



University of  
Lancashire

# Data Science: Fundamentals

Karen O'Shea

Where opportunity creates success

# Learning Objectives



Discuss data analysis design



Evaluate tools for Data Science

# Data Analysis



Terms used interchangeably



Asking intelligent questions



Finding patterns



Deriving insight

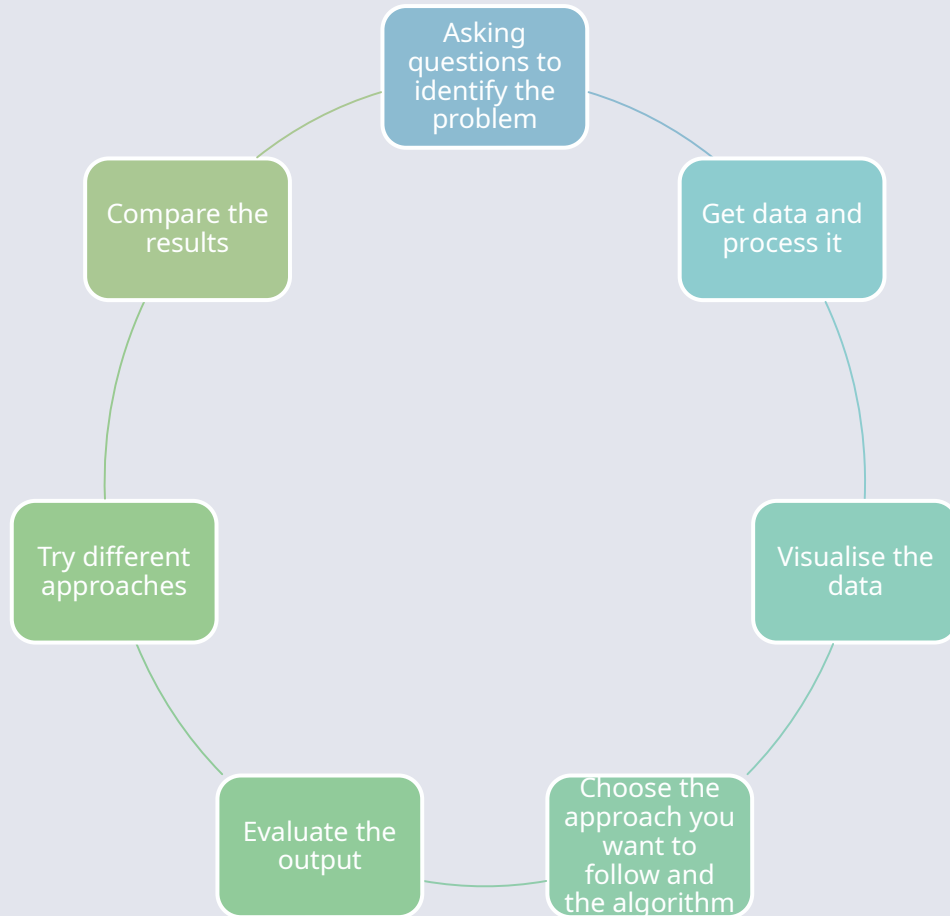


Decision making

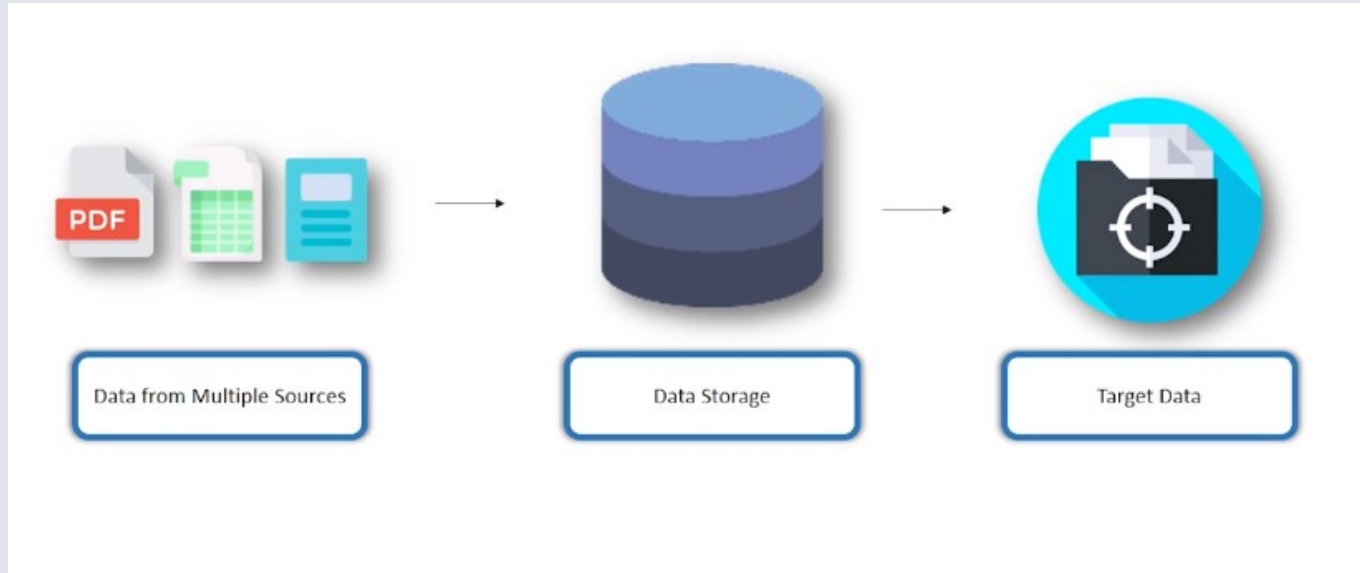


Machine Learning...

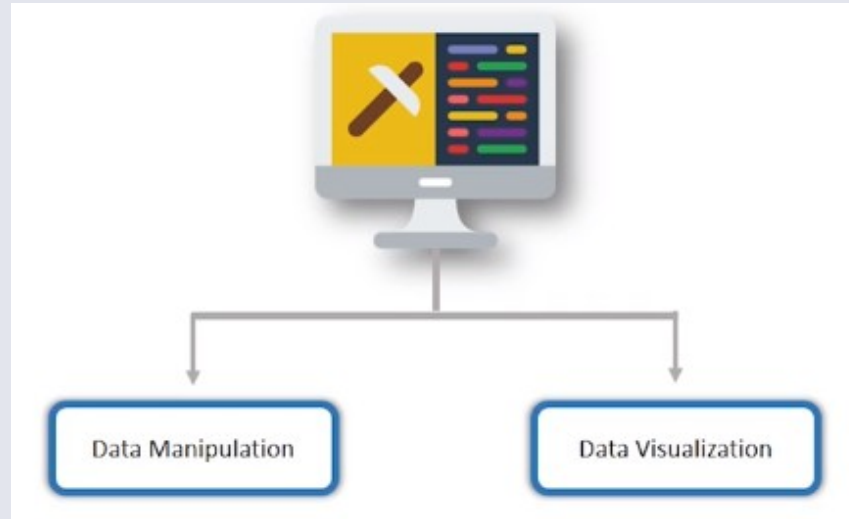
## Data Science: Workflow Model



# Data Science Workflow (1)



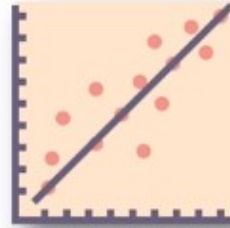
# Data Science Workflow (2)



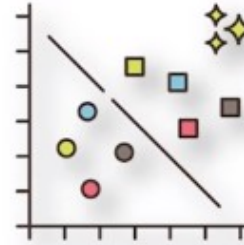
# Data Science Workflow (3)



Classification



Regression



Clustering

# Python for Data Science



Versatile platform



Easy to learn



Open-source



Widely used language

## Dataset examples (last week's lab)

- Dataset available on Blackboard (download from 'Kaggle')

<https://www.kaggle.com/altruistdelhite04/loan-prediction-problem-dataset>

- Examine dataset first. Task:
  - What questions could be asked of the data?
- We will use a Jupyter Notebook or Google Colab. Set up a folder to store dataset and notebook

# Data Analysis Python Libraries

## Many libraries:

- NumPy
- Matplotlib
- SciPy
- Scikit Learn
- Statmodels
- Seaborn
- Blaze
- Scrapy
- SymPy
- Bokeh
- **PANDAS**

# Data Science using PANDAS



Data Analysis and  
Manipulation Tool



One of the most  
useful libraries



Extensive means  
of data analysis



Methods for data  
filtering



Fast, flexible, and  
user friendly

*Data Science algorithms provide benefit to increase company profit, gain insight into business intelligence and enhance processing ability.*

# PANDAS – Data Structures

DataFrame and Series form the basic data model in Pandas:

- Series – one-dimensional indexed array. Easy to access elements of array
- DataFrame – similar to an Excel workbook. Column names (indexes) and row numbers.

# Import Dataset and Libraries

```
import pandas as pd
```

```
df=pd.read_csv("dataset.csv")
```

\*Reads dataset into DataFrame ready for viewing

# Data Exploration

## **df.head(15)**

\* Will display 15 rows

## **df.describe()**

\* Will provide you with mean, quartile, count, min, max and standard deviation as well as their outputs.

## **df['Property\_Area'].value\_counts()**

\* Accessing a particular column in the data frame

```
df = pd.read_csv("LoanPredictionDataset.csv")
df.head()
```



	Loan_ID	Gender	Married	Dependents	Education	Sel
0	LP001002	Male	No	0	Graduate	
1	LP001003	Male	Yes	1	Graduate	
2	LP001005	Male	Yes	0	Graduate	
3	LP001006	Male	Yes	0	Not Graduate	
4	LP001008	Male	No	0	Graduate	

```
df.describe()
```

	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History
<b>count</b>	614.000000	614.000000	592.000000	600.000000	564.000000
<b>mean</b>	5403.459283	1621.245798	146.412162	342.000000	0.842199
<b>std</b>	6109.041673	2926.248369	85.587325	65.12041	0.364878
<b>min</b>	150.000000	0.000000	9.000000	12.000000	0.000000
<b>25%</b>	2877.500000	0.000000	100.000000	360.000000	1.000000
<b>50%</b>	3812.500000	1188.500000	128.000000	360.000000	1.000000
<b>75%</b>	5795.000000	2297.250000	168.000000	360.000000	1.000000
<b>max</b>	81000.000000	41667.000000	700.000000	480.000000	1.000000

```
df['Property_Area'].value_counts()
```

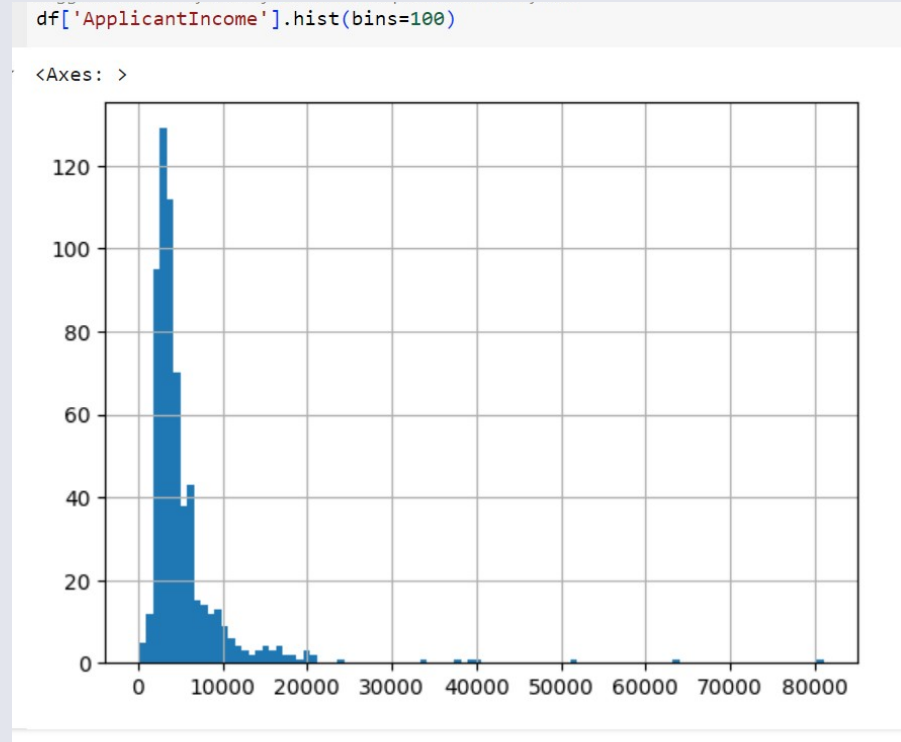


	count
<b>Property_Area</b>	
<b>Semiurban</b>	233
<b>Urban</b>	202
<b>Rural</b>	179

# Distribution Analysis

```
df['ApplicantIncome'].hist(bins=100)
```

\* Using a histogram to observe extreme values



# Learning Objectives



Discuss data analysis design



Evaluate tools for Data Science