# Machine Learning

# Dr Tamara Clelford

To download a copy of the slides go to https://tamaraclelford.co.uk/iop.html



# INTRODUCTION TO DATA SCIENCE

• What is data science

• Examples from my work

Data science workflow walkthrough



#### IS DATA SCIENCE: "THE SEXIEST JOB OF THE 21<sup>ST</sup> CENTURY"?

 In 2012 the Harvard Business Review described data science as being 'The sexiest job of the 21<sup>st</sup> Century'

• 'Data will be the raw material of the 21<sup>st</sup> Century' Angela Merkel, German Chancellor, Davos 2018



# WHAT IS THIS THING CALLED DATA SCIENCE?





#### WIKIPEDIA The Free Encyclopedia

Main page Contents Featured content Current events Random article Article Talk

#### Data science

From Wikipedia, the free encyclopedia

Not to be confused with information science.

**Data science** is a multi-disciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from data in various forms, both structured and unstructured,<sup>[1][2]</sup> similar to data mining.





Dr TAMARA CLELFORD

Read

Edit

💄 Not I

View

### WHAT IS THIS THING CALLED DATA SCIENCE?

- Any data can be looked at, multiple data sets used
  - Structured, unstructured, numbers, words, clicks, actions, images
- Multi-disciplinary
  - Part physics, statistics, maths, computer science, coding, big data, business
- Clean data before use
  - Removes incorrect and partial data
- Uses scientific methods and algorithms
  - Coding, decision trees, machine learning, neural networks
- Extracts knowledge and insights from data
  - The output is verified as being plausible

#### THE VENN DIAGRAM OF DATA SCIENCE













https://www.needpix.com/photo/920796/mountain-telescopehawaii-summit-astronomy-astrophysics-mauna-kea-kecktelescope-subaru-telescope Dr TA



#### HOUSE OCCUPANCY CALCULATIONS

House Occupancy (red=asleep, green=in, orange=out)



SWAMPHEN ENTERPRISES

### COMMERCIAL EXAMPLES OF DATA SCIENCE

- Business analytics
- Recommender engines
- Online advertisements
- Financial predictions
- Internet of Things
- Diagnostics
- Kaggle competitions



### THE RISE OF DATA SCIENCE

- As more data is collected in all walks of life the need to analyse this data and find it's hidden secrets becomes increasingly important
- Blurred line between engineer, statistician, data analyst, data scientist, data engineer, business analyst
- People from different backgrounds are re-training and becoming data scientists



- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results



SWAMPHEN ENTERPRISE

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
  - 5. Need more data? -
- → 6. Clean data set
  - 7. Setup model
  - 8. Evaluate model
  - 9. Calculate results
  - 10. Present results

SWAMPHEN ENTERPRISE

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results



SWAMPHEN ENTERPRISE

#### SURVIVAL OF THE BEST DATA



By F.G.O. Stuart (1843-1923) - http://www.uwants.com/viewthread.php?tid=3817223&extra=page%3D1, Public Domain, https://commons.wikimedia.org/w/index.php?curid=2990792

SWAMPHEN ENTERPRISES

**SWAMPHE** 

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results



## 2) SET UP DATA SCIENCE TOOLS

Python programming language
Toolboxes such as NumPy, pandas and scikit-learn

Integrated Development Environment (IDE)
Jupyter notebook style
Google Colab



# GO TO GOOGLE COLAB https://colab.research.google.com

	Google	google colab							<b>پ</b> Q	
		All	News	Images	Videos	Books	More	Settings	Tools	
		About	40,100,000	) results (0.34	4 seconds)					
		Hello, Colaboratory - Colaboratory - Google								
		https://colab.research.google.com/ ▼								
		To load a specific notebook from github, append the github path to http://colab.								
		research.google.com/github/. For example to load this notebook in Colab:								

You've visited this page 2 times. Last visit: 15/03/19

#### Welcome To Colaboratory

Colaboratory is a free Jupyter notebook environment that ...

#### Colaboratory – Google

SWAMPHEN ENTERPRISES

What is Colaboratory? Colaboratory is a research tool ...

#### Overview of Colaboratory

Colaboratory is built on top of Jupyter Notebook. Below are ...

#### Run in Google Colab

View on TensorFlow.org, Run in Google Colab, View source on ...

#### OPEN A NEW PYTHON 3 NOTEBOOK



# 2) SETUP DATA SCIENCE TOOLS

- *# open up google colab*
- 2 # open a python 3 notebook and rename to 'Titanic\_analysis'
- *# import packages going to use*
- **import** numpy **as** np
- **import** pandas **as** pd
- **import** matplotlib.pyplot **as** plt



**SWAMPHE** 

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results



# 3) GET DATA

1 # import the data to google colab

- 2 from google.colab import files
- 3 uploaded = files.upload()



**SWAMPHE** 

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results



### 4) LOOK AT YOUR DATA SET

- Important step to get an understanding of the data
  - If you don't understand your data set it will be hard to get insights from it

• Check that is has been read in correctly

Look at raw data and plot things to help understanding



### TITANIC DATA INVESTIGATION

1 # how many rows in the data set 2 len(raw\_data)

1 # import data set
2 raw\_data = pd.read\_csv('titanic.csv')

```
1 # print out the headder
2 raw_data.head()
3 # try putting a number in the brackets
4 # what does .tail() do?
```

1 # get basic shape information

2 raw\_data.shape

# print out a selection of the data
raw\_data[20:25][['Name']]

- 1 # sort wrt one column
- 2 raw\_data.sort\_values('Fare')
- 1 # calculate average fare
- np.mean(raw\_data['Fare'])



```
1 # describe data set - nothing with letters in the result
2 raw_data.describe()
```

```
1 # get the range of ages on board
2 age_female = raw_data.loc[raw_data['Sex'] == 'female', 'Age']
3 age_male = raw_data.loc[raw_data['Sex'] == 'male', 'Age']
4 print(age_female)
5 print(age_male)
```

```
1 # Look at the range of ages on board, often scatter graph, but not in this case
2 fig, graph = plt.subplots()
3 graph.hist([age_female, age_male])
```

```
1 # graph options
2 fig, graph = plt.subplots()
3 graph.hist([age_female, age_male], color = ['b','g'], label = ['female', 'male'])
4 # pick any colours to use from b, g, k, r, c, m, y, w
5 # print out the axis labels
6 graph.set_ylabel("number of people")
7 graph.set_xlabel("age (y)")
8 # add the legend
9 plt.legend(loc = 'best')
10 # move the legend around, options are:
11 # upper left, upper right, lower left, lower right, best, upper center, lower center,
12 # center left, center right
```

SWAMPHEN ENTERPRISES

```
Dr TAMARA CLELFORD
```

**SWAMPHE** 

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results



```
5) NEED MORE DATA?
```

```
1 # we want to be able to include gender information
2 # change female = 1, male = 0
3 raw_data['GenderNumerical'] = raw_data.apply(lambda row:1 if row.Sex == 'female' else 0, axis = 1)
4 raw_data.head()
```

```
1 # describe the data set again
2 raw_data.describe()
```



**SWAMPHE** 

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results



# 6) CLEAN DATA SET

- This is where most of your time is spent as a data scientist
- Need to ensure you do not introduce bias or errors to the data set during cleaning
- Examples of cleaning practices:
  - Remove lines with missing data
  - Use an average to fill missing data
  - Use a different data set to fill missing data
  - Interpolate between points

#### 6) CLEAN DATA SET

1 raw\_data.head(7)

1 # pick data want to include 2 include = ['Survived', 'Pclass', 'Age', 'GenderNumerical', 'Fare']

1 # check include 2 raw\_data[include]

1 # remove nans

2 data\_set = raw\_data[include].dropna()

```
1 # check removal of nans
```

- 2 data\_set.head(7)
- 3 data\_set[-5:]

SWAMPHEN ENTERPRISES

**SWAMPHE** 

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results



# MACHINE LEARNING AND ARTIFICIAL INTELLIGENCE

- Al is getting computers to behave like humans
- ML is an application of AI
  - gives the ML code access to data and lets it learn for itself
- Two categories of ML:
  - Supervised
    - decision trees, random forests, linear regression, Gaussian regression, Bayesian statistics
  - Unsupervised

SWAMPHEN ENTERPRISES

 clustering, k-means, anomaly detection, neural networks, method of moments

TAMARA CI FI FO



SWAMPHEN ENTERPRISES

#### 7) SET UP MODEL



1 *# import decision tree* 

2 **from** sklearn **import** tree

#### 1 *# set up the model*

2 tree\_model = tree.DecisionTreeClassifier(random\_state = 42)

1 *# run the model on our data* 

```
2 titanic_tree = tree_model.fit(data, answer)
```

```
3 print(titanic_tree)
```

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results





#### 8) EVALUATE MODEL

- 1 *# see how good model is at predicting survival*
- 2 **from** sklearn.model\_selection **import** cross\_val\_score **as** CVS
- 3 scores\_tree = CVS(titanic\_tree, data, answer)
- 4 print(scores\_tree)
- 5 print(np.mean(scores\_tree))



### DECISION TREE



SWAMPHEN ENTERPRISES

**SWAMPHE** 

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results





#### 7) SET UP MODEL

- 1 # select data for random forest
- 2 include
- 3 include[1:]
- 4 data = data\_set[include[1:]]
- 5 answer = data\_set['Survived']

1 *# import random forest* 

2 **from** sklearn.ensemble **import** RandomForestClassifier **as** RFC

Dr TAMARA CLELFORD

1 *# create the model* 

- 2 forest\_model = RFC(n\_estimators=100)
- 1 # run the model on the data
- 2 titanic\_forest = forest\_model.fit(data, answer)
- 3 print(titanic\_forest)



- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results





#### **CROSS VALIDATION SCORE**

- Splits data randomly into k sections (fold)
- Model is fitted on k-1 of sections training data
- Model is evaluated on remaining section testing data

Dr TAMARA CLELFORD

• Repeated k times



#### CROSS VALIDATION SCORE

	DATA SET SPLIT				
FOLD 1	1	2	3		
FOLD 2	1	2	3		
FOLD 3	1	2	3		

#### TRAINING DATA TESTING DATA



### 8) EVALUATE

1 # see how good model is at predicting survival

- 2 from sklearn.model\_selection import cross\_val\_score as CVS
- 3 scores\_forest = CVS(titanic\_forest, data, answer)
- 4 print(scores\_forest)
- 5 print(np.mean(scores\_forest))



- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results





#### 9) CALCULATE RESULTS

1 # choose a person to see if they survived 2 pclass = np.mean(data\_set['Pclass']) # 2.5 3 age = np.mean(data\_set['Age']) # 29.7 years 4 gender = np.mean(data\_set['GenderNumerical']) # 0.37 5 fare = np.mean(data\_set['Fare']) # £34.69

1 # put person data into an array 2 person = np.array([[pclass, age, gender, fare]]) 3 print(person)

- 1 # calculate survival score
- 2 survival\_score = titanic\_forest.predict\_proba(person)
- 3 print(survival\_score)

SWAMPHEN ENTERPRISES

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results
- 10. Present results



#### **10) PRESENT RESULTS**

1 # did the person survive the sinking of the titanic? 2 print({'survival chances': survival\_score[0,1]\*100, 'death chances': survival\_score[0,0]\*100})



- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results ←
   10. Present results ←



#### 9) CALCULATE RESULTS

```
1 # how does age and gender alter survival chances
2 # create a list covering all ages on the titanic
3 age_max = int(np.max(data_set['Age']))
4 print(age_max)
5 age_list = list(range(0, age_max + 1))
6 print(age list)
```

```
1 # to look at survival probability
2 survival = []
3 selected_class = 1
4 average_fare = np.mean(data_set['Fare'])
5 gender = 1
6 for i in age_list:
7 person = np.array([[selected_class, i, gender, average_fare]])
8 survival_chances = titanic_forest.predict_proba(person)
9 survival.append(survival_chances[0,1]*100)
```

```
1 # to look at male and female survival probability
 2 survival_female = []
 3 survival_male = []
 4 selected_class = 2
 5 average_fare = np.mean(data_set['Fare'])
 6 gender_list = [0,1]
7 for i in age list:
       for j in gender list:
 8
 9
           if j == 1:
               array = survival_female
10
11
            else:
                array = survival male
12
           person = np.array([[selected_class, i, j, average_fare]])
13
            survival_chances = titanic_forest.predict_proba(person)
14
            array.append(survival_chances[0,1]*100)
15
16 len(survival female)
```

SWAMPHEN ENTERPRISES

- 1. Identify requirement
- 2. Setup data science tools
- 3. Get data
- 4. Look at data
- 5. Need more data?
- 6. Clean data set
- 7. Setup model
- 8. Evaluate model
- 9. Calculate results -
- 10. Present results



#### **10) PRESENT RESULTS**

```
1 # plot just survival
2 plt.plot(age_list, survival)
3 plt.title('survival in ' + str(selected_class) + ' class')
4 plt.xlabel('age (y)')
5 plt.ylabel('survival chances (%)')
```

```
1 # plot male and female survival chances
2 plt.plot(age_list, survival_female, label = 'female')
3 plt.plot(age_list, survival_male, label = 'male')
4 plt.title('survival in ' + str(selected_class) + ' class')
5 plt.xlabel('age (y)')
6 plt.ylabel('survival chances (%)')
7 plt.legend()
```

Dr TAMARA CI FI FORD



# SO, IS DATA SCIENCE: "THE SEXIEST JOB OF THE 21<sup>ST</sup> CENTURY"?

- It is certainly a job with its profile on the rise
- Analytical role with varied and interesting work
- Data Science is a well paid and sought-after skill set
  a lot of people are re-training or re-branding
- It covers a skill set desperately needed in this new digital world
- It needs people with good analytical skills to champion
   its evolution as a numerical science

#### HOW DO I LEARN MORE?

- Learn to program in Python
- Get to grips with the data science packages
  - NumPy, pandas and scikit-learn
- Brush up on your stats
- Understand the business you want to work in
- Practice!!

#### QUESTIONS!



#### https://tamaraclelford.co.uk/online\_courses.html

