

# **DATA SCIENCE**

**11 WEEK PART TIME COURSE**

**Week 3 – Logistic Regression**  
**Wednesday 6th April 2016**

1. Motivation
2. What is Logistic Regression?
3. Why use Logistic Regression
4. Lab
5. Homework Review

# scikit-learn algorithm cheat-sheet

START

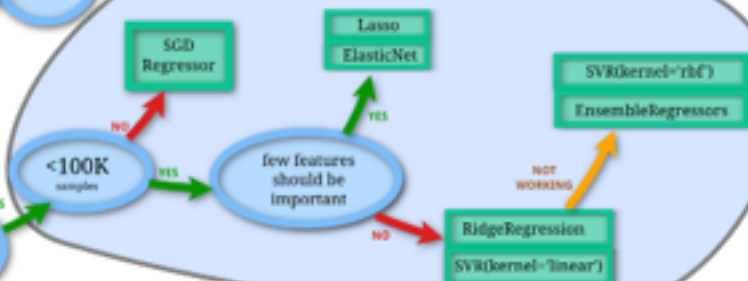
## classification



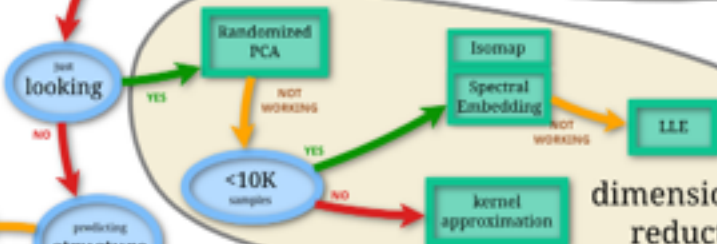
## clustering



## regression



## dimensionality reduction



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) then we have a classification problem - we are trying to classify what group that  $y$  belongs to.

**DATA SCIENCE PART TIME COURSE**

---

# **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?

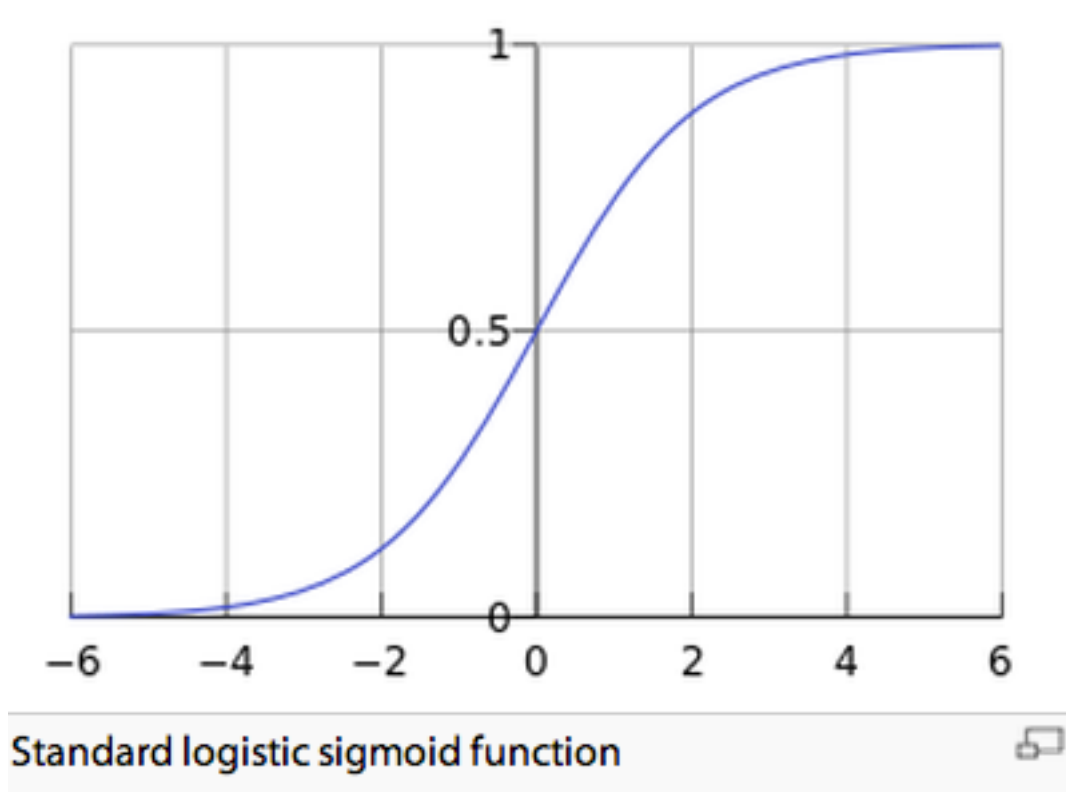
$$y = X\beta + \epsilon$$

- 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





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

- 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’

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

---

**DATA SCIENCE PART TIME COURSE**

---

A young man and woman are riding a roller coaster. The woman is in the foreground, wearing a dark jacket over a white lace top, with her arms outstretched and a joyful expression. The man is behind her, also with his arms outstretched. They are both looking towards the right. The background shows a sunset over a body of water, with a warm orange and yellow glow. The word "LAB" is overlaid in large white letters on the right side of the image.

**LAB**

# **DISCUSSION TIME**

- **Review of last week**
- **Further Reading for Logistic Regression**
- **Check in with homework/course project**

## WEEK 3

Wednesday 16<sup>th</sup> December

- Classification ☒ Understand Supervised VS. Unsupervised Learning
- Regression (numeric) ☒ Describe process of Linear Regression
- variables ☒ Build a Linear Regression Model
- 1 tests ☒ List of Resources to Review

# DISCUSSION TIME

## **An Introduction to Statistical Learning**

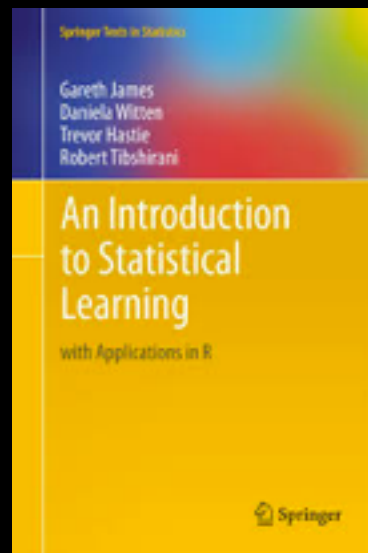
### ‣ **Chapter 4 – Logistic Regression**

## **Logistic Regression applied to loan applications**

### ‣ **<https://github.com/nborwankar/LearnDataScience>**

## **Odds Ratio in Logistic Regression**

### ‣ **[http://www.ats.ucla.edu/stat/mult\\_pkq/faq/general/odds\\_ratio.htm](http://www.ats.ucla.edu/stat/mult_pkq/faq/general/odds_ratio.htm)**



---

**DATA SCIENCE - Week 3 Day 2**

---

# **DISCUSSION TIME**

**Homework/Course Project**

‣ **How's Homework 1 going ?**