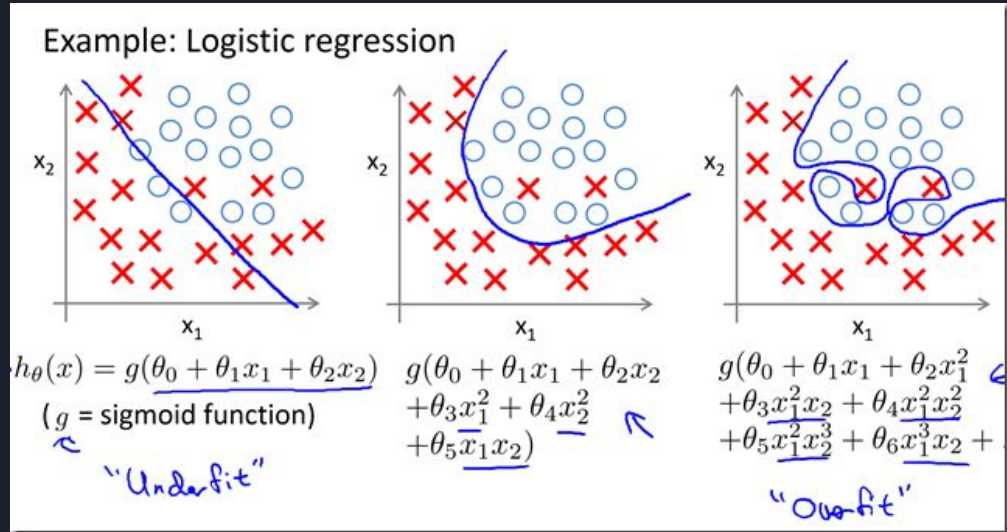# Decision Trees & Random Forests

Data Science Curriculum

# Last week we learnt..

- 3x Lessons - Logistic Regression for Classification tasks
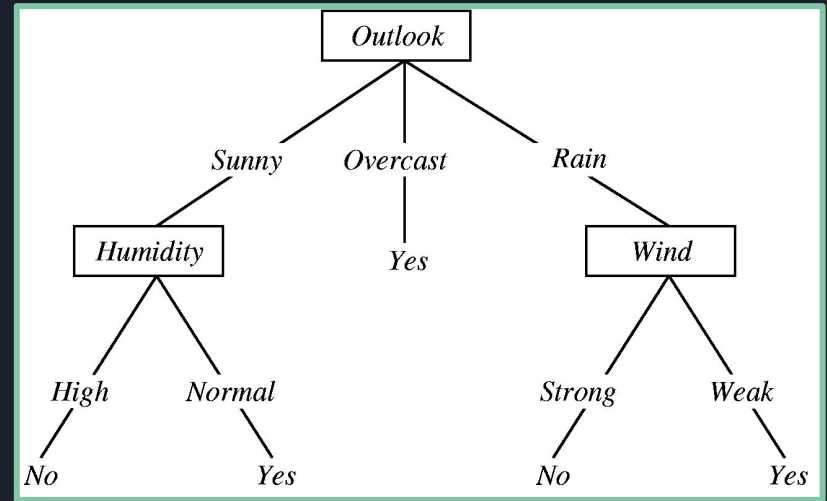  - Classifies 2 class problems, e.g. will respond to marketing campaign or will not respond

- Advantages?

- Disadvantages?



Example: Logistic regression

$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$

($g$ = sigmoid function)

"Underfit"

$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$

$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \ldots$

"Overfit"

# Lesson Objectives

1. Understand how a decision tree functions and how to build one
2. Learn to visually represent a decision tree
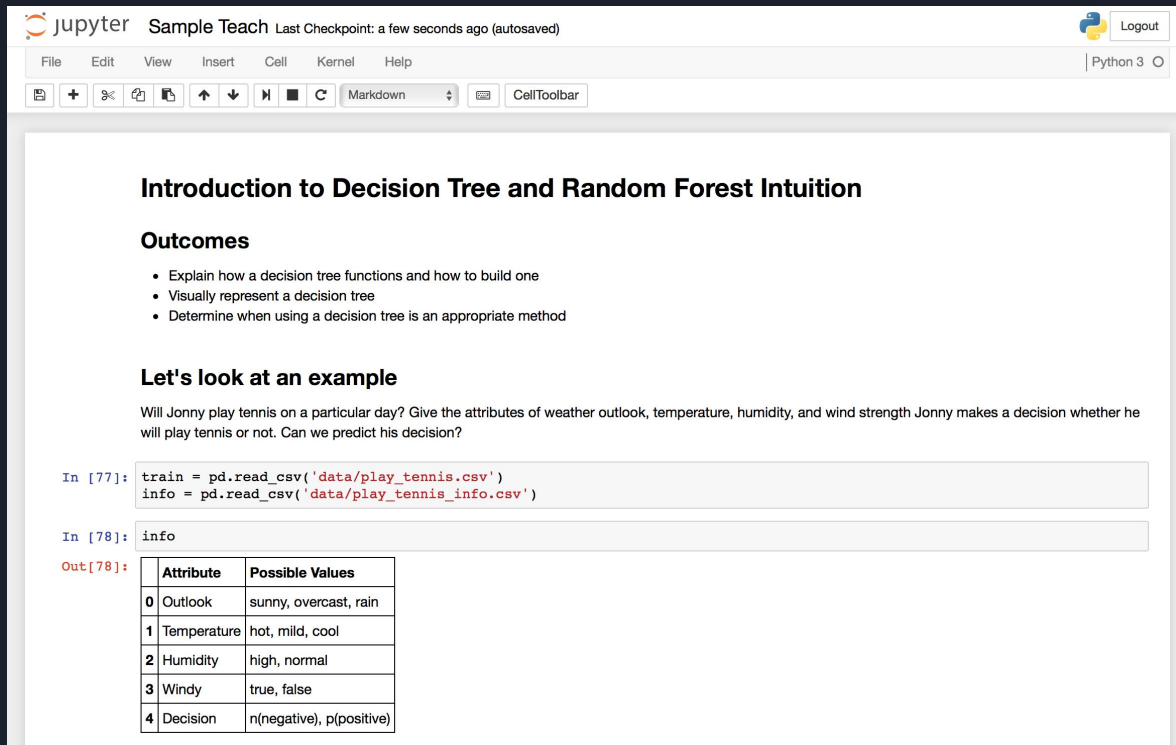3. Understand when to use a decision tree

# Agenda

1. To Play or Not to Play…. Tennis Example of decision trees

2. Worked example of Iris Decision Trees in Python using Scikit-Learn

3. When to use decision trees

4. Intro to Random Forests

# Jupyter Notebook

# The End

By Andrew Szwec

# Split data on column labels then calculate number of outcomes

https://youtu.be/eKD5gxPPeY0

# Decision Trees for Classification

- Decision trees for classification - classes like dog/cat, will purchase/won't purchase
- Decision trees for regression - predicts a value like house price
- Make a decision tree for taxi dataset
- Visualise it
- Now use a toy dataset so they can do it themselves

Prune tree using validation set. Take away each node separately and see how performance improve/degrades against the validation set

# Decision Tree