

# Data Mining

---

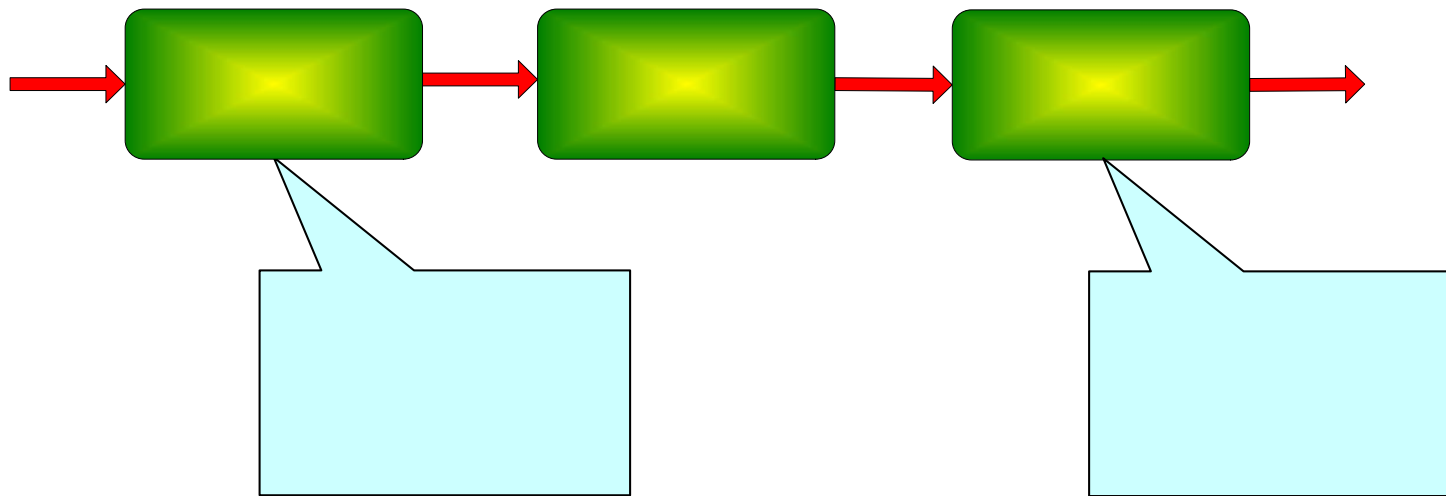
---

## Lecture 1: Introduction

# What is Data Mining?

---

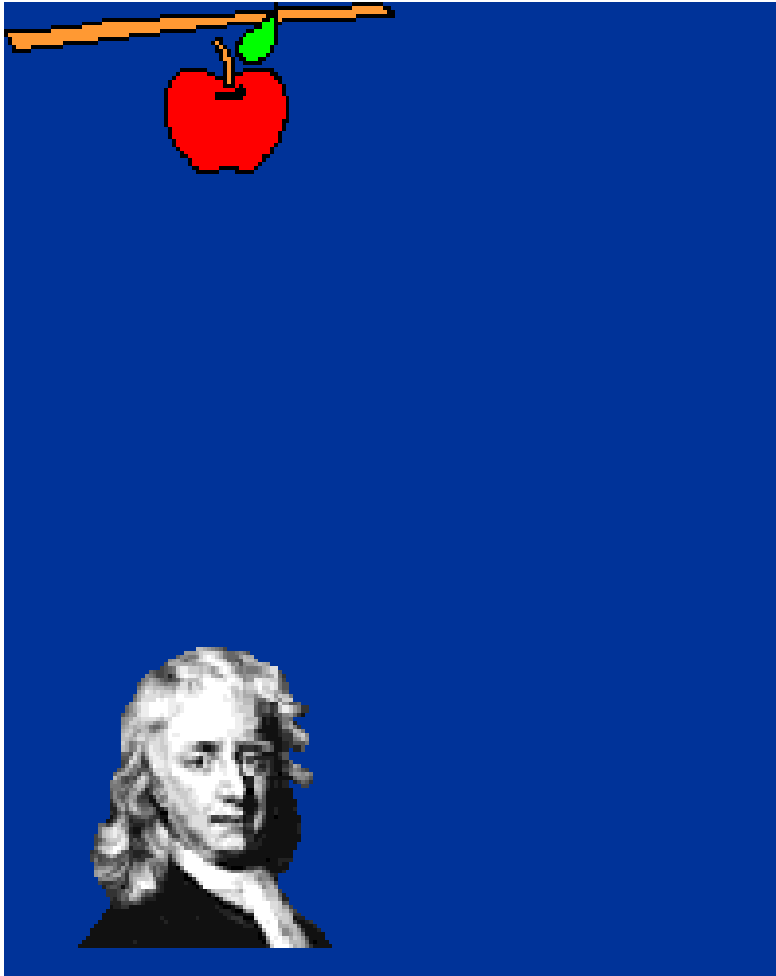
- **Non-trivial extraction of implicit, previously unknown and potentially useful information from data**



**Knowledge Discovery in Databases (KDD)**

# Knowledge Discovery in the old days

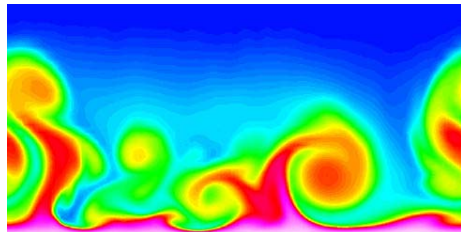
---



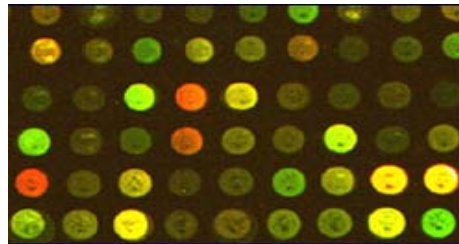
1. Observe phenomena
2. Formulate theory
3. Validate theory via experimentation

From: <http://csep10.phys.utk.edu/astr161/lect/history/newtongrav.html>

# Nowadays...



**Data,  
Data  
Data  
Everywhere**



1. Lots of data
2. Lack of theory
3. Data mining can help to generate new hypothesis or help analysts to make sense out of the data

# 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’Rourke, O’Reilly... in Boston area)
- Group together similar documents returned by search engine according to their context (e.g. Amazon rainforest, Amazon.com,)

# Why Mine Data? Commercial Viewpoint

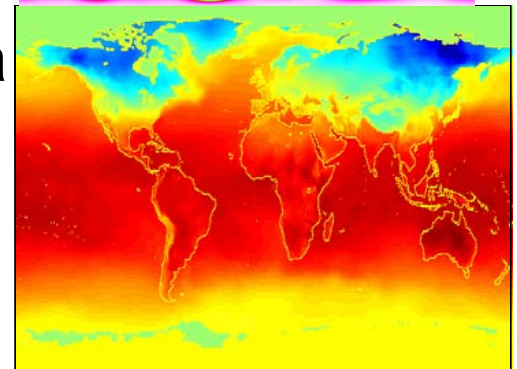
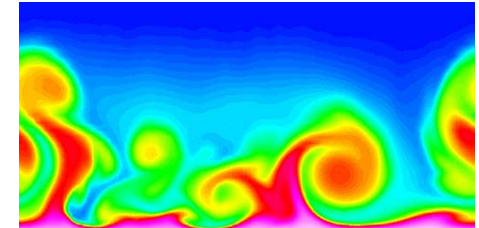
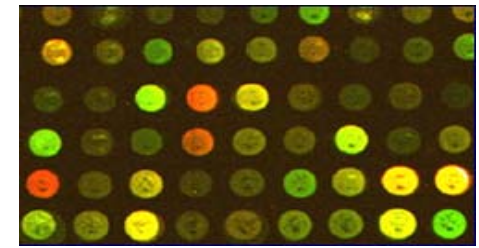
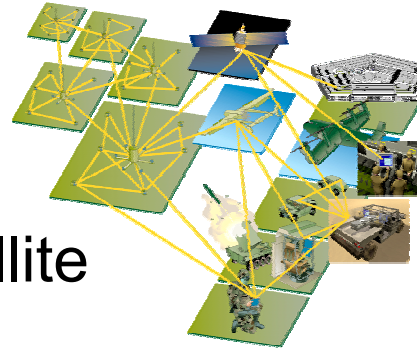
---

- Lots of data is being collected and warehoused
  - Walmart: ~20 million txn/day
  - Google: > 3 billion Web pages
  - Yahoo: ~10 GB Web data/hr
- Computers have become cheaper and more powerful
- Competitive Pressure is Strong
  - Provide better, customized services for an *edge* (e.g. in Customer Relationship Management)



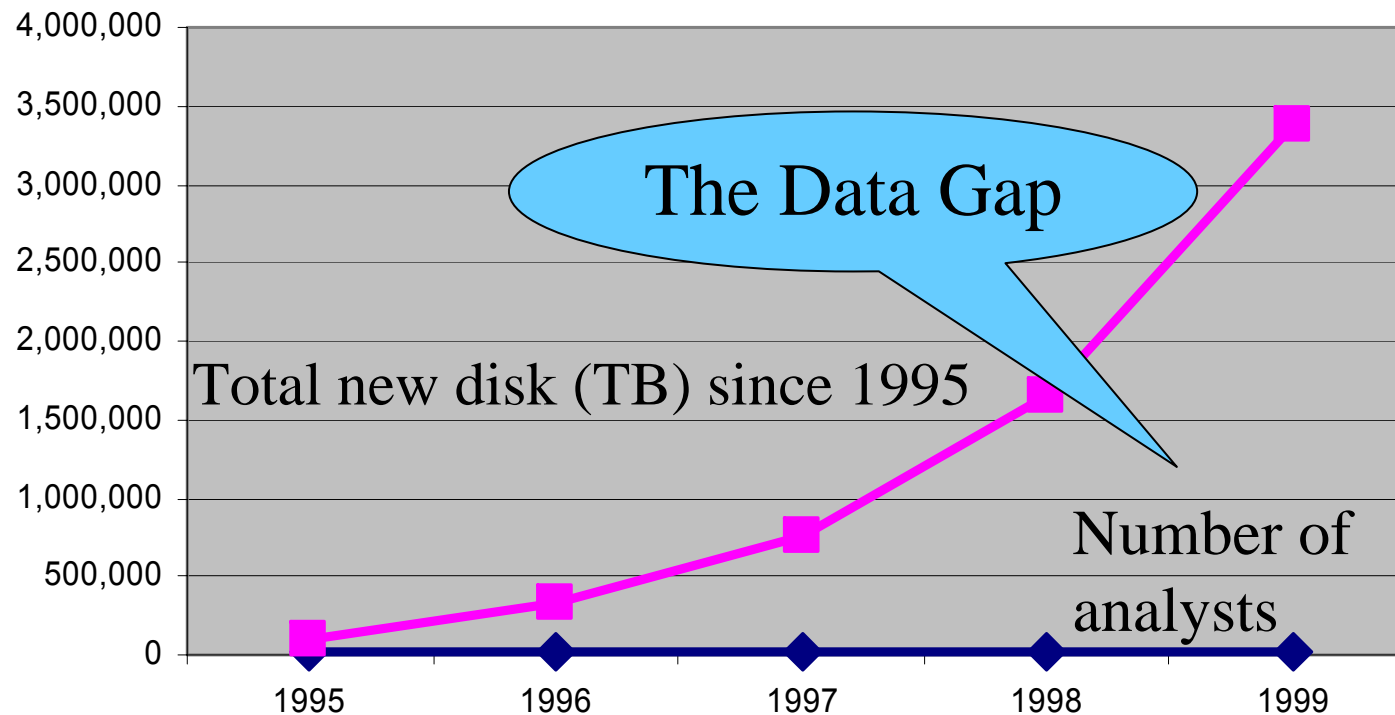
# Why Mine Data? Scientific Viewpoint

- 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
- Need techniques that can automatically analyze the data and form new hypotheses for further evaluation by scientists



# 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



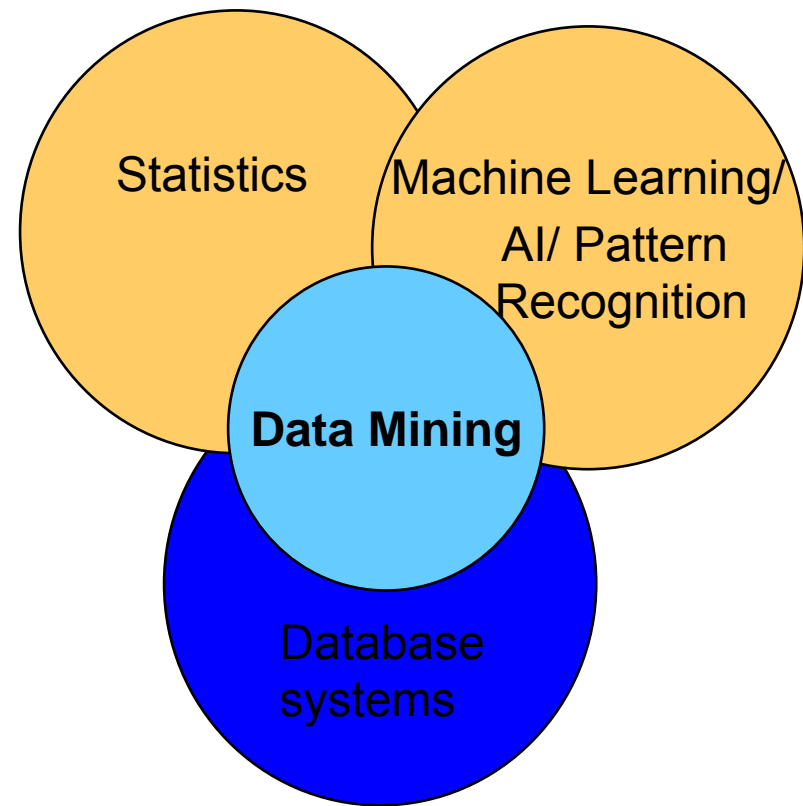
From: R. Grossman, C. Kamath, V. Kumar, “Data Mining for Scientific and Engineering Applications”



# Origins of Data Mining

---

- Draws ideas from machine learning/AI, pattern recognition, statistics, and database systems
- Traditional Techniques may be unsuitable due to
  - Enormity of data
  - High dimensionality of data
  - Heterogeneous, distributed nature of data



# Data Mining Tasks

---

- Predictive Methods
  - Use some variables to predict unknown or values of other variables.
- Descriptive Methods
  - Find human-interpretable patterns that describe the data.

From [Fayyad, et.al.] Advances in Knowledge Discovery and Data Mining, 1996

# Data Mining Tasks...

---

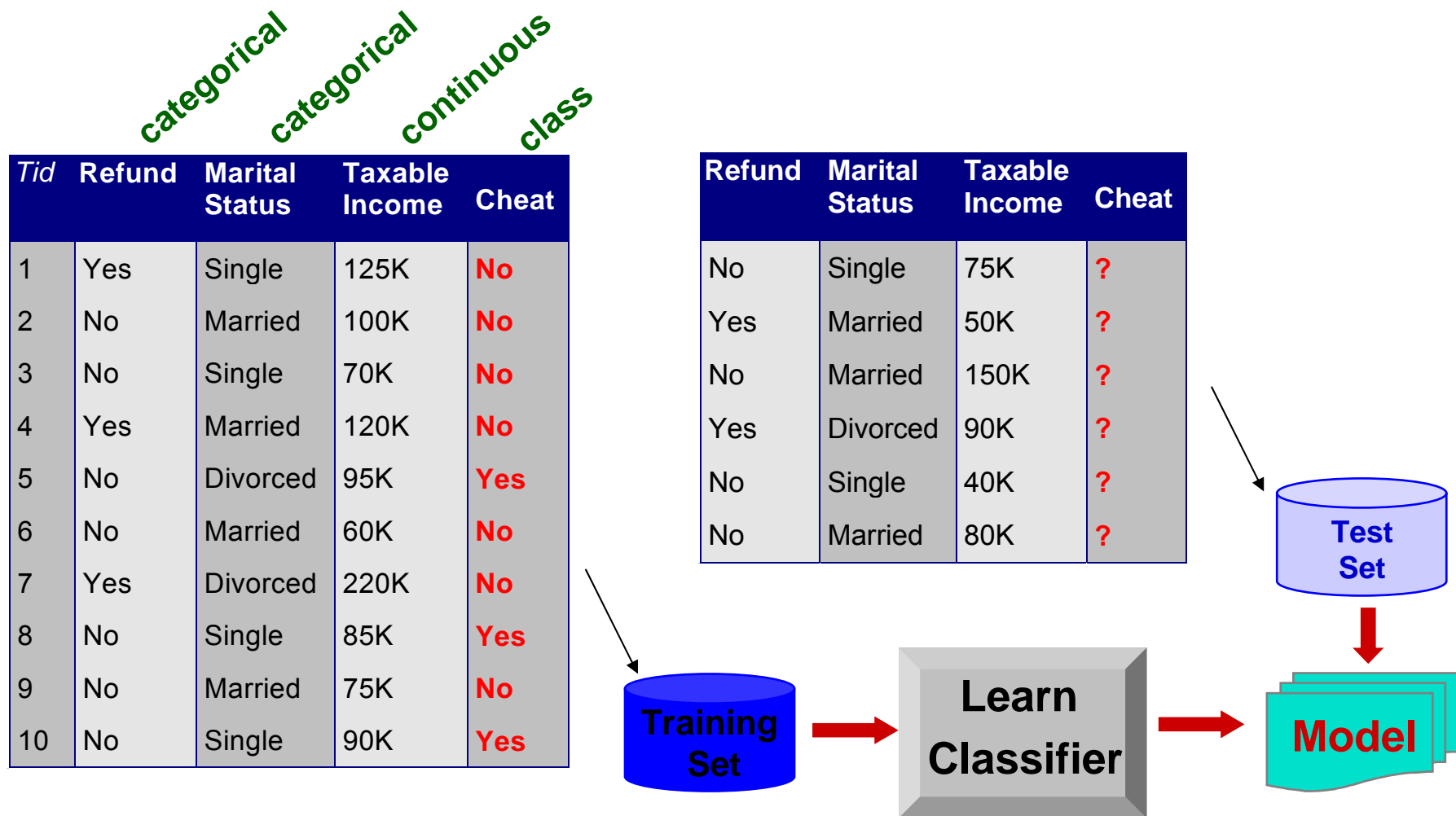
- Classification [Predictive]
- Clustering [Descriptive]
- Association Rule Discovery [Descriptive]
- Sequential Pattern Discovery [Descriptive]
- Regression [Predictive]
- Anomaly Detection [Predictive]

# Classification: Definition

---

- Given a collection of records (*training set*)
  - Each record contains a set of *attributes*, one of the attributes is the *class*.
- Task: 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.

# Classification Example



# 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 4

---

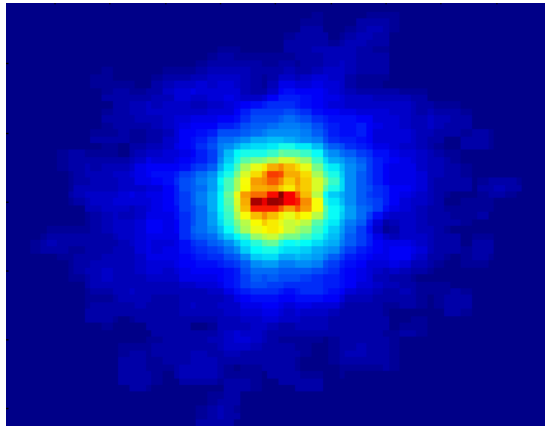
- Sky Survey Cataloging
  - Goal: To predict class (star or galaxy) of sky objects, especially visually faint ones, based on the telescopic survey images (from Palomar Observatory).
    - 3000 images with 23,040 x 23,040 pixels per image.
  - Approach:
    - ◆ Segment the image.
    - ◆ Measure image 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!

From [Fayyad, et.al.] Advances in Knowledge Discovery and Data Mining, 1996

# Classifying Galaxies

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

*Early*



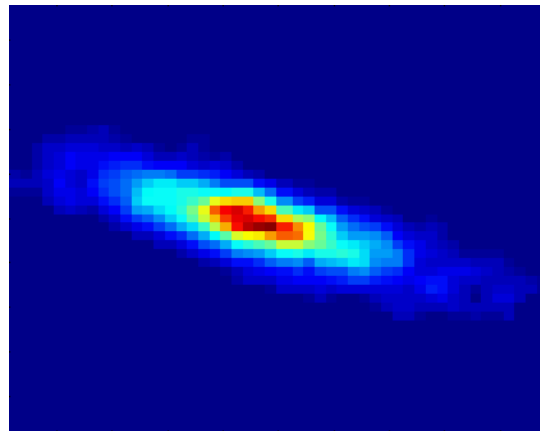
**Class:**

- Stages of Formation

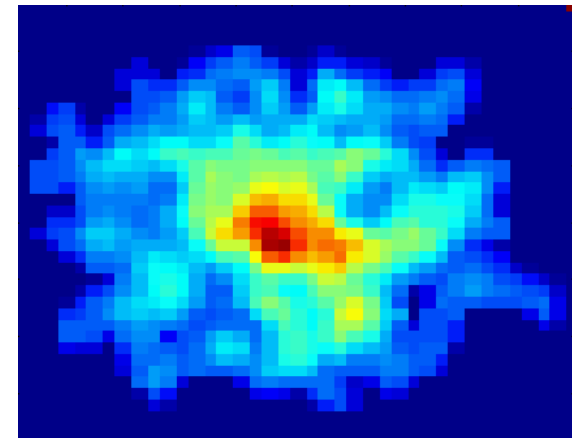
**Attributes:**

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

*Intermediate*



*Late*



**Data Size:**

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



# 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

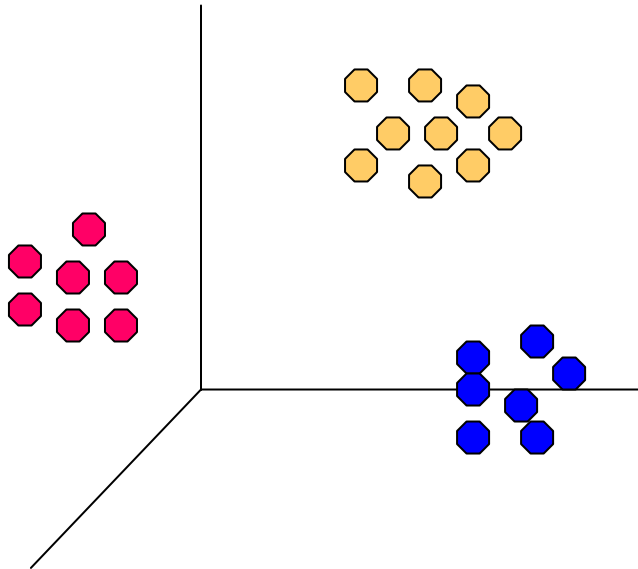
---

---

☒ Euclidean Distance Based Clustering in 3-D space.

Intracluster distances  
are minimized

Intercluster distances  
are maximized



# Clustering 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>
<b>1</b>	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, Orac1-DOWN, SGI-DOWN, Sun-DOWN	Technology1-DOWN
<b>2</b>	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
<b>3</b>	Fannie-Mae-DOWN, Fed-Home-Loan-DOWN, MBNA-Corp-DOWN, Morgan-Stanley-DOWN	Financial-DOWN
<b>4</b>	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.

<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}**

**{Diaper, Milk} --> {Beer}**

# 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.

# Regression

---

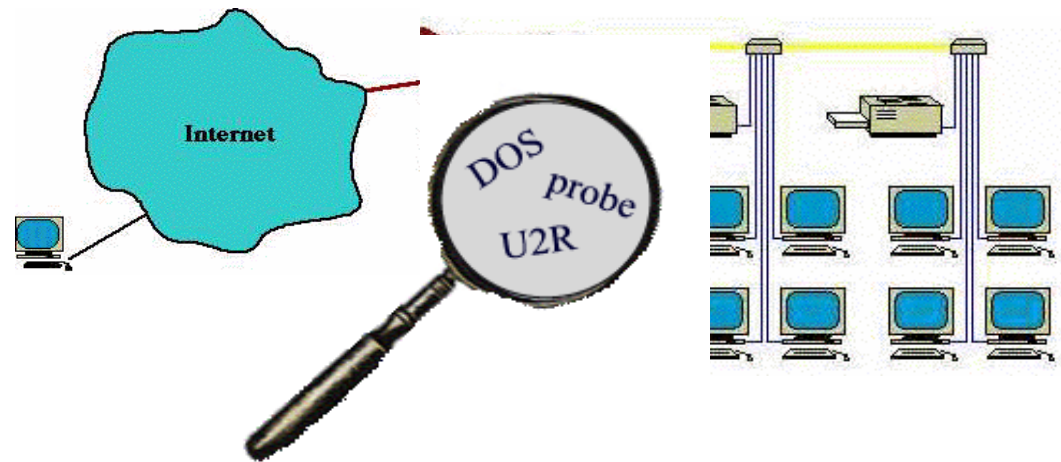
- Predict a value of a given continuous valued variable based on the values of other variables, assuming a linear or nonlinear model of dependency.
- Greatly studied in statistics, neural network fields.
- Examples:
  - Predicting sales amounts of new product based on advertising expenditure.
  - Predicting 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*