

Data Mining: Introduction

Lecture Notes for Chapter 1

Introduction to Data Mining

by

Tan, Steinbach, Kumar

(modified by Predrag Radivojac, 2018)

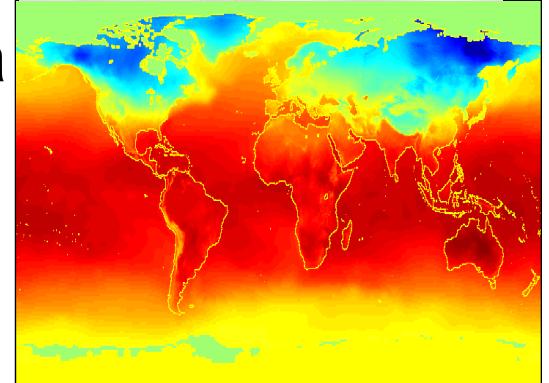
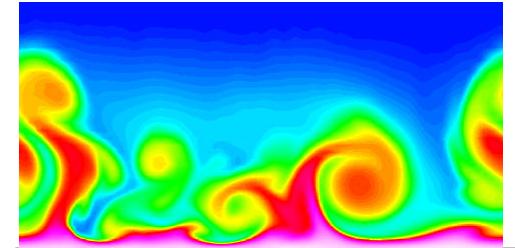
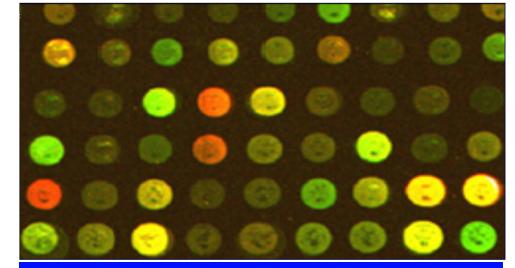
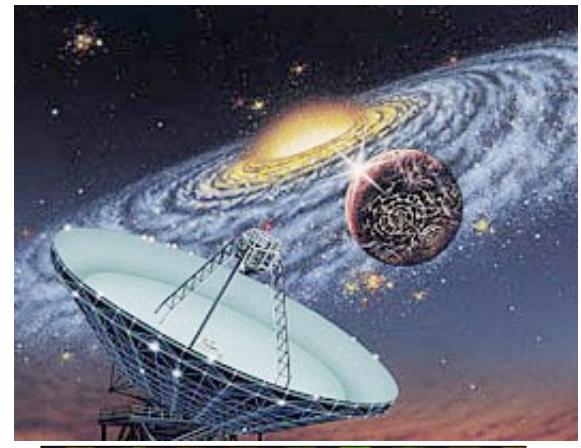
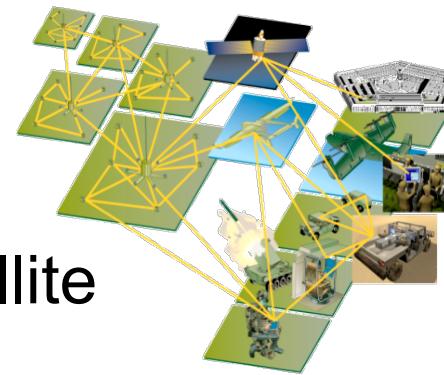
Why Mine Data? Commercial Viewpoint

- Lots of data is being collected and warehoused
 - Web data, e-commerce
 - purchases at department/grocery stores
 - 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)



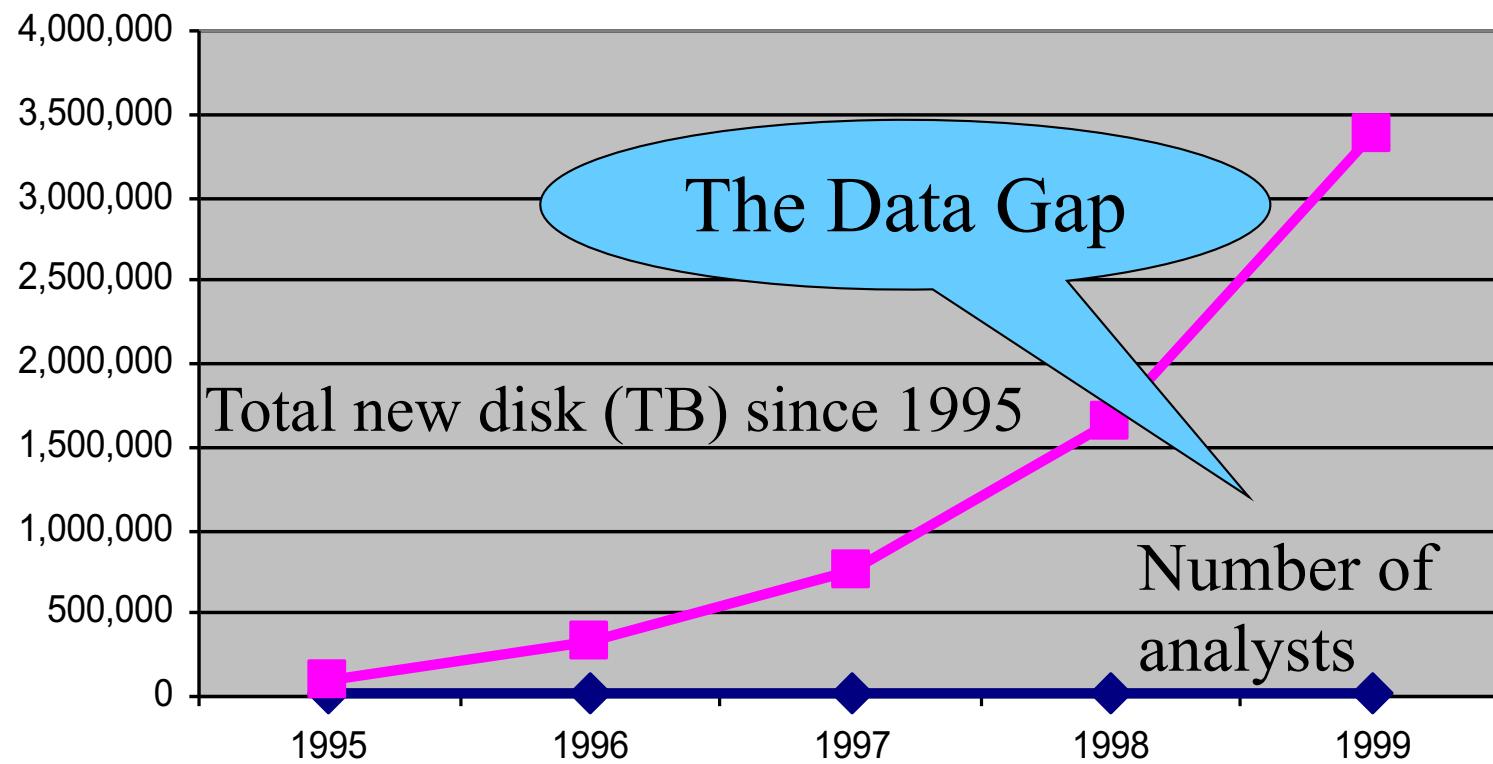
Why Mine Data? Scientific Viewpoint

- Data collected and stored at enormous speeds (GB/hour)
 - remote sensors on a satellite
 - telescopes scanning the skies
 - microarrays generating gene expression data
 - scientific simulations generating terabytes of data
- Traditional techniques infeasible for raw data
- Data mining may help scientists
 - in classifying and segmenting data
 - in hypothesis formation



Mining Large Data Sets - Motivation

- There is often information “hidden” in the data
- Human analysts may take weeks to discover useful information
- Much of the data is never analyzed at all

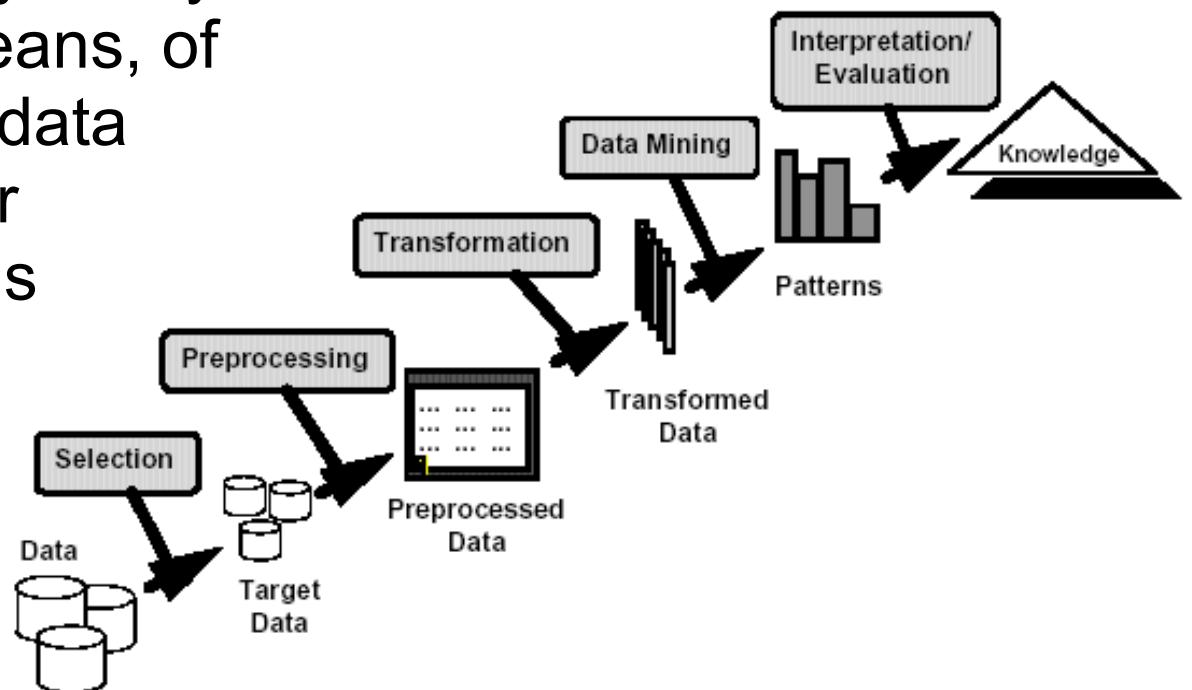


From: R. Grossman, C. Kamath, V. Kumar, “Data Mining for Scientific and Engineering Applications”

What is Data Mining?

● Many definitions

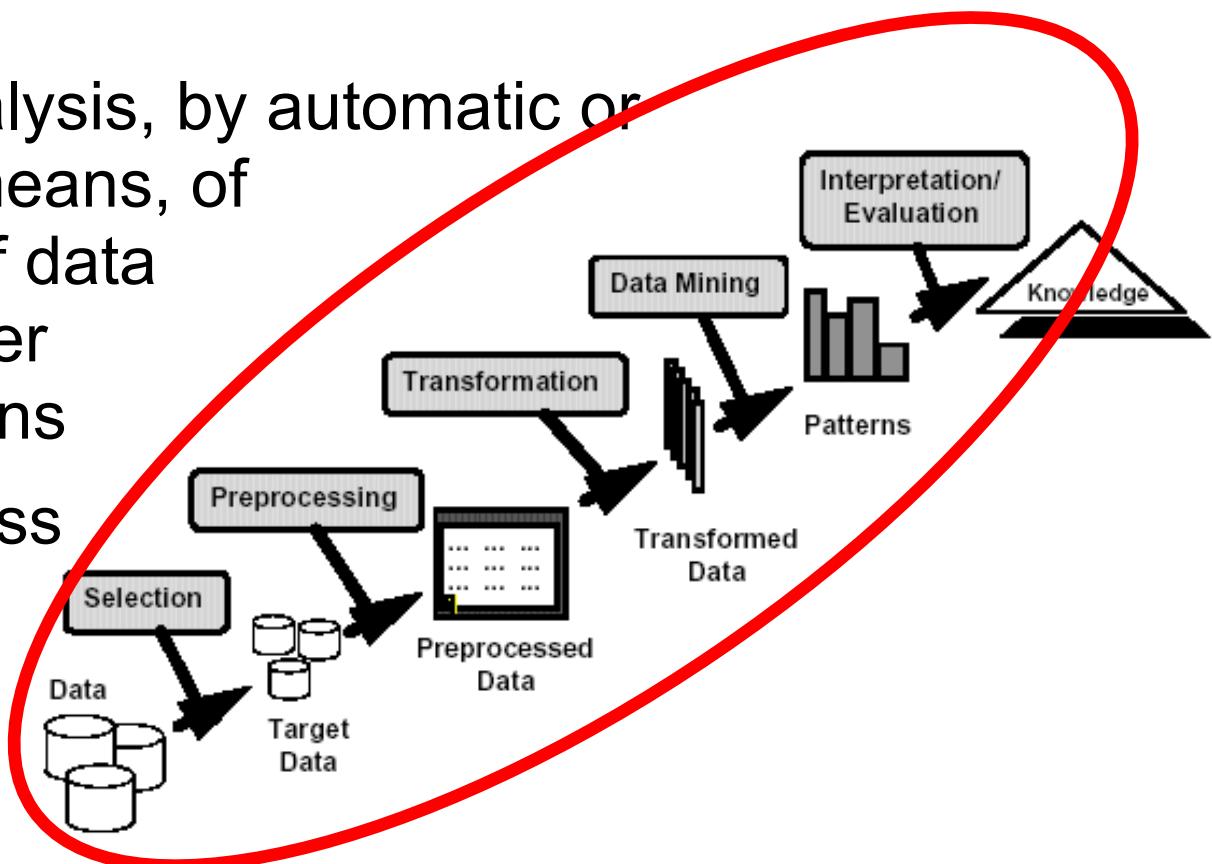
- Non-trivial extraction of implicit, previously unknown and potentially useful information from data
- Exploration & analysis, by automatic or semi-automatic means, of large quantities of data in order to discover meaningful patterns



What is Data Mining?

● Many definitions

- Non-trivial extraction of implicit, previously unknown and potentially useful information from data
- Exploration & analysis, by automatic or semi-automatic means, of large quantities of data in order to discover meaningful patterns
- End-to-end process of pattern discovery



What is (not) Data Mining?

- What is not Data Mining?

- Look up phone number in phone directory
- Query a database to retrieve particular information

- 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 (not so useful) patterns...

- “rules” for American presidents (before 2004 elections)

Some (not so useful) patterns...

- “rules” for American presidents (before 2004 elections)
 - if the Washington Redskins win their last home game before the election, the incumbent’s party will be re-elected

Some (not so useful) patterns...

- “rules” for American presidents (before 2004 elections)
 - if the Washington Redskins win their last home game before the election, the incumbent’s party will be re-elected
 - no Republican has ever won a presidential election without carrying Ohio

Some (not so useful) patterns...

- “rules” for American presidents (before 2004 elections)
 - if the Washington Redskins win their last home game before the election, the incumbent’s party will be re-elected
 - no Republican has ever won a presidential election without carrying Ohio
 - no incumbent with a four-letter last name has ever been re-elected (Polk, Taft, Ford, Bush Sr.)

Some (not so useful) patterns...

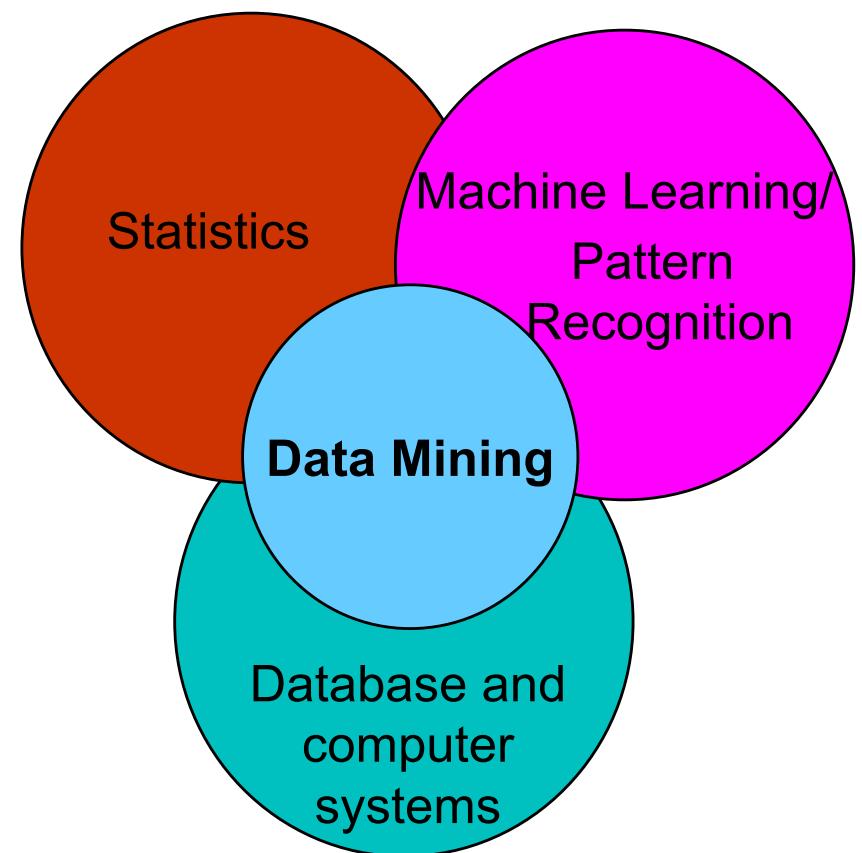
- “rules” for American presidents (before 2004 elections)
 - if the Washington Redskins win their last home game before the election, the incumbent’s party will be re-elected
 - no Republican has ever won a presidential election without carrying Ohio
 - no incumbent with a four-letter last name has ever been re-elected (Polk, Taft, Ford, Bush Sr.)
 - Americans won’t unseat a wartime President

Some (not so useful) patterns...

- “rules” for American presidents (before 2004 elections)
 - if the Washington Redskins win their last home game before the election, the incumbent’s party will be re-elected (**Redskins vs. Panthers: 13-21**)
 - no Republican has ever won a presidential election without carrying Ohio
 - no incumbent with a four-letter last name has ever been re-elected (Polk, Taft, Ford, Bush Sr.) (**GWB**)
 - Americans won’t unseat a wartime President

Origins of Data Mining

- Draws ideas from machine learning/AI, pattern recognition, statistics, and systems
- Traditional techniques may be unsuitable due to
 - enormity of data
 - high dimensionality of data
 - heterogeneous, distributed nature of data



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.

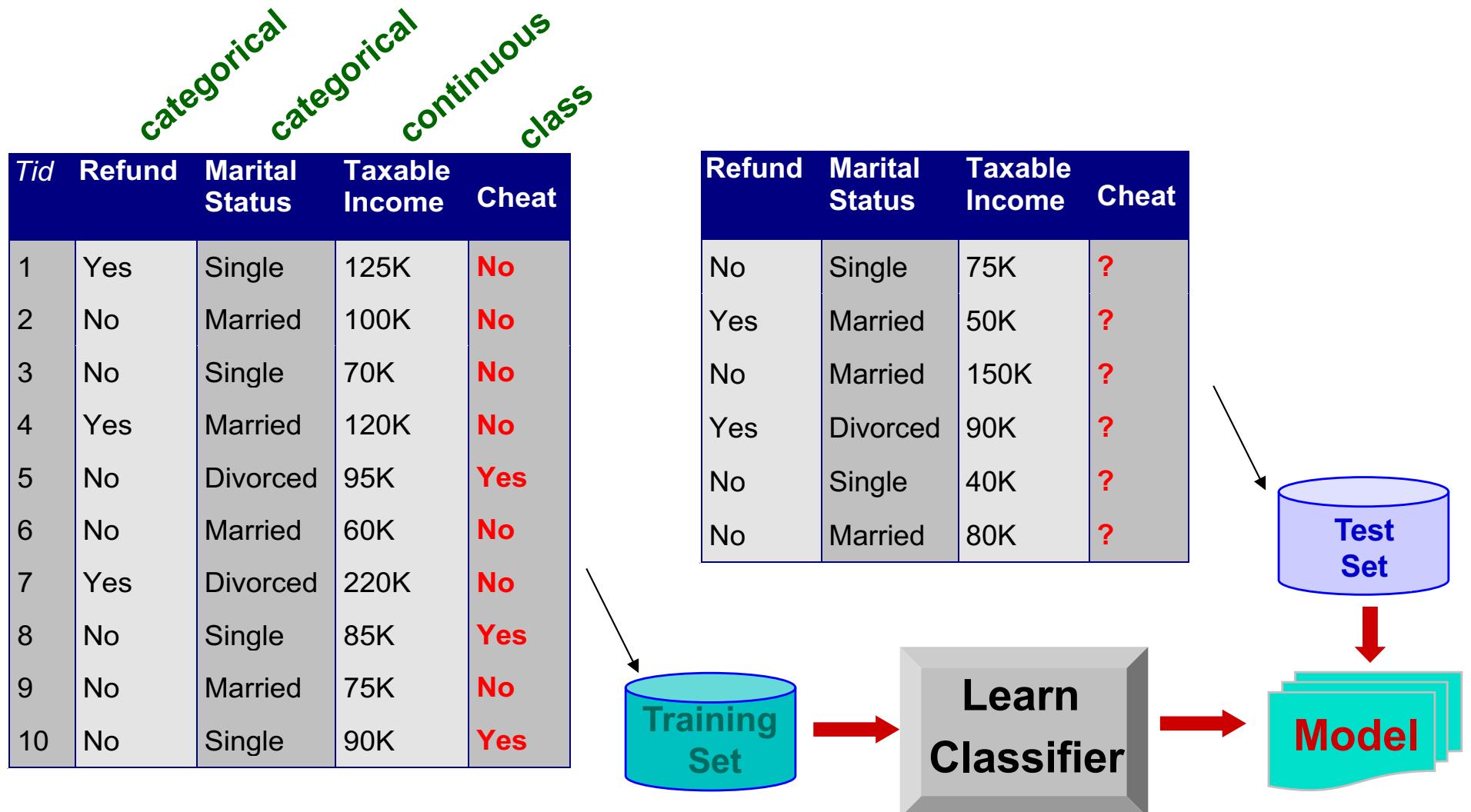
Data Mining Tasks...

- Classification [Predictive]
 - Regression [Predictive]
 - Deviation Detection [Predictive]
-
- Clustering [Descriptive]
 - Association Rule Discovery [Descriptive]
 - Sequential Pattern Discovery [Descriptive]

Classification: Definition

- Given a collection of objects (*training set*)
 - Objects may be records that contain a set of *attributes*; one of the attributes is the *class*. Objects may also be defined by a similarity measure between them.
- Find a *model* that maps an object into the class value; e.g. as a function of the values of other attributes.
- Goal: previously unseen objects 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.

Classification: Example



Classification: Application 1

● 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 stay, how much they earn, etc.
 - ◆ Use this information as input attributes to learn a classifier.

From [Berry & Linoff] Data Mining Techniques, 1997

Classification: Application 2

● Fraud Detection

- Goal: Predict 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: Application 3

- **Customer Attrition/Churn:**
 - Goal: To predict whether a customer is likely to be lost to a competitor.
 - Approach:
 - ◆ Use detailed record of transactions with each of the past and present customers, to find attributes.
 - How often the customer calls, where he calls, what time-of-the day he calls most, his financial status, marital status, etc.
 - ◆ Label the customers as loyal or disloyal.
 - ◆ Find a model for loyalty.

From [Berry & Linoff] Data Mining Techniques, 1997

Classification: Application 4

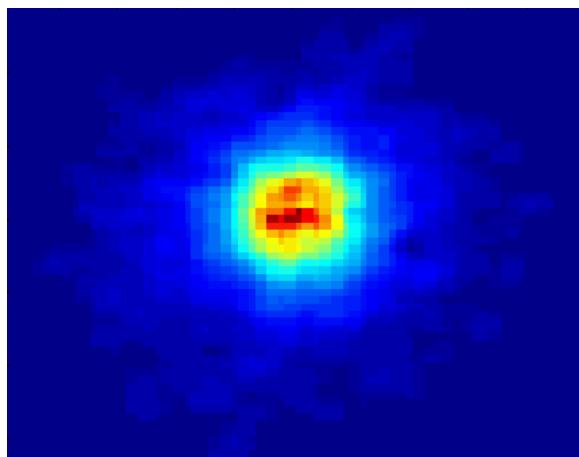
● Sky Survey Cataloging

- Goal: To predict class (star or galaxy) of sky objects, especially visually faint ones, based on the telescopic survey images (from Palomar Observatory).
 - 3000 images with $23,040 \times 23,040$ pixels per image.
- Approach:
 - ◆ Segment the image.
 - ◆ Measure image attributes (features) - 40 of them per object.
 - ◆ Model the class based on these features.
 - ◆ Success Story: Could find 16 new high red-shift quasars, some of the farthest objects that are difficult to find!

Classifying Galaxies

Courtesy: <http://aps.umn.edu>

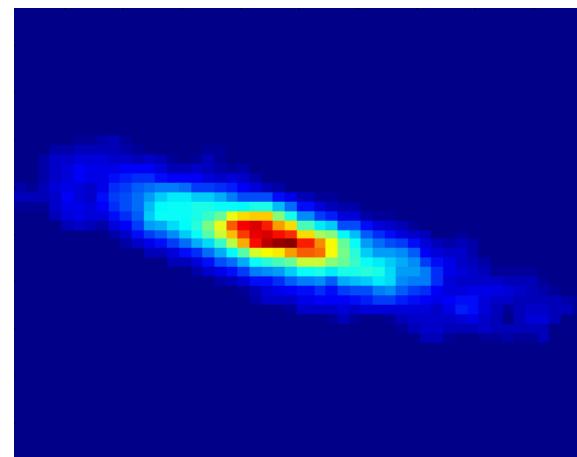
Early



Class:

- **Stages of Formation**

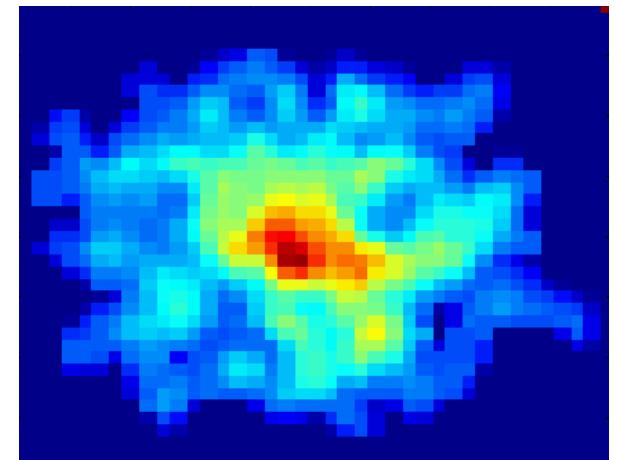
Intermediate



Attributes:

- Image features,
- Characteristics of light waves received, etc.

Late

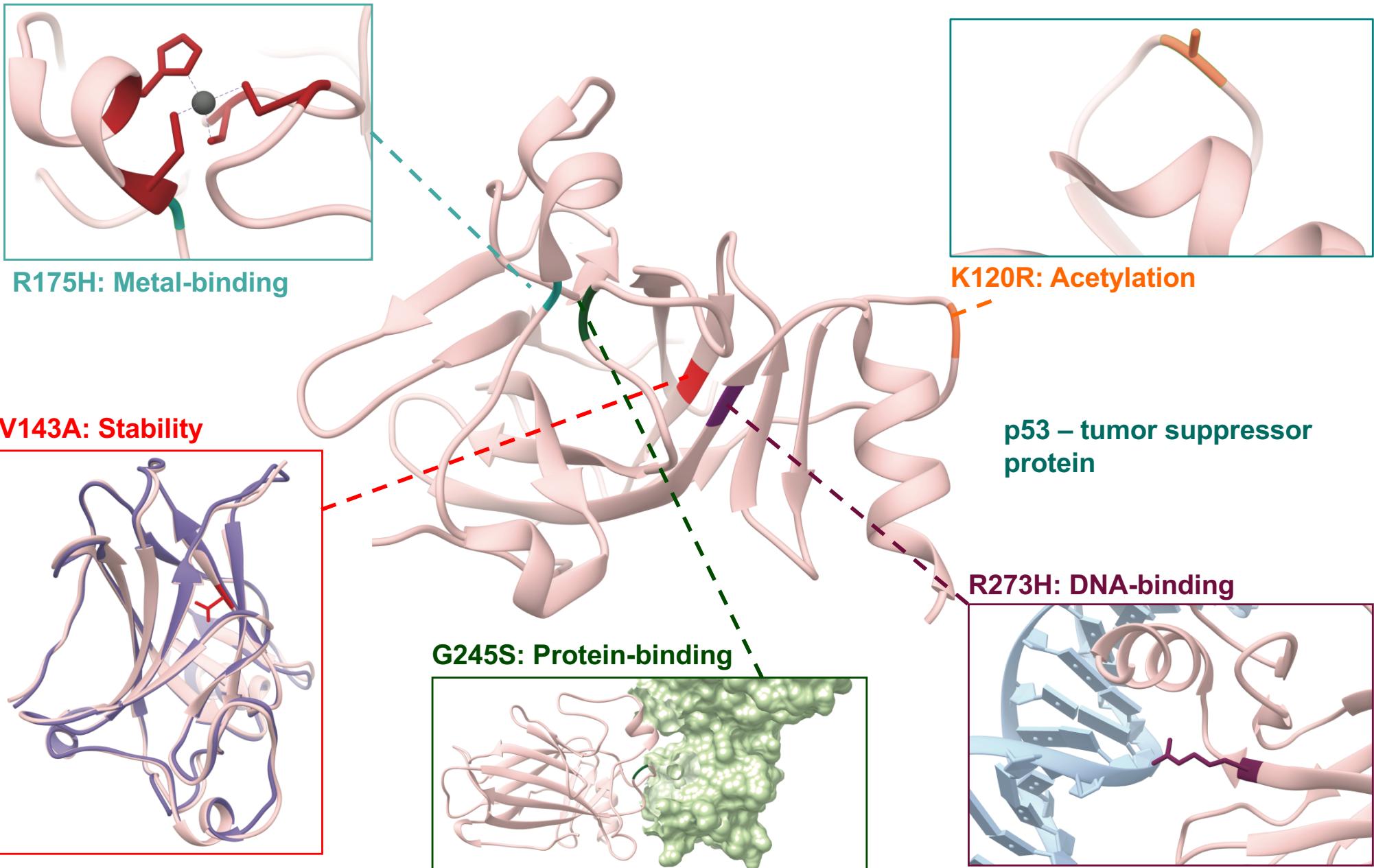


Data Size:

- 72 million stars, 20 million galaxies
- Object Catalog: 9 GB
- Image Database: 150 GB

Classification: Application 5

Disease mechanisms

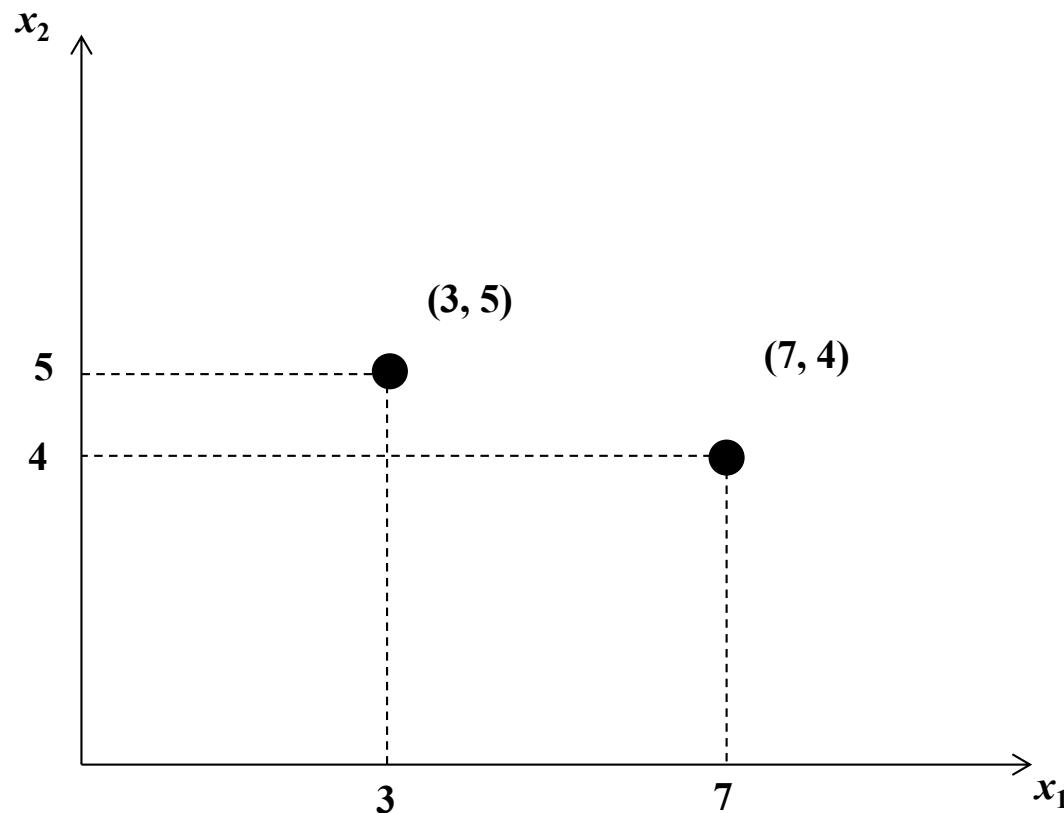


Toy classification problem

Assume medical situation:

x_1 is blood pressure

x_2 is temperature



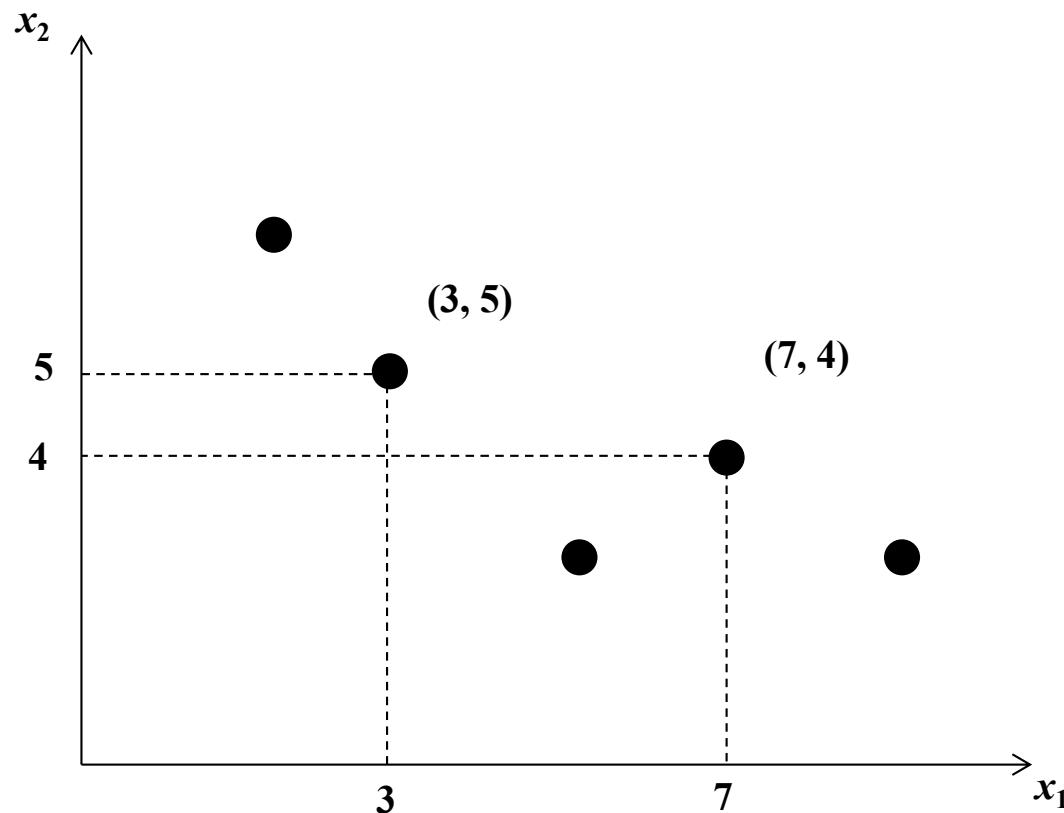
Then, each vector (x_1, x_2) corresponds to one patient

Toy classification problem

Assume medical situation:

x_1 is blood pressure

x_2 is temperature



Then, each vector (x_1, x_2) corresponds to one patient

Dataset

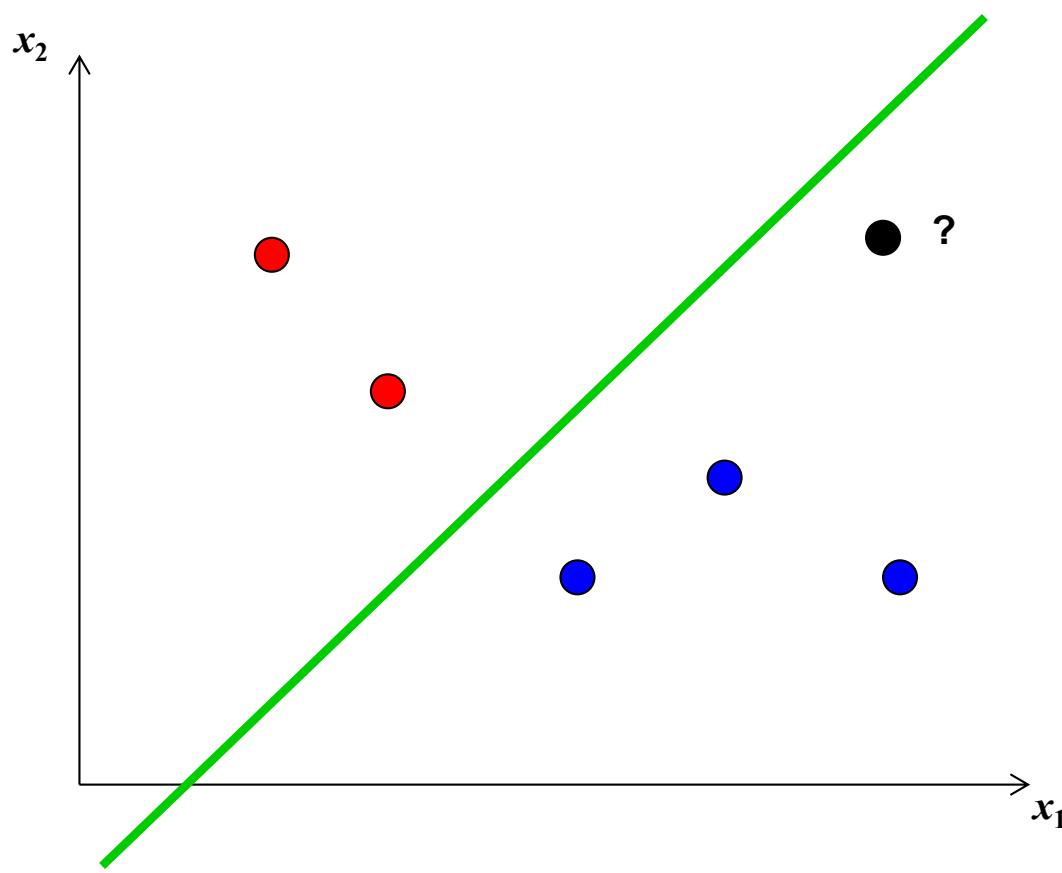
3	5
7	4
2	7
5	4
9	4

Toy classification problem

2 patients with class 0

3 patients with class 1

1 patient with unknown class (we need to predict)



Assume medical situation:

x_1 is blood pressure

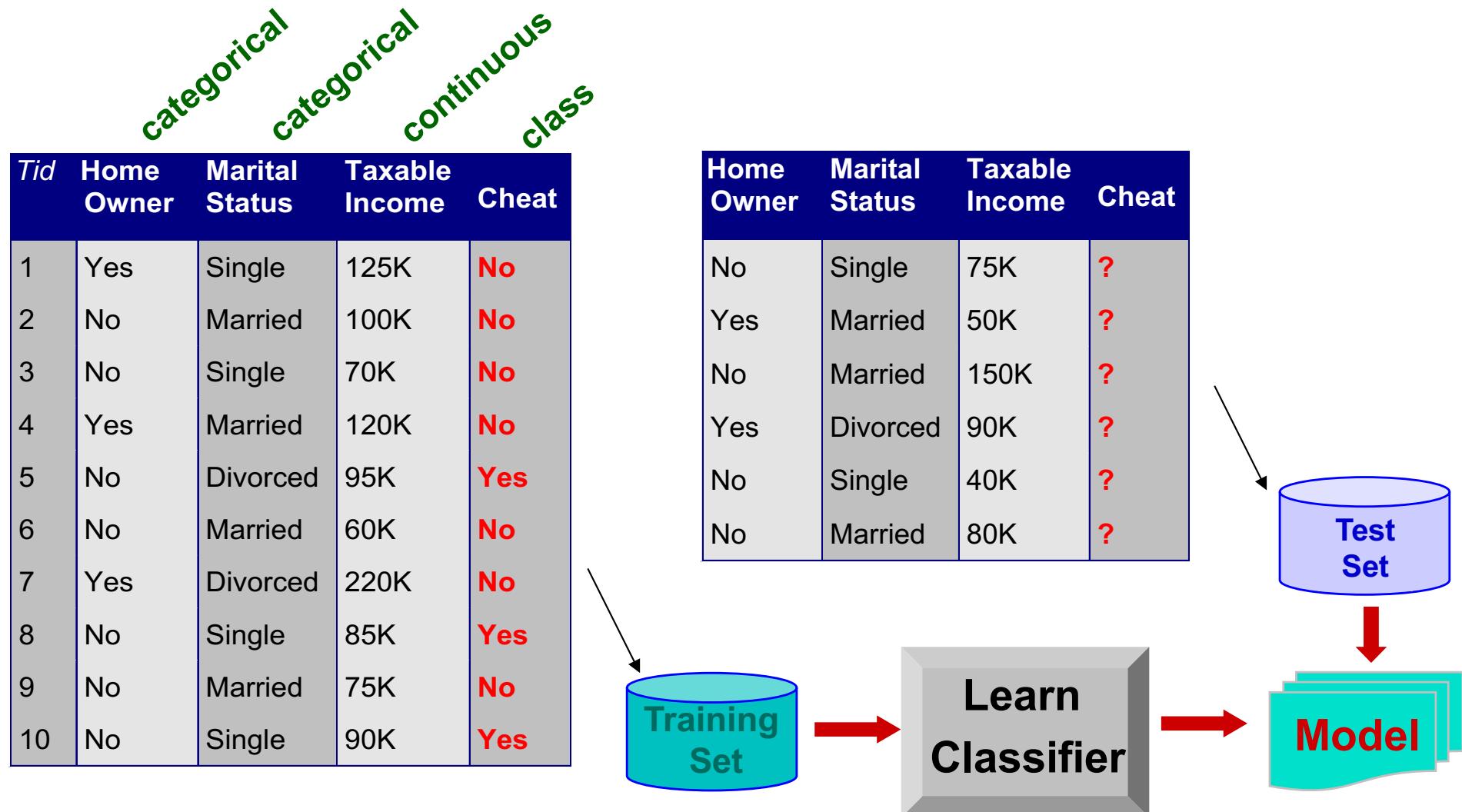
x_2 is temperature

Then, each vector (x_1, x_2) corresponds to one patient

Dataset

3	5	0
7	4	1
2	7	0
5	4	1
9	4	1

Do you see it differently now?



Clustering: Definition

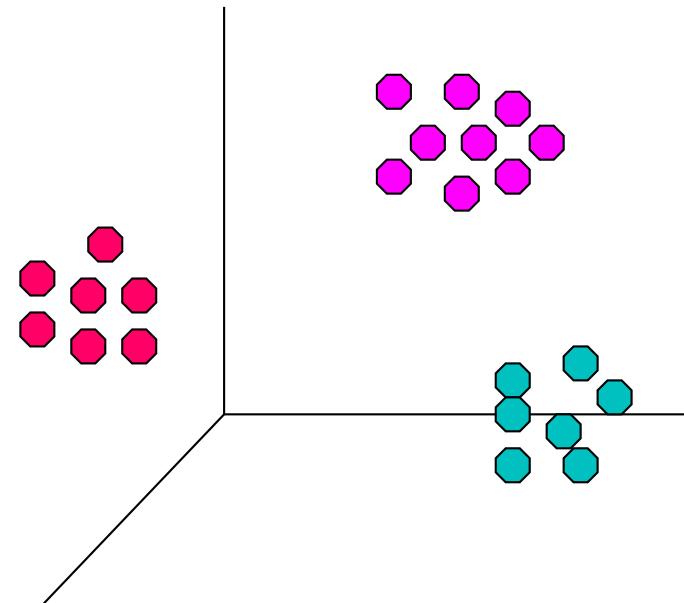
- Given a set of data points, each having a set of attributes, and a similarity measure among them, find clusters such that
 - Data points in one cluster are more similar to one another.
 - Data points in separate clusters are less similar to one another.
- Similarity Measures:
 - Euclidean distance if attributes are continuous.
 - Other problem-specific measures.

Illustrating Clustering

| Euclidean Distance Based Clustering in 3-D space.

Intracluster distances
are minimized

Intercluster distances
are maximized



Clustering: Application 1

- **Market Segmentation:**

- Goal: subdivide a market into distinct subsets of customers where any subset may conceivably be selected as a market target to be reached with a distinct marketing mix.
 - Approach:
 - ◆ Collect different attributes of customers based on their geographical and lifestyle related information.
 - ◆ Find clusters of similar customers.
 - ◆ Measure the clustering quality by observing buying patterns of customers in the same cluster vs. those from different clusters.

Clustering: Application 2

- **Document Clustering:**
 - Goal: To find groups of documents that are similar to each other based on the important terms appearing in them.
 - Approach:
 - ◆ To identify frequently occurring terms in each document, form a similarity measure based on the frequencies of different terms. Use it to cluster.
 - Gain: Information retrieval can utilize the clusters to relate a new document or search term to clustered documents.

Illustrating Document Clustering

- Clustering Points: 3204 Articles of Los Angeles Times.
- Similarity Measure: How many words are common in these documents (after some word filtering).

<i>Category</i>	<i>Total Articles</i>	<i>Correctly Placed</i>
<i>Financial</i>	555	364
<i>Foreign</i>	341	260
<i>National</i>	273	36
<i>Metro</i>	943	746
<i>Sports</i>	738	573
<i>Entertainment</i>	354	278

Clustering of S&P 500 Stock Data

- Observe Stock Movements every day.
- Clustering points: Stock-{UP/DOWN}
- Similarity Measure: Two points are more similar if the events described by them frequently happen together on the same day.
 - We used association rules to quantify a similarity measure.

	<i>Discovered Clusters</i>	<i>Industry Group</i>
1	Applied-Matl-DOWN,Bay-Network-Down,3-COM-DOWN, Cabletron-Sys-DOWN,CISCO-DOWN,HP-DOWN, DSC-Comm-DOWN,INTEL-DOWN,LSI-Logic-DOWN, Micron-Tech-DOWN,Texas-Inst-Down,Tellabs-Inc-Down, Natl-Semiconduct-DOWN,Oracl-DOWN,SGI-DOWN, Sun-DOWN	Technology1-DOWN
2	Apple-Comp-DOWN,Autodesk-DOWN,DEC-DOWN, ADV-Micro-Device-DOWN,Andrew-Corp-DOWN, Computer-Assoc-DOWN,Circuit-City-DOWN, Compaq-DOWN, EMC-Corp-DOWN, Gen-Inst-DOWN, Motorola-DOWN,Microsoft-DOWN,Scientific-Atl-DOWN	Technology2-DOWN
3	Fannie-Mae-DOWN,Fed-Home-Loan-DOWN, MBNA-Corp-DOWN,Morgan-Stanley-DOWN	Financial-DOWN
4	Baker-Hughes-UP,Dresser-Inds-UP,Halliburton-HLD-UP, Louisiana-Land-UP,Phillips-Petro-UP,Unocal-UP, Schlumberger-UP	Oil-UP

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.

TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

Rules Discovered:

$\{\text{Milk}\} \rightarrow \{\text{Coke}\}$

$\{\text{Diaper}, \text{Milk}\} \rightarrow \{\text{Beer}\}$

Association Rule Discovery: Application 1

- **Marketing and Sales Promotion**

- Let the rule discovered be
 $\{Bagels, \dots\} \rightarrow \{Potato\ Chips\}$
 - Potato Chips as consequent => Can be used to determine what should be done to boost its sales.
 - Bagels in the antecedent => Can be used to see which products would be affected if the store discontinues selling bagels.
 - Bagels in antecedent and Potato chips in consequent => 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.
 - A classic rule --
 - ◆ If a customer buys diaper and milk, then he is very likely to buy beer.
 - ◆ So, don't be surprised if you find six-packs stacked next to diapers!

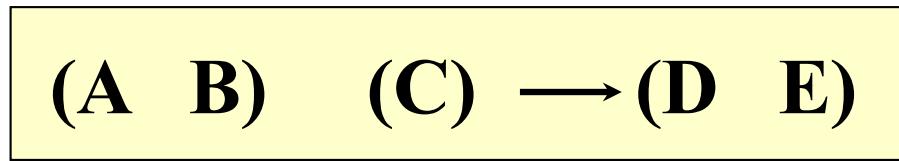
Association Rule Discovery: Application 3

● Inventory Management

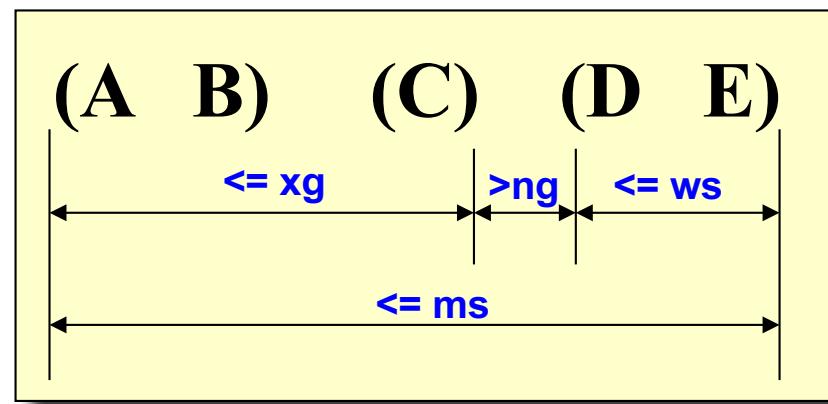
- Goal: A consumer appliance repair company wants to anticipate the nature of repairs on its consumer products and keep the service vehicles equipped with right parts to reduce on number of visits to consumer households.
- Approach: Process the data on tools and parts required in previous repairs at different consumer locations and discover the co-occurrence patterns.

Sequential Pattern Discovery: Definition

- Given is a set of *objects*, with each object associated with its own *timeline of events*, find rules that predict strong **sequential dependencies** among different events.



- Rules are formed by first discovering patterns. Event occurrences in the patterns are governed by timing constraints.



Sequential Pattern Discovery: Examples

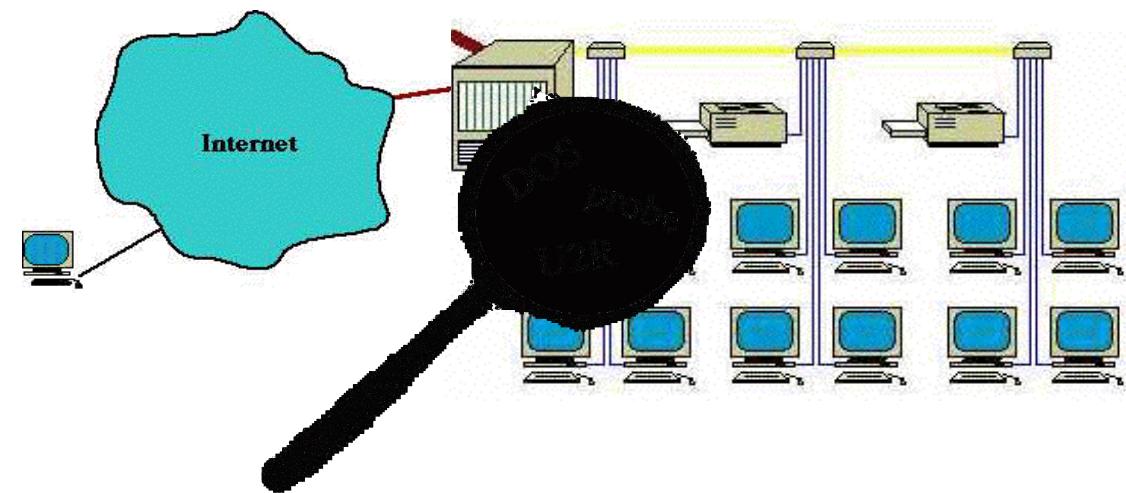
- In telecommunications alarm logs,
 - (Inverter_Problem Excessive_Line_Current)
(Rectifier_Alarm) --> (Fire_Alarm)
- In point-of-sale transaction sequences,
 - Computer Bookstore:
(Intro_To_Visual_C) (C++_Primer) -->
(Perl_for_dummies,Tcl_Tk)
 - Athletic Apparel Store:
(Shoes) (Racket, Racketball) --> (Sports_Jacket)

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.
- Greatly studied in statistics, machine learning, etc.
- Examples:
 - Predicting sales amounts of new product based on advertising expenditure.
 - Predicting wind velocities as a function of temperature, humidity, air pressure, etc.
 - Time series prediction of stock market indices.

Deviation/Anomaly Detection

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



Typical network traffic at University level may reach over 100 million connections per day

(Some) Challenges of Data Mining

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