

Data Mining

Yannis Velegrakis

Administrative Issues

- Schedule Wed & Thu @ 17:30
 - Lectures (in the class) – A212
 - Not online.
- Evaluation: Project
 - Topic: Will be announced later (not a free choice)
 - Development of a system and analysis of some dataset

Course Syllabus

- <https://www.esse3.unitn.it/Guide/PaginaRicercaInse.do>
- Enter in the prof field: “Velegrakis”
- Select “Data Mining”

- Also <https://www.disi.unitn.it/it/bachelor-student>

Administrative Issues

- Material Online (didattica Online)
 - <https://didatticaonline.unitn.it/dol/course/view.php?id=34838>

The screenshot shows a Moodle course page for "Data mining 2020/2021". The page includes a header with the university logo and navigation links for "Back", "UniTrento", "Didattica Online", "Contatti", "Moodle", and "Italiano (i)". A red banner at the top displays the course title and a Zoom meeting link: "Data mining 2020/2021 Join through Zoom: https://unitn.zoom.us/j/92957613136 with passcode: 1821". The main content area features a sidebar with links for "Partecipanti", "Badge", "Valutazioni", "Back", "miei corsi - docente", "Calendario", "File personali", and "My Media". The central area contains sections for "Announcements" (empty), "NEWS: INFORMATIVA LEZIONI ED ESAMI A DISTANZA" (with a link to "Consulto le nuove informative sullo svolgimento di lezioni a distanza ed esami a distanza per il periodo dell'emergenza epidemiologica Covid19"), "Descrizione del corso" (with a "Join Zoom Meeting" button and the Zoom link), and "Linee guida MOODLE" (with a link to "Clicca su LINEE GUIDA e scarica i tutori in fondo alla pagina"). The right side of the page shows a user profile for "Yannis Velegrakis".

Administrative Issues

■ Textbook

- Mining Massive Datasets, Ullman et. al.
- Available Online at <http://www.mmds.org/>

The screenshot shows the homepage of the Infolab website for the book "Mining of Massive Datasets". The header includes the Stanford University logo and the text "Infolab" and "Stanford University". The main title "Mining of Massive Datasets" is in orange, followed by the authors' names: Jure Leskovec, Anand Rajaraman, Jeff Ullman. A green box contains the text: "Big-data is transforming the world. Here you will learn data mining and machine learning techniques to process large datasets and extract valuable knowledge from them." On the left, there's a sidebar with links: Home, Book & Slides, Stanford Courses, and Supporting Materials. Below the sidebar are three versions of the book covers: the first edition (purple), the second edition (green), and the third edition (blue). The third edition cover features a character holding a large gold nugget. To the right of the book covers, sections include "The book" (describing the course connection and publisher), "The 3rd edition of the book" (describing new material and deep learning chapter), and a table of contents with links to PDF and PPT files for each chapter.

Chapter	Title	Book	Slides	Videos
	Preface and Table of Contents	PDF		
Chapter 1	Data Mining	PDF	PDF PPT	

Questions?

Course Evaluation

- Project
- You get a real world problem. You get no instructions on what method to use. You have to select from those you have learned in the course.
- It is not a straight forward application of a method.
- You need to prove that your solution is correct and good.
 - Basically you have to convince us to “buy” your solution

**What is Data Mining?
Knowledge discovery from data**

2019 every **MINUTE** of **the** **DAY**

PRESENTED BY DOMO

1,000,000
VIDEOS

BLR
ERS PUBLISH

2,340
POSTS

0,030
ARE DOWNLOADED

100,000
ARE SENT

GOOGLE

4,500,000
VIDEOS

TW

511,2

TWEETS

188,000,
EMAILS ARE

SKYPE
USERS
231,8

INSTAGRAM



\$600 to buy a disk drive that can store all of the world's music

5 billion mobile phones in use in 2010

30 billion pieces of content shared on Facebook every month

40% projected growth in global data generated per year vs.

5% growth in global IT spending

\$5 million vs. \$400

Price of the fastest supercomputer in 1975¹ and an iPhone 4 with equal performance

235 terabytes data collected by the US Library of Congress by April 2011

15 out of 17 sectors in the United States have more data stored per company than the US Library of Congress



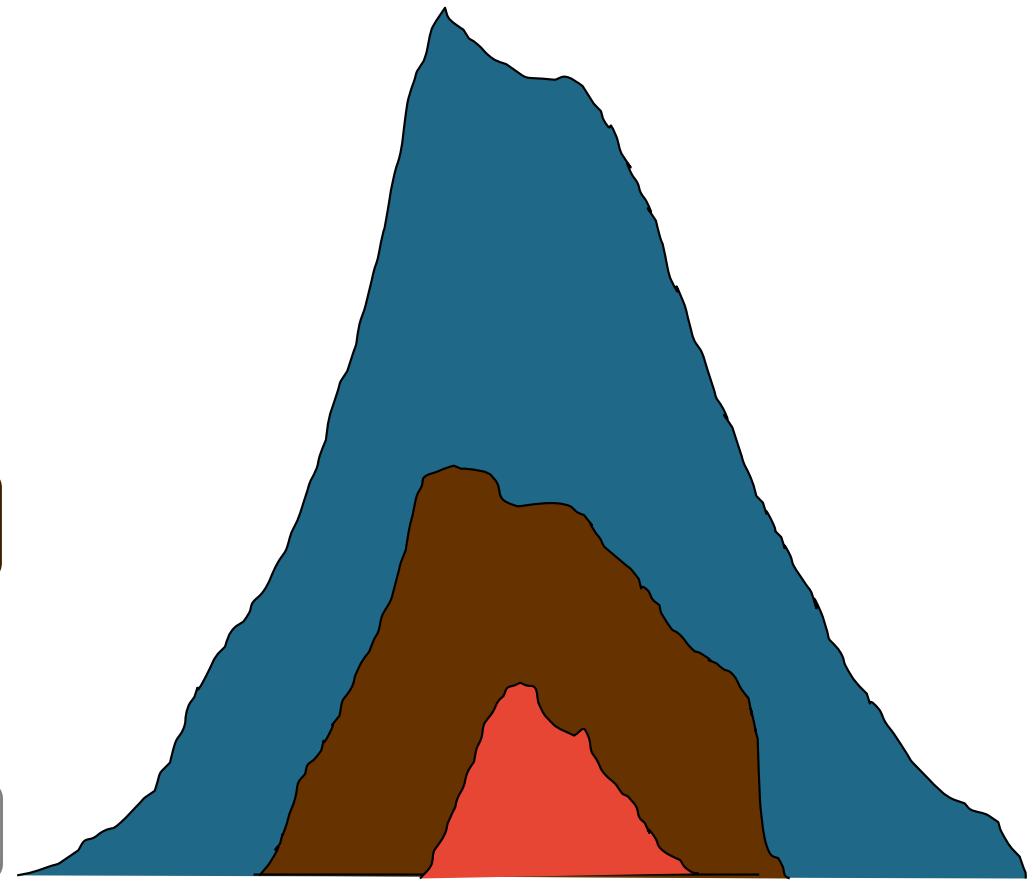
“DATA IS THE NEW GOLD”

Machine Generated Data

Internet Users

Employees / Data Engineers

More Data





Data contains value and knowledge

Data Mining

- But to extract the knowledge data needs to be
 - Stored
 - Managed
 - And ANALYZED ← this class

**Data Mining ≈ Big Data ≈
Predictive Analytics ≈ Data Science**

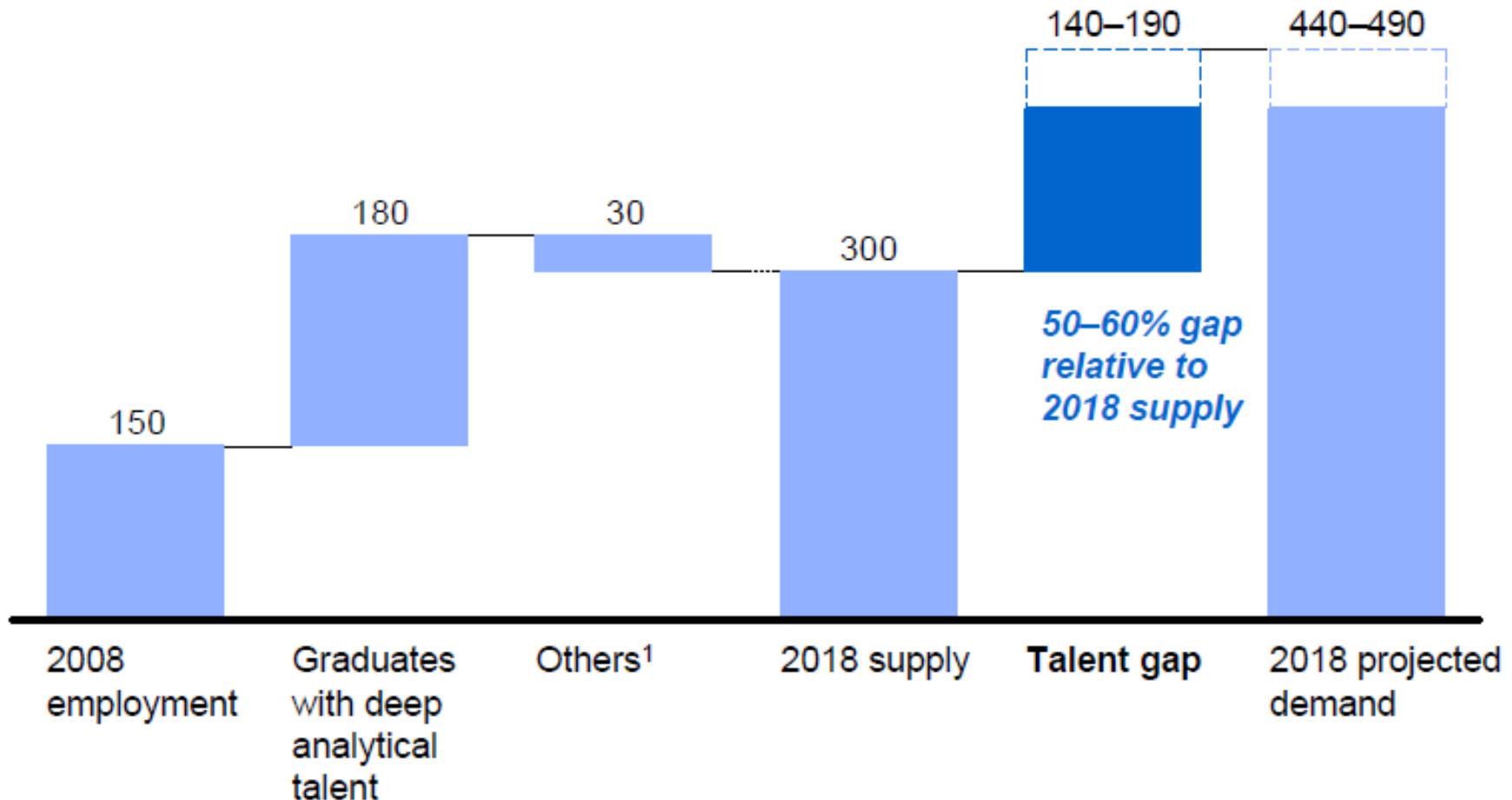
Or not?

Good news: Demand for Data Mining

Demand for deep analytical talent in the United States could be 50 to 60 percent greater than its projected supply by 2018

Supply and demand of deep analytical talent by 2018

Thousand people



¹ Other supply drivers include attrition (-), immigration (+), and reemploying previously unemployed deep analytical talent (+).

SOURCE: US Bureau of Labor Statistics; US Census; Dun & Bradstreet; company interviews; McKinsey Global Institute analysis

What is Data Mining?

- Given lots of data
- Discover patterns and models that are:
 - Valid: hold on new data with some certainty
 - Useful: should be possible to act on the item
 - Unexpected: non-obvious to the system
 - Understandable: humans should be able to interpret the pattern

Machine Learning

- **Good when we do not know what we are looking for.**
- **Not good if you know characteristics of the data you are looking**
- **Supervised vs unsupervised.**

Data Mining

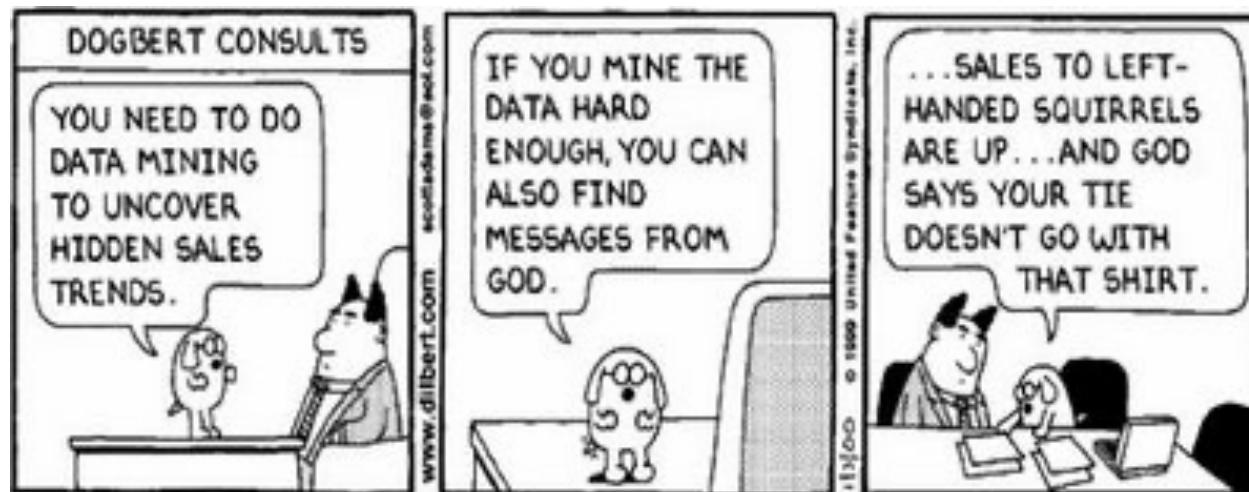
- Algorithmic
 - 1. Summarization
 - e.g. Page Rank
 - 2. Feature Extraction
 - The Cholera case
 - Frequent Itemsers
 - Similar Items

Data Mining Tasks

- **Descriptive methods**
 - Find human-interpretable patterns that describe the data
 - **Example:** Clustering
- **Predictive methods**
 - Use some variables to predict unknown or future values of other variables
 - **Example:** Recommender systems

Meaningfulness of Analytic Answers

- A risk with “Data mining” is that an analyst can “discover” patterns that are meaningless
- Statisticians call it **Bonferroni’s principle:**
 - Roughly, if you look in more places for interesting patterns than your amount of data will support, you are bound to find crap



Bonferroni by example

- We need to find evil-doers
 - they periodically gather in a hotel to plan an attack
- $1,000,000,000 = 10^9$ people
- Everyone goes to a hotel once every 100 days
- A hotel holds 100 people

- In 100 days almost everyone has been to a hotel. So, 10^9 people in 100 days,
- hence $10^9 / 100 = 10,000,000 = 10^7$ people in a hotel per day
- $10^7 / (100 \text{ people per hotel}) = 100,000$ hotels we have
- We examine the hotel records for 1000 days
- We try to find two people that in 2 different days they met in different hotels.
- Suppose there are no evil-doers
- Probability 2 people on the same day to visit a hotel on the same day is $1/100 * 1/100 = 10^{-4}$
- Chance that it is the same hotel is $V * 10^{-4} = 10^{-4}$ We
 - Probability they both select the 1st hotel = $\frac{1}{10^5} * \frac{1}{10^5} = \frac{1}{10^{10}}$
 - Probability they both select the 2nd hotel = $\frac{1}{10^5} * \frac{1}{10^5} = \frac{1}{10^{10}}$
 - Probability they both select the 3rd hotel = $\frac{1}{10^5} * \frac{1}{10^5} = \frac{1}{10^{10}}$
 - ...
 - Probability they both select the 10^5 th hotel = $\frac{1}{10^5} * \frac{1}{10^5} = \frac{1}{10^{10}}$
 - -----
 - Probability they choose the same hotel = prob they meet in the 1st + prob they meet in the 2nd etc $= \frac{1}{10^{10}} + \frac{1}{10^{10}} + \dots + \frac{1}{10^{10}} = 10^5 * \frac{1}{10^{10}} = \frac{1}{10^5}$
- Probability they meet = prob decide to go on same day and choose the same hotel = $10^{-4} * \frac{1}{10^5} = 10^{-9}$

Bonferroni by example

- Probability they meet = prob decide to go on same day and choose the same hotel = $10^{-4} * \frac{1}{10^5} = 10^{-9}$
- Probability they meet at the same hotel in two different days = $10^{-9} * 10^{-9} = 10^{-18}$
- Pairs of two people meeting on the same hotel in two different days:
- Number of pairs of people is $\binom{10^9}{2} \cong \frac{10^{9^2}}{2} = 5*10^{17}$
- Number of pairs of days is $\binom{1000}{2} \cong \frac{1000^2}{2} = 5*10^5$
- The number of cases of 2 random people visiting the same hotel in two different days is:
 - $5*10^{17} * 5*10^5 * 10^{-18} = 250,000$
- Assuming that there are 10 evil doers... then you need to do 250,000 checks to catch them. ... too many combinations to check – we need to have some additional evidence to find “suspicious” pairs of people in some more efficient way

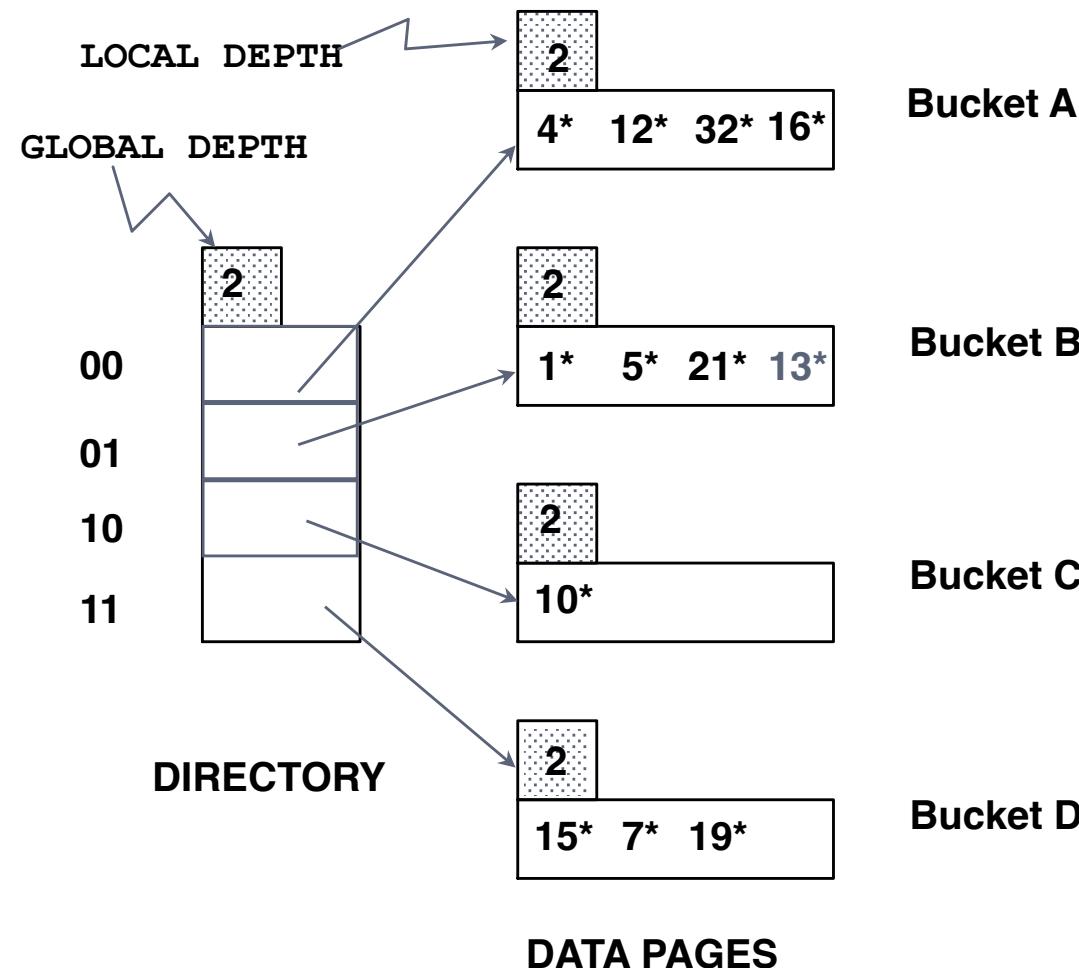
TF-IDF

- Collection of N documents
 - Need to find documents related to COVID
 - (or decide if a document is related to COVID)
- Frequency of term word i in document j:

$$TF_{ij} = \frac{f_{ij}}{\max_k f_{kj}}$$

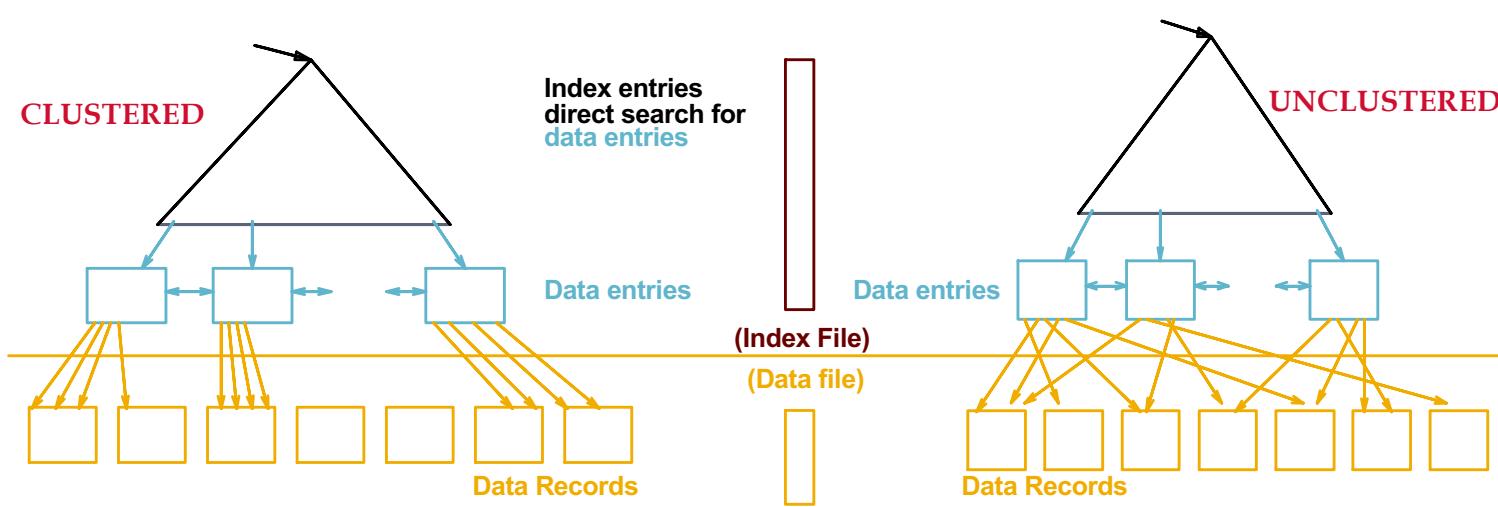
- Inverse Term Frequency
- $IDF_i = \log_2 \left(\frac{N}{n_i} \right)$
- $TF_{ij} * IDF_i$

Hash Functions

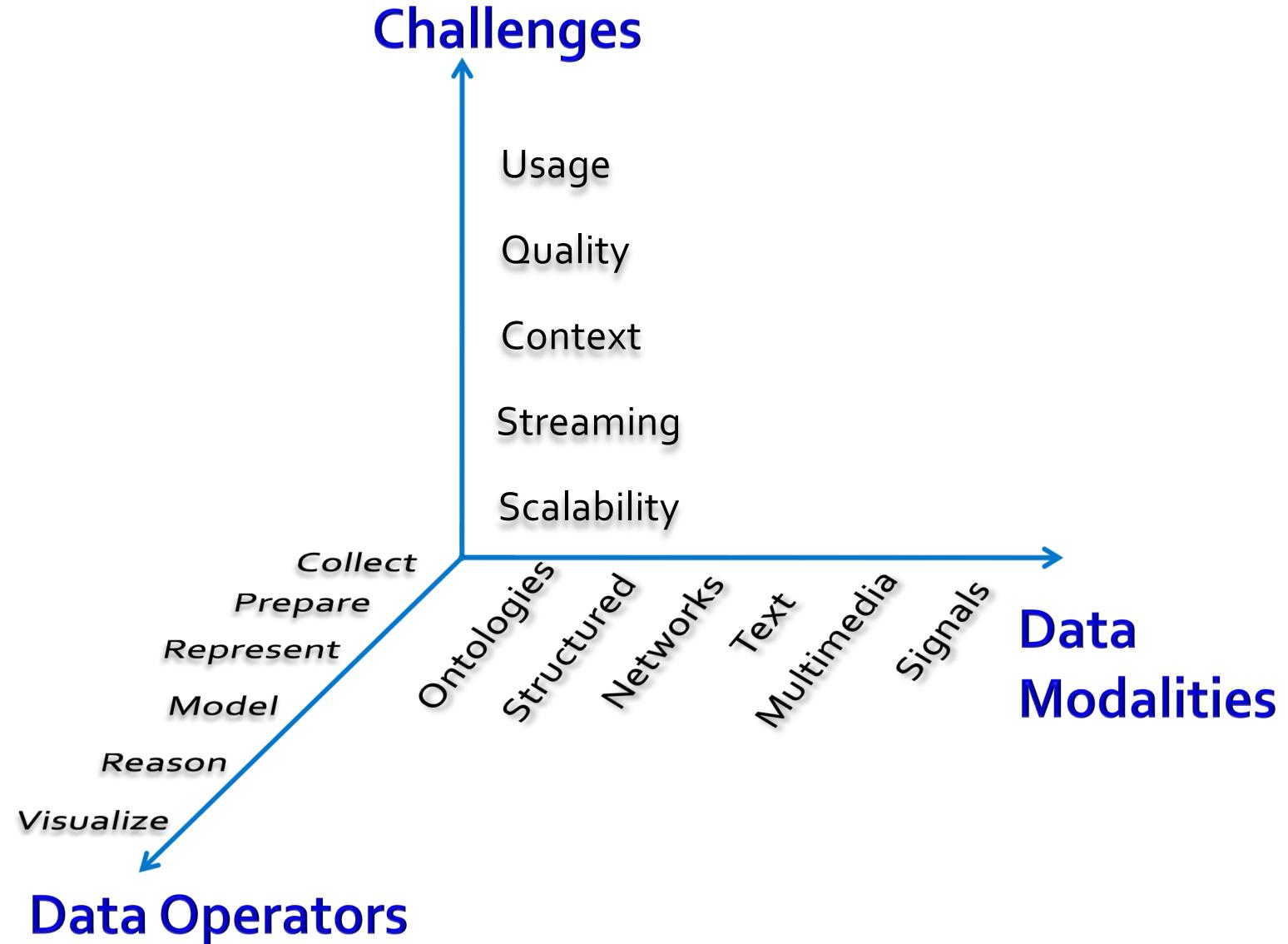


Hash function = last 2 bits in its binary representation

Clustered vs. Unclustered Index

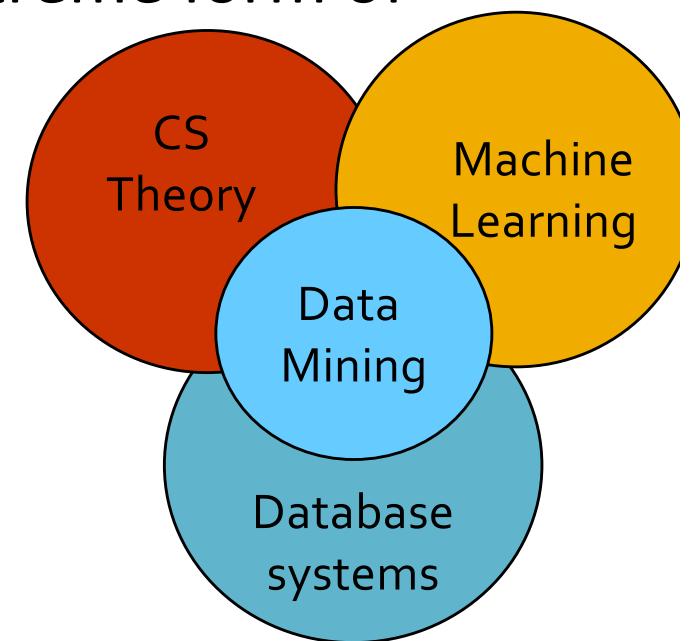


What matters when dealing with data?



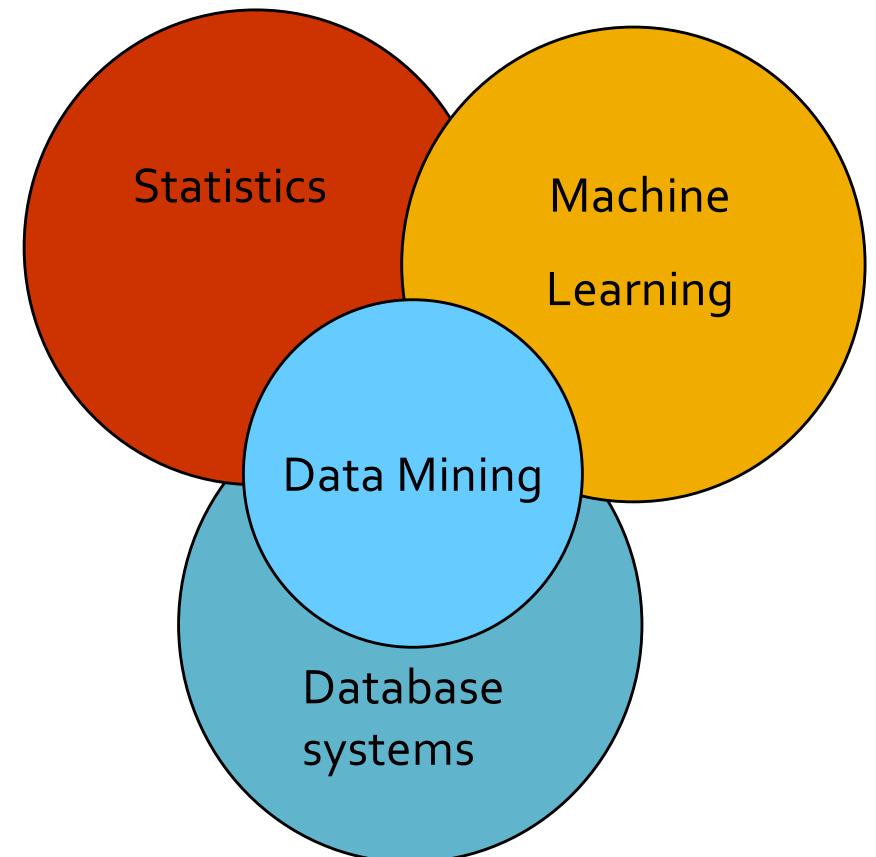
Data Mining: Cultures

- Data mining overlaps with:
 - **Databases:** Large-scale data, simple queries
 - **Machine learning:** Small data, Complex models
 - **CS Theory:** (Randomized) Algorithms
- Different cultures:
 - To a DB person, data mining is an extreme form of **analytic processing** – queries that examine large amounts of data
 - Result is the query answer
 - To a ML person, data-mining is the **inference of models**
 - Result is the parameters of the model
- In this class we will do both!



This Class

- This class overlaps with machine learning, statistics, artificial intelligence, databases but more stress on
 - Scalability (big data)
 - Algorithms
 - Computing architectures
 - Automation for handling large data



What will we learn?

- **We will learn to mine different types of data:**
 - Data is high dimensional
 - Data is a graph
 - Data is infinite/never-ending
 - Data is labeled
- **We will learn to use different models of computation:**
 - MapReduce (high level. See Big Data Course for details)
 - Streams and online algorithms
 - Single machine in-memory

What will we learn?

- **We will learn to solve real-world problems:**
 - Recommender systems
 - Market Basket Analysis
 - Spam detection
 - Duplicate document detection
- **We will learn various “tools”:**
 - Linear algebra (SVD, Rec. Sys., Communities)
 - Optimization (stochastic gradient descent)
 - Dynamic programming (frequent itemsets)
 - Hashing (LSH, Bloom filters)

How It All Fits Together

High dim. data

Locality
sensitive
hashing

Clustering

Dimensional
ity
reduction

Graph data

PageRank,
SimRank

Community
Detection

Spam
Detection

Infinite data

Filtering
data
streams

Web
advertising

Queries on
streams

Machine learning

SVM

Decision
Trees

Perceptron,
kNN

Apps

Recommen
der systems

Association
Rules

Duplicate
document
detection