i>
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

Rules Discovered:
{Milk} --> {Coke}

Urban Legend

- Classic Association Rule Example:
 - If a customer buys diaper and milk, then he is very likely to buy beer.
 - Any plausible explanations ? 😊

Association Rule Discovery: Application 1

- Marketing and Sales Promotion:
 - Let the rule discovered be
 $\{\text{Bagels, ...}\} \rightarrow \{\text{Potato Chips}\}$
 - **Potato Chips as consequent** \Rightarrow Can be used to determine what should be done to boost its sales.
 - **Bagels in the antecedent** \Rightarrow Can be used to see which products would be affected if the store discontinues selling bagels.
 - **Bagels in antecedent and Potato chips in consequent** \Rightarrow Can be used to see what products should be sold with Bagels to promote sale of Potato chips!

Association Rule Discovery: Application 2

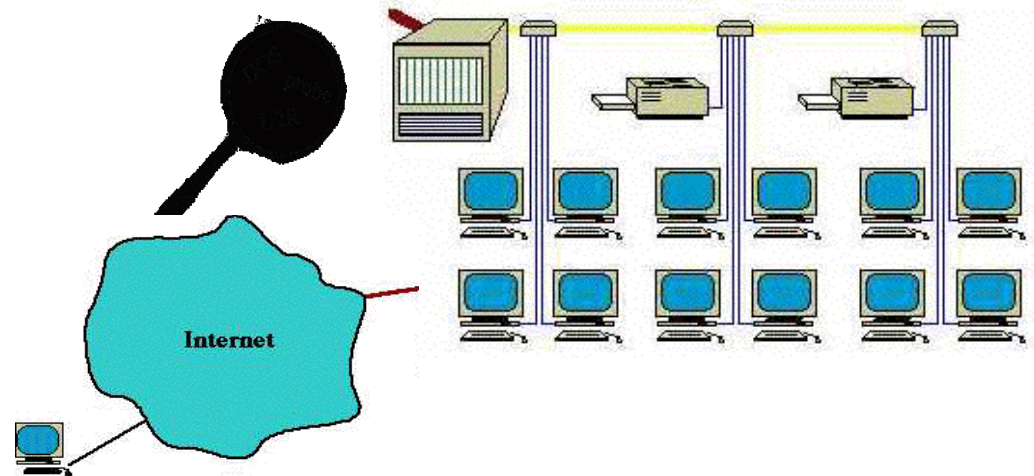
- Supermarket shelf management.
 - Goal: To identify items that are bought together by sufficiently many customers.
 - Approach: Process the point-of-sale data collected with barcode scanners to find dependencies among items.
 - Wal-mart, Target, and departmental store managers are big into this.
 - All your transactions gets processed & analyzed in a warehouse.

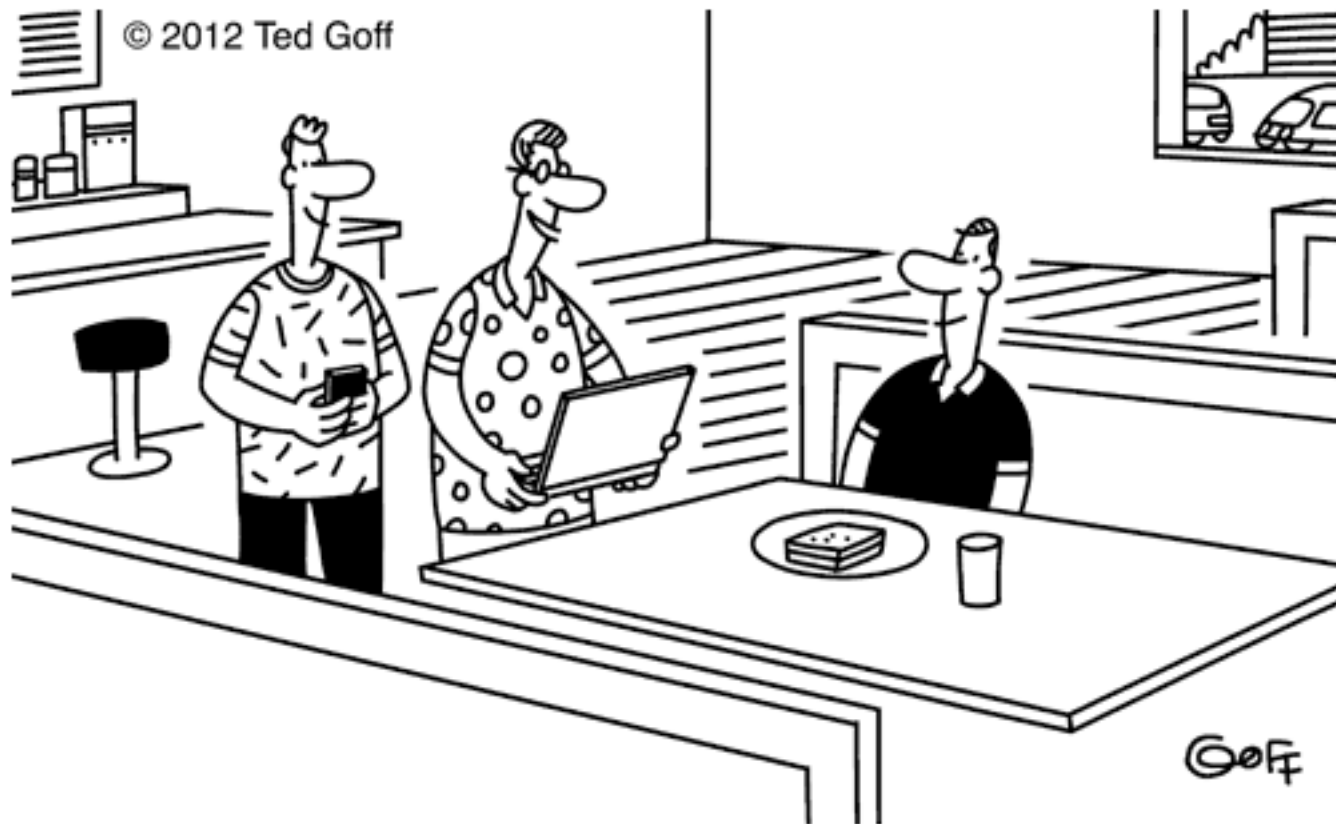
Regression

- Predict a value of a given continuous valued variable based on the values of other variables, assuming a linear or nonlinear model of dependency.
- Also called density estimation.
- Greatly studied in statistics, neural network fields.
- Examples:
 - Predicting wind velocities as a function of temperature, humidity, air pressure, etc.
 - Use personal income to predict auto sales

Deviation/Anomaly Detection

- Detect significant deviations from normal behavior
- Applications:
 - Credit Card Fraud Detection
 - Network Intrusion Detection





“Twitter and Facebook can’t predict the election, but they did predict what you’re going to have for lunch: a tuna salad sandwich. You’re having the wrong sandwich.”

What else can Data Mining do ?



Dilbert

Challenges of Data Mining

- Scalability
- Dimensionality
- Complex and Heterogeneous Data
- Data Quality
- Data Ownership and Distribution
- Privacy Preservation
- Streaming Data