# DATA SCIENCE
## 10 WEEK PART TIME COURSE

## Week 6 – Decision Trees
## Tuesday 27th June 2017

1. What are decision trees?
2. How decision trees work
3. Visual example
4. Lab
5. Discussion

# DECISION TREES

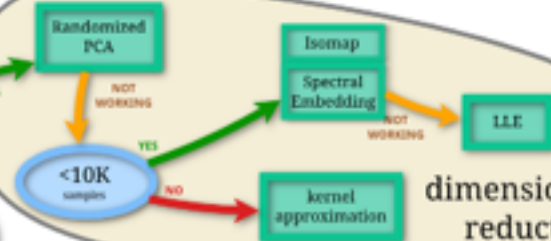scikit-learn
algorithm cheat-sheet

**START**

**classification**

- kernel approximation
- SVC
- Ensemble Classifiers
- KNeighbors Classifier
- SGD Classifier
- Naive Bayes
- Text Data
- Linear SVC
- <100K samples

**regression**

- SGD Regressor
- Lasso ElasticNet
- SVR(kernel='rbf') EnsembleRegressors
- <100K samples
- few features should be important
- RidgeRegression SVR(kernel='linear')

**clustering**

- Spectral Clustering
- GMM
- KMeans
- number of categories known
- <10K samples
- MiniBatch KMeans
- MeanShift
- VBGMM

**dimensionality reduction**

- Randomized PCA
- Isomap
- Spectral Embedding
- LLE
- <10K samples
- kernel approximation

- get more data
- >50 samples
- predicting a category
- do you have labeled data
- predicting a quantity
- just looking
- predicting structure
- tough luck

Back

scikit learn

‣ A supervised learning technique that can be used for classification or regression.

‣ A supervised learning technique that can be used for classification or regression.

‣ Visually engaging and easy to interpret.

‣ A supervised learning technique that can be used for classification or regression.

‣ Visually engaging and easy to interpret.

‣ Foundation for getting into very powerful techniques.

- A supervised learning technique that can be used for classification or regression.

- Visually engaging and easy to interpret.

- Foundation for getting into very powerful techniques.

- Great for explaining to people!
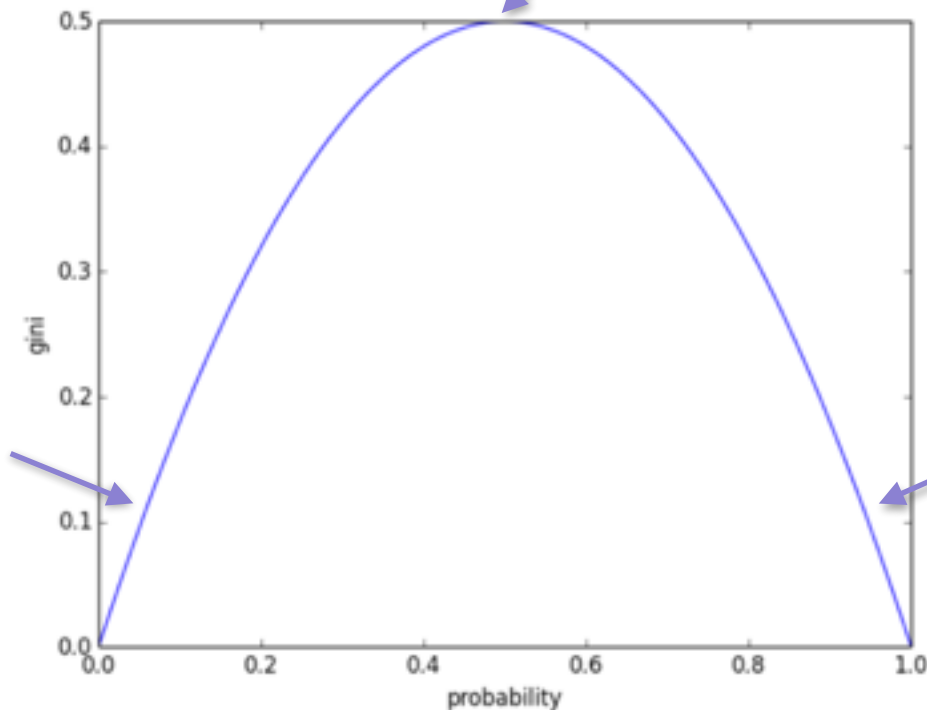
‣ Prone to overfitting.

‣ Prone to overfitting.

‣ Predictive power is lower in comparison to many other modern techniques.

‣ Scans for a feature to split on that results in the greatest separation between classes in the resulting nodes.
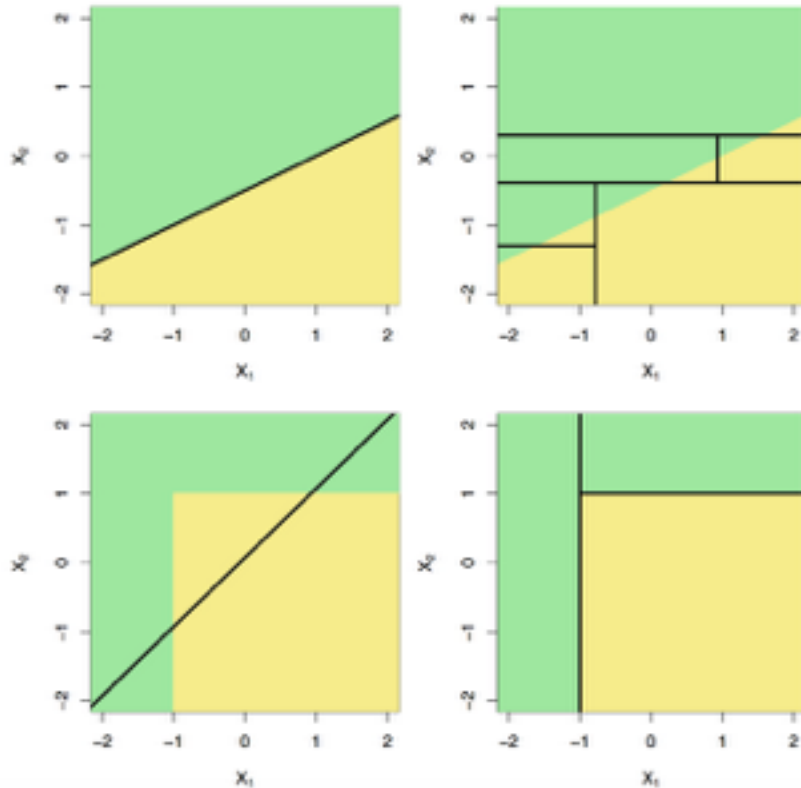
**The Gini Index**



Equal ratio of
target classes
50:50

High purity
of class 0

High purity
of class 1

‣ Scans for a feature to split on that results in the greatest separation between classes in the resulting nodes.

‣ Non-linear.

Linear decision boundary

Non-linear decision boundary

‣ Scans for a feature to split on that results in the greatest separation between classes in the resulting nodes.

‣ Non-linear

‣ Greedy process

‣ Splits within splits

‣ Scans for a feature to split on that results in the greatest separation between classes in the resulting nodes.

‣ Non-linear

‣ Greedy process

‣ Splits within splits

‣ For a classification tree, we predict that each observation belongs to the most commonly occurring class of training observations in the region to which it belongs.

‣ Scans for a feature to split on that results in the greatest separation between classes in the resulting nodes.

‣ Non-linear

‣ Greedy process

‣ Splits within splits

‣ For a classification tree, we predict that each observation belongs to the most commonly occurring class of training observations in the region to which it belongs.

‣ We naturally get combinations of features used for our prediction.

http://www.r2d3.us/visual-intro-to-machine-learning-part-1/

# LAB

1.  re-name your labs with lab_name.<yourname>.ipynb  (to prevent a conflict)
2.  cd <path to the root of your SYD_DAT_8 local repo>
3.  commit your changes ahead of sync
    - git status
    - git add .
    - git commit -m "descriptive label for the commit"
    - git status
4.  download new material from official course repo (upstream) and merge it
    - git checkout master  (ensures you are in the master branch)
    - git fetch upstream
    - git merge upstream/master

# DISCUSSION TIME

- ‣ Review of last week
- ‣ Further Reading for Decision Trees
- ‣ Check in with homework/course project
- ‣ Pre-Reading

# DISCUSSION TIME

‣ **Recommendations**

‣ **SQL & Productivity Tools**
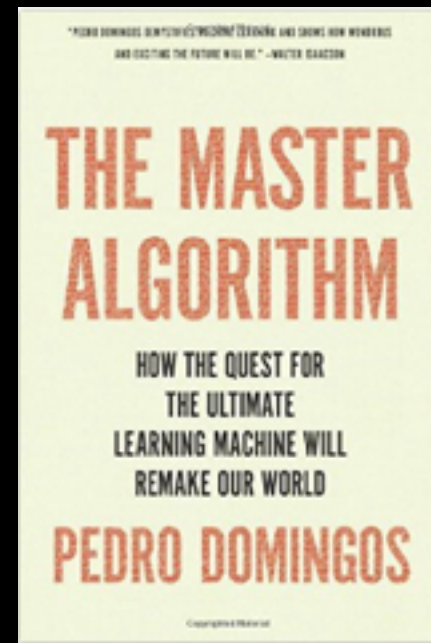
‣ **Linear Regression**

# DISCUSSION TIME

**CIA using Decision Trees**

‣ https://www.cia.gov/library/center-for-the-study-of-intelligence/kent-csi/
vol18no4/html/v18i4a03p_0001.htm

**Overview of difference approaches to building models**

‣ https://www.amazon.com/Master-Algorithm-Ultimate-Learning-Machine/dp/
0465065708

# DISCUSSION TIME

**Homework/Course Project**

‣ **How's Homework 2 going ?**

‣ **How are the projects going?**

# PRE-READING

**An Introduction to Statistical Learning**

‣ **Chapter 8 – Tree-Based Methods**

**Signup to Google Cloud Platform**

‣ **https://cloud.google.com/**