# BST 219
# Core Principles of Data Science

Lecture 30: Machine Learning continued, Data Science next steps
December 19, 2024

# Recipe of the Day!
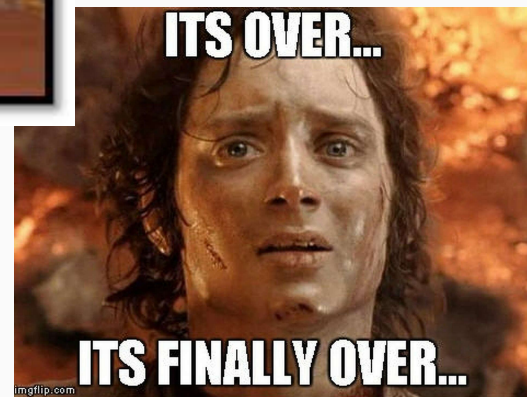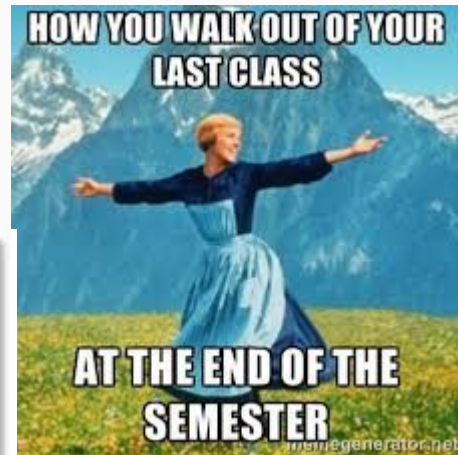
Holiday Cocktails

Holiday Mocktails

# Agenda

- Announcements
  - No lab this week!

  - No office hour on Thursday

  - Please complete the course evaluation -  we value your feedback!

- Continue Machine Learning module
  - Regularization
    - Nice LASSO tutorial video
  - Cross-validation explanation

- Next steps in Data Science

# Course Highlight Reel
## What did we learn?
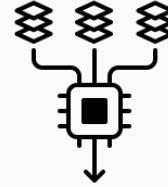
# Course Roadmap
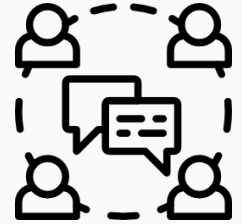


**Importing (loading) the data**

**Processing (cleaning, wrangling) the data**

**Visualizing and summarizing the data**

**Building models (statistical and ML)**

**Interpretation and communication of results**

# Topics we covered

| Module 1<br>R Basics | Module 2<br>Git and GitHub | Module 3<br>Data Visualization | Module 4<br>Advanced Data Wrangling | Module 5<br>Machine Learning |
|---|---|---|---|---|
| -R Markdown<br><br>-Data types<br><br>-Vectors<br><br>-Sorting<br><br>-Vector arithmetic<br><br>-Indexing<br><br>-Data wrangling<br><br>-Importing data<br><br>-Functions<br><br>-For loops<br><br>-If else statements | -Cloning repositories via RStudio<br><br>-Committing and pushing assignments via RStudio<br><br>-Pulling course notes via RStudio | -ggplot2<br><br>-Data visualization principles<br>-Faceting<br><br>-Fixing scales<br><br>-Time series plots<br><br>-Transformations<br><br>-Ordering by a value<br><br>-Making comparisons<br><br>-Maps | -Tidy data format<br><br>-Importing data<br><br>-Reshaping data<br><br>-Join functions and combining tables<br><br>-Dates and times<br><br>-TableOne | -Fundamentals of ML including train/test split and the process<br><br>-Logistic regression, Naive Bayes, kNN, QDA, LDA, Decision trees, Random Forests, LASSO, Ridge regression, PCA<br><br>-ROC and AUROC<br><br>-Confusion matrix<br><br>-Performance metrics |

Data Science Next Steps

# Spring Courses

| BST 263 - Statistical Learning (Full Spring) | BST 261 - Data Science II (Deep Learning, Spring 2) | BST 221 - Applied Data Structures and Algorithms (Full Spring) |
|---|---|---|
| Language: R or Python<br>● Material: An Introduction to Statistical Learning<br>● Additional Reading: The Elements of Statistical Learning<br><br>Probability Basics<br>Assessing Model Accuracy<br>Linear and Logistic Regression<br>Classification<br>Cross-validation and bootstrap<br>Subset selection<br>Penalty-based methods<br>Ridge regression and LASSO<br>Dimension Reduction<br>Polynomial Regression<br>Step functions and Basis functions<br>Generalized Additive Models (GAMs)<br>Classification and Regression Trees (CART)<br>Bagging and Random Forests<br>Ensembles<br>Support Vector Machines<br>Unsupervised Learning<br>Bayesian Methods | Language: Python<br>● Material: Deep Learning, Deep Learning with Python<br><br>Brief review of Python 3<br>Brief review of Linear Algebra and Probability<br>Brief review of Machine Learning<br>Feedforward Neural Networks<br>Convolutional Neural Networks<br>Recurrent Neural Networks<br>Generative Adversarial Networks (GANs)<br>Reinforcement Learning<br>Transfer Learning<br>Hyperparameter Tuning<br>Model Selection | Language: Python<br><br>Overview of numerical analysis material<br>Concepts of Algorithms, Complexity and Sorting Algorithms<br>Data Structures and Heapsort<br>Parallel Programming<br>Greedy Algorithms and Dynamic Programming<br>Numerical Stability<br>(Pseudo) Random Number Generation<br>Efficient Algorithms for Linear Algebra<br>Least-Squares Problem, Eigenvalue Decompositions<br>Algorithms for Numerical Integration, MC-Integration, Importance Sampling<br>Graphs and Network Algorithms |

# Other Courses

- BST 262: Computing for Big Data (only offered in Fall 2)

- BST 267: Introduction to Social and Biological Networks (only offered in Fall 2)

- APMTH 120: Applied Linear Algebra and Big Data (offered in Spring, main campus)

- Other machine learning and NLP courses at Harvard and MIT
  - Caution: these tend to be VERY difficult and a ton of work

# Skills Development

- Keep using R for research or other projects!

- Learn Linear and Matrix Algebra

- Diversify your Programming Languages (Python, SQL, C++, Java, etc.)

- Kaggle competitions

- Online courses on Coursera, edX, Data Camp, etc.

- Data Science Meetups in Boston

- Data science podcasts and social media

- Data science blogs

- Data science tutorials and videos

# Become a TF!

We are always in need of teaching fellows for our Biostats courses, including this one!