Update anything and share with people in 474
Student Project Presentations:

| Group Number and Project | Algorithm used | Notes |
|---|---|---|
| 9 - Price prediction of used cars | Linear Regression | Kijiji -DM-> their interface -> $ price of car<br>SQL db<br>Clean data - remove noise - add 0's |
| 15 - Github issue label prediction | Text classification<br>Clustering | 5 main types of issues |
| 4 - Violence in the News | Random forest | Categorize appropriate for kids<br>Vectorizors<br>Tried: bayes, trees, regression |
| 1 - NBA Expected shots<br><br>GitHub repo:<br>https://github.com/erikreppel/seng474-nba-shots<br><br>Presentation slides:<br>https://goo.gl/WVwciJ | - Multilayer perceptron (final decision)<br>- Linear regression, naive bayes (bernoulli and gaussian), AdaBoost all tested | - Classification > regression<br>- Calculating expected value of shot, i.e. 75% chance to hit a 3 pointer, 0.75 * 3 = 2.25 expected value. |
| 20 - NHL Points Projection | Clustering<br>Euclidean distance<br>Neural network | Age regression - in prime?<br>Accuracy - 91.6% - 43.4%<br>Good for average - bad for outliers<br>Neural net was having major problems giving stars a 0 prediction |
| 12 - Predictive Policing (Crime Prediction) | Lots used:<br>Gaussian NB, decision tree, multinomial NB, bernoulli np, perceptron, logistic regression, SVM | Government data site: from Victoria - 2006-2016 categorizing |
| 13 - MLB MVP | Multiple Regression models:<br>Linear regression, ridge, lasso<br><br>Proposal started with ID3, had to change due to working with continuous values not discreet. | Scraped data from 1961-2016<br>Used player stats for position players and pitchers and evaluated both separately using multiple regression models to create a weighted voting scheme, and apply votes to players for stats they achieved.<br><br>Tried to mimic BBWAA voting mathematically. |

| | | |
|---|---|---|
| 14 - Music Popularity ("hotttness" prediction) | Classification = good<br><br>Logistic regression, linear regression, support vector machines(support vector regression, RBD vs polynomial kernel, coefficient $r^2=0.34$), gaussian naive bayes | Million song db<br>12 statistics<br>4 algs |
| 6 - Edible Mushrooms (we focused on precision for TP having low FP) | Linear Svc (linear svm)<br>Svm<br>Gaussian Naive bayes | -Very accurate 95% or more<br>-Over fit/100%<br>-Did well 85%-70% |
| 10 - Oscars | | Used 2 to predict top 10 then took intersection<br>Social Media |
| 7 - Bike share load balancing | Decision tree | Predicted when bike station will be empty, full or none at any hour of day<br>Used for days in advance<br>Allows for<br>Bay Area Bike Share data for San Francisco |
| 11 - Trump Tweets | Neural Network | Fooled people 17% of time<br>Had website for guessing real trump tweet or theirs |
| 2 - Academy Awards Prediction | | |
| 3 - Soccer Premier League Prediction | Regression Models | Didn't handle relegation data well, chopped all the new teams data. |
| 5 - Sentiment Analysis | Vectorizer gsvc | Yelp texts review -> /5 stars<br>Maps were bad |
| 17 - Eve Online Market Prediction | Linear Regression, Support Vector Regression, AdaBoost Regression | Calculated momentum attribute from simple market data; found a high return of 1.5x (50%) on initial over 3.5 months of trading. |

Guest Lectures:

| | | |
|---|---|---|
| Daniel German | | Demonstration using emacs and R on the Iris dataset |
| Brian Ziebart | a supervised machine learning framework that adversarially approximates the training data and uses the exact performance measure | Supervised machine learning as an adversarial game. Provides flexibility for addressing sample selection bias and for inductively Optimizing multivar performance measures like F-measure discounted Cumulative gain from information retrieval & ranking tasks. |
| George Tzanetakis | Markov Logic Networks | Automatically classifying the kind of music by sample. Tagging a song > bags of words > |
| David Johnson | Support Vector Machine | Kinect image processing edge detection in piano learners - computer vision |

Pros and Cons of various models:

| Model | When to use | When not to use |
|---|---|---|
| Decision tree | Discrete data | Continuous data |
| Regression tree | Discrete values<br>Continuous prediction | Anything else |
| Naïve bayes | Text classification<br>Few parameters, large data set | Lots of parameters |
| HMM | Sequence or temporal models | Anything else |
| Logistic regression | Lots of parameters<br>Continuous values | Few parameters<br>Discrete |
| K-means | Clustering<br>Glob forms (spheres)<br>Know how many clusters | Differing sizes, densities, shapes |
| Bisecting K-means | Clustering<br>Don't know how many clusters spherical clusters | Assignment 3 Q3 |
| Linear Regression | Continuous data and prediction | Discrete, lots of parameters |
| Neural Networks | Anything graphical<br>Human brain replication | Very complicated<br>Blackbox<br>Slow to make and run |
| SVM | Linearly separable data | Non-binary set of data |
| Min/Max/Avg | Min: non-globular shapes<br>Max: Good with noise | Min: Sensitive to noise and outliers<br>Max: Tends to break large clusters<br>Avg: Slow – big calculations |
| Euclidean distance | Finding similarities | Different scales |
| Pearson correlation | Finding similarities with different scales | Doesn't handle outliers well |
| | | |