# DATA SCIENCE
## 10 WEEK PART TIME COURSE

## Week 3 – Logistic Regression
## Tuesday 6th June 2017

1. Motivation
2. What is Logistic Regression?
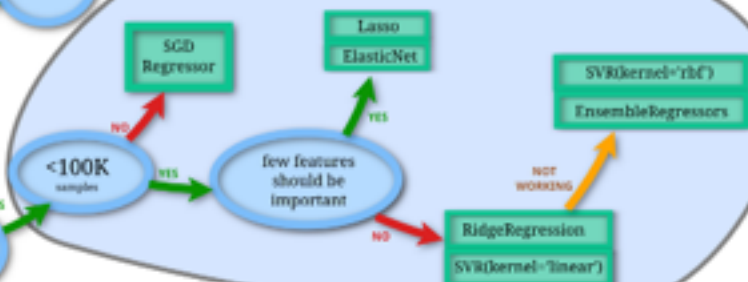3. Evaluating Logistic Regression
4. Lab
5. Homework Review
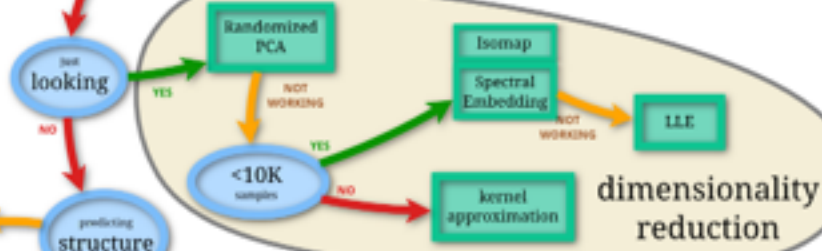
scikit-learn algorithm cheat-sheet

**classification**

kernel approximation

SVC

Ensemble Classifiers

KNeighbors Classifier

SGD Classifier

Naive Bayes

Text Data

Linear SVC

<100K samples

**START**

get more data

>50 samples

predicting a category

do you have labeled data

predicting a quantity

just looking

predicting structure

tough luck

**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

<10K samples

MiniBatch KMeans

MeanShift

VBGMM

**dimensionality reduction**

Randomized PCA

Isomap

Spectral Embedding

LLE

<10K samples

kernel approximation

Back

scikit learn

If the y variable is numeric then we have a regression problem - we are trying to predict a continuous number

If the y variable is a category (for example trying to predict a type of flower) the we have a classification problem - we are trying to classify what group that y belongs to.

# WHAT IS LOGISTIC REGRESSION?

We want to build a classifier that correctly identifies which class our target variable y belongs to given our input variable x.
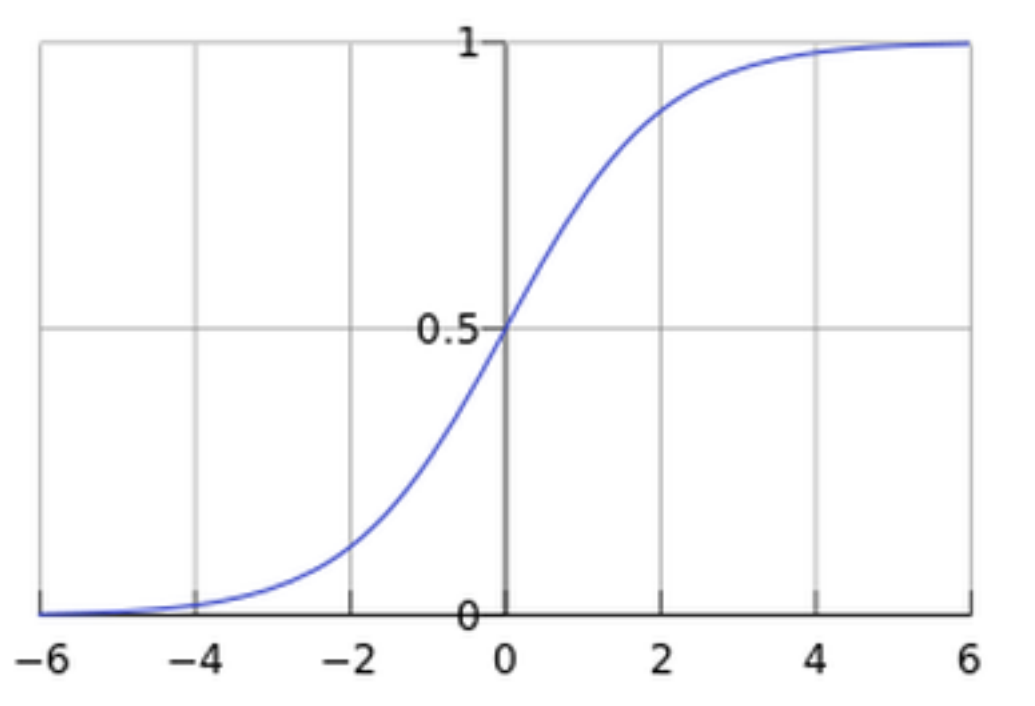
Why not use the linear regression model?

$$y = X\beta + \epsilon$$

‣ If we only have a binary response variable (0 or 1) it might make sense… BUT we can have our estimated value of y > 1 or y < 0 … which doesn't make sense.

‣ What of the case where we have more than one class? Linear regression cannot easily handle these cases.

‣ We want a classification method that can handle these cases and give us results we can easily interpret.

$$p(Y=1|X) = \beta_0 + \beta_1 X.$$

‣ This is a good starting point but we still have the problem of $p(Y)$ being outside the 0,1 range.

‣ We need to model $p(Y=1|X)$ using a function that gives outputs between 0 and 1.

‣ Basically we want something that looks like the following

Standard logistic sigmoid function

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x$$

- This is the logit function,
- We can see that it this function is linear in X
- $\frac{p}{1-p}$ is called the 'odds' and can be any value from 0 to $\infty$
- $\log\left(\frac{p}{1-p}\right)$ is called the 'log-odds' or 'logit'

# EVALUATING LOGISTIC REGRESSION

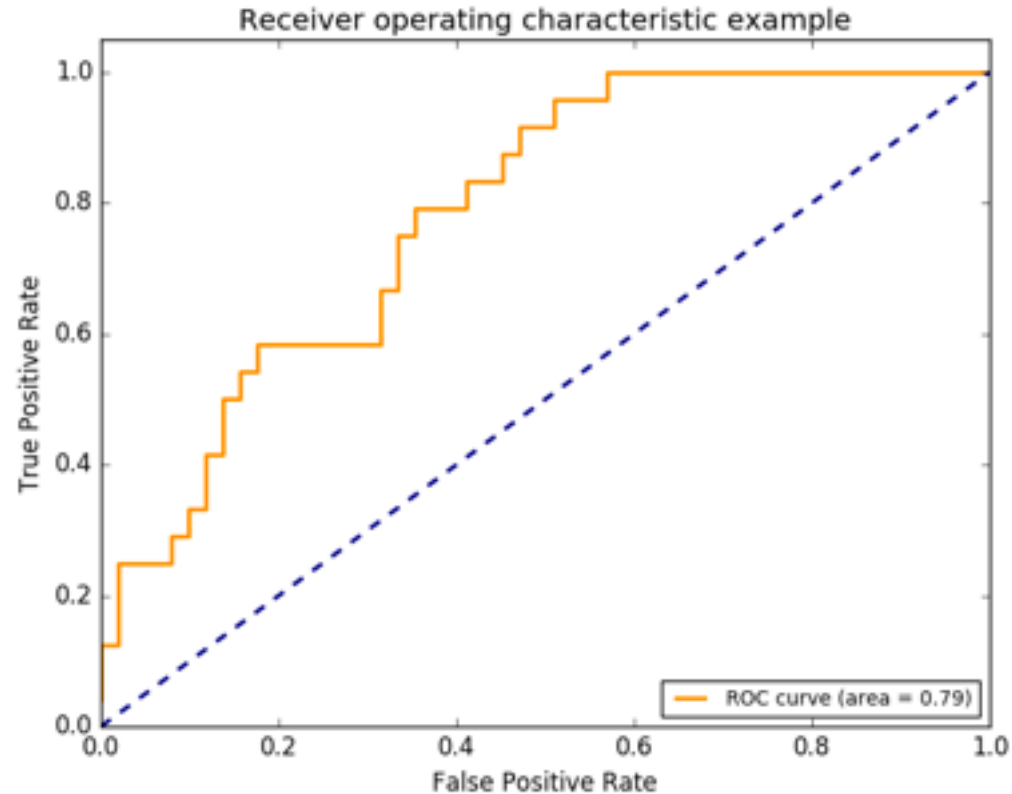This is simply the fraction of correct predictions from the model.

So it is the number of correct predictions divided by the number of observations in our dataset.

A confusion matrix shows us what the predicted class was against what the actual class was. The true class makes up the rows or the vertical axes and the predicted class makes up the columns or the horizontal axis.

Any entries in the diagonal of the matrix are those that are correctly classified.

|  |  | Predicted class | |
|---|---|---|---|
|  |  | P | N |
| Actual Class | P | True Positives (TP) | False Negatives (FN) |
|  | N | False Positives (FP) | True Negatives (TN) |

The Receiver Operating Characteristic or ROC curve shows the performance of a binary classifier system as its discrimination threshold is varied. It is created by plotting the fraction of true positives out of the positives (TPR = true positive rate) vs. the fraction of false positives out of the negatives (FPR = false positive rate), at various threshold settings.

Receiver operating characteristic example

By computing the Area Under the Curve of the ROC curve we get a single number summary of accuracy.

The closer that number is to 1 the more accurate our model is.

‣ We will step through a notebook together and cover these concepts in a more tangible way.

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 Logistic Regression
- ‣ Check in with homework/course project

# DISCUSSION TIME

- ‣ Visualisation
- ‣ Supervised vs Unsupervised Learning
- ‣ Linear Regression

# DISCUSSION TIME

**Logistic Regression applied to loan applications**

‣ **https://github.com/nborwankar/LearnDataScience**

**Odds Ratio in Logistic Regression**

‣ **http://www.ats.ucla.edu/stat/mult_pkg/faq/general/odds_ratio.htm**

# DISCUSSION TIME

**Homework/Course Project**

‣ **How's Homework 1 going ?**

‣ **How are the projects going?**

# PRE-READING

**An Introduction to Statistical Learning**

▸ **Chapter 5 – Resampling Methods**