# Data Mining

**SU 5050**

**LECTURE 1**

**JESSICA L. MCCARTY, PH.D.**

**AND MIKE BILLMIRE, M.S.**

# What is Data Mining?

- Recently* coined term for:
  - Confluence of ideas from statistics and computer science
    - Machine learning and database methods
  - Applied to large databases in science, engineering, business

  - *First International Workshop on Knowledge Discovery and Data Mining was in 1995
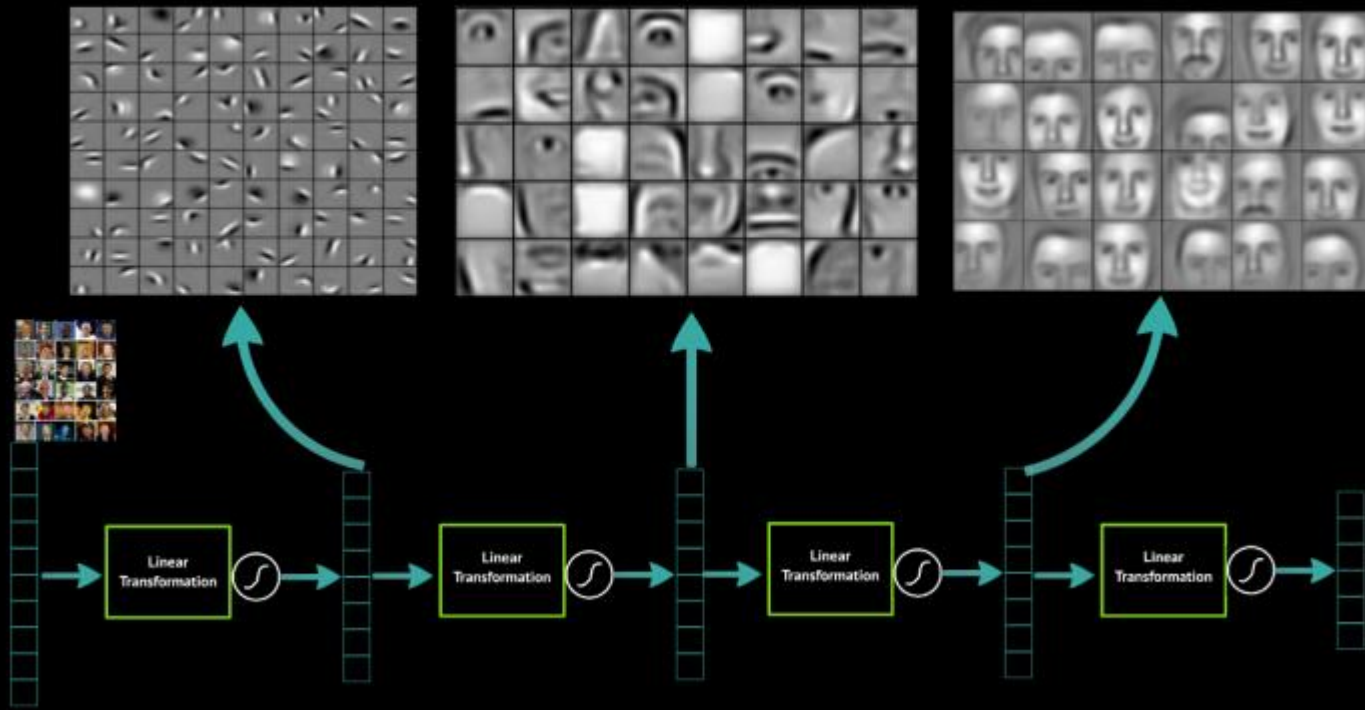
# What is Data Mining?

- As less than 20-year-old discipline, is in state of flux
  - Debate over what it is and what it is not

- Terminology is not standard

- Bias, classification, prediction, feature = independent variable

- Target = dependent variable

- Case = exemplar = row

# Data Mining's Cousins



Deep Learning learns layers of features
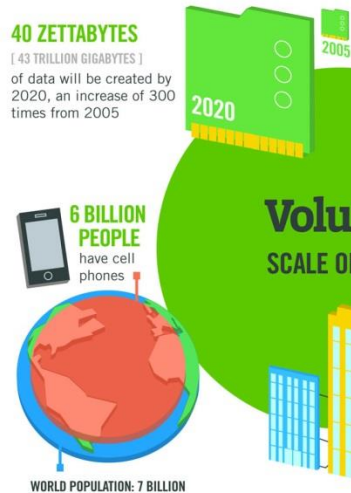
# Data Mining's Cousins

**40 ZETTABYTES**
[ 43 TRILLION GIGABYTES ] of data will be created by 2020, an increase of 300 times from 2005

**6 BILLION PEOPLE** have cell phones

WORLD POPULATION: 7 BILLION

## Volume
### SCALE OF DATA

It's estimated that **2.5 QUINTILLION BYTES** [ 2.3 TRILLION GIGABYTES ] of data are created each day

Most companies in the U.S. have at least **100 TERABYTES** [ 100,000 GIGABYTES ] of data stored

## The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015
**4.4 MILLION IT JOBS** will be created globally to support big data, with 1.9 million in the United States

As of 2011, the global size of data in healthcare was estimated to be **150 EXABYTES** [ 161 BILLION GIGABYTES ]

**30 BILLION PIECES OF CONTENT** are shared on Facebook every month

By 2014, it's anticipated there will be **420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

**4 BILLION+ HOURS OF VIDEO** are watched on YouTube each month

## Variety
### DIFFERENT FORMS OF DATA

**400 MILLION TWEETS** are sent per day by about 200 million monthly active users

The New York Stock Exchange captures **1 TB OF TRADE INFORMATION** during each trading session
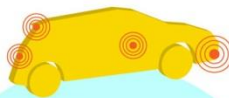
By 2016, it is projected there will be **18.9 BILLION NETWORK CONNECTIONS** – almost 2.5 connections per person on earth

Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure

## Velocity
### ANALYSIS OF STREAMING DATA

**1 IN 3 BUSINESS LEADERS** don't trust the information they use to make decisions

**27% OF RESPONDENTS** in one survey were unsure of how much of their data was inaccurate

Poor data quality costs the US economy around **$3.1 TRILLION A YEAR**

## Veracity
### UNCERTAINTY OF DATA

IBM

# Data Mining and Deep Learning's Evil Spawn

- Google's Deep Dream Neural Network



- http://googleresearch.blogspot.com/2015/07/deepdream-code-example-for-visualizing.html

# Broad and Narrow Definition

- Broad -> traditional statistical methods

- Narrow - > automated and heuristic methods

- Heuristics – exploratory problem-solving techniques that give a non-optimal solution (exhaustive search impractical)

  - Rule of thumb, educated guess, intuitive judgment, stereotyping, common sense
  - Computer Science – technique used when classic methods are too slow; approximate solution after exact solution not found.
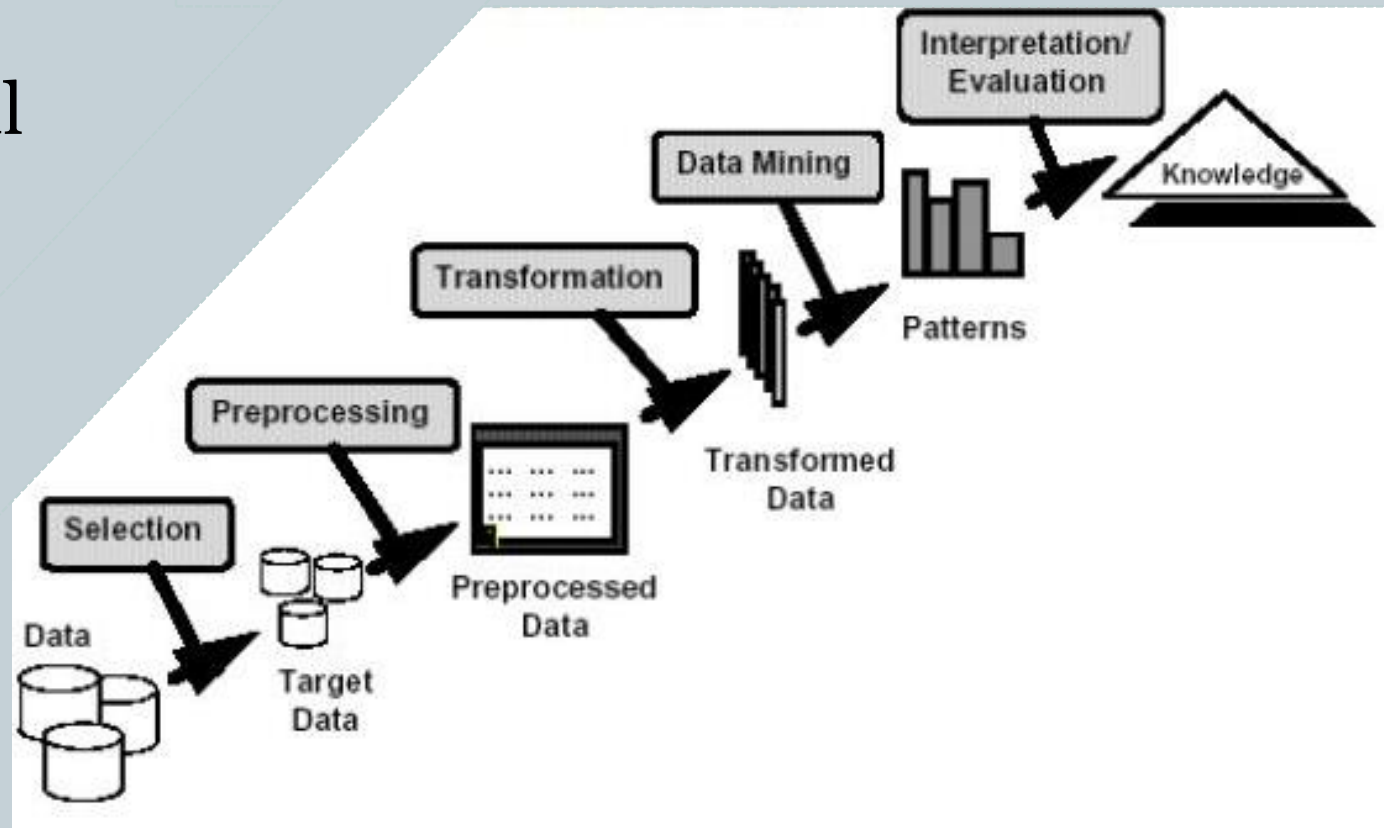
# Broad and Narrow Definitions

- Data mining, data dredging, fishing expeditions

- Knowledge Discovery in Databases (KDD)
  - Interactive and iterative
  - Many interactions and feedback loops between steps
  - http://www.usc.edu/dept/ancntr/Paris-in-LA/Analysis/discovery.html

# Schema for Data Mining

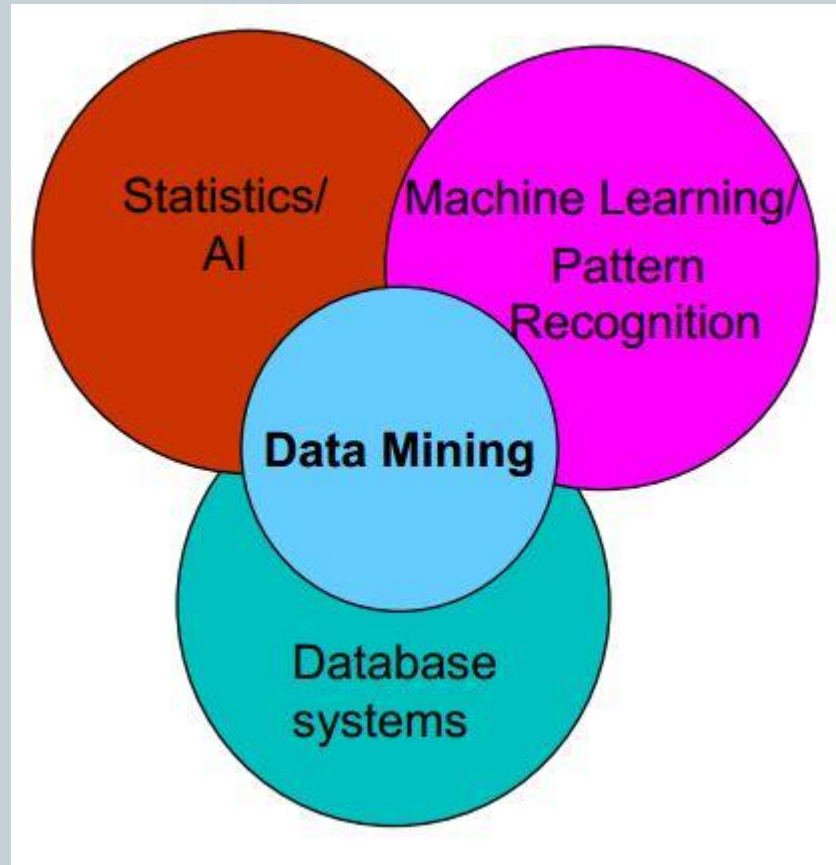- Search for meaningful patterns

# Drivers

- Market -> From focus on products/service to focus on customers

- IT -> From focus on up-to-date balance to focus on patterns in transactions (Data warehouses, cloud)

- Automatic Data Capture of Transactions (bar codes, POS devices, mouse clicks, GPS/locational data)

- Internet -> personalized interactions, longitudinal data

# Core Disciplines

- Statistics: Visualization (Descriptive Stats) & Regression, Cluster Analysis (Models)

- Machine Learning: Neural Nets

- Database Retrievals: Association Rules

- Parallel developments: Decision trees, k means, nearest neighbors, Online Analytical Processing (OLAP) Exploratory Data Analysis (EDA)

# Core Disciplines

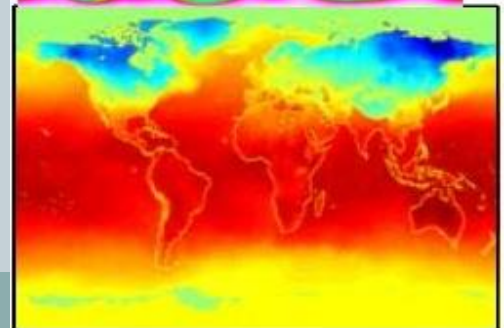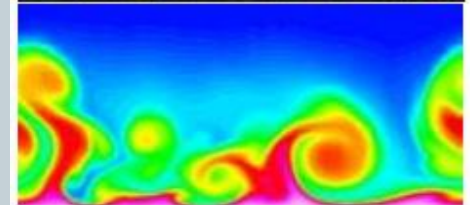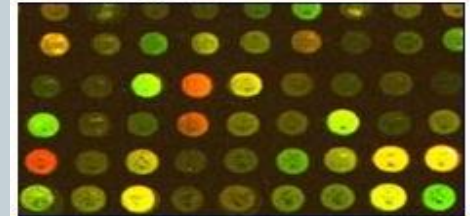# Why Mine Data? Commercial Viewpoint



- Loads of data collected and stored

- Computers are cheaper and more powerful

- Competitive pressure is strong
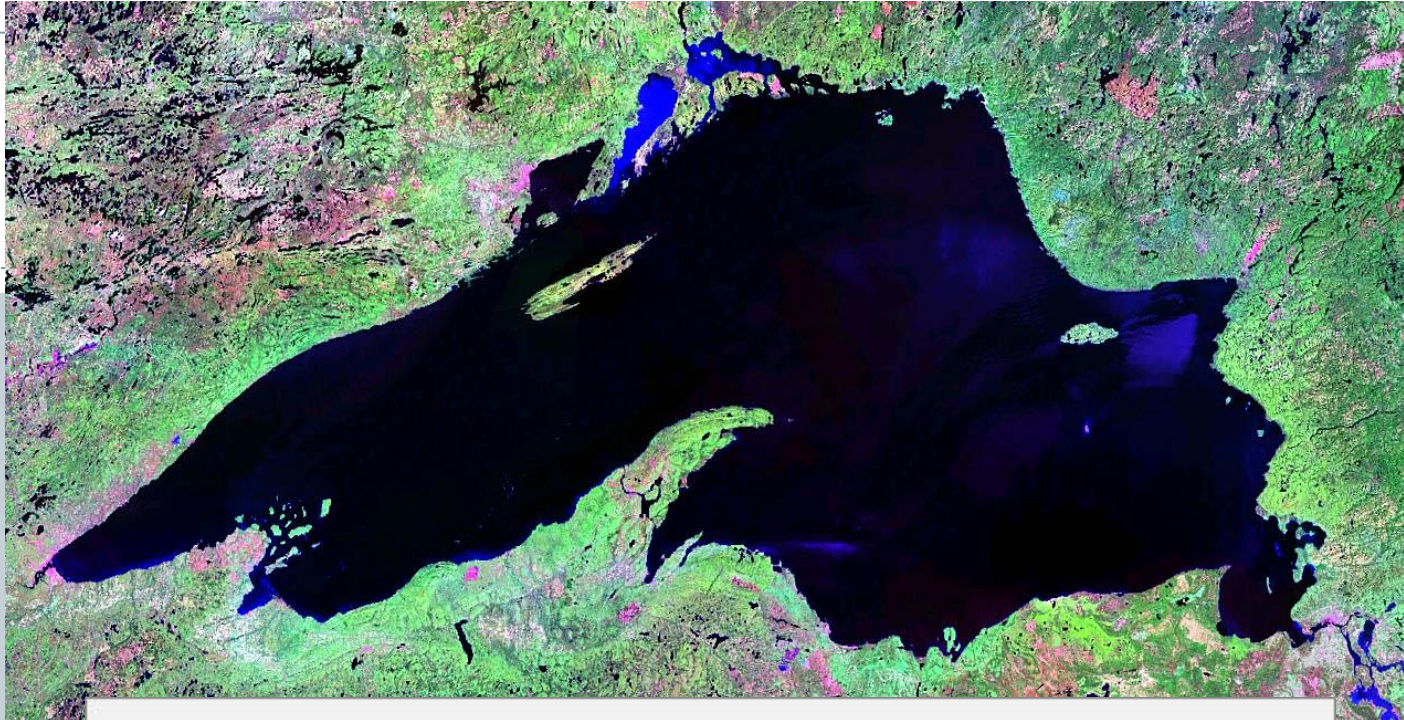  - Provide better, customized service for an **edge**
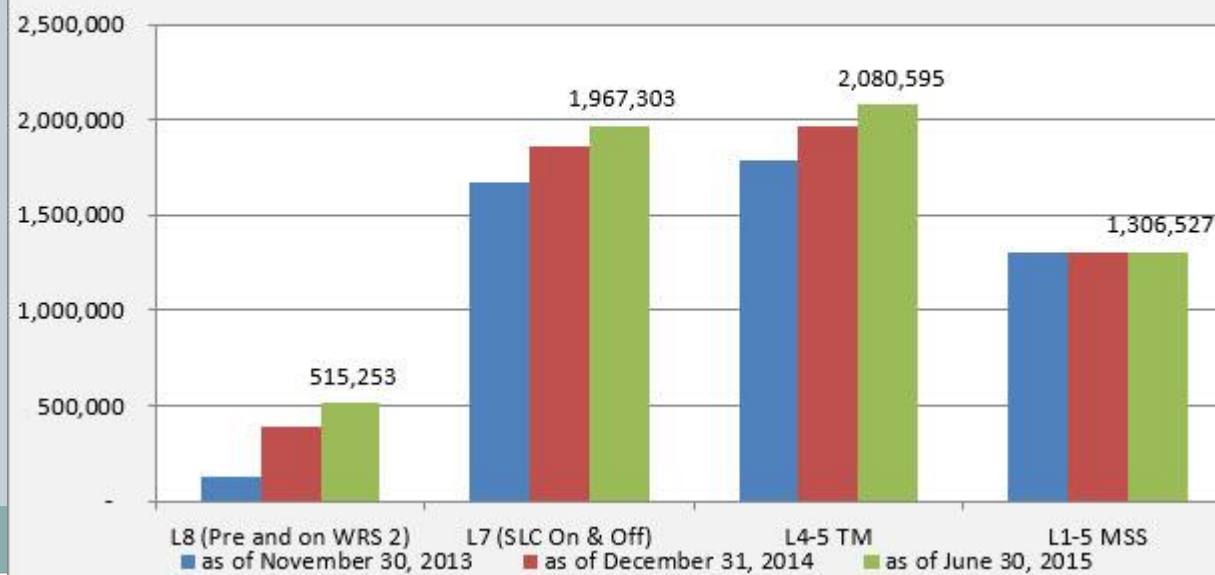
# Why Mine Data? Scientific Viewpoint

- Data collected and stored at enormous speeds and quantities (TB/hour globally)
  - Remote sensors on a satellite
  - High powered telescopes
  - Microarrays replicating genome
  - Scientific simulations
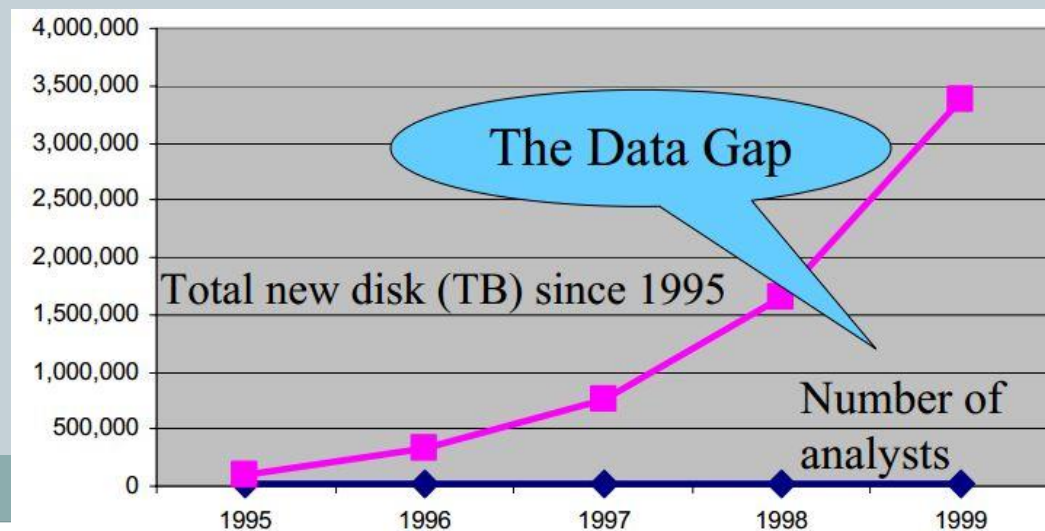- Classifying data
- Hypothesis formation
- Visualizations
- http://usdaapps.devpost.com/

Landsat Scenes Visible in EarthExplorer

# Motivation for Mining Large Data Sets

- Often information "hidden" in the data that is not readily evident
- Human analysts may take weeks to discover useful information
- Scope of analyst-based methods different
- Much of the data never analyzed at all

# What is (not) Data Mining?

## What is not Data Mining?

– Look up phone number in phone directory

– Query a Web search engine for information about "Amazon"

## What is Data Mining?

– Certain names are more prevalent in certain US locations (O'Brien, O'Rurke, O'Reilly… in Boston area)

– Group together similar documents returned by search engine according to their context (e.g. Amazon rainforest, Amazon.com,)

# Process

1. Develop understanding of application, goals
2. Create dataset
3. Data cleaning and preprocessing
4. Data reduction and projection
5. Choose data mining task
6. Choose data mining algorithms
7. Use algorithms to perform task
8. Interpret and iterate thru 1-7 *if necessary*
9. Deploy: integrate into/create new operational system
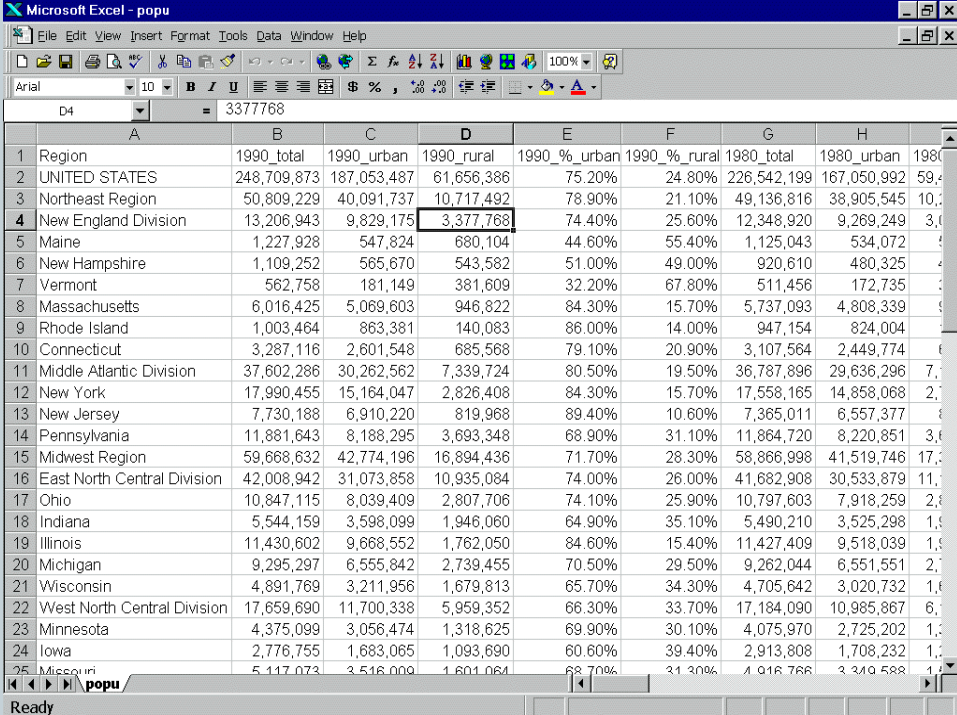
Data Mining

# Challenges of Data Mining

- Scalability
- Dimensionality
- Complex and Heterogeneous Data
- Data Quality
- Data Ownership and Distribution
- Privacy Preservation
- Streaming Data

# Typical characteristics of mining data

- "Standard format is spreadsheet
  - Row = observation unit
  - Column = variable

- Many rows, many columns

- Many rows, few columns

- Few rows, many columns

- Opportunistic data collect

# What is Data?

- Collection of data objects and their attributes

- An attribute is a property or characteristic of an object
  - Examples: eye color of a person, temperature, etc.
  - Attribute is also known as variable, field, characteristic, or feature

- A collection of attributes describe an object
  - Object is also known as record, point, case, sample, entity, or instance

**Attributes**

**Objects**

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

# Attribute Values

- **Attribute values are numbers or symbols assigned to an attribute**

- **Distinction between attributes and attribute values**
  - Same attribute can be mapped to different attribute values
    - Example: height can be measured in feet or meters

  - Different attributes can be mapped to the same set of values
    - Example: Attribute values for ID and age are integers
    - But properties of attribute values can be different
      - ID has no limit but age has a maximum and minimum value

# Measurement of Length

- **The way you measure an attribute is somewhat may not match the attributes properties.**

| | | |
|---|---|---|
| 5 | A | 1 |
| 7 | B | 2 |
| 8 | C | 3 |
| 10 | D | 4 |
| 15 | E | 5 |

# Types of Attributes

- **There are different types of attributes**
  - Nominal
    - Examples: ID numbers, eye color, zip codes
  - Ordinal
    - Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height in {tall, medium, short}
  - Interval
    - Examples: calendar dates, temperatures in Celsius or Fahrenheit.
  - Ratio
    - Examples: temperature in Kelvin, length, time, counts

# Properties of Attribute Values

- The type of an attribute depends on which of the following properties it possesses:

    - Distinctness:          $=$ $\neq$

    - Order:                 $<$ $>$

    - Addition:              $+$ $-$

    - Multiplication:        $*$ $/$

    - Nominal attribute: distinctness

    - Ordinal attribute: distinctness & order

    - Interval attribute: distinctness, order & addition

    - Ratio attribute: all 4 properties

| Attribute Type | Description | Examples | Operations |
|---|---|---|---|
| Nominal | The values of a nominal attribute are just different names, i.e., nominal attributes provide only enough information to distinguish one object from another. $(=, \neq)$ | zip codes, employee ID numbers, eye color, sex: {*male, female*} | mode, entropy, contingency correlation, $\chi^2$ test |
| Ordinal | The values of an ordinal attribute provide enough information to order objects. $(<, >)$ | hardness of minerals, {*good, better, best*}, grades, street numbers | median, percentiles, rank correlation, run tests, sign tests |
| Interval | For interval attributes, the differences between values are meaningful, i.e., a unit of measurement exists. $(+, -)$ | calendar dates, temperature in Celsius or Fahrenheit | mean, standard deviation, Pearson's correlation, $t$ and $F$ tests |
| Ratio | For ratio variables, both differences and ratios are meaningful. $(*, /)$ | temperature in Kelvin, monetary quantities, counts, age, mass, length, electrical current | geometric mean, harmonic mean, percent variation |

| Attribute Level | Transformation | Comments |
|---|---|---|
| Nominal | Any permutation of values | If all employee ID numbers were reassigned, would it make any difference? |
| Ordinal | An order preserving change of values, i.e., $new\_value = f(old\_value)$ where $f$ is a monotonic function. | An attribute encompassing the notion of good, better best can be represented equally well by the values $\{1, 2, 3\}$ or by $\{0.5, 1, 10\}$. |
| Interval | $new\_value = a * old\_value + b$ where a and b are constants | Thus, the Fahrenheit and Celsius temperature scales differ in terms of where their zero value is and the size of a unit (degree). |
| Ratio | $new\_value = a * old\_value$ | Length can be measured in meters or feet. |

# Discrete and Continuous Attributes

- **Discrete Attribute**
  - Has only a finite or countably infinite set of values
  - Examples: zip codes, counts, or the set of words in a collection of documents
  - Often represented as integer variables.
  - Note: binary attributes are a special case of discrete attributes

- **Continuous Attribute**
  - Has real numbers as attribute values
  - Examples: temperature, height, or weight.
  - Practically, real values can only be measured and represented using a finite number of digits.
  - Continuous attributes are typically represented as floating-point variables.

# Types of Data Sets

- **Record**
  - Data Matrix
  - Document Data
  - Transaction Data

- **Graph**
  - World Wide Web
  - Molecular Structures

- **Ordered**
  - Spatial Data
  - Temporal Data
  - Sequential Data
  - Genetic Sequence Data

# Syllabus

- Class Goals

- Assignments

- Expectations

- Schedule

- Canvas Site



*"Drowning in data, yet starving for knowledge."*
-- **Anonymous**

*"Where is the knowledge we have lost in information?"* -- **T.S. Eliot**

## Data Mining
### SU 5050

| Instructors | Jessica L. McCarty, PhD<br>Adjunct Faculty, School of Technology, MTU<br>Research Scientist,<br>Michigan Tech Research Institute<br>mtri.org | Michael Billmire, MS<br>and CMS-GIS/LIS<br>Research Scientist,<br>Michigan Tech<br>Research Institute |
|---|---|---|
| Contact | jmccarty@mtu.edu<br>Cell: 502.415.1628<br>Work: 734.994.7236 | mgbillmi@mtu.edu<br>Work: 734.913.6853 |
| Office Hours | Thurs 10 am – 12 pm, Online via Adobe Connect, via Google Hangout, Email or Phone (work, then cell) *Any communication will answered immediately.* | * Please discuss lab issues with Billmire during lab times; email only if you have tried 10 times and can NOT make it work. |
| Class Meets | Online Lecture: Mondays and Wednesdays 12:00 to 12:55 pm via http://mtu.adobeconnect.com/datamining/<br><br>Online Lab Instruction: Fridays 12:00 to 12:55 pm via http://mtu.adobeconnect.com/datamining/ | |
| Canvas | The Canvas site will be used to distribute pdf copies of lecture slides and lab assignments, online midterm and final, and for online submissions. Lectures will be made available after each class, including links to video recordings (requires Adobe and Flash Player). | |
| Objectives | This course will be taught in three modules:<br>1. Overview of current techniques, including theory and applications of data mining and big data for geospatial techniques;<br>2. Application focuses on open source programming and library development (Python);<br>3. Writing a research plan suitable for research submission and proof-of-concept study. | |
| Prerequisites | This course is a lot of work. Lab assignments usually require work outside of class and lab times. The course is designed so that students without a programming and geospatial background can succeed, but previous experience will no doubt be helpful. Although **not required for successful completion** of this course, courses in the following areas can be a helpful background: computer programming, statistics, surveying, remote sensing/GIS. | |
| Required Readings | 1. **Textbook** – Russell, M.A.. 2013. *Mining the Social Web*. Second Edition. Sebastapol, CA: O'Reilly Media, Inc. Available at as Ebook, Print & Ebook, Print: http://shop.oreilly.com/product/0636920030195.do . Instructors have a copy of the Ebook. NOTE: You can save money on your online purchase: http://www.retailmenot.com/view/oreilly.com?c=5659596<br><br>2. **Miscellaneous Readings** - from various sources will be available on Canvas. | |

# Who is this Prof McCarty Person?



- PhD in Geography, University of Maryland (2009)
- Research Scientists and Adjunct Professor at Michigan Tech
- mtri.org
- @jmccarty_geo
- Climate and Carbon, Fire, Air Quality, Land Cover/Land Use Change, Food Security, Regional & Natural Planning, Data Mining, Remote Sensing, GIS

# On a personal note





- Native of Eastern Kentucky (Appalachia)
- Izzy Dawg
- Recently moved to Houghton