

# CAS CS 565, Data Mining

# Course logistics

- Course webpage:
  - <http://www.cs.bu.edu/~evimaria/cs565-13.html>
- Schedule: Mon – Wed, 4:00–5:30
- Instructor: Evimaria Terzi,  
[evimaria@cs.bu.edu](mailto:evimaria@cs.bu.edu)
- Office hours: Tues 5pm–6:30pm, Mon 2:00pm–3:30pm (or by appointment)

# Topics to be covered

- What is data mining?
- Distance functions
- Finding similar entities
- Dimensionality reduction
- Clustering
- Classification
- Link analysis ranking
- Covering problems and submodular function optimization
- Applications: Web advertising, recommendation systems

# Course workload

- Two programming assignments (25%)
- Three problem sets (25%)
- Midterm exam (20%)
- Final exam (30%)
- **Late assignment policy:** 10% per day up to three days; credit will be not given after that
- Incompletes will not be given

# Learn what you (don't)know

The main goal of the class is for you to get to know what you know and what you don't know (20% rule)

# Textbooks

- A. Rajaraman and J. Ullman: Mining of Massive Datasets. Cambridge University Press, 2012.
- P.-N. Tan, M. Steinbach, V. Kumar: Introduction to Data Mining. Addison-Wesley, 2006.

# Prerequisites

- **Basic algorithms**: sorting, set manipulation, hashing
- **Analysis of algorithms**:  $O$ -notation and its variants, perhaps some recursion equations, NP-hardness
- **Programming**: some programming language, ability to do small experiments reasonably quickly
- **Probability**: concepts of probability and conditional probability, expectations, binomial and other simple distributions
- Some **linear algebra**: e.g., eigenvector and eigenvalue

# Above all

- The goal of the course is to learn and enjoy
- The basic principle is to ask questions when you don't understand
- Say when things are unclear; not everything can be clear from the beginning
- Participate in the class as much as possible
- We will do a lot of thinking together...better to think with company



# Introduction to data mining

- Why do we need data analysis?
- What is data mining?
- Examples where data mining has been useful
- Data mining and other areas of computer science and statistics
- Some (basic) data-mining tasks

# There are lots of data around

# There are lots of data around

- Web (3.7 billion pages)
- Online social networks (Facebook has 1.2 million users)
- Recommendation systems (33 million subscribers on Netflix)
- Wikipedia has 4.4 million articles and counting
- Genomic sequences:  $3 \times 10^9$  nucleotides per individual for 1000 people  $\rightarrow 3 \times 10^{12}$  nucleotides... + medical history + census information

# Example: environmental data

- Climate data (just an example)  
<http://www.ncdc.gov/oa/climate/ghcn-monthly/index.php>
- “a database of temperature, precipitation and pressure records managed by the National Climatic Data Center, Arizona State University and the Carbon Dioxide Information Analysis Center”
- “6000 temperature stations, 7500 precipitation stations, 2000 pressure stations”

# Data complexity

- Multiple types of data: tables, time series, images, graphs, etc
- Spatial and temporal aspects
- Large number of different variables
- Lots of observations → large datasets

# We have large datasets...so

- **Goal:** obtain useful knowledge from large masses of data
- “Data mining is the analysis of (often large) observational data sets to find unsuspected relationships and to summarize the data in novel ways that are both understandable and useful to the data analyst”
- Tell me something interesting about the data; describe the data
- Exploratory analysis on large datasets

# What can data-mining methods do?

# What can data-mining methods do?

- **Rank** web-query results
  - What are the most relevant web-pages to the query: “Student housing BU”?



# What can data-mining methods do?

- **Rank** web-query results
  - What are the most relevant web-pages to the query: “Student housing BU”?
- Find **groups** of entities that are similar (clustering)
  - Find groups of facebook users that have similar friends/interests
  - Find groups amazon users that buy similar products
  - Find groups of walmart customers that buy similar products

# What can data-mining methods do?

- **Rank** web-query results
  - What are the most relevant web-pages to the query: “Student housing BU”?
- Find **groups** of entities that are similar (clustering)
  - Find groups of facebook users that have similar friends/interests
  - Find groups amazon users that buy similar products
  - Find groups of walmart customers that buy similar products
- Find good **recommendations** for users
  - Recommend amazon customers new books
  - Recommend facebook users new friends/groups

# Goal of this course

- Describe some **problems** that can be solved using data-mining methods
- Discuss the **intuition** behind data-mining methods that solve these problems
- Illustrate the **theoretical underpinnings** of these methods
- Show how these methods can be **useful in**

# Data mining and related

- How does data mining relate to machine learning?
- How does data mining relate to statistics?
- Other related areas?

# Data mining vs machine learning

# Data mining vs machine learning

- Machine learning methods are used for data mining
  - Classification, clustering

# Data mining vs machine learning

- Machine learning methods are used for data mining
  - Classification, clustering
- Amount of data makes the difference
  - Data mining deals with much larger datasets and scalability becomes an issue

# Data mining vs machine learning

- Machine learning methods are used for data mining
  - Classification, clustering
- Amount of data makes the difference
  - Data mining deals with much larger datasets and scalability becomes an issue
- Data mining has more modest goals
  - Automating tedious discovery tasks, not aiming at human performance in real discovery
  - Helping users, not replacing them



# Data mining vs. statistics

# Data mining vs. statistics

- “tell me something interesting about this data” – what else is this than statistics?

# Data mining vs. statistics

- “tell me something interesting about this data” – what else is this than statistics?
  - The goal is similar
  - Different types of methods
  - In data mining one investigates lot of possible hypotheses
  - Data mining is more exploratory data analysis
  - In data mining there are much larger datasets → algorithmics/scalability is an issue

# Data mining and algorithms

- Lots of nice connections
- A wealth of interesting research questions
- We will focus on some of these questions later in the course

# Some simple data-analysis tasks

- Given a stream or set of numbers (identifiers, etc)
- How many numbers are there?
- How many distinct numbers are there?
- What are the most frequent numbers?
- How many numbers appear at least  $K$  times?
- How many numbers appear only once?
- etc

# Finding the majority element

- A neat problem
- A stream of identifiers; one of them occurs more than 50% of the time
- How can you find it using no more than a few memory locations?
- Suggestions?

# Finding the majority element

- $A =$  first item you see;  $\text{count} = 1$
  - **for** each subsequent item  $B$ 
    - if**  $(A==B)$   $\text{count} = \text{count} + 1$
    - else**
      - $\text{count} = \text{count} - 1$
      - if**  $(\text{count} == 0)$   $A=B$ ;  $\text{count} = 1$
  - endfor**
  - return**  $A$
- Why does this work correctly?

# Finding the majority element

- A = first item you see;
  - count = 1
  - **for** each subsequent item B
    - if** (A==B)
      - count = count + 1
    - else**
      - count = count - 1
      - if** (count == 0)
        - A=B;
        - count = 1
  - endfor**
  - return** A
- **Basic observation:**  
Whenever we discard element **u** we also discard a unique element **v** different from **u**



# Finding a number in the top half

- Given a set of  $N$  numbers ( $N$  is very large)
- Find a number  $x$  such that  $x$  is **\*likely\*** to be larger than the **median** of the numbers
- Simple solution
  - Sort the numbers and store them in sorted array  $A$
  - Any value larger than  $A[N/2]$  is a solution

# Finding a number in the top half

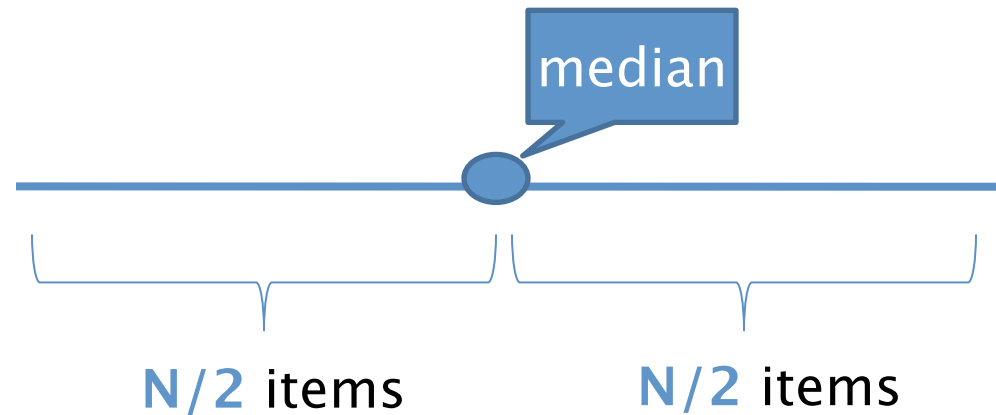
- Given a set of  $N$  numbers ( $N$  is very large)
- Find a number  $x$  such that  $x$  is **\*likely\*** to be larger than the **median** of the numbers
- Simple solution
  - Sort the numbers and store them in sorted array  $A$
  - Any value larger than  $A[N/2]$  is a solution

# Finding a number in the top half

- Given a set of  $N$  numbers ( $N$  is very large)
- Find a number  $x$  such that  $x$  is **\*likely\*** to be larger than the **median** of the numbers
- Simple solution
  - Sort the numbers and store them in sorted array  $A$
  - Any value larger than  $A[N/2]$  is a solution
- Other solutions?

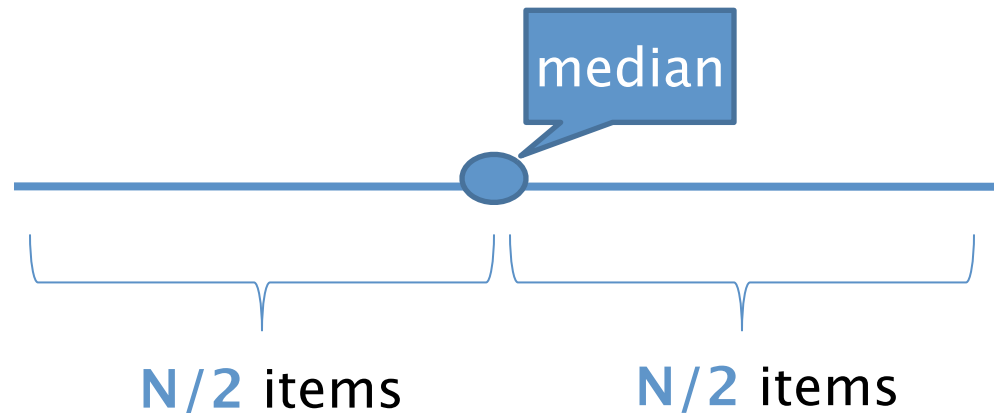
# Finding a number in the top half **efficiently**

# Finding a number in the top half efficiently



# Finding a number in the top half efficiently

- A solution that uses small number of operations
  - Randomly sample **K** numbers from the file
  - Output their maximum



- Failure probability  $(1/2)^K$

# Sampling a sequence of items

- **Problem:** Given a sequence of items  $P$  of size  $N$  form a **random sample**  $S$  of  $P$  that has size  $n$  ( $n < N$ )  $\rightarrow$  sampling without replacement
- What does random sample mean?
  - Every element in  $P$  appears in  $S$  with probability  $n/N$
  - Equivalent as if you generate a **random permutation** of the  $N$  elements and take the **first**  $n$  elements of the permutation

# Sampling algorithm v.0.

- $R = \{\}$  // empty set
- **for**  $i=1$  **to**  $n$ 
  - $\text{rnd} = \text{Random}([1\dots N])$
  - while** ( $\text{rnd}$  in  $R$ )
    - $\text{rnd} = \text{Random}([1\dots N])$
  - endwhile**
  - $R = R \cup \{\text{rnd}\}$
  - $S[i] = P[\text{rnd}]$
- endfor**
- return**  $S$



# Sampling algorithm v.0.

- $R = \{\}$  // empty set
  - **for**  $i=1$  **to**  $n$ 
    - $\text{rnd} = \text{Random}([1\dots N])$
    - while** ( $\text{rnd}$  in  $R$ )
      - $\text{rnd} = \text{Random}([1\dots N])$
    - endwhile**
    - $R = R \cup \{\text{rnd}\}$
    - $S[i] = P[\text{rnd}]$
  - endfor**
  - return**  $S$
- 
- Running time?

# Sampling algorithm v.0.

- $R = \{\}$  // empty set
  - **for**  $i=1$  **to**  $n$ 
    - $\text{rnd} = \text{Random}([1\dots N])$
    - while** ( $\text{rnd}$  in  $R$ )
      - $\text{rnd} = \text{Random}([1\dots N])$
    - endwhile**
    - $R = R \cup \{\text{rnd}\}$
    - $S[i] = P[\text{rnd}]$
  - endfor**
  - return**  $S$
- Running time?

# Sampling algorithm v.0.

- $R = \{\}$  // empty set
- **for**  $i=1$  **to**  $n$ 
  - $\text{rnd} = \text{Random}([1\dots N])$
  - while** ( $\text{rnd}$  in  $R$ )
    - $\text{rnd} = \text{Random}([1\dots N])$
  - endwhile**
  - $R = R \cup \{\text{rnd}\}$
  - $S[i] = P[\text{rnd}]$
- **endfor**
- **return**  $S$
- Running time?
- The algorithm assumes that  $S$  and its size are known in advance!

# Sampling algorithm v.1.

- **Step 1:** Create a random permutation  $\pi$  of the elements in  $P$
- **Step 2:** Return the first  $n$  elements of the permutation,  $S[i] = \pi[i]$ , for  $(1 \leq i \leq n)$

# Sampling algorithm v.1.

- **Step 1:** Create a random permutation  $\pi$  of the elements in  $P$
- **Step 2:** Return the first  $n$  elements of the permutation,  $S[i] = \pi[i]$ , for  $(1 \leq i \leq n)$

You can do Step 2 in linear time 😊

# Sampling algorithm v.1.

- **Step 1:** Create a random permutation  $\pi$  of the elements.

Can you do Step 1 in linear time?

- **Step 2:** Return the first  $n$  elements of the permutation,  $S[i] = \pi[i]$ , for  $(1 \leq i \leq n)$ .

You can do Step 2 in linear time 😊

# Creating a random permutation in linear time

- **for**  $i=1\dots N$  **do**
  - $j = \text{Random}([1\dots i-1])$
  - swap  $P[i]$  with  $P[j]$
- **endfor**
- Is this really a random permutation?  
(see CLR for the proof)
- It runs in linear time

# Sampling algorithm v.1.

- **Step 1:** Create a random permutation  $\pi$  of the elements in  $P$
- **Step 2:** Return the first  $n$  elements of the permutation,  $S[i] = \pi[i]$ , for  $(1 \leq i \leq n)$



# Sampling algorithm v.1.

- **Step 1:** Create a **random permutation**  $\pi$  of the elements in  $P$
- **Step 2:** Return the first  $n$  elements of the permutation,  $S[i] = \pi[i]$ , for  $(1 \leq i \leq n)$
- The algorithm works in **linear time**  $O(N)$

# Sampling algorithm v.1.

- **Step 1:** Create a **random permutation**  $\pi$  of the elements in  $P$
- **Step 2:** Return the first  $n$  elements of the permutation,  $S[i] = \pi[i]$ , for  $(1 \leq i \leq n)$
- The algorithm works in **linear time**  $O(N)$
- The algorithm assumes that  $P$  is **known in advance**

# Sampling algorithm v.1.

- **Step 1:** Create a **random permutation**  $\pi$  of the elements in  $P$
- **Step 2:** Return the first  $n$  elements of the permutation,  $S[i] = \pi[i]$ , for  $(1 \leq i \leq n)$
- The algorithm works in **linear time**  $O(N)$
- The algorithm assumes that  $P$  is **known in advance**
- The algorithm makes **2 passes** over the data

# Sampling algorithm v.2.

- **for**  $i = 1$  to  $n$   
     $S[i] = P[i]$   
**endfor**
- $t = n + 1$
- **while**  $P$  has more elements  
     $\text{rnd} = \text{Random}([1\dots t])$   
    if ( $\text{rnd} \leq n$ )  
         $\{S[\text{rnd}] = P[t]\}$   
     $t = t + 1$   
**endwhile**

# Sampling algorithm v.2.

- **for**  $i = 1$  to  $n$  **Correctness proof**  
     $S[i] = P[i]$   
**endfor**
- $t = n + 1$
- **while**  $P$  has more elements  
     $\text{rnd} = \text{Random}([1\dots t])$   
    if ( $\text{rnd} \leq n$ )  
         $\{S[\text{rnd}] = P[t]\}$   
     $t = t + 1$   
**endwhile**

# Sampling algorithm v.2.

- **for**  $i = 1$  to  $n$   
     $S[i] = P[i]$   
**endfor**
- $t = n + 1$
- **while**  $P$  has more elements  
     $\text{rnd} = \text{Random}([1\dots t])$   
    if ( $\text{rnd} \leq n$ )  
         $\{S[\text{rnd}] = P[t]\}$   
     $t = t + 1$   
**endwhile**

## Correctness proof

- At iteration  $t+1$  a **new** item is included in the sample with probability  $n/(t+1)$

# Sampling algorithm v.2.

- **for**  $i = 1$  to  $n$   
     $S[i] = P[i]$   
**endfor**
- $t = n + 1$
- **while**  $P$  has more elements  
     $\text{rnd} = \text{Random}([1\dots t])$   
    if ( $\text{rnd} \leq n$ )  
         $\{S[\text{rnd}] = P[t]\}$   
     $t = t + 1$   
**endwhile**

## Correctness proof

- At iteration  $t+1$  a **new** item is included in the sample with probability  $n/(t+1)$
- At iteration  $(t+1)$  an **old** item is kept in the sample with probability  $n/(t+1)$ 
  - **Inductive argument:** at iteration  $t$  the old item was in the sample with probability  $n/t$

# Sampling algorithm v.2.

- **for**  $i = 1$  to  $n$   
     $S[i] = P[i]$   
**endfor**
- $t = n + 1$
- **while**  $P$  has more elements  
     $\text{rnd} = \text{Random}([1\dots t])$   
    if ( $\text{rnd} \leq n$ )  
         $\{S[\text{rnd}] = P[t]\}$   
     $t = t + 1$   
**endwhile**

## Correctness proof

- At iteration  $t+1$  a **new** item is included in the sample with probability  $n/(t+1)$
- At iteration  $(t+1)$  an **old** item is kept in the sample with probability  $n/(t+1)$ 
  - **Inductive argument:** at iteration  $t$  the old item was in the sample with probability  $n/t$
  - $\text{Pr}(\text{old item in sample at } t+1) = \text{Pr}(\text{old item was in sample at } t) \times (\text{Pr}(\text{rnd} > n) + \text{Pr}(\text{rnd} \leq n) \times \text{Pr}(\text{old item was not chosen for eviction}))$



# Sampling algorithm v.2.

- **for**  $i = 1$  to  $n$   
     $S[i] = P[i]$   
**endfor**
- $t = n + 1$
- **while**  $P$  has more elements {  
     $\text{rnd} = \text{Random}([1\dots t])$   
    if ( $\text{rnd} \leq n$ )  
         $\{S[\text{rnd}] = P[t]\}$   
     $t = t + 1$   
**endwhile**

# Sampling algorithm v.2.

- **for**  $i = 1$  to  $n$   
     $S[i] = P[i]$   
**endfor**
- $t = n + 1$
- **while**  $P$  has more elements {  
     $\text{rnd} = \text{Random}([1\dots t])$   
    if ( $\text{rnd} \leq n$ )  
         $\{S[\text{rnd}] = P[t]\}$   
     $t = t + 1$   
**endwhile**

## Advantages

# Sampling algorithm v.2.

- **for**  $i = 1$  to  $n$   
     $S[i] = P[i]$   
**endfor**
- $t = n + 1$
- **while**  $P$  has more elements {  
     $\text{rnd} = \text{Random}([1\dots t])$   
    if ( $\text{rnd} \leq n$ )  
         $\{S[\text{rnd}] = P[t]\}$   
     $t = t + 1$   
**endwhile**

## Advantages

- Linear time

# Sampling algorithm v.2.

- **for**  $i = 1$  to  $n$   
     $S[i] = P[i]$   
**endfor**
- $t = n + 1$
- **while**  $P$  has more elements {  
     $\text{rnd} = \text{Random}([1\dots t])$   
    if ( $\text{rnd} \leq n$ )  
         $\{S[\text{rnd}] = P[t]\}$   
     $t = t + 1$   
**endwhile**

## Advantages

- Linear time
- **Single pass** over the data

# Sampling algorithm v.2.

- **for**  $i = 1$  to  $n$   
     $S[i] = P[i]$   
**endfor**
- $t = n + 1$
- **while**  $P$  has more elements {  
     $\text{rnd} = \text{Random}([1\dots t])$   
    if ( $\text{rnd} \leq n$ )  
         $\{S[\text{rnd}] = P[t]\}$   
     $t = t + 1$   
**endwhile**

## Advantages

- Linear time
- **Single pass** over the data
- **Any time**; the length of the sequence need not be known in advance