# Data Mining

CS4821-CS5831-s24

Laura E. Brown

Jan. 10, 2024

Some slides from G. Piatetsky-Shapiro; Han, Kamber, & Pei; M. Hahsler
P. Smyth; C. Volinsky; Tan, Steinbach, & Kumar; J. Taylor; G. Dong; L. Hannah
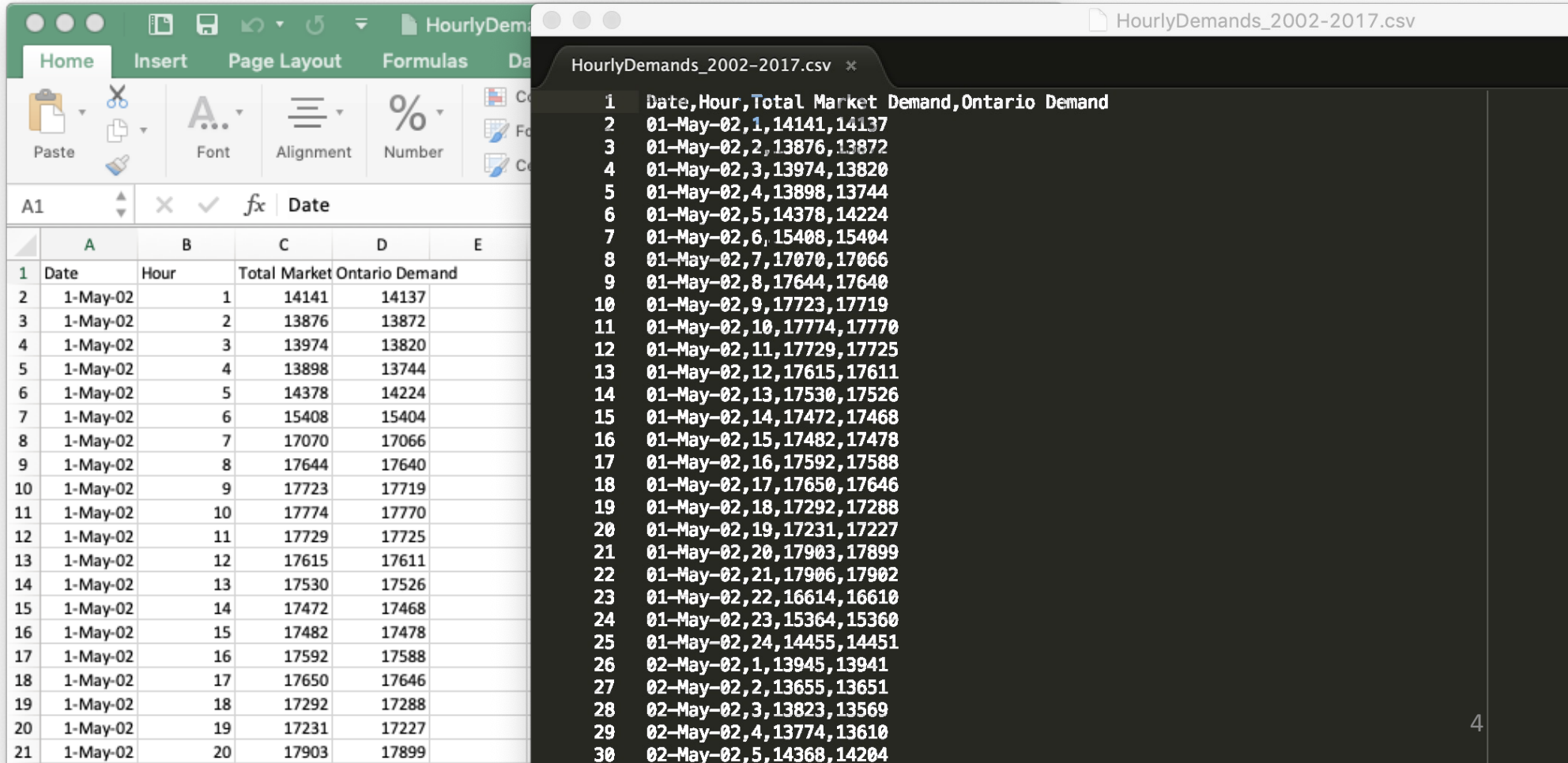
# Today's Agenda

- Examples of Data Mining (continue from 01.intro.part1)

- Multidimensional View of Data Mining

- Tools for Data Mining

- Data Mining vs. Privacy

# Multi-dimensional View of Data Mining

- What kinds of **data** can be mined?

- What kinds of **patterns** can be mined?

- What kinds of **techniques** are used?

- What kinds of **applications** are targeted?

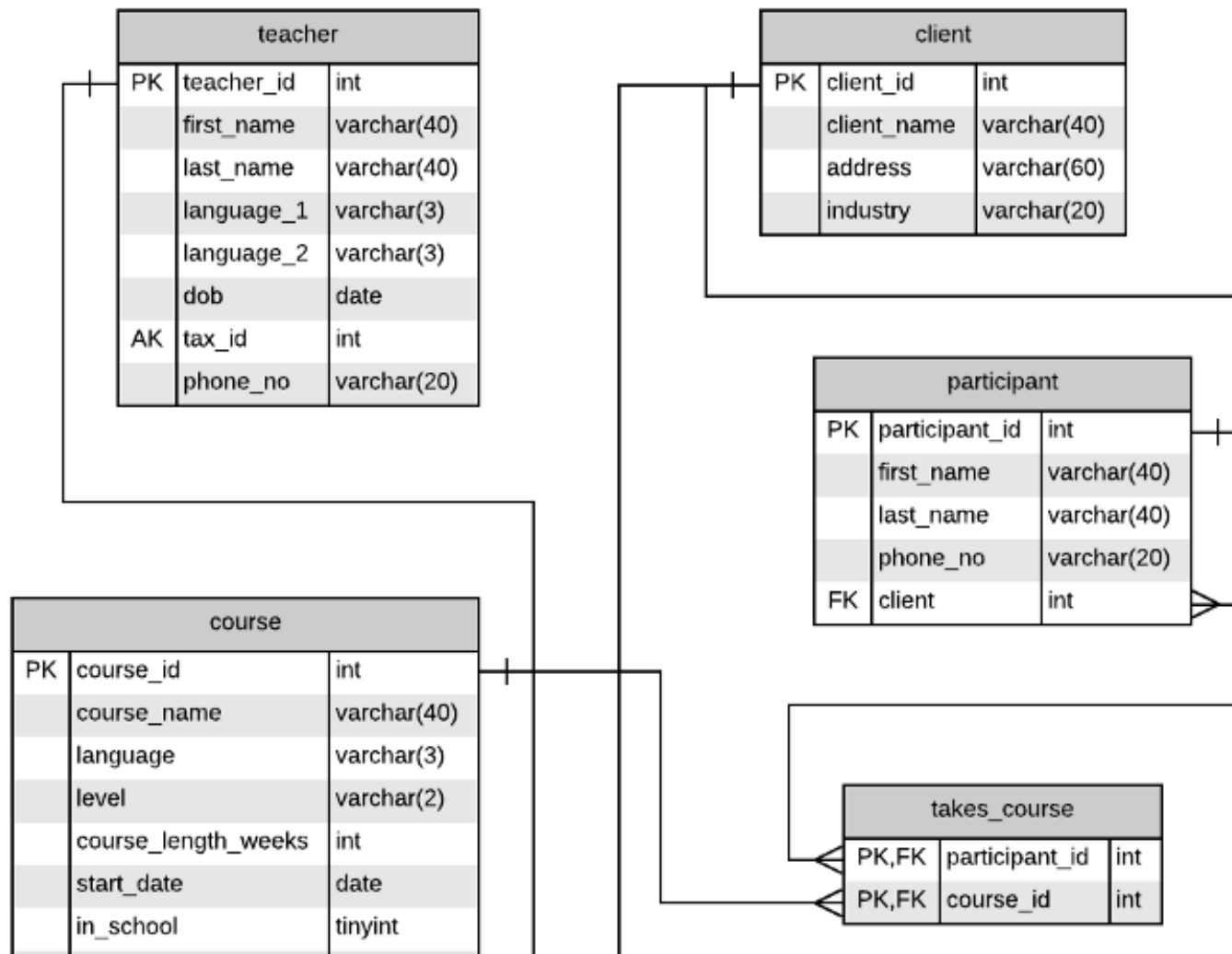# What kind of data?

- Spreadsheets
- Flat file, vector data

# What kind of data?

- Relational data - Databases

# What kind of data?

- Time Series Data

# What kind of data?

● Text files, document collections

# What kind of data?

- Image data

# What kind of data?

- Spatial Data

# What kind of data?

● Transaction Data

| TID | Items |
|---|---|
| 1 | Bread, Coke, Milk |
| 2 | Beer, Bread |
| 3 | Beer, Coke, Diaper, Milk |
| 4 | Beer, Bread, Diaper, Milk |
| 5 | Coke, Diaper, Milk |



item frequency (relative)

# What kind of data?

- Ratings Data



| | SHERLOCK | HOUSE of CARDS | AVENGERS | ARRESTED DEVELOPMENT | Breaking Bad | WALKING DEAD |
|---|---|---|---|---|---|---|
| | 2 | | | 4 | 5 | |
| | (5) | | 4 | | | (1) |
| | | | 5 | | 2 | |
| | | 1 | | 5 | | 4 |
| | | | 4 | | | 2 |
| | 4 | 5 | | 1 | | |

# What kinds of <u>data</u>?

- Flat File, Vector data
- Relational data
- Text files, document collections
- Time series data
- Spatial data
- Spatio-temporal data
- Transactional data

- Ratings data
- Network data
- Image data
- Custom data for particular application
- …, and many more

# What types of <u>patterns</u>?

- Prediction Methods
  - Use some variables to predict unknown or future values of other variables

- Descriptive Methods
  - Find human-interpretable patterns that describe the data

# What type of <u>patterns</u>?

- Generalization / Characterization
- Classification
- Association and Correlation Analysis
- Cluster Analysis
- Recommender Systems
- Structure / Network Analysis
- Outlier Analysis
- Sequential Pattern Analysis
- …

# What kinds of techniques?

- Confluence of techniques and disciplines
  - Statistics
  - Machine Learning
  - Database technology
  - Algorithms
  - Pattern Recognition
  - Parallel Computing
  - Visualization

# What kinds of <u>applications</u>?

- Web page analysis: from web page classification, clustering to PageRank & HITS algorithms

- Collaborative analysis & recommender systems

- Basket data analysis to targeted marketing

- Biological and medical data analysis: classification, cluster analysis (microarray data analysis),  biological sequence analysis, biological network analysis

- Data mining and software engineering (e.g., IEEE Computer, Aug. 2009 issue)

- From major dedicated data mining systems/tools (e.g., SAS, MS SQL-Server Analysis Manager, Oracle Data Mining Tools) to invisible data mining

# Example: Fraud Detection

- Which credit card transactions are fraudulent?
- Goal: Predict fraudulent cases in credit card transactions.
  - Data: credit card transactions, information on account-holder
  - Pattern: classification
    - Labeled transaction data: fair or fraud
    - Learn a model to prediction this label
  - Technique: machine learning

# Example: Fraud Detection (2)

- Which credit card transactions are fraudulent?
- Goal: Predict fraudulent cases in credit card transactions.
  - Data: credit card transactions, information on account-holder
  - Pattern: outlier analysis
    - Fraud events are rare (hopefully!)
    - Detect transaction as unusual from prior historical data
  - Technique: machine learning, statistics

# Example: Clustering

- How should a set of images be placed into groups?
  - Data: image files
  - Pattern: cluster analysis
    - Group images by similarity
    - How to measure similarity?
  - Techniques: machine learning

# In this class …

- You will study algorithms that exploit and reveal patterns in data.

- Goal:
  - You will learn how to think about problems in data mining
  - You will learn about a set of data analysis tools:
    - How to use them
    - What their assumptions are
    - The capabilities and limitations

# Methods to Be Examined

- Supervised learning – Classification
- Unsupervised learning – Data Reduction
- Text Mining
- Unsupervised learning – Clustering
- Association Mining
- Recommendation Systems
- Web Mining (if time permits)

# Popular Tools for Data Mining



Top Analytics, Data Science, Machine Learning Software 2017-2019, KDnuggets Poll

https://www.kdnuggets.com/polls/

# Types of Tools

Simple graphical user interface

Process oriented

Programming oriented

# Tools: Simple GUI

- Weka: Waikato Environment for Knowledge Analysis (Java API)

- Rattle: GUI for Data Mining using R

# Tools: Process oriented

- SAS Enterprise Miner
- IBM SPSS Modeler
- RapidMiner
- Knime
- Orange

# Tools: Programming oriented

- R
  - RStudio IDE
    - Caret
    - Tidyverse

- Python
  - Jupyter notebooks
    - Numpy
    - Pandas
    - Scikit-learn

# Data Mining vs. Privacy

- Tension between data mining and personal privacy

# Data Mining vs. Privacy

- How can we leverage sensitive personal data for research / commercial purposes?

- 3 cases
  - AOL search data set
  - Netflix prize
  - Barabasi mobile study

Examples from C. Volinsky

# Case 1: AOL Search Data

- Aug. 4, 2006 – AOL releases 20M search terms by anonymized users "for research purposes"
- Within hours, uproar on blogs
  - "The utter stupidity of this is staggering" – TechCrunch
- Aug. 7, 2006 – AOL removes data, issues apology
  - "this was a screw-up, and we are angry"
  - "an innocent enough attempt to reach out the the research community"
- Aug. 9, 2006 – NYT front page story
  - Identifies user
- Aug. 21, 2006 – CTO resigns

# Case 1: AOL Search Data

- ## What's the big deal?
  - ### How and why people search is often personal and may contain information they do not want released to public

- ## What went wrong?
  - ### Not well thought out
  - ### Poor internal controls
  - ### Lack of understanding on anonymizing

- ## Fallout
  - ### CTO + at least two others fired
  - ### Data is still out there

# Case 2: Netflix Prize

- Oct. 2006: Netflix released anonymized movie ratings from its customer database
    - 100M ratings, 500K customers (<10% of data)
    - Random integer for user ID
    - "some of the rating data for some customers in the training and qualifying sets have been deliberately perturbed in one or more of the following ways: deleting ratings; inserting alternative ratings and dates; and modifying rating dates"
- 2007, Paper claiming to de-anonymize Netflix data

# Case 2: Netflix Prize

- Narayanan and Shmatikov
  - "The adversary with a small amount of background knowledge about an individual ... can identify with high probability that individual's record in the data and learn ... sensitive attributes"
  - Claim Netflix's data sanitization not relevant
  - Basic Idea:
    - With aux info on 8 movies, where 2 can be wrong, and dates are known within 14 days, 99% de-anonymization
    - Aux info can come from other web-sites (IMDB), personal contact, etc.

# Case 2: Netflix Prize

- Much ado about nothing
  - Paper is technically correct, but dates are key
  - Without dates, you must know 8 movies, all outside the top 500 to get over 80% chance of de-anonymization
  - Aux info is not easy to come by for many people
  - No identities released
- Netflix did it right
  - Consulted with top machine learning experts
  - Gained new knowledge in machine learning and also privacy fields
- Fallout
  - Netflix was planning another challenge was canceled to due privacy concerns

# Case 3: Barabasi Mobile Study

Gonzalez, Hidalgo, and Barabasi (2008)

- Article in Nature outlines study on human mobility patterns
  - 100000 individuals selected randomly from dataset of 6 million
  - Unidentified country (unclear if researchers knew)
  - Cell tower location at start of call
  - 206 individuals were "pinged" every two hours for a week
- Findings
  - "humans follow simple, reproducible patterns"
  - Sample finding: Nearly three-quarters of those studied mainly stayed in a 20-mile circle for half a year.
  - Results "could impact all phenomena driven by human mobility, from epidemic prevention to emergency response and urban planning"

# Case 3: Barabasi Mobile Study

- Uproar over "secret tracking" of cell phone users
  - Blowback of negative feedback to Nature and scientists
  - Study would be "illegal in the US"
  - Approval from ONR review board and Northeastern review board. Barabasi did not check with an "ethics panel"

- Response
  - Hidalgo: "the data could be misused", but we were "not trying to do evil things. We are trying to make the world a little better."
  - Northeastern and Nature backed the research
  - Continues to be referenced as an example of dangerous research
  - Risk and reward both very high

# Data Mining and Ethics

- Privacy is not the only issue in data mining
  - Selling of Data
  - Transparency of Data Collection
  - Security
  - Discrimination and Bias

- We will explore these topics periodically throughout the semester

# Data Mining Introduction Summary

| Data Mining is interdisciplinary and overlaps significantly with many fields | Data Mining requires a team effort with members who have expertise in several areas |
|---|---|
| • Statistics<br>• CS (machine learning, AI, data bases)<br>• Optimization | • Data management<br>• Statistics<br>• Programming<br>• Communication<br>• + Application domain |