

CSC-321: Data Mining & Machine Learning

Prof. Nick Webb
webbn@union.edu

Data Mining: What's it all about?

- It's a buzzword!
 - Good for me: Get you interested?
 - Bad for me: Maybe it's not what you think?
- A practical definition:
 - Finding *useful* patterns in large amounts of data
 - What certain agencies call *actionable intelligence*

Data Mining

- Two 'modes' of thought
- Looking for patterns in large amount of data
 - These might not be wildly complex patterns
- Looking for patterns or relationships that are too hard to describe
 - These might only exist in small amounts of data

The data fire hose...



...is not new

- Many people used to be employed at Lloyds of London to scour newspapers for reports of ship wrecks

No. 24. **LLOYD'S LIST.** April 15, 1-27 p.m.
CASUALTY REPORT

T I T A N I C (a).

In reply to enquiry signal station at Cape Race cable:-
10-25 p.m. Titanic reports by wireless struck iceberg and calls for immediate assistance at 11 p.m. she reported sinking by head women being put off in boats gave position as 41.40 N. 50.14 W. Baltic Olympic and Virginian all making towards scene disaster latter was last to hear Titanic signals at 12-27 a.m. reported them then blurred and ending abruptly believed Virginian will be first ship to reach.

LLOYD'S LIST. Apr. 15, 2-34 p.m.
CASUALTY REPORT

T I T A N I C (a)

Re Tel. Co's telegram dated New York Apr. 15 states:-
less message received at Halifax at 4-30 this morning stated
out of the passengers from the TITANIC had been put in the life-
boats and that the sea was calm.

An Exchange Tel. Co's telegram dated New York Apr. 15 states:-
The White Star officials here state that the Virginian is standing
by the TITANIC and that there is no danger of loss of life.

Machine Learning: What's it all about?

- We could find these patterns by hand
- Indeed, we did find these patterns by hand!
- As data increases, so too does the time, and the complexity
- Maybe our intuition is a good place to start, but maybe machines can do better.

Machine Learning

- Automated (or semi-automated) way for machine to ‘discover’ patterns in data
- CAUTION: There is NO magic
- You should know what you’re looking for
- You should have an idea how to find it



NO MAGIC

Machine Learning

- Automated (or semi-automated) way for machine to ‘discover’ patterns in data
- CAUTION: There is NO magic
- You should (mostly) know what you’re looking for
- You should (mostly) have an idea how to find it

ML \neq AI

- Machine Learning and Artificial Intelligence are NOT the same
- ML – find patterns in data
- AI – give machines wisdom and intelligence
- ML will find an answer, but will not know if it makes sense
- AI will know if an answer is sensible

How is it used?

- Academic Examples
 - Speech: Recognition, Production
 - Image: Recognition, Emotions
 - Language: Translation, Analysis
 - Bioinformatics
 - Robotics
 - Financial markets
 - Sports
 - Politics
 - And so on and so on...

How is it used?

- Real world
 - Amazon
 - Netflix
 - Google
 - Intelligence Community
 - Supermarkets
 - Self driving cars
 - Dating websites
 - And so on and so on...

ML of the day

- <https://www.wired.com/story/machine-learning-march-madness/>



So what will you learn?

- What ML is:
 - Key algorithms
 - How we implement them
 - How we use them
 - The theory behind them

Key topics

- Basic algorithms
 - Linear regression
 - Naïve Bayes
 - KNN
 - Decision Trees
- How we evaluate algorithms
- Practical application
- Ethics and data

Skills you need

- Python
 - Implementation of algorithms
 - Acquisition and manipulation of data
- Statistics
 - Data description
 - Model evaluation

Learning Objectives

- Theoretical and practical aspects of machine learning as a tool for mining data
- Knowledge of methods and techniques for learning, application and evaluation
- Familiarity with open source machine learning tools

Assignments

- This is a developing course
 - TRANSLATION: I can and WILL change my mind
- Homework Assignments
- Midterm
- Challenge Problem
- Final Project
 - Literature survey
 - Final research paper

Homework

- Getting familiar with specific Python functionality
- Implementing algorithms
- Exercises using including data sets
- Explore the features of open source tool kits and start to work with real data

Challenge Problem

- Working on a data set I give you
- WHAT can you do to it, how well can you classify the data?
- A chance to get familiar with the data pipeline

Final Project

- Find a data set
- Perform experiments
- Create 8 page research paper
 - Intro, method, results, conclusion
 - Tell me about important attributes
 - What learning experiments you used

Course Policies

- There are ZERO late assignments accepted in this course
- Academic Integrity
 - DO NOT use code from ANYWHERE

Tools

- Python 3
- scipy / numpy / matplotlib / scikit-learn
- Jupyter notebooks
- Easiest way is to use anaconda
 - <https://www.anaconda.com/>