

Data Mining

CS 5140 / CS 6140

Instructor: Jeff M. Phillips
Presenter: TA Xingyuan Pan

January 3, 2018

Data Mining

Instructor : [Jeff Phillips \(email\)](#) | Office hours: Thursday morning 10-11am @ MEB 3442 (and directly after class in WEB L104)

TAs: [Sunipa Dev \(email\)](#) | Office hours: Monday 11am-1pm, MEB 3115

+ [Maryam Baryouti \(email\)](#) | Office Hours: TBA, MEB 3115

+ [Yang Gao \(email\)](#) | Office Hours: TBA, (online, TBD)

+ [Xingyuan Pan \(email\)](#) | Office Hours: TBA, MEB 3115

+ [Trang Tran \(email\)](#) | Office Hours: TBA, MEB 3115

Spring 2018 | 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, programming, probability, and linear algebra are required 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:

The "book" for this course will be [my own course notes](#) 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.

MADA: *Math for Data Analysis* by Jeff M. Phillips. This is a gradual introduction to many of the topics this course builds on.

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](#).

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. A great primer on the [Mathematics of Data Analysis](#) can be found in the [linked book](#).

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 formal prerequisites are CS 3500 and CS 3130 and MATH 2270 (or equivalent), and CS 4150 is a corequisite. We recommend undergraduates take a new course [CS 4964 \(Foundations of Data Analysis\)](#) before this course, but it is not currently required, and many students have done well without having taken this course. I will grant exceptions to the pre-requisites for students with (a reasonable grade in) [Foundations of Data Analysis](#).

For graduate students, there are no enforced pre-requisites. Still it may be useful to review material in the [Math for Data book](#)

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.

For an example of what sort of mathematical material I **expect you to be to be familiar with**, see these notes on [probability and linear algebra](#).

Schedule: (subject to change)

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

Data Mining

Instructor : **Jeff Phillips** ([email](#)) | Office hours: Thursday morning 10-11am @ MEB 3442 (and direct)

TAs: **Sunipa Dev** ([email](#)) | Office hours: Monday 11am-1pm, MEB 3115

+ **Maryam Baryouti** ([email](#)) | Office Hours: TBA, MEB 3115

+ **Yang Gao** ([email](#)) | Office Hours: TBA, (online, TBD)

+ **Xingyuan Pan** ([email](#)) | Office Hours: TBA, MEB 3115

+ **Trang Tran** ([email](#)) | Office Hours: TBA, MEB 3115

Spring 2018 | 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 topics: (1) understanding the theory of data mining, and abstract one, (2) applying scalable and probabilistic algorithms to these well-structured abstract data sets. The course is divided into two parts (1) and (2), including choice of data representation and trade-offs between accuracy and scalability. Topics include: Algorithms, programming, probability, and linear algebra are required tools for understanding these approaches. Topics will include: similarity search, clustering, regression/dimensionality reduction, graph analysis, PageRank, 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 papers.

Books:

The "book" for this course will be [my own course notes](#) serve as the defacto book. However, the following are recommended reading for this course.

MMDS(v1.3): [Mining Massive Data Sets](#) by Anand Rajaraman, Jure Leskovec, and Jeff Ullman. The digital

FoDS: [Foundations of Data Science](#) by Avrim Blum, John Hopcroft and Ravindran Kannan. This provides

M4DA: [Math for Data Analysis](#) by Jeff M. Phillips. This is a gradual introduction to many of the topics this

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, programming, probability, and linear algebra are required 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 answers**. Such posts will be immediately removed, and not answered.

Data Mining

Instructor : **Jeff Phillips** ([email](#)) | Office hours: Thursday morning 10-11am @ MEB 3442 (and direct)

TAs: **Sunipa Dev** ([email](#)) | Office hours: Monday 11am-1pm, MEB 3115

+ **Maryam Baryouti** ([email](#)) | Office Hours: TBA, MEB 3115

+ **Yang Gao** ([email](#)) | Office Hours: TBA, (online, TBD)

+ **Xingyuan Pan** ([email](#)) | Office Hours: TBA, MEB 3115

+ **Trang Tran** ([email](#)) | Office Hours: TBA, MEB 3115

Spring 2018 | 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 topics: (1) understanding the theory of data mining, and abstract one, (2) applying scalable and probabilistic algorithms to these well-structured abstract data sets. The course is divided into two parts (1) and (2), including choice of data representation and trade-offs between accuracy and scalability. Topics include: Algorithms, programming, probability, and linear algebra are required tools for understanding these approaches. Topics will include: similarity search, clustering, regression/dimensionality reduction, graph analysis, PageRank, 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 papers.

Books:

The "book" for this course will be [my own course notes](#) serve as the defacto book. However, the following are recommended reading for this course.

MMDS(v1.3): Mining Massive Data Sets by Anand Rajaraman, Jure Leskovec, and Jeff Ullman. The digital

FoDS: Foundations of Data Science by Avrim Blum, John Hopcroft and Ravindran Kannan. This provides

M4DA: Math for Data Analysis by Jeff M. Phillips. This is a gradual introduction to many of the topics this

Date	Topic (+ Notes)	Video	Link	Assignment (latex)	Project
Mon 1.08	Class Overview	Vid	MMDS 1.1		
Wed 1.10	Statistics Principles	Vid	M4DA 3 MMDS 1.2 FoDS 12.4		
Mon 1.15	MLK DAY				
Wed 1.17	Similarity : Jaccard + k-Grams (S)	Vid	MMDS 3.1 + 3.2 FoDS 7.3		
Mon 1.22	Similarity : Min Hashing	Vid	MMDS 3.3		
Wed 1.24	Similarity : LSH	Vid	MMDS 3.4	Statistical Principles	
Mon 1.29	Similarity : Distances	Vid	MMDS 3.5 + 7.1 FoDS 8.1		
Wed 1.31	Similarity : SIFT and ANN vs. LSH	Vid	MMDS 3.7 + 7.1.3		Proposal
Mon 2.05	Clustering : Hierarchical	Vid	MMDS 7.2 FoDS 8.7		
Wed 2.07	Clustering : K-Means	Vid	M4DA 8 MMDS 7.3 FoDS 8.3		
Mon 2.12	Clustering : Spectral (S)	Vid	MMDS 10.4 FoDS 8.4	Document Hash	
Wed 2.14	Streaming : Misra-Greis and Frugal	Vid	MMDS 4.1 FoDS 7.1.3		
Mon 2.19	PRESIDENTS DAY				
Wed 2.21	Streaming : Count-Min + Apriori Algorithm	Vid	MMDS 6+4.3 BF Analysis		Data Collection Report
Mon 2.26	Regression : Basics in 2-dimensions	Vid	M4DA 5 ESL 3.2 and 3.4		
Wed 2.28	Regression : SVD + PCA	Vid	M4DA 4 and 7 FoDS 4	Clustering	
Mon 3.05	Regression : Matrix Sketching	Vid	MMDS 9.4 FoDS 2.7 + 7.2.2 arXiv		
Wed 3.07	MIDTERM TEST				
Mon 3.12	Regression : Random Projections	Vid	FoDS 2.9		
Wed 3.14	Regression : Compressed Sensing and OMP	Vid	FoDS 10.2 Tropp + Gilbert	Frequent	
Mon 3.19	SPRING BREAK				
Wed 3.21	SPRING BREAK				
Mon 3.26	Regression : L1 Regression and Lasso	Vid	Davenport ESL 3.8		Intermediate Report
Wed 3.28	Noise : Noise in Data	Vid	MMDS 9.1 Tutorial		
Mon 4.02	Lecture on Ethics/Fairness -- By Suresh	Vid	10 Simple Rules		
Wed 4.04	Noise : Privacy	Vid	McSherry Dwork	Regression	
Mon 4.09	Graph Analysis : Markov Chains (S)	Vid	MMDS 10.1 + 5.1 FoDS 5 Weckesser		
Wed 4.11	Graph Analysis : PageRank	Vid	MMDS 5.1 + 5.4		
Mon 4.16	Graph Analysis : MapReduce	Vid	MMDS 2		Final Report
Wed 4.18	Graph Analysis : Communities	Vid	MMDS 10.2 + 5.5 FoDS 8.8 + 3.4		Poster Outline
Mon 4.23	ENDTERM TEST				
Mon 4.30				Graphs	
Wed 5.02	Poster Day !!! (3:30-5:30pm)				Poster Presentation

4.2 Min Hashing

Last time we saw how to compare document sets via sets. This discussion how to compare sets, specifically using the Jaccard similarity. Specifically, for two sets $A = \{0, 1, 2, 3, 4\}$ and $B = \{0, 2, 3, 5, 7, 9\}$. The Jaccard similarity is defined

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{3}{10 + 1, 2, 3, 5, 6, 7, 9} = \frac{3}{13} \approx 0.231$$

Although this gives us a single numeric score to compare similarity (or distance) it is not easy to compare, and will be especially cumbersome if the sets are quite large.

This leads us to a technique called *min hashing* that uses a randomized algorithm to quickly estimate the Jaccard similarity. Furthermore, we can show how accurate it is through the Chernoff bounding bound.

To achieve these results we consider a new abstract data type, a matrix. This format is incredibly useful conceptually, but often extremely wasteful if full written out.

4.1 Matrix Representation

Here we see how to convert a series of sets (e.g. a set of sets) to be represented as a single matrix. Consider sets:

- $S_1 = \{1, 2, 5\}$
- $S_2 = \{3\}$
- $S_3 = \{2, 3, 4, 6\}$
- $S_4 = \{1, 4, 6\}$

For instance $J(S_1, S_3) = |\{2\}| / |\{1, 2, 3, 4, 5, 6\}| = 1/6$.

We can represent these four sets as a single matrix:

Element	S_1	S_2	S_3	S_4
1	1	0	0	1
2	1	0	1	0
3	0	1	1	0
4	0	0	1	1
5	1	0	0	0
6	0	0	1	1

represents matrix $M = \begin{pmatrix} 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix}$

The element in the i th row and the j th column denotes if element i is in set S_j . It is 1 if the element is in the set, and 0 otherwise. This captures each of the same data set as the set representation, but may take much more space. If the matrix is sparse, meaning that most entries (e.g. $> 90\%$ or maybe $> 99\%$...) are zero conceptually, as the matrix becomes \times of the non-zero entries grows as roughly \times^2 , but the space grows as \times (i.e. it wastes a lot of space). But still it is very useful to think about. There are also sparse matrix representations built into many languages such as Matlab which do not store all of the 0s, they just store the locations of the non-zeroes.

4.2 Min Hashing

The next approach, called *min hashing*, initially seems even simpler than the clustering approach. It will need to evolve through several steps to become a useful trick.

Step 1: Randomly permute the items by permuting the rows of the matrix.

Element	S_1	S_2	S_3	S_4
5	1	0	0	0
6	0	0	1	1
1	1	0	0	1
4	0	0	1	1
3	0	1	1	0
2	1	0	1	0

Step 2: Record the first 1 in each column, using a min function m_i . That is, given a permutation, applied to set S , the function $m(S)$ records the element from S which appears earliest in the permutation.

- $m(S_1) = 2$
- $m(S_2) = 3$
- $m(S_3) = 2$
- $m(S_4) = 6$

Step 3: Estimate the Jaccard similarity $J(S_1, S_3)$ as

$$J(S_1, S_3) = \begin{cases} 1 & \text{if } m(S_1) = m(S_3) \\ 0 & \text{otherwise.} \end{cases}$$

Lemma 4.2.1. $\Pr\{m(S_1) = m(S_3)\} = J(S_1, S_3)$.

Proof: There are three types of rows.

- (T1) There are x rows with 1 in both columns.
- (T2) There are y rows with 1 in one column and 0 in the other.
- (T3) There are z rows with 0 in both columns.

The total number of rows is $x + y + z$. The Jaccard similarity is precisely $J(S_1, S_3) = x / (x + y)$. (Note that usually $x \gg y$ (mostly empty) and we can ignore them.)
Let row i be the i th $m(S_1) = m(S_3)$. It is either type (T1) or (T2), and it is (T1) with probability exactly $x / (x + y)$, since the permutation is random. This is the only case that $m(S_1) = m(S_3)$, otherwise S_1 or S_3 has 1, but not both. \square

Thus this approach only gives 0 or 1, but has the right expectation. To get a better estimate, we need to repeat this several (k) times. Consider k random permutations $\{m_1, m_2, \dots, m_k\}$ and also k random variables $\{X_1, X_2, \dots, X_k\}$ where

$$X_i = \begin{cases} 1 & \text{if } m_i(S_1) = m_i(S_3) \\ 0 & \text{otherwise.} \end{cases}$$

Now we can estimate $J(S_1, S_3) = J(S_1, S_3) = \frac{1}{k} \sum_{i=1}^k X_i$, the average of the k simple random variables.

So how large should we set k so that this gives us an accurate measure? Since it is a randomized algorithm, we will have an error tolerance $\epsilon \in (0, 1)$ (e.g. we want $J(S_1, S_3) - J(S_1, S_3) \leq \epsilon$) and a probability of failure δ (e.g. the probability we have more than ϵ error). We will use our Theorem 2.4.2 where $M = \sum_{i=1}^k X_i$ and hence $\mathbb{E}[M] = k \cdot J(S_1, S_3)$. We have $0 \leq X_i \leq 1$ so each $\Delta_i = 1$. Now we can write for some value:

$$\Pr\{J(S_1, S_3) - J(S_1, S_3) \geq 2\epsilon\} = \Pr\{J(S_1, S_3) - J(S_1, S_3) \geq \epsilon\} \\ = \Pr\{M - \mathbb{E}[M] \geq \epsilon\} \leq 2 \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^k \Delta_i^2}\right) = 2 \exp(-2\epsilon^2/k)$$

Setting $\epsilon = \delta$ and $k = (1/(2\epsilon^2)) \ln(2/\delta)$ we obtain

$$\Pr\{J(S_1, S_3) - J(S_1, S_3) \geq \epsilon\} \leq 2 \exp(-2\epsilon^2/(1/\delta)) = 2 \exp(-2\delta \epsilon^2) = \delta$$

Or in other words, if we set $k = (1/(2\epsilon^2)) \ln(2/\delta)$, then the probability that our estimate $J(S_1, S_3)$ is within ϵ of $J(S_1, S_3)$ is at least $1 - \delta$.

Say for instance we want error at most $\epsilon = 0.01$ and can tolerate a failure δ of the time (if $\delta = 0.01$), then we need $k = (1/(2 \cdot 0.01^2)) \ln(2/0.01) = 200 \ln(200) \approx 1000$. Note that the rounding error of converting a structure into a set may be more than $\epsilon = 0.05$, so this should be an acceptable loss in accuracy.

Tip 6. It is sometimes more efficient to use the top- k (the same small number $k > 1$) hash values for each hash function, than just the top one. For instance, see Cohen and Euphrates (Stamford, Data using Bloom & Sketches, PODC 2007). This approach requires a bit more intricate analysis, as well as a bit more careful implementation.

4.2.1 Fast Min Hashing Algorithm

This is efficient since we need to construct the full matrix, and we need to permute it k times. A faster way is the min hash algorithm.

Make one pass over the data. Let $n = |S|$. Maintain k random hash functions $\{h_1, h_2, \dots, h_k\}$ chosen from a hash family at random so $h_i: U \rightarrow [n]$ (one can use a larger range of n' so values of $n' \geq 2n$ for a power of two). An initial k values at $\{r_1, r_2, \dots, r_k\}$ so $r_i \in [n]$.

Algorithm 4.2.1 Min Hash on S

for $i = 1$ to k

 for $j = 1$ to n do

$h_j(i) = r_j$ then

$r_j = h_j(i)$

 end for

end for

On output $m_j(S) = r_j$. The algorithm runs in $O(kn)$ steps, for a set of size $|S|$. Note this is independent of the size of all possible elements. And the output space of a single set is only $k = (1/(2\epsilon^2)) \ln(2/\delta)$ which is independent of the size of the original set. The space for N sets is only $O(kN)$.

Finally, we can use our estimate $J(S_1, S_3)$ for two sets S_1 and S_3 if

when $|S_1| = 1$ if $r_j = \text{True}$ and 0 otherwise. This only takes $O(k)$ time, again independent of n (sets and δ).

Date	Topic (+ Notes)	Video	Link	Assignment (latex)	Project
Mon 1.08	Class Overview	Vid	MMDS 1.1		
Wed 1.10	Statistics Principles	Vid	M4DA 3 MMDS 1.2 FoDS 12.4		
Mon 1.15	MLK DAY				
Wed 1.17	Similarity : Jaccard + k-Grams (S)	Vid	MMDS 3.1 + 3.2 FoDS 7.3		
Mon 1.22	Similarity : Min Hashing	Vid	MMDS 3.3		
Wed 1.24	Similarity : LSH	Vid	MMDS 3.4	Statistical Principles	
Mon 1.29	Similarity : Distances	Vid	MMDS 3.5 + 7.1 FoDS 8.1		
Wed 1.31	Similarity : SIFT and ANN vs. LSH	Vid	MMDS 3.7 + 7.1.3		Proposal
Mon 2.05	Clustering : Hierarchical	Vid	MMDS 7.2 FoDS 8.7		
Wed 2.07	Clustering : K-Means	Vid	M4DA 8 MMDS 7.3 FoDS 8.3		
Mon 2.12	Clustering : Spectral (S)	Vid	MMDS 10.4 FoDS 8.4	Document Hash	
Wed 2.14	Streaming : Misra-Greis and Frugal	Vid	MMDS 4.1 FoDS 7.1.3		
Mon 2.19	PRESIDENTS DAY				
Wed 2.21	Streaming : Count-Min + Apriori Algorithm	Vid	MMDS 6+4.3 BF Analysis		Data Collection Report
Mon 2.26	Regression : Basics in 2-dimensions	Vid	M4DA 5 ESL 3.2 and 3.4		
Wed 2.28	Regression : SVD + PCA	Vid	M4DA 4 and 7 FoDS 4	Clustering	
Mon 3.05	Regression : Matrix Sketching	Vid	MMDS 9.4 FoDS 2.7 + 7.2.2 arXiv		
Wed 3.07	MIDTERM TEST				
Mon 3.12	Regression : Random Projections	Vid	FoDS 2.9		
Wed 3.14	Regression : Compressed Sensing and OMP	Vid	FoDS 10.2 Tropp + Gilbert	Frequent	
Mon 3.19	SPRING BREAK				
Wed 3.21	SPRING BREAK				
Mon 3.26	Regression : L1 Regression and Lasso	Vid	Davenport ESL 3.8		Intermediate Report
Wed 3.28	Noise : Noise in Data	Vid	MMDS 9.1 Tutorial		
Mon 4.02	Lecture on Ethics/Fairness -- By Suresh	Vid	10 Simple Rules		
Wed 4.04	Noise : Privacy	Vid	McSherry Dwork	Regression	
Mon 4.09	Graph Analysis : Markov Chains (S)	Vid	MMDS 10.1 + 5.1 FoDS 5 Weckesser		
Wed 4.11	Graph Analysis : PageRank	Vid	MMDS 5.1 + 5.4		
Mon 4.16	Graph Analysis : MapReduce	Vid	MMDS 2		Final Report
Wed 4.18	Graph Analysis : Communities	Vid	MMDS 10.2 + 5.5 FoDS 8.8 + 3.4		Poster Outline
Mon 4.23	ENDTERM TEST				
Mon 4.30				Graphs	
Wed 5.02	Poster Day !!! (3:30-5:30pm)				Poster Presentation

Project*

Final Report Due: Monday, April 16
Turn in report by 2:45pm (through Canvas).

1 Overview

Your project will consist of five elements.

- Project Proposal : Due January 31
- Data Collection Report : Due February 21
- Intermediate Report : Due March 26
- Final Report : Due April 16
- 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 led to much stronger projects without considerably more work required.

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.

CS 6955 Data Mining; Spring 2012

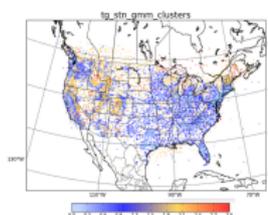
Clustering

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

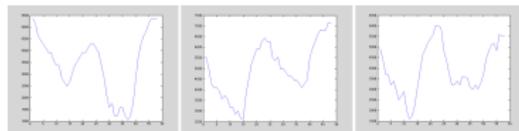


Curve Matching

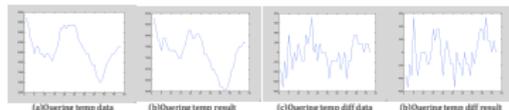
LSH family: Pick a random projection of R^d onto a 1-dimensional line and chop the line into segments of length w , shifted by a random value $b \in [0, w]$.

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

Temperature difference data querying



Temperature difference data querying



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.

Instructor: Jeff M. Phillips, University of Utah

Data Mining

What is Data Mining (in this course)?

Data Mining

What is Data Mining (in this course)?

- ▶ How to think about data analytics.

Data Mining

What is Data Mining (in this course)?

- ▶ How to think about data analytics.

What are course goals?

Data Mining

What is Data Mining (in this course)?

- ▶ How to think about data analytics.

What are course goals?

- ▶ Intuition and principals for data analytics
- ▶ How to model data (convert to abstract data types)
- ▶ How to process data efficiently (balance models with algorithms)
- ▶ To be able to read understand data mining research papers

Data Mining

What is Data Mining (in this course)?

- ▶ How to think about data analytics.

What are course goals?

- ▶ Intuition and principals for data analytics
- ▶ How to model data (convert to abstract data types)
- ▶ How to process data efficiently (balance models with algorithms)
- ▶ To be able to read understand data mining research papers

- ▶ **Not** how to use software toolkits (e.g., scikit-learn).
- ▶ **Not** how to program.
- ▶ We will not cover everything, but cover basics, exposures to many cool modern approaches

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.

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

Videos

Class will be video-recorded and live-streamed.

- ▶ <https://www.youtube.com/channel/UCDUS80bdunpmvWVPyFRPqFQ>
- ▶ links off of webpage to live stream and playlist

Come to class if you can.

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

Data Group

Interested in Research?

(1) Get feedback from Jeff on your projects!

or

(2)

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>