# Data Mining
## CS 5140 / CS 6140

Jeff M. Phillips

January 9, 2017

# Data Mining

**Instructor :** Jeff Phillips (email) | Office hours: TBA @ MEB 3442 (and directly after class in WEB L104)
**TAs:** WaiMing Tai (email) | Office hours: 11am-noon Friday, location TBA
    + Deb Paul (email) | Office hours: 3-4pm Thursday, location TBA
    + Yang Gao (email) | Office Hours: 9-11am Tuesdays, location TBA
    + Shweta Singhal (email) | Office Hours: TBA
    others TBA
**Spring 2017** | Mondays, Wednesdays 3:00 pm - 4:20 pm
**WEB L104**
**Catalog number: CS 5140 01 or CS 6140 01**

---

**Syllabus**

**Description:**
Data mining is the study of efficiently finding structures and patterns in large data sets. We will focus on several aspects of this: (1) converting from a messy and noisy raw data set to a structured and abstract one, (2) applying scalable and probabilistic algorithms to these well-structured abstract data sets, and (3) formally modeling and understanding the error and other consequences of parts (1) and (2), including choice of data representation and trade-offs between accuracy and scalability. These steps are essential for training as a data scientist.
Algorithms, probability, and linear algebra are required mathematical tools for understanding these approaches.
Topics will include: similarity search, clustering, regression/dimensionality reduction, graph analysis, PageRank, and small space summaries. We will also cover several recent developments, and the application of these topics to modern applications, often relating to large internet-based companies.
Upon completion, students should be able to read, understand, and implement ideas from many data mining research papers.

**Books:**
We will in general not follow any book. My own course notes (linked below) serve as the defacto book. However, the following two free online books may serve as useful references that have good overlap with the course.
**MMDS**(v1.3): Mining Massive Data Sets by Anand Rajaraman, Jure Leskovec, and Jeff Ullman. The digital version of the book is free, but you may wish to purchase a hard copy.
**FoDS**: Foundations of Data Science by Avrim Blum, John Hopcroft and Ravindran Kannan. This provide some proofs and formalisms not explicitly covered in lecture.

**Videos:** We plan to videotape all lectures, and make them available online. They will appear on this playlist on our YouTube Channel.
Videos will also **livestream here**.
Lectures will also be live-streamed and available through Luum. More information to come.

**Prerequisites:** A student who is comfortable with basic probability, basic linear algebra, basic big-O analysis, and basic programming and data structures should be qualified for the class. There is no specific language we will use. However, programming assignments will often (intentionally) not be as specific as in lower-level classes. This will partially simulate real-world settings where one is given a data set and asked to analyze it; in such settings even less direction is provided.
For undergrads, the prerequisits are CS 3500 and CS 3130 and MATH 2270 (or equivalent), and CS 4150 is a corequisite. I will grant exceptions for those with (a reasonable grade in) CS 4964 (Fall 2016).
In the past, this class has had undergraduates, masters, and PhD students, including many from outside of Computer Science. Most (but not all) have kept up fine, and still most have been challenged. If you are unsure if the class is right for you, contact the instructor.

---

**Schedule:** (subject to change - some linked material is from the previous iteration of the class)

| Date | Topic (+ Notes) | Video | Link | Assignment (latex) | Project |
|------|-----------------|-------|------|--------------------|---------|

# Data Mining

**Instructor :** Jeff Phillips (email) | **Office hours: TBA @ MEB 3442 (and directly after class in WEB L1**
**TAs: WaiMing Tai** (email) | **Office hours: 11am-noon Friday, location TBA**
    **+ Deb Paul** (email) | **Office Hours: 3-4pm Thursday, location TBA**
    **+ Yang Gao** (email) | **Office Hours: 9-11am Tuesdays, location TBA**
    **+ Shweta Singhal** (email) | **Office Hours: TBA**
    **others TBA**
**Spring 2017 | Mondays, Wednesdays 3:00 pm - 4:20 pm**
**WEB L104**
**Catalog number: CS 5140 01 or CS 6140 01**

---

## Syllabus

**Description:**
Data mining is the study of efficiently finding structures and patterns in large data sets. We will focus on se
and abstract one, (2) applying scalable and probabilistic algorithms to these well-structured abstract data s
parts (1) and (2), including choice of data representation and trade-offs between accuracy and scalability. T
Algorithms, probability, and linear algebra are required mathematical tools for understanding these approac
Topics will include: similarity search, clustering, regression/dimensionality reduction, graph analysis, PageR
the application of these topics to modern applications, often relating to large internet-based companies.
Upon completion, students should be able to read, understand, and implement ideas from many data mini

**Books:**
We will in general not follow any book. My own course notes (linked below) serve as the defacto book. How
good overlap with the course.
**MMDS**(v1.3): Mining Massive Data Sets by Anand Rajaraman, Jure Leskovec, and Jeff Ullman. The digita
**FoDS**: Foundations of Data Science by Avrim Blum, John Hopcroft and Ravindran Kannan. This provide

**Videos:** We plan to videotape all lectures, and make them available online. They will appear on this playlis

## Syllabus

**Instructor:** Jeff M. Phillips. | 3442 MEB | http://www.cs.utah.edu/~jeffp

**Class Meetings:** Mondays and Wednesdays, 3:00pm – 4:20pm, WEB L104.

**Course Web Page:** http://www.cs.utah.edu/~jeffp/teaching/cs5140.html

Data mining is the study of efficiently finding structures and patterns in large data sets. We will focus on several aspects of this: (1) converting from a messy and noisy raw data set to a structured and abstract one, (2) applying scalable and probabilistic algorithms to these well-structured abstract data sets, and (3) formally modeling and understanding the error and other consequences of parts (1) and (2), including choice of data representation and trade-offs between accuracy and scalability. These steps are essential for training as a data scientist.

Algorithms, probability, and linear algebra are required mathematical tools for understanding these approaches.

Topics will include: similarity search, clustering, regression/dimensionality reduction, graph analysis, PageRank, and small space summaries. We will also cover several recent developments, and the application of these topics to modern applications, often relating to large internet-based companies.

Upon completion, students should be able to read, understand, and implement many data mining research papers.

### Getting Help

Take advantage of the instructor and TA office hours (posted on course web page). We will work hard to be accessible to students. Please send us email if you need to meet outside of office hours. Don't be shy if you don't understand something: come to office hours, send email, or speak up in class!

Students are encouraged to use a discussion group for additional questions outside of class and office hours. The class will rely on the Canvas discussion group. Feel free to post questions regarding any questions related to class: homeworks, schedule, material covered in class. Also feel free to answer questions, the instructors and TAs will also actively be answering questions. But, **do not post potential homework**

# Data Mining

**Instructor :** Jeff Phillips (email) | **Office hours: TBA @ MEB 3442 (and directly after class in WEB L1**
**TAs: WaiMing Tai** (email) | **Office hours: 11am-noon Friday, location TBA**
     **+ Deb Paul** (email) | **Office Hours: 3-4pm Thursday, location TBA**
     **+ Yang Gao** (email) | **Office Hours: 9-11am Tuesdays, location TBA**
     **+ Shweta Singhal** (email) | **Office Hours: TBA**
     **others TBA**
**Spring 2017 | Mondays, Wednesdays 3:00 pm - 4:20 pm**
**WEB L104**
**Catalog number: CS 5140 01 or CS 6140 01**

---

## Syllabus

**Description:**
Data mining is the study of efficiently finding structures and patterns in large data sets. We will focus on se
and abstract one, (2) applying scalable and probabilistic algorithms to these well-structured abstract data s
parts (1) and (2), including choice of data representation and trade-offs between accuracy and scalability. T
Algorithms, probability, and linear algebra are required mathematical tools for understanding these approa
Topics will include: similarity search, clustering, regression/dimensionality reduction, graph analysis, PageI
the application of these topics to modern applications, often relating to large internet-based companies.
Upon completion, students should be able to read, understand, and implement ideas from many data minir

**Books:**
We will in general not follow any book. My own course notes (linked below) serve as the defacto book. Hov
good overlap with the course.
**MMDS**(v1.3): Mining Massive Data Sets by Anand Rajaraman, Jure Leskovec, and Jeff Ullman. The digita
**FoDS**: Foundations of Data Science by Avrim Blum, John Hopcroft and Ravindran Kannan. This provide

**Videos:** We plan to videotape all lectures, and make them available online. They will appear on this playlis

# www.cs.utah.edu/~jeffp/teaching/cs5140.html

| Date | Topic (+ Notes) | Video | Link | Assignment (latex) | Project |
|------|-----------------|-------|------|-------------------|---------|
| Mon 1.09 | Class Overview | Vid | MMDS 1.1 | | |
| Wed 1.11 | Statistics Principles + Chernoff Bounds | Vid | MMDS 1.2 | | |
| Mon 1.16 | (MLK Day - No Class) | | | | |
| Wed 1.18 | Similarity : Jaccard + k-Grams | Vid | MMDS 3.1 + 3.2 \| FoDS 7.3 | | |
| Mon 1.23 | Similarity : Min Hashing | Vid | MMDS 3.3 | | |
| Wed 1.25 | Similarity : LSH | Vid | MMDS 3.4 | Statistical Principles | |
| Mon 1.30 | Similarity : Distances | Vid | MMDS 3.5 + 7.1 \| FoDS 8.1 | | Proposal |
| Wed 2.01 | Similarity : SIFT and ANN vs. LSH | Vid | MMDS 3.7 + 7.1.3 | | |
| Mon 2.06 | Clustering : Hierarchical | Vid | MMDS 7.2 \| FoDS 8.7 | | |
| Wed 2.08 | Clustering : K-Means | Vid | MMDS 7.3 \| FoDS 8.3 | | |
| Mon 2.13 | Clustering : Spectral (S) | Vid | MMDS 10.4 \| FoDS 8.4 \| Speilman \| Gleich | Document Hash | |
| Wed 2.15 | Streaming : Misra-Greis and Frugal | Vid | MMDS 4.1 \| FoDS 7.1.3 \| Min-Count Sketch \| Misra-Gries | | |
| Mon 2.20 | (Presidents Day - No Class) | | | | |
| Wed 2.22 | Streaming : Count-Min + Apriori Algorithm | Vid | MMDS 6+4.3 \| Careful Bloom Filter Analysis | | Data Collection Report |
| Mon 2.27 | Regression : Basics in 2-dimensions | Vid | ESL 3.2 and 3.4 | | |
| Wed 3.01 | Regression : SVD + PCA | Vid | Geometry of SVD - Chap 3 \| FoDS 4 | Clustering | |

# Lecture Notes

| Date | Topic | Vid | Reading | | |
|------|-------|-----|---------|--|--|
| Mon 3.06 | Regression : Matrix Sketching | Vid | MMDS 9.4 | FoDS 2.7 + 7.2.2 | arXiv | | |
| Wed 3.08 | **MIDTERM TEST** | | | | |
| Mon 3.13 | (Spring Break - No Class) | | | | |
| Wed 3.15 | (Spring Break - No Class) | | | | |
| Mon 3.20 | Regression : Random Projections | Vid | FoDS 2.9 | | Intermediate Report |
| Wed 3.22 | Regression : Compressed Sensing and OMP | Vid | FoDS 10.3 | Tropp + Gilbert | Frequent | |
| Mon 3.27 | Regression : L1 Regression and Lasso | Vid | Davenport | ESL 3.8 | bias-variance example | | |
| Wed 3.29 | Noise : Noise in Data | Vid | MMDS 9.1 | Tutorial | | |
| Mon 4.03 | Noise : Privacy | Vid | McSherry | Dwork | | |
| Wed 4.05 | Graph Analysis : Markov Chains (S) | Vid | MMDS 10.1 + 5.1 | FoDS 5 | Weckesser notes | | |
| Mon 4.10 | Graph Analysis : PageRank | Vid | MMDS 5.1 + 5.4 | Regression | |
| Wed 4.12 | Graph Analysis : MapReduce | Vid | MMDS 2 | | | |
| Mon 4.17 | Graph Analysis : Communities | Vid | MMDS 10.2 + 5.5 | FoDS 8.8 + 3.4 | | Final Report |
| Wed 4.19 | Graph Analysis : Graph Sparsification | Vid1,2 | MMDS 4.1 | | Poster Outline |
| Mon 4.24 | **ENDTERM TEST** | | | | |
| Mon 5.01 | | | | Graphs | |
| Tue 5.02 | Poster Day !!! (3:30-5:30pm) | | | | Poster Presentation |

## Project[*]

Final Report Due: Monday, April 17
Turn in report by 2:45pm (through Canvas).

### 1 Overview

Your project will consist of five elements.

- Project Proposal : Due January 30
- Data Collection Report : Due February 22
- Intermediate Report : Due March 20
- Final Report : Due April 17
- Poster Presentation : May 2 | (3:30pm - 5:30pm or 6:00pm)

As in any research in order to get people to pay attention, you will need to be able to present your work efficiently in written and oral form.

You may work in teams of 2 or 3, but the amount of work you perform will need to scale accordingly. Teams of size 1 might be allowed under unusual circumstances with special permission from the instructor. All students will need to have clearly defined roles as demonstrated in the final report and presentation. I highly recommend groups of size 3. Although the project work will scale with students, the administrative parts will remain constant, so having a large group will make it easier for you.

Note that some topics will not be covered before many elements of the project are due. I realize this is not ideal. However, typically, most work on a project is crammed in the last week or two of the semester, which is also not ideal. In the past this has lead to much stronger projects without considerably more work required.

### 1.1 Scale of Project

# Example Posters



## Station Evaluation and Time-Series Curve Matching for Meteorological Observation
### Yan Zheng

### Introduction

A meteorological observation at a given place can be inaccurate for a variety of reasons. Quality control can help spot which meteorological observations are inaccurate.

The project data is mainly from MesoWest group of Atmosphere Science Department, which are the results of UU2DVAR analysis(bias, Impact) for clustering and weather observations from 100 stations of six-year data for curve matching.

### Key Idea

Based on long-term statistical information with widely neighbor stations and the pattern of a specific day of a station, QC methods are explored to distinguish high impact stations using clustering algorithm and to find a weather pattern by time-series curve matching using nearest neighbor search based LSH algorithm. Euclidean distance is used to measure the distance of two curves.

### Clustering

K-Mean++ and Gaussian mixture modeling clustering algorithm have been applied and the cluster index is used as the score to evaluate the quality of a station.

**Result of k-mean++ clustering**



**Result of Gaussian Mixture Modeling**



### Curve Matching

LSH family: Pick a random projection of $R^d$ onto a 1-dimensional line and chop the line into segments of length w, shifted by a random value $b \in [0,w)$.

Choose L functions $g_j$, $j=1...L$, by setting $g_j=(h_{1,j}, h_{2,j}, ... h_{k,j})$, where $h_{1,j}, h_{2,j}, ... h_{k,j}$ are chosen at random from the LSH family, $H$. Then construct L hash tables.

**Tempreture difference data quering**



**Tempreture difference data quering**



### Conclusion

- Understanding the data, key to data mining.
- Finding the right algorithm, need to explore many options.
- Correctly use the data, do experiments and compare the results.

# Data Mining

What is Data Mining?

# Data Mining

What is Data Mining?

- ▶ Finding structure in data?
- ▶ Machine learning on large data?
- ▶ Unsupervised learning?
- ▶ Large scale computational statistics?

# Data Mining

What is Data Mining?

- ▶ Finding structure in data?
- ▶ Machine learning on large data?
- ▶ Unsupervised learning?
- ▶ Large scale computational statistics?

- ▶ How to think about data analytics.

# Data Mining

What is Data Mining?

- ▶ Finding structure in data?
- ▶ Machine learning on large data?
- ▶ Unsupervised learning?
- ▶ Large scale computational statistics?

- ▶ How to think about data analytics.

- ▶ *Principals* of converting from messy raw data to abstract representations.
- ▶ Algorithms of how to analyze data in abstract representations.
- ▶ Addressing challenges in scalability, error, and modeling.

# Methods for Data Analytics

Machine Learning (CS 5350/6350)

- ▶ Classification: Given labeled data $\ell(x) \in \{\text{TRUE or FALSE}\}$, build model so given new data, you can guess a label.
- ▶ More continuous optimization (DM more discrete)

# Methods for Data Analytics

Machine Learning (CS 5350/6350)

- ▶ Classification: Given labeled data $\ell(x) \in \{\text{TRUE or FALSE}\}$, build model so given new data, you can guess a label.
- ▶ More continuous optimization (DM more discrete)

Artificial Intelligence (CS 4300 / CS 6300)

- ▶ Interaction with World/Data: Observe, Learn, Act; repeat.

# Methods for Data Analytics

Machine Learning (CS 5350/6350)

- ▶ Classification: Given labeled data $\ell(x) \in \{\text{TRUE or FALSE}\}$, build model so given new data, you can guess a label.
- ▶ More continuous optimization (DM more discrete)

Artificial Intelligence (CS 4300 / CS 6300)

- ▶ Interaction with World/Data: Observe, Learn, Act; repeat.

More advanced Topics:

- ▶ Probabilistic Learning
- ▶ Structured Prediction
- ▶ Natural Language Processing
- ▶ Clustering

# Methods for Data Analytics

Machine Learning (CS 5350/6350)

- ▶ Classification: Given labeled data $\ell(x) \in \{\text{True or False}\}$, build model so given new data, you can guess a label.
- ▶ More continuous optimization (DM more discrete)

Artificial Intelligence (CS 4300 / CS 6300)

- ▶ Interaction with World/Data: Observe, Learn, Act; repeat.

More advanced Topics:

- ▶ Probabilistic Learning
- ▶ Structured Prediction
- ▶ Natural Language Processing
- ▶ Clustering

Data Mining has some ($< 10\%$) overlap with each of these.

# Modeling versus Efficiency

Two Intertwined (and often competing) Objectives:

- Model Data Correctly
- Process Data Efficiently

# Other Data Mining Courses

Every university teaches data mining differently!

# Other Data Mining Courses

Every university teaches data mining differently!

What flavor is offered in this class:

- Focus on techniques for *very* large scale data
- Broad coverage ... with recent developments
- Formally and generally presented (proof sketches)
- ... but useful in practice (e.g. internet companies)
- Probabilistic algorithms: connections to CS and Stat

# Other Data Mining Courses

Every university teaches data mining differently!

What flavor is offered in this class:
- Focus on techniques for *very* large scale data
- Broad coverage ... with recent developments
- Formally and generally presented (proof sketches)
- ... but useful in practice (e.g. internet companies)
- Probabilistic algorithms: connections to CS and Stat
- *no specific software tools / programming languages*

# Other Data Mining Courses

Every university teaches data mining differently!

What flavor is offered in this class:

- Focus on techniques for *very* large scale data
- Broad coverage ... with recent developments
- Formally and generally presented (proof sketches)
- ... but useful in practice (e.g. internet companies)
- Probabilistic algorithms: connections to CS and Stat
- *no specific software tools / programming languages*

Maths: Linear Algebra, Probability, High-dimensional geometry

# Outline

Statistical Principals:

- ► 1. **Hashing, Concentration of Measure**

Data and Distances:

- ► 2. **Similarity** (find duplicates and similar items)
- ► 3. **Clustering** (aggregate close items)

Structure in Data:

- ► 3. **Clustering** (aggregate close items)
- ► 4. **Regression** (linearity of (high-d) data)
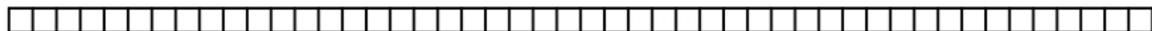- ► 5. **Noisy Data** (anomalies in data)

Controlling for Noise and Uncertainty:

- ► 5. **Noisy Data** (anomalies in data)
- ► 6. **Link Analysis** (prominent structure in large graphs)

# Statistical Principals

What happens as data is generated with replacement
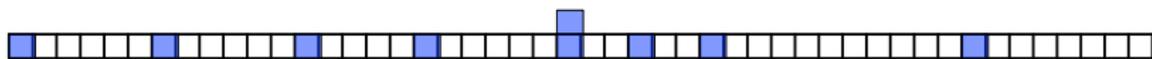{IP addresses, words in dictionary, edges in graph, hash table}
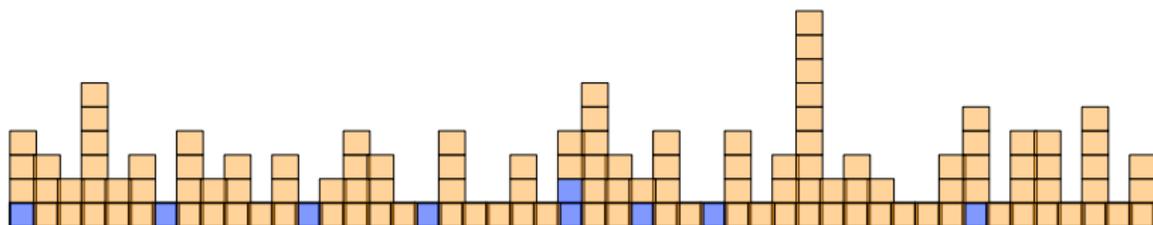
- ▶ When do items collide?
- ▶ When do you see all items?
- ▶ When is the distribution almost uniform?

# Statistical Principals

What happens as data is generated with replacement
{IP addresses, words in dictionary, edges in graph, hash table}

- ► When do items collide?
- ► When do you see all items?
- ► When is the distribution almost uniform?

# Statistical Principals

What happens as data is generated with replacement
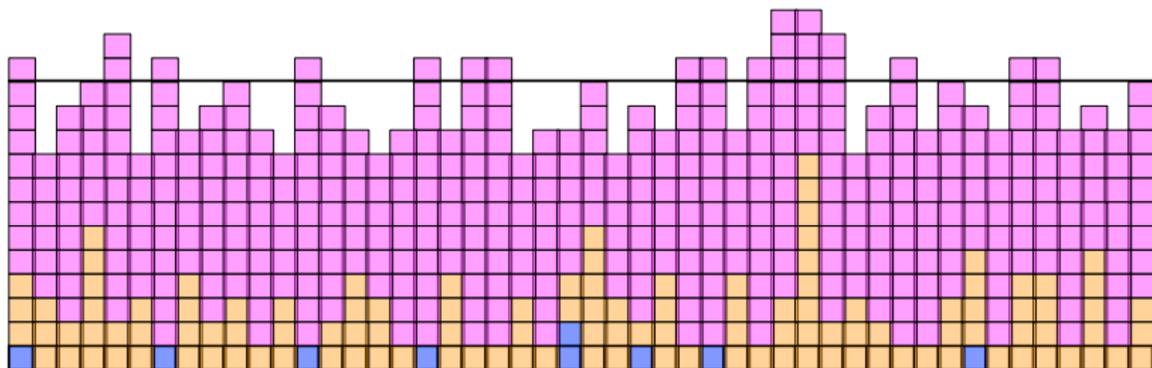{IP addresses, words in dictionary, edges in graph, hash table}

- ► When do items collide?
- ► When do you see all items?
- ► When is the distribution almost uniform?

# Statistical Principals

What happens as data is generated with replacement
{IP addresses, words in dictionary, edges in graph, hash table}

- ▶ When do items collide?
- ▶ When do you see all items?
- ▶ When is the distribution almost uniform?

# Raw Data to Abstract Representations

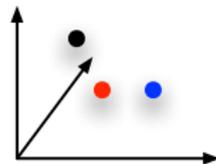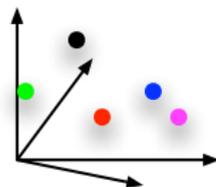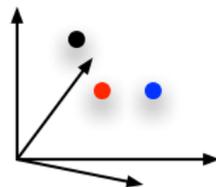How to measure similarity between data?
Key idea: data $\rightarrow$ point



a quick brown fox jumped ...

|     | age | income | height |
|-----|-----|--------|--------|
| joe | 25  | 90K    | 1.85   |
| bob | 32  | 45K    | 1.52   |
| sue | 28  | 38K    | 1.61   |

# Similarity

Given a large set of data $P$.
Given new point $q$, is $q$ in $P$?

Given a large set of data $P$.
Given new point $q$, what is the *closest* point in $P$ to $q$?
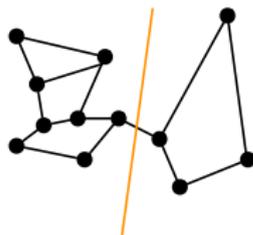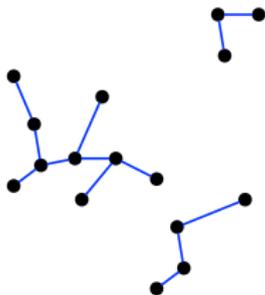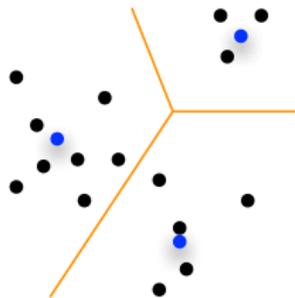
# Clustering

How to find groups of similar data.

- do we need a representative?
- can groups overlap?
- what is structure of data/distance?

# Clustering

How to find groups of similar data.

- ▶ do we need a representative?
- ▶ can groups overlap?
- ▶ what is structure of data/distance?

- ▶ **Hierarchical clustering** : When to combine groups?
- ▶ *k*-**means clustering** : $k$-median, $k$-center, $k$-means++
- ▶ **Graph clustering** : modularity, spectral

# Regression

Consider a data set $P \in \mathbb{R}^d$, where $d$ is BIG!

Want to find representation of $P$ in some $\mathbb{R}^k$

$\mu(P) \to Q \in \mathbb{R}^k$ so $\|p_i - p_j\| \approx \|q_i - q_j\|$
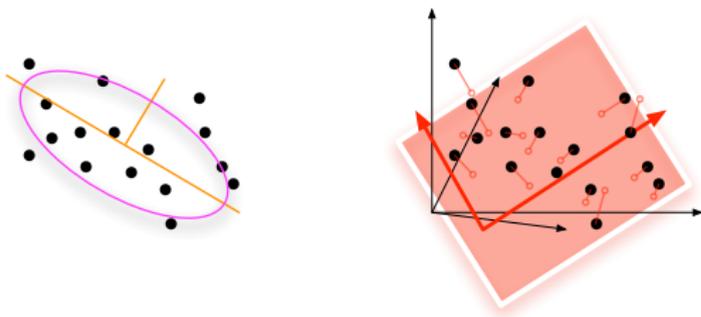
$Q \in \mathbb{R}^k$ should capture most data in $P$

# Regression

Consider a data set $P \in \mathbb{R}^d$, where $d$ is BIG!

Want to find representation of $P$ in some $\mathbb{R}^k$

$\mu(P) \to Q \in \mathbb{R}^k$ so $\|p_i - p_j\| \approx \|q_i - q_j\|$
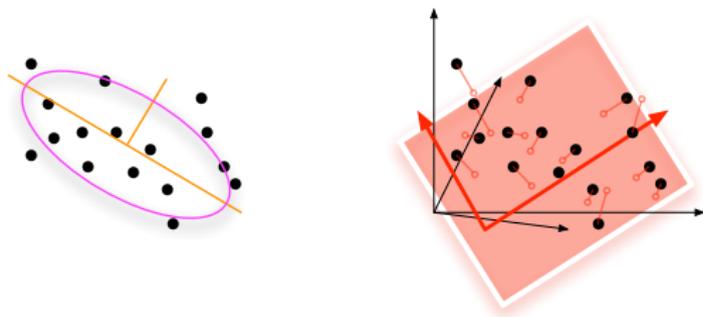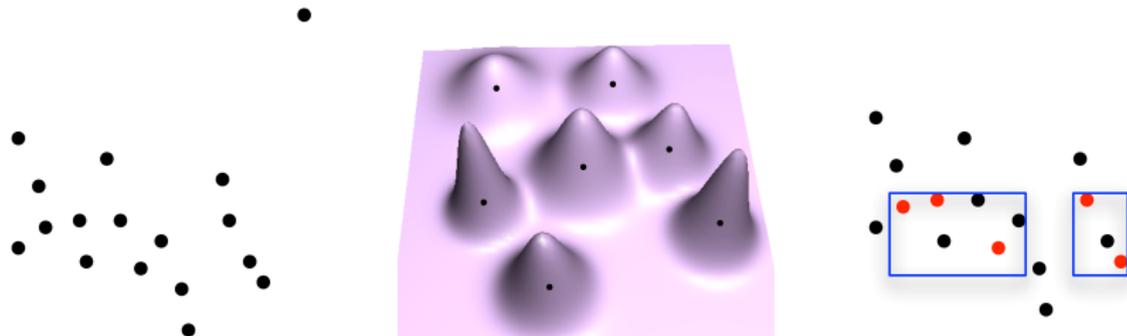
$Q \in \mathbb{R}^k$ should capture most data in $P$



- $L_2$ **Regression + PCA** : Common easy approach
- **Multidimensional Scaling** : Fits in $\mathbb{R}^k$ with $k$ small
- **Matrix Sketching**: Random Projections, Sampling, FD
- $L_1$ **Regression** : "Better", Orthogonal Matching Pursuit
- **Info Recovery** : Compressed Sensing

# Noisy Data

What to do when data is noisy?

- **Identify it** : Find and remove outliers
- **Model it** : It may be real, affect answer
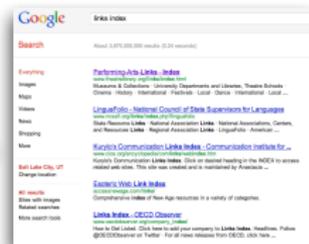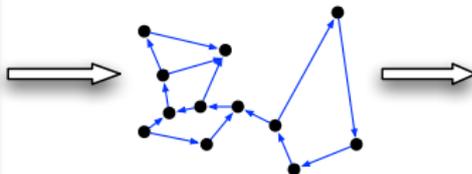- **Exploit it** : Differential privacy *(ethics in data)*

# Link Analysis

How does Google Search work?
Converts webpage links into directed graph.

- **Markov Chains** : Models movement in a graph
- **PageRank** : How to convert graph into important nodes
- **MapReduce** : How to scale up PageRank
- **Communities** : Other important nodes in graphs

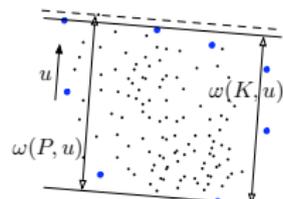# Summaries

Reducing *massive* data to small space.
Want to retain as much as possible (not specific structure)
error guarantees

- **OnePass Sampling** : Reservoir Sampling
- **MinCount Hash** : Sketching data, $\rightarrow$ abstract features
- **Density Approximation** : Quantiles
- **Matrix Sketching** : Preprocessing complex data
- **Spanners** : graph approximations

# Themes

What are course goals?

- ► Intuition for data analytics
- ► How to model data (convert to abstract data types)
- ► How to process data efficiently (balance models with algorithms)

# Themes

What are course goals?

- ▶ Intuition for data analytics
- ▶ How to model data (convert to abstract data types)
- ▶ How to process data efficiently (balance models with algorithms)

Work Plan:

- ▶ 2-3 weeks each topic.
  - ▶ Overview classic techniques
  - ▶ Focus on modeling / efficiency tradeoff
  - ▶ Special topics
  - ▶ Short homework for each (analysis + with data) **(45% grade)**
- ▶ 2 Tests **(10% grade)**
- ▶ Course Project **(45% grade)**.
  - ▶ Focus on specific data set
  - ▶ Deep exploration with technique
  - ▶ Ongoing refinement of presentation + approach

# On Homeworks

Managed through Canvas (should be up)

- ▶ No restriction on programming language.
- ▶ Some designed for matlab, others better in python or C++.
- ▶ Programming assignments with not too many specifications.
- ▶ Bonus Questions!

# On Canvas

Class management communication through Canvas

- All homework turn ins (typically as pdfs).
- Grades assigned
- Announcements
- Discussion (emails to instructor may not be responded)
  no posting potential solutions

# Videos

Class will be video-recorded and live-streamed.

- ▶ `https://www.youtube.com/channel/UCDUS80bdunpmvWVPyFRPqFQ`
- ▶ links off of webpage to live stream and playlist
- ▶ Experiment with Luum.io

# Videos

Class will be video-recorded and live-streamed.

- ▶ https://www.youtube.com/channel/
  UCDUS80bdunpmvWVPyFRPqFQ
- ▶ links off of webpage to live stream and playlist
- ▶ Experiment with Luum.io

Come to class if you can.

- ▶ Easier to ask questions, interact
  (mechanism through video, with delay)
- ▶ Talk to me before/after class!
- ▶ Attendance required for MIDTERM, FINAL, Poster Day
- ▶ Help your grade, and understanding.

# Data Group

Data Group Meeting
Thursdays @ 12:15-1:30 in MEB 3147 (LCR)

CS 7941 *Data Reading Group*
requires one presentation if taken for credit

`http://datagroup.cs.utah.edu`

# Utah Data Science Day

# DATA SCIENCE
# DAY 2017
## UTAH

http://datascience.utah.edu/dataday
Friday, January 13 : 11:30 - 6pm Union Ballroom

| Time | Event |
|---|---|
| 11:30 AM - 1:00 PM | Data Science Job Fair |
| 1:00 PM - 1:10 PM | Welcome: Data Science at Utah |
| 1:10 PM - 2:00 PM | Panel: Data Science in Industry |
| 2:00 PM - 3:30 PM | Posters and Demos |
| 3:30 PM - 4:50 PM | Data Science + X Talks |
| 5:00 PM - 6:00 PM | Keynote |
| 6:00 PM - 6:15 PM | Poster Awards !! |