

---

# Data Science

---

---

# Large-scale Data is Everywhere!

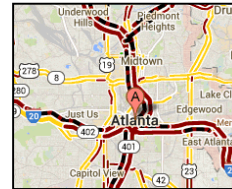
- There has been enormous data growth in both commercial and scientific databases due to advances in data generation and collection technologies
- New mantra
  - Gather whatever data you can whenever and wherever possible.
- Expectations
  - Gathered data will have value either for the purpose collected or for a purpose not envisioned.



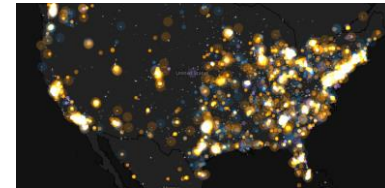
**Cyber Security**



**E-Commerce**



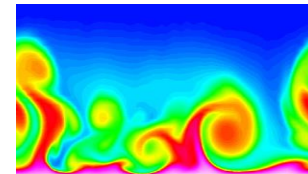
**Traffic Patterns**



**Social Networking: Twitter**



**Sensor Networks**



**Computational Simulations**

# Why Data Science? Commercial Viewpoint

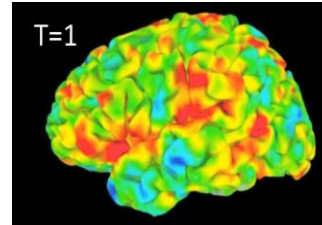
---

- Lots of data is being collected and warehoused
  - Web data
    - ◆ Google has Peta Bytes of web data
    - ◆ Facebook has billions of active users
  - purchases at department/grocery stores, e-commerce
    - ◆ Amazon handles millions of visits/day
  - 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)



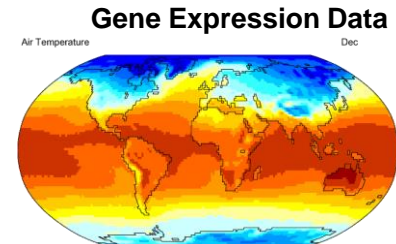
# Why Data Science? Scientific Viewpoint

- Data collected and stored at enormous speeds
  - remote sensors on a satellite
    - ◆ NASA EOSDIS archives over petabytes of earth science data / year
  - telescopes scanning the skies
    - ◆ Sky survey data
  - High-throughput biological data
  - scientific simulations
    - ◆ terabytes of data generated in a few hours
- Data science helps scientists
  - in automated analysis of massive datasets
  - In hypothesis formation



fMRI Data from Brain

Sky Survey Data

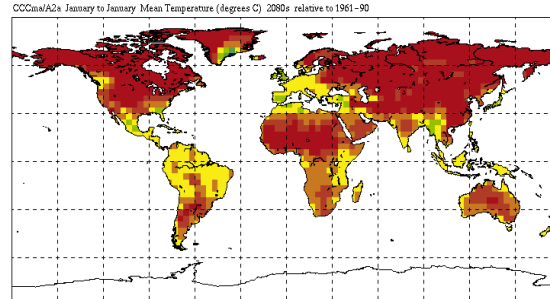


Surface Temperature of Earth

# Great Opportunities to Solve Society's Major Problems



Improving health care and reducing costs



Predicting the impact of climate change



Finding alternative/ green energy sources



Reducing hunger and poverty by increasing agriculture production

# What is Data Science?

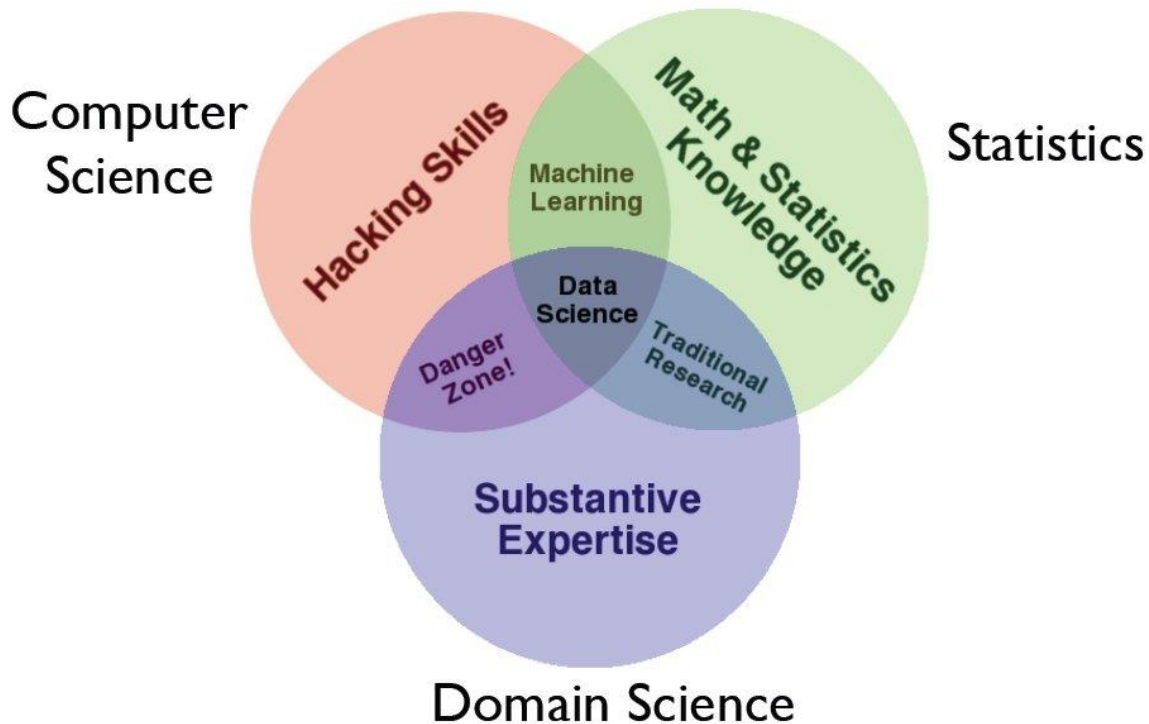
---

Like any emerging field, it isn't yet well defined, but incorporates elements of:

- Exploratory Data Analysis and Visualization
  - Machine Learning and Statistics
  - High-Performance Computing technologies for dealing with scale.
-

# Skill Sets for Data Science

---



# Appreciating Data

---

Computer Scientists do not naturally appreciate data: it's just stuff to run through a program.

The usual way to test algorithm performance is to run the implementation on “random data”.

But interesting data sets are a scarce resource, which requires hard work and imagination to obtain.

---



# Computer vs. Real Scientists (1)

---

- Scientists strive to understand the complicated and messy natural world, while computer scientists build their own clean and organized virtual worlds. Thus:
  - Nothing is ever completely true or false in science, while everything is either true or false in Computer Science / Mathematics.
-

# Computer vs. Real Scientists (2)

---

- Scientists are data-driven, while computer scientists are algorithm-driven.
  - Scientists obsess about discovering things, which computer scientists invent rather than discover.
  - Scientists are comfortable with the idea that data has errors; computer scientists are not.
-

# Genius vs. Wisdom

---

Software developers are hired to produce code.

Data Scientists are hired to produce insights.

Genius shows in finding the right answer!!!

Wisdom shows in avoiding the wrong answers.

Data science (like most things) benefits more from wisdom than from genius.

---

# Developing Wisdom

---

- Wisdom comes from experience.
- Wisdom comes from general knowledge.
- Wisdom comes from listening to others.
- Wisdom comes from humility, observing how often you have been wrong and why/how.

I seek pass on wisdom, through my experience on the difficulty of making good predictions.

---

# Developing Curiosity

---

- The good data scientist develops a curiosity about the domain/application they are working in.
  - They talk shop with the people whose data they are working on.
  - They read the newspaper every day, to get a broader perspective on the world.
-

# Asking Good Questions

---

Software developers are not encouraged to ask questions, but data scientists are:

- What exciting things might you be able to learn from a given data set?
  - What things do you/your people really want to know?
  - What data sets might get you there?
-

# Let's Practice Asking Questions!

---

Who, What, Where, When, and Why on the following datasets:

- [Baseball-reference.com](http://baseball-reference.com)
  - Google ngrams
  - NYC taxi cab records
-

# Baseball-Reference.com: biosketch



play index **players** teams seasons managers leaders awards postseason boxes japan nlb minors draft

Mobile Site You Are Here > Home > Encyclopedia of Players > R Listing > Babe Ruth Statistics and

News: s-r blog:KBO Stats back to 1999 - Baseball-Reference.com

Babe Ruth Player Page > Batting Pitching Fielding Minors News Archive (1456) Bullpen Oracle



## Babe Ruth

Like 1,213 people like this. +25 Recommend this

George Herman Ruth (Babe, The Bambino or The Sultan Of Swat)

**Positions:** Outfielder and Pitcher  
**Bats:** Left, **Throws:** Left  
**Height:** 6' 2", **Weight:** 215 lb.

**Born:** February 6, 1895 in Baltimore, MD   
**High School:** St. Mary's HS (Baltimore, MD) (All Transactions)  
**Debut:** July 11, 1914 (Age 19.155)  
**Rookie Status:** Exceeded rookie limits during 1915 season [\*]  
**Teams** (by GP): Yankees/RedSox/Braves 1914-1935

**Final Game:** May 30, 1935 (Age 40.113)  
**Inducted** into the Hall of Fame by BBWAA as Player in 1936 (215/226 ballots). Induction ceremony in [View Babe Ruth Page](#) at the Baseball Hall of Fame (plaque, photos, videos).  
**Died:** August 16, 1948 in New York, NY (Aged 53.192)  
**Buried:** Gate of Heaven Cemetery, Hawthorne, NY  
**View Player Bio** from the [SABR BioProject](#)  
[About biographical information](#)



S-R: M

## Transactions

**July 9, 1914:** Purchased with [Ernie Shore](#) and [Ben Egan](#) by the [Boston Red Sox](#) from Baltimore (International) for more than \$25000. more than \$25000  
**December 26, 1919:** Purchased by the [New York Yankees](#) from the [Boston Red Sox](#) for \$100,000.  
**February 26, 1935:** Released by the [New York Yankees](#).  
**February 26, 1935:** Signed as a Free Agent with the [Boston Braves](#).

The transaction information used here was obtained free of charge from and is copyrighted by [RetroSheet](#). We attempt to update transactions throughout the season.

## Salaries

Convert to YYYY \$'s Salaries may not be complete (especially pre-1985) and may not include some earned bonuses

Year	Age	Team	Salary	ServTm (OpnDay)	Sources	Notes/Other Sources
1914	19	Boston Red Sox	\$2,500	?	Bill James Historical Abstract	Annualized rate; came up late in season
1915	20	Boston Red Sox	\$3,500	?	Bill James Historical Abstract	
1916	21	Boston Red Sox	\$3,500	?	Contract at HOF	
1917	22	Boston Red Sox	\$3,500	?	Contract at HOF	BJHA: \$5,000; Baseball Timeline \$7,000
1918	23	Boston Red Sox	\$9,000	?	Allan Wood, 1918, at 183	Includes \$1,000 midseason raise, \$1,000 WS bonus
1919	24	New York Yankees	\$10,000*	?	Michael Haupert research of HOF contracts	Contract at HOF:10000.00,
1920	25	New York Yankees	\$20,000*	?	Michael Haupert research of HOF contracts	Bill James Historical Abstract:20000.00,
1921	26	New York Yankees	\$20,000*	?	Michael Haupert research of HOF contracts	Bill James Historical Abstract:30000.00,Plus \$5K for '20 and '21 exhibitions, \$50/HR (\$9)m
1922	27	New York Yankees	\$52,000*	?	Michael Haupert research of HOF contracts	Bill James Historical Abstract:52000.00,
1923	28	New York Yankees	\$52,000*	?	Michael Haupert research of HOF contracts	Bill James Historical Abstract:52000.00,
1924	29	New York Yankees	\$52,000*	?	Michael Haupert research of HOF contracts	Bill James Historical Abstract:52000.00,
1925	30	New York Yankees	\$52,000*	?	Michael Haupert research of HOF contracts	Bill James Historical Abstract:52000.00,
1926	31	New York Yankees	\$52,000*	?	Michael Haupert research of HOF contracts	Bill James Historical Abstract:52000.00,
1927	32	New York Yankees	\$52,000*	?	Michael Haupert research of HOF contracts	S/23/27 AL letter:70000.00,
1928	33	New York Yankees	\$52,000*	?	Michael Haupert research of HOF contracts	S/23/27 AL letter:70000.00,
1929	34	New York Yankees	\$52,000*	?	Michael Haupert research of HOF contracts	S/23/27 AL letter:70000.00,
1930	35	New York Yankees	\$70,000*	?	Michael Haupert research of HOF contracts	Bill James Historical Abstract:80000.00,
1931	36	New York Yankees	\$70,000*	?	Michael Haupert research of HOF contracts	Bill James Historical Abstract:80000.00,
1932	37	New York Yankees	\$70,000*	?	Michael Haupert research of HOF contracts	M. Smelser, Life That Ruth Built, p. 441:75000.00,Plus 25% of all exhibition-game profits
1933	38	New York Yankees	\$80,000*	?	Michael Haupert research of HOF contracts	M. Smelser, Life That Ruth Built, p. 456:52000.00,Plus 25% of revenue from in-season exhibitions
1934	39	New York Yankees	\$80,000*	?	Michael Haupert research of HOF contracts	1/16/36 TSN, per government report:36696.00,\$35,000 salary plus 25% of exhibition profits
1935	40	New York Yankees	\$75,000*	?	Michael Haupert research of HOF contracts	Bill James Historical Abstract:35000.00,Annualized rate; retired early in season
1936	41	New York Yankees	\$52,000*	?	Michael Haupert research of HOF contracts	
1937	42	New York Yankees	\$35,000	?	Michael Haupert research of HOF contracts	

Career to date (may be incomplete) **\$1,020,000**





# Baseball Questions

---

- How to best measure individual player's skill, value or performance?
  - How fair do trades between teams work out?
  - What is the trajectory of player's performances as they mature and age?
  - To what extent does batting performance correlate with the position played?
-

# Demographic Questions

---

- Do left-handed people have shorter lifespans than right-handers?
  - How often do people return to where they were born?
  - Do player salaries reflect past, present, or future performance?
  - Are heights and weights increasing in the population?
-

# Google Ngrams

---

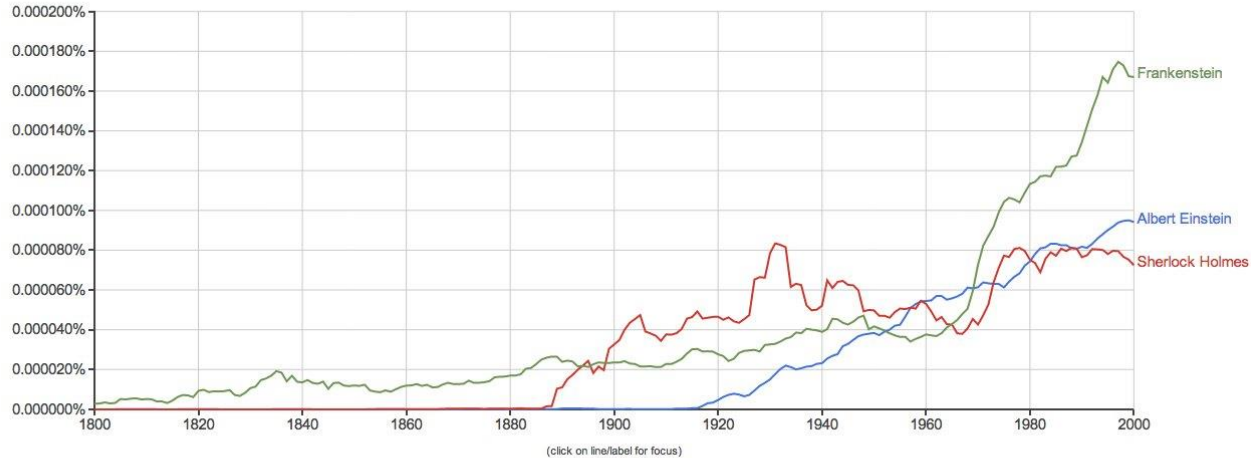
- Presents an annual time series of the frequency of every “popular” word/phrase with 1 to 5 words occurs in scanned books.
  - ‘Popular’ means appears >40 times in total.
  - Google has scanned about 15% of all books ever published, making this resource quite comprehensive.
-

# Google Ngram Viewer

Google books Ngram Viewer

Graph these comma-separated phrases:   case-insensitive

between  and  from the corpus  with smoothing of  [Search lots of books](#)



Run your own experiment! Raw data is available for download [here](#).

# Ngram Questions

---

- How has the amount of cursing changed over time?
  - What is the lifespan of fame and technologies? Is it increasing/decreasing?
  - How often do new words emerge? Do they stay in common usage?
  - What words are associated with other words, i.e. can you build a language model?
-



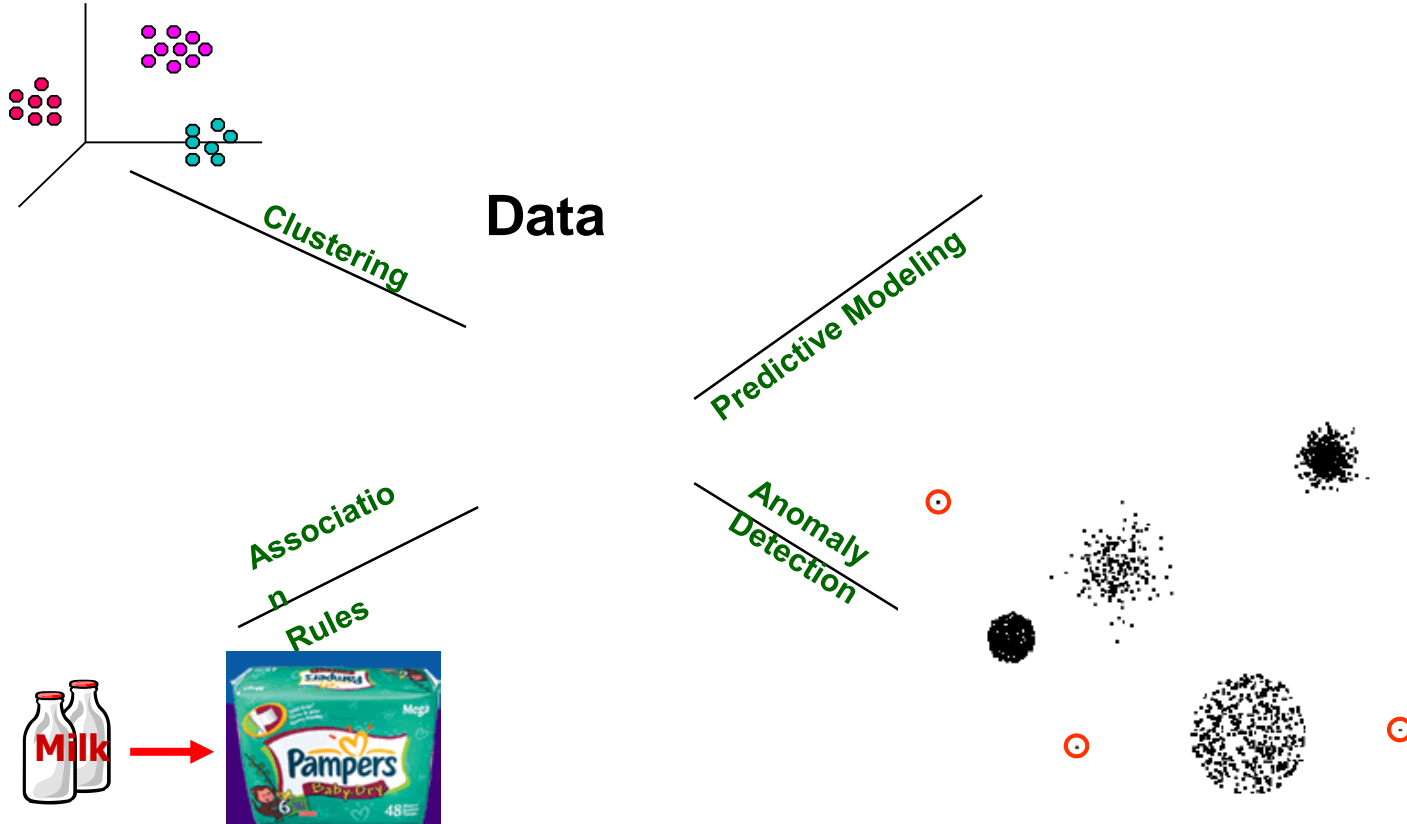
# Taxicab Questions

---

- How much do drivers make each night?
  - How far do they travel?
  - How much slower is traffic during rush hour?
  - Where are people traveling to/from at different times of the day?
  - Do faster drivers get tipped better?
  - Where should drivers go to pick up their next fare?
-



# Machine Learning Tasks ...



# Predictive Modeling: Classification

---

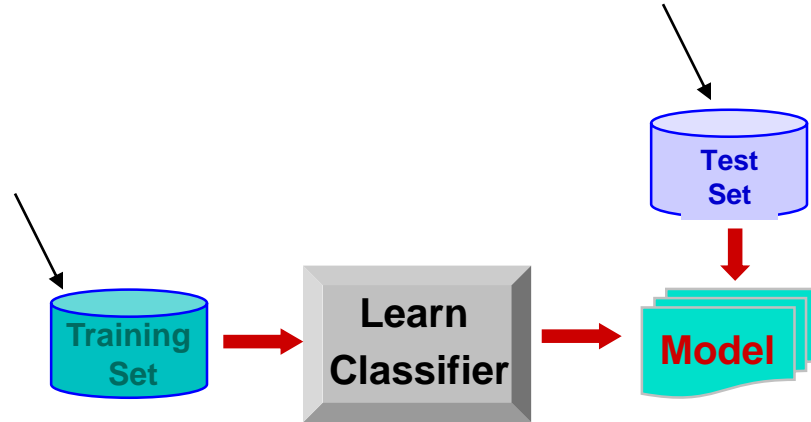
- Find a model for class attribute as a function of the values of other attributes

**Model for predicting credit worthiness**

**Class**

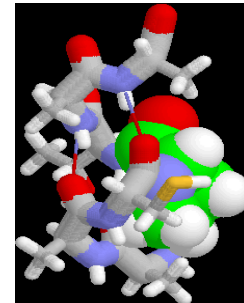
# Classification Example

categorical  
categorical  
quantitative  
class



# Examples of Classification Task

- Classifying credit card transactions as legitimate or fraudulent
- Classifying land covers (water bodies, urban areas, forests, etc.) using satellite data
- Categorizing news stories as finance, weather, entertainment, sports, etc
- Identifying intruders in the cyberspace
- Predicting tumor cells as benign or malignant
- Classifying secondary structures of protein as alpha-helix, beta-sheet, or random coil



# Classification: Application 1

---

- 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 2

---

- Churn prediction for telephone customers
  - **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.

# Classification: Application 3

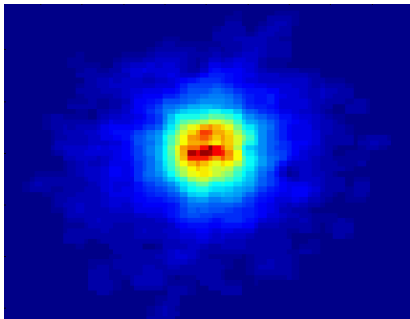
---

- 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!

# Classifying Galaxies

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

*Early*



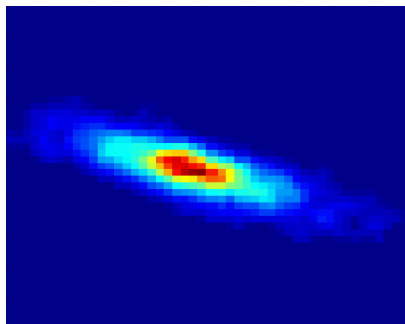
**Class:**

- Stages of Formation

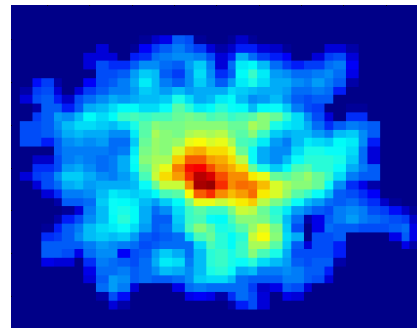
**Attributes:**

- 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



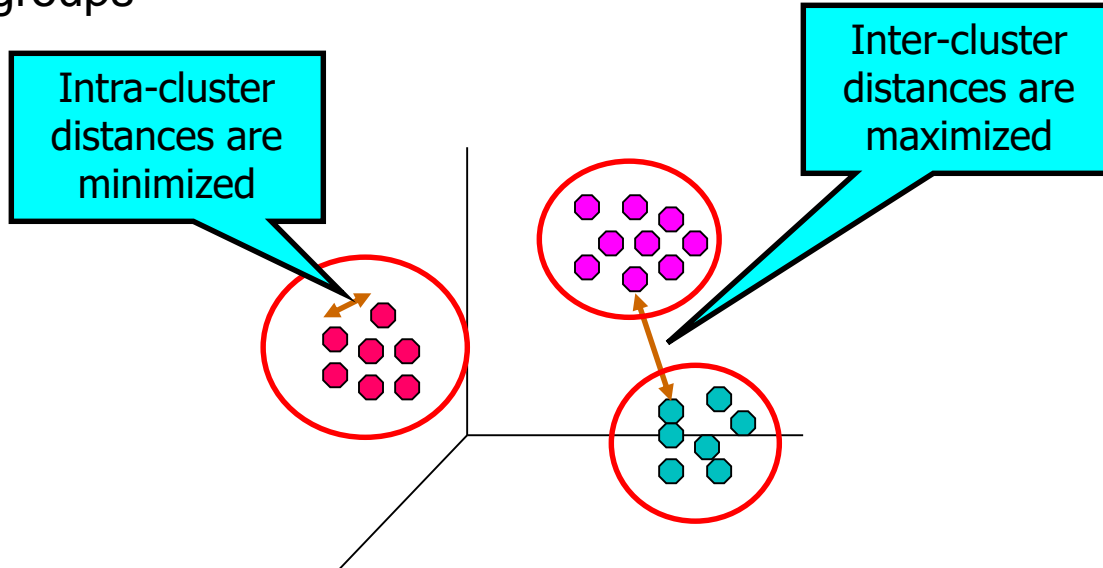
# Regression

---

- Predict a value of a given continuous valued variable based on the values of other variables, assuming a linear or nonlinear model of dependency.
- Extensively 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.

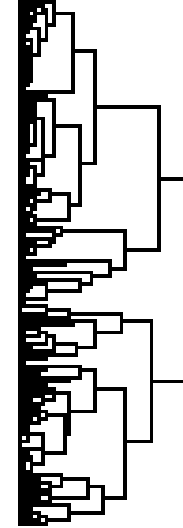
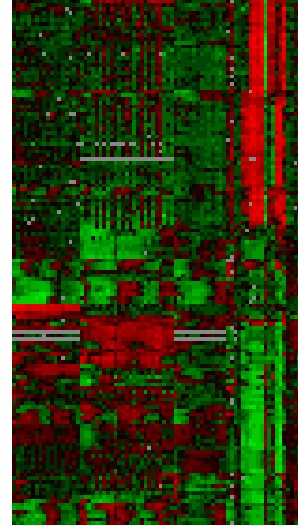
# Clustering

- Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups

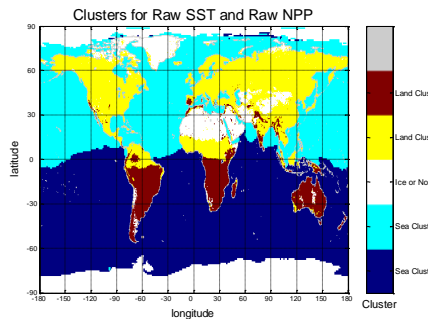


# Applications of Cluster Analysis

- **Understanding**
  - Custom profiling for targeted marketing
  - Group related documents for browsing
  - Group genes and proteins that have similar functionality
  - Group stocks with similar price fluctuations
- **Summarization**
  - Reduce the size of large data sets



Courtesy: Michael Eisen



**Use of K-means to partition Sea Surface Temperature (SST) and Net Primary Production (NPP) into clusters that reflect the Northern and Southern Hemispheres.**

# Clustering: Application 1

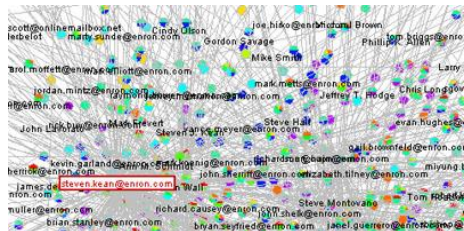
---

- 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 same cluster vs. those from different clusters.

# Clustering: Application 2

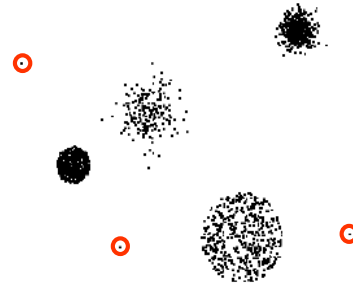
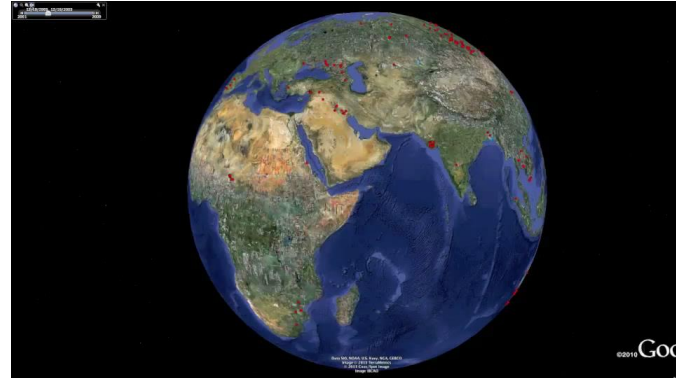
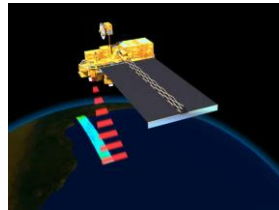
- Document Clustering:
  - **Goal:** To find groups of documents that are similar to each other based on the important terms appearing in them.
  - **Approach:** To identify frequently occurring terms in each document. Form a similarity measure based on the frequencies of different terms. Use it to cluster.

Enron email dataset



# Deviation/Anomaly/Change Detection

- Detect significant deviations from normal behavior
- Applications:
  - Credit Card Fraud Detection
  - Network Intrusion Detection
  - Identify anomalous behavior from sensor networks for monitoring and surveillance.
  - Detecting changes in the global forest cover.



# Motivating Challenges

---

- Scalability
- High Dimensionality
- Heterogeneous and Complex Data
- Data Ownership and Distribution
- Non-traditional Analysis



# DS Career path



---

# Introduction

- Graduates of data science program will mostly, and preferably, work as Data Scientists
- Data Scientists can work in any type of organization:
  - Private
  - Governmental
  - Non-for-Profit



# Industries

- Any organization can benefit from the data they have, so data scientists can work in any industry:
  - Financial Institutions (E.g., Banks)
  - Government agencies (E.g., Civil Status and Passports Department and Police Department)
  - Healthcare (E.g., Hospitals)
  - Online platforms (E.g., Uber)
  - Large Retailers (E.g., Carrefour and Amazon)
  - Agricultural Companies
  - And much more ...



# Data Scientist Responsibilities

- Data scientists usually need to build models of verified and validated data sets
- These models will be used by the employer to predict, recommend, or evaluate any future business decision



# Data Scientist Responsibilities

- For example, a data scientist, working for a hospital, can build a data model that predicts the best treatment for a specific patient
- The data scientist will use the data that was collected by the hospital about the patients and the treatments that worked and did not work for them in the past.



# Data Scientist Responsibilities

- Another example could be a data scientist, working for the police department, can build a data model that predicts the location and time of the next crime before it happens
- The data scientist will use the data that was collected by the police department about the previous crimes to build the proposed model



# Data Scientist Responsibilities

- Another example could be a data scientist, working for a large retailer, can build a data model that predicts the demand for certain products and services
- The data scientist will use the data that was collected by the retailer about the previous purchasing transactions
- The data scientist may use data that is provided by external entities



# Data Scientist Responsibilities

- Before building the model, data scientist usually need to clean and normalize the data
- Data could be collected from internal sources or/and external sources
- Data scientists need to communicate with data management guys to make sure that necessary data is being collected
  - Data compliance department should be involved to make sure that data collection is properly handled from a legal perspective





---

# More Opportunities

- In addition to working as data scientists, graduates of data science program can work as software development engineers
- In this field, they will mostly specialize in developing platforms that help data scientists in their jobs
- They also can develop dashboards that present business intelligence charts and reports to users





---

# **CIS Career path**

---

# Introduction

- Graduates of Computer Information Systems (CIS) program can pursue a job in of the following fields:
  - Business Analysis
  - Software Development
  - System Implementation



---

# Introduction

- CIS is an interdisciplinary program that encompasses technology and business courses
- This makes the graduates of this program knowledgeable about how business works and how technology can make businesses more efficient and more effective



# Introduction

- People who have knowledge about the technology only will have the following issues while working in the software development field:
  - Difficulty in developing a software that satisfies the business requirements
  - Difficulty in architecting the software systems according to the international standards
  - Difficulty in maintaining existing systems due to lack of knowledge about the business behind them



---

# Example

- CIS program exposes students to healthcare information systems
- When a CIS graduate joins a software development team that is responsible for developing an electronic health record (EHR), he/she will be already aware of the features and functionality of the proposed system



---

# You as a Business Analyst

- You will help customers define their requirements of any proposed software system
- Because you are already aware of how existing systems work, you can make notes and suggestions on how the proposed software system should look like
- Also, It is less likely you will misinterpret the requirements provided by customers



# You as a Software Developer

- You will write code to make a software system
- Because you are already aware of how business works, you will be able to choose the right architecture for the system
- The right architecture is one that supports any future improvements without making radical changes to the existing architecture



# You as a System Implementer

- You will help users use the software system the right way
- Because you are already aware of how business works, you will be able to provide a very helpful advice on how the software should be used and utilized

