

# Data Science Workshop

*Practical . Python*

By: Ali Hamdi

# Introduction to Data Science

01

Life cycle

02

Data  
science  
team

03

Applications

04

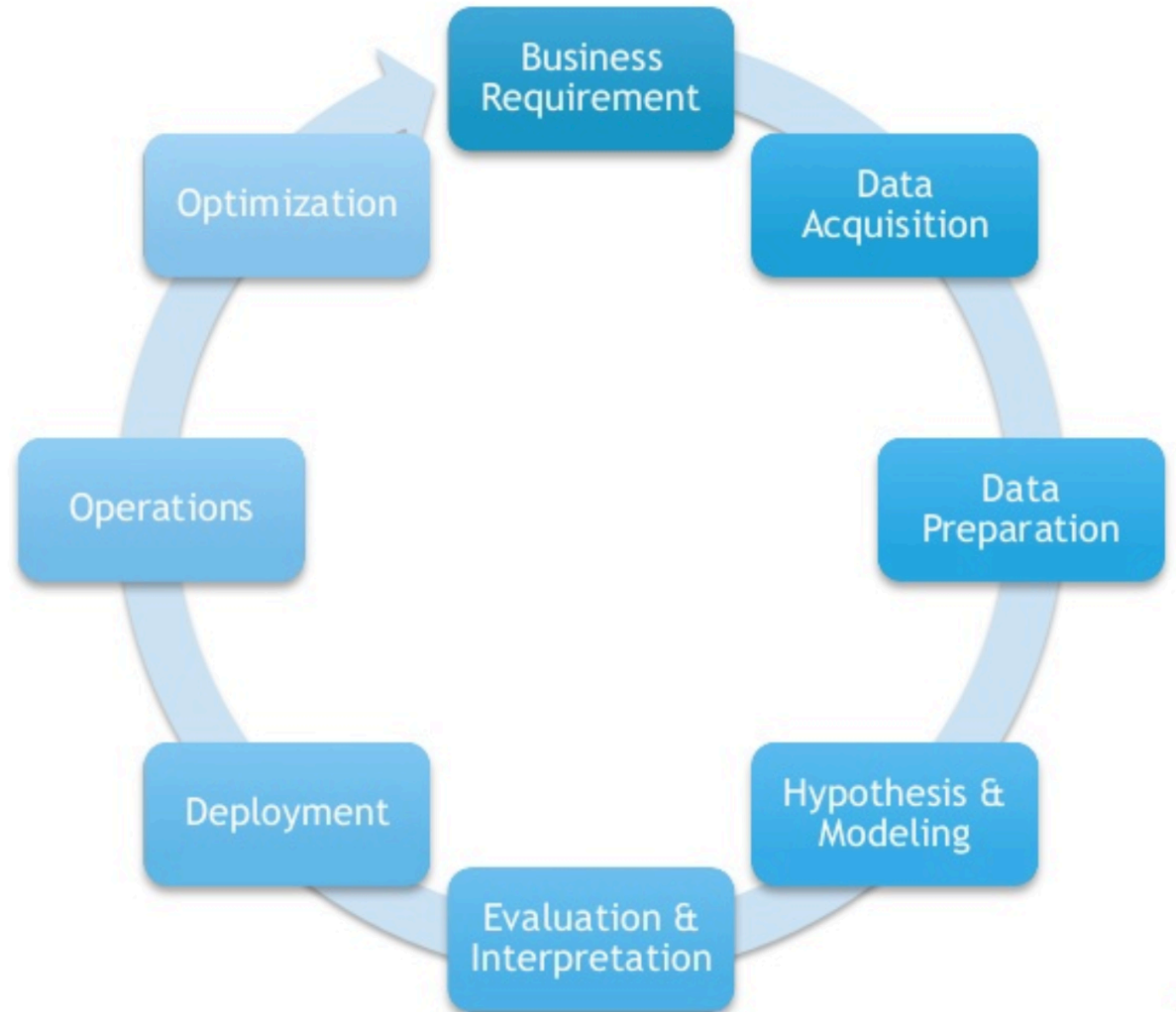
Research  
directions

05

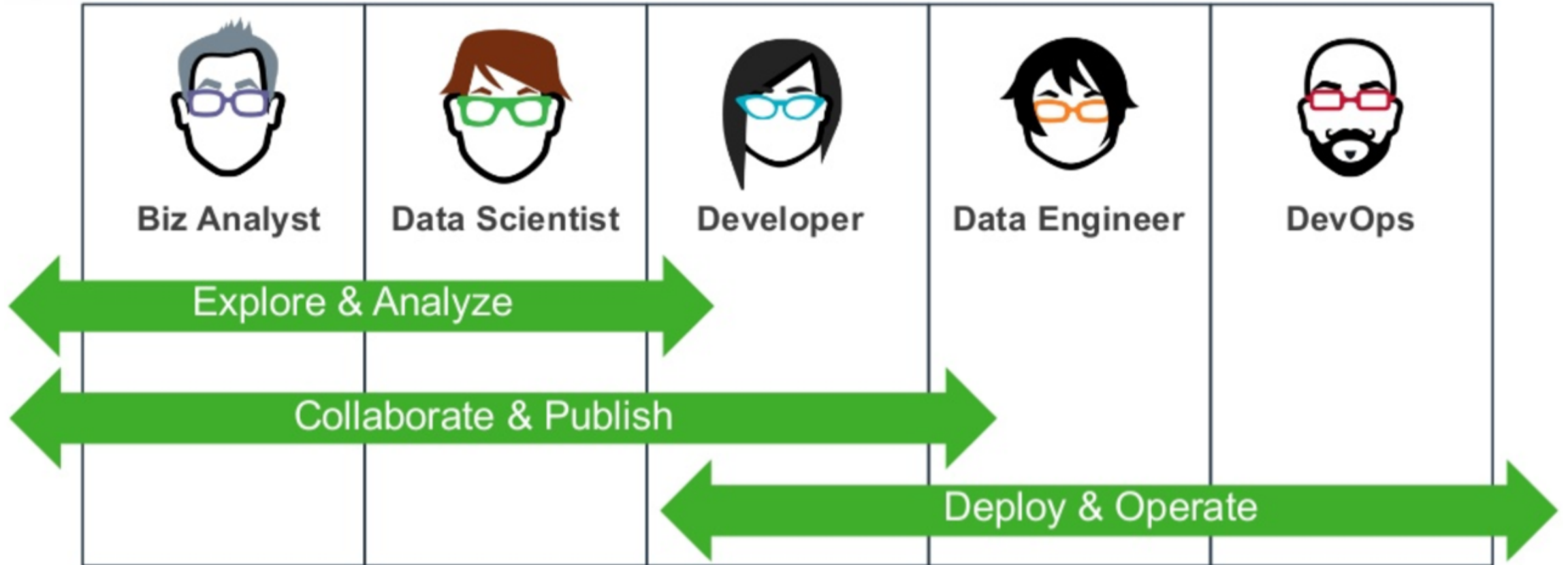
Jobs

# Data Science Project Life Cycle

---



# Data science team



# Data Science Application

---

- Internet Search
- Digital Advertisements
- Recommender Systems
- Image Recognition
- Speech Recognition
- Gaming
- Airline Route Planning
- Fraud and Risk Detection
- Self Driving Cars
- Robots

# Research directions

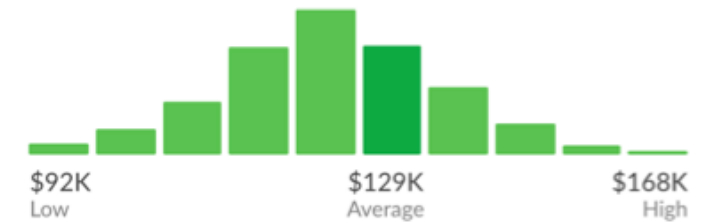
---








- Knowledge Discovery & Data Mining.
- Machine Learning & Deep Learning.
- Visual Computing & Multimedia Analytics.
- Spatial & Context-Aware Data Management.
- Recommender Systems & Preference Analytics.
- Natural Language Processing & Text Mining.








# Data Scientist Salary

Average Base Pay

**\$128,549** /yr



	<b>Data Scientist</b> Facebook 65 salaries	\$135,117/yr	\$100K	\$180K
	<b>Data Scientist</b> Microsoft 44 salaries	\$123,556/yr	\$94K	\$152K
	<b>Data Scientist</b> IBM 40 salaries	\$109,177/yr	\$81K	\$144K
	<b>Data Scientist</b> Booz Allen Hamilton 38 salaries	\$84,450/yr	\$58K	\$140K
	<b>Data Scientist</b> Capital One 28 salaries	\$106,750/yr	\$79K	\$131K
	<b>Data Scientist</b> Nielsen 23 salaries	\$73,725/yr	\$63K	\$85K
	<b>Data Scientist</b> Airbnb 21 salaries	\$126,287/yr	\$93K	\$170K

	<b>Data Scientist</b> Twitter 19 salaries	\$135,360/yr	\$120K	\$155K
	<b>Data Scientist</b> Uber 18 salaries	\$123,686/yr	\$103K	\$152K
	<b>Data Scientist</b> KPMG 16 salaries	\$96,922/yr	\$80K	\$122K
	<b>Data Scientist</b> Civis Analytics 16 salaries	\$76,284/yr	\$61K	\$95K
	<b>Data Scientist</b> LinkedIn 16 salaries	\$132,059/yr	\$111K	\$261K
	<b>Data Scientist</b> Rang Technologies 15 salaries	\$104,480/yr	\$92K	\$133K
	<b>Data Scientist</b> Apple 14 salaries	\$144,833/yr	\$112K	\$180K

# Python Programming

01

Environment Setup

02

Python Basics



# Data Exploratory and Analysis

## Manipulation

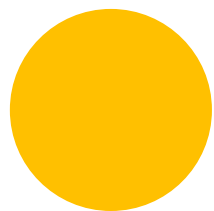
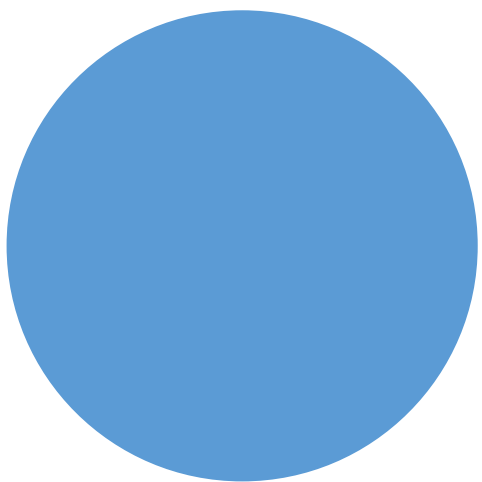
- Numpy
- Pandas

## Visualization

- Matplotlib

## Retrieval

- BeautifulSoup



Break



# Machine Learning

---

## Regression

---

## Classification

---

## Clustering

---

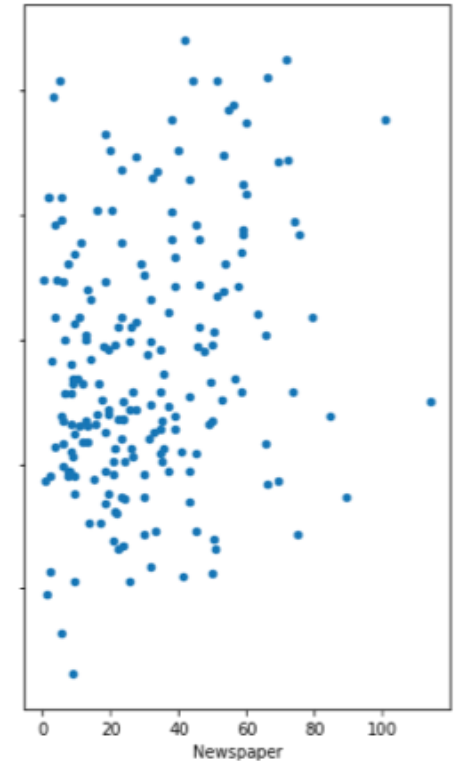
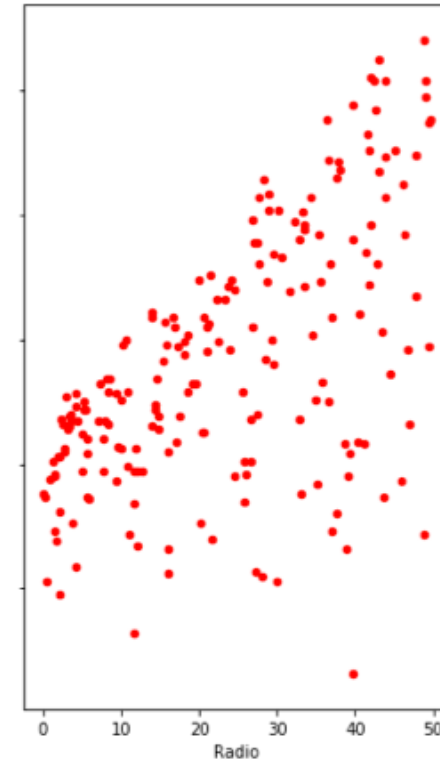
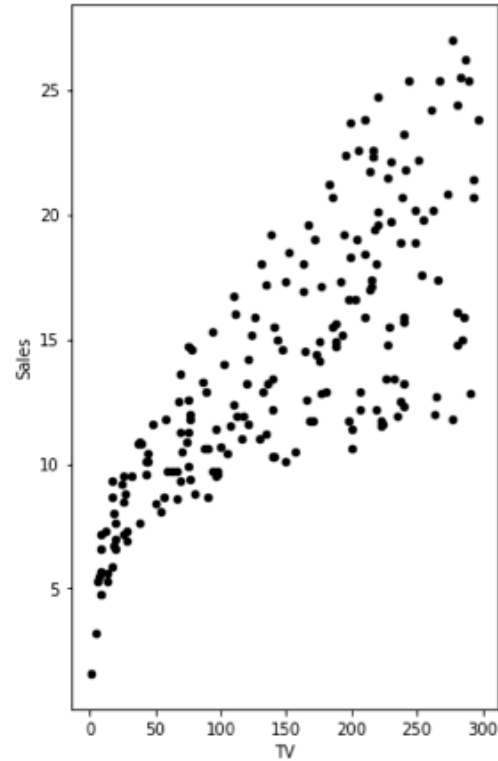
## Association Rules

# Supervised & Unsupervised Learning

	Continuous	Categorical
Supervised	Regression	Classification
Unsupervised	Dimensionality Reduction	Clustering

# Regression

	TV	Radio	Newspaper	Sales
1	230.1	37.8	69.2	22.1
2	44.5	39.3	45.1	10.4
3	17.2	45.9	69.3	9.3
4	151.5	41.3	58.5	18.5
5	180.8	10.8	58.4	12.9
6	8.7	48.9	75	7.2
7	57.5	32.8	23.5	11.8
8	120.2	19.6	11.6	13.2

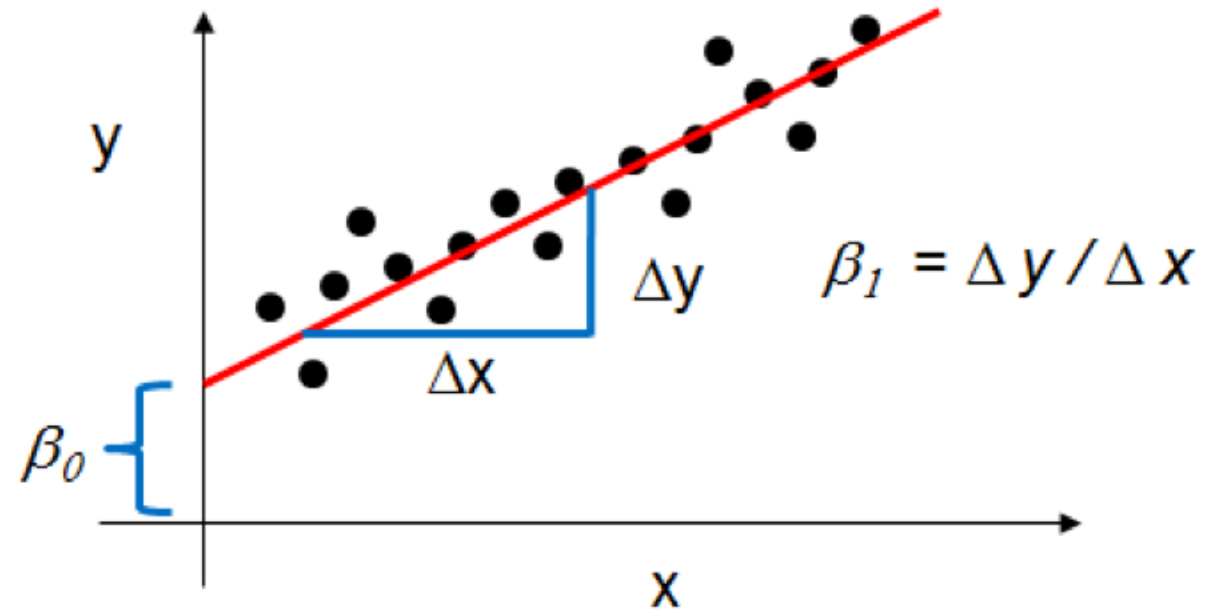


# Simple Linear Regression

- To predict a **quantitative response** using a **single feature**.

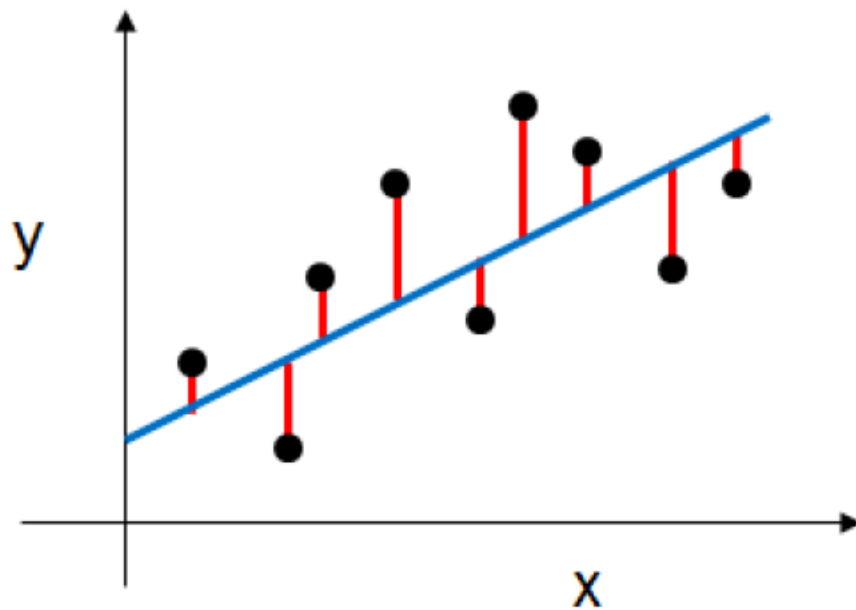
$$y = \beta_0 + \beta_1 x$$

- $y$  is the response (Dependent Variable)
- $x$  is the feature (Independent variable)
- $\beta_0$  is the **intercept** (the value of  $y$  when  $x=0$ )
- $\beta_1$  is the **slope** (the change in  $y$  divided by change in  $x$ )



# Estimating Model Coefficients

- Coefficients are estimated using the **least squares criterion**
  - Minimizes the **sum of squared residuals** or "sum of squared errors"



$$SS_{residuals} = \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

Model Prediction

Observed Result

The diagram shows the formula for the sum of squared residuals. A red arrow points from the label 'Model Prediction' to the term  $\hat{y}_i$  in the formula. Another red arrow points from the label 'Observed Result' to the term  $y_i$  in the formula.

# Multiple Linear Regression

- Using multiple features:

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$

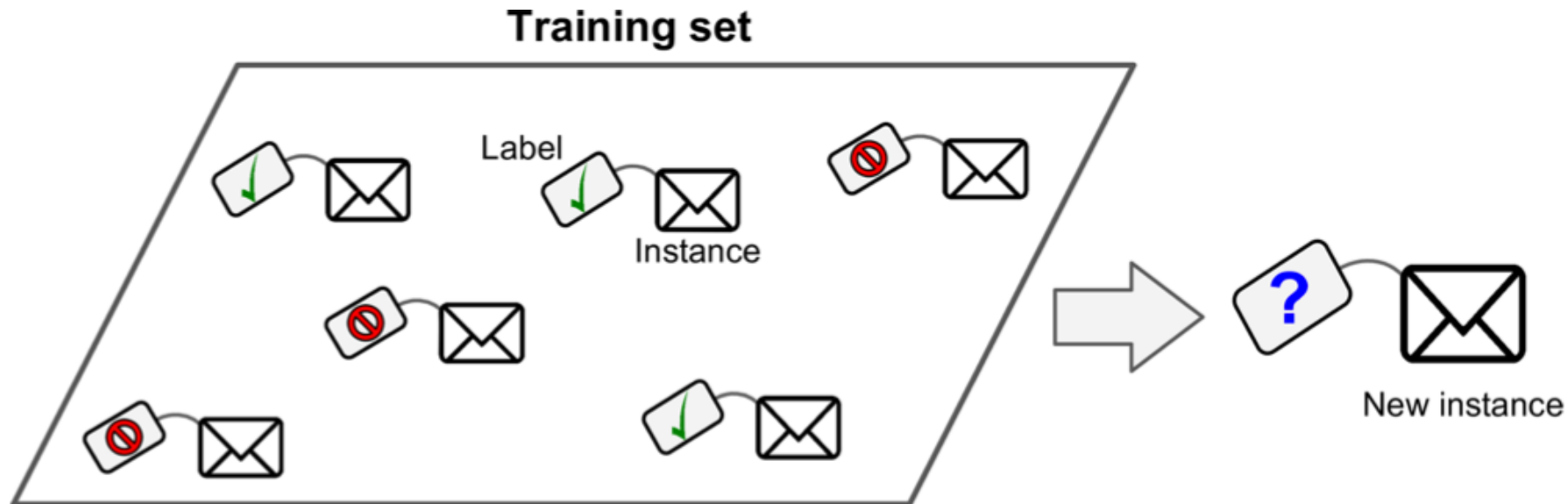
- $x$  represents different feature and  $\beta$  is the feature coefficient :

$$y = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{Radio} + \beta_3 \times \text{Newspaper}$$



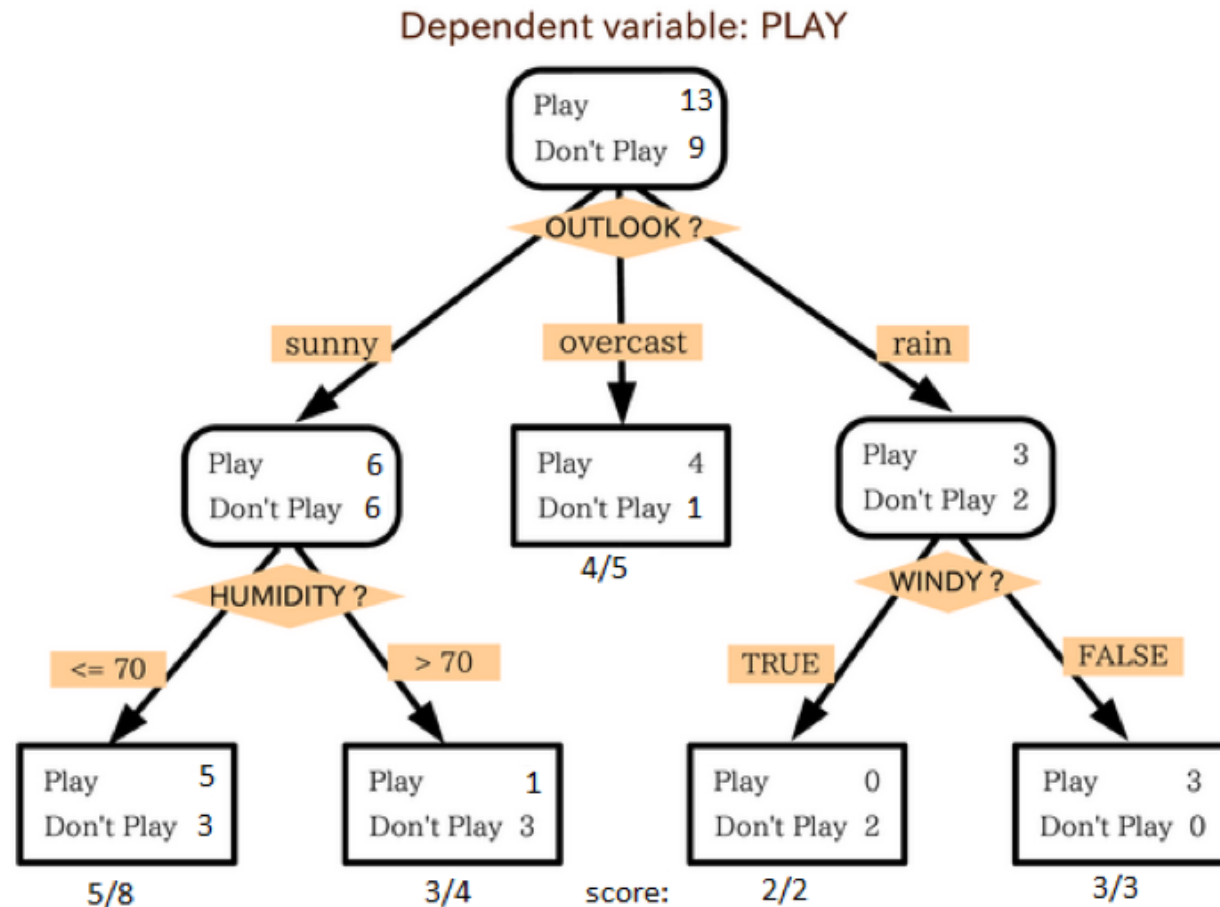
# Machine Learning / Statistical classification

To identify a new observation belongs to which category, based on training data containing observations whose category membership is known.

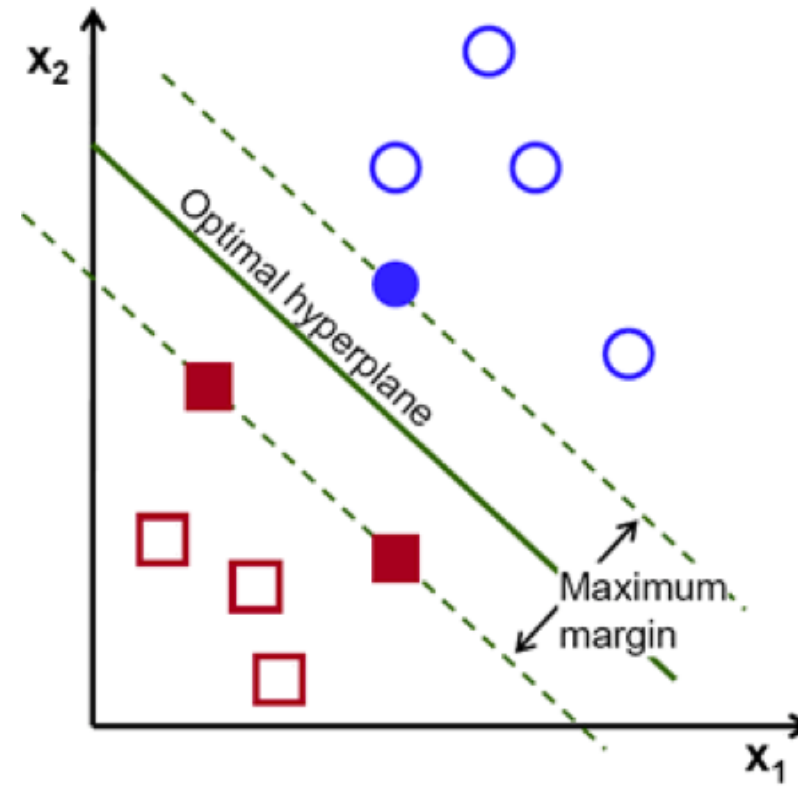
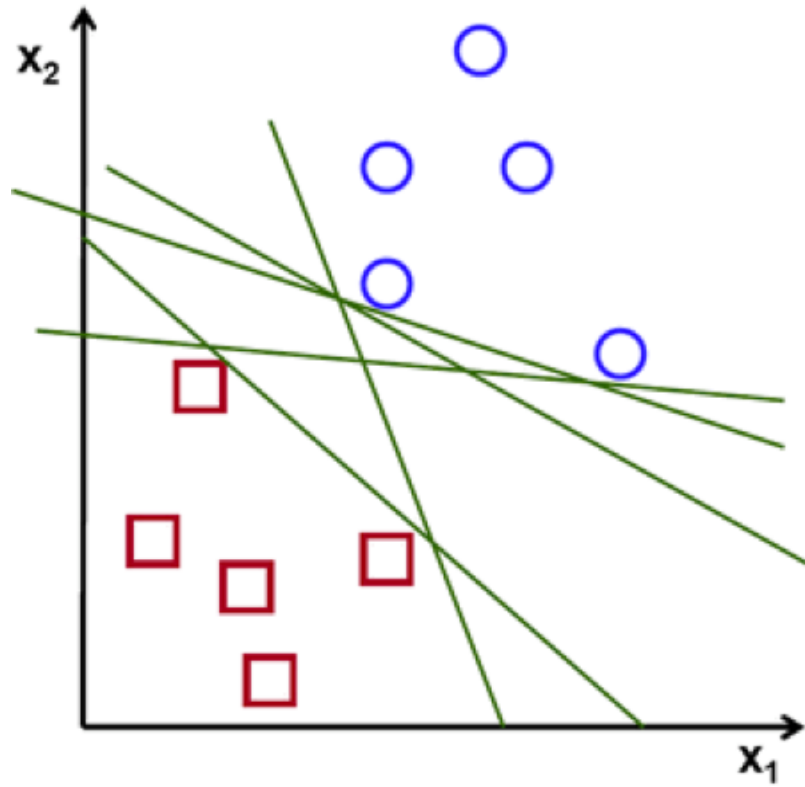


# Decision Tree

Split the data into two or more homogeneous sets, based on most significant attributes to make as distinct groups as possible.



# Support Vector Machine (SVM)



# Naive Bayes

**Will** players will pay if weather is sunny?

$$P(\text{Yes} \mid \text{Sunny}) = P(\text{Sunny} \mid \text{Yes}) * P(\text{Yes}) / P(\text{Sunny})$$

$$P(\text{Sunny} \mid \text{Yes}) = 3/9 = 0.33,$$

$$P(\text{Sunny}) = 5/14 = 0.36,$$

$$P(\text{Yes}) = 9/14 = 0.64$$

$$P(\text{Yes} \mid \text{Sunny}) = 0.33 * 0.64 / 0.36 = 0.60, \text{ which has higher probability.}$$

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No

14 data instances

Frequency Table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9

Likelihood table				
Weather	No	Yes		
Overcast		4	=4/14	0.29
Rainy	3	2	=5/14	0.36
Sunny	2	3	=5/14	0.36
All	5	9		
	=5/14	=9/14		
	0.36	0.64		

$$P(c \mid x) = \frac{P(x \mid c)P(c)}{P(x)}$$

Likelihood
Class Prior Probability

Posterior Probability
Predictor Prior Probability

$$P(c \mid X) = P(x_1 \mid c) \times P(x_2 \mid c) \times \dots \times P(x_n \mid c) \times P(c)$$

$P(c/x)$  is the posterior probability of class ( $c$ ) given predictor ( $x$ ).

$P(c)$  is the prior probability of class.

$P(x/c)$  is the likelihood which is the probability of predictor given class.

$P(x)$  is the prior probability of predictor.

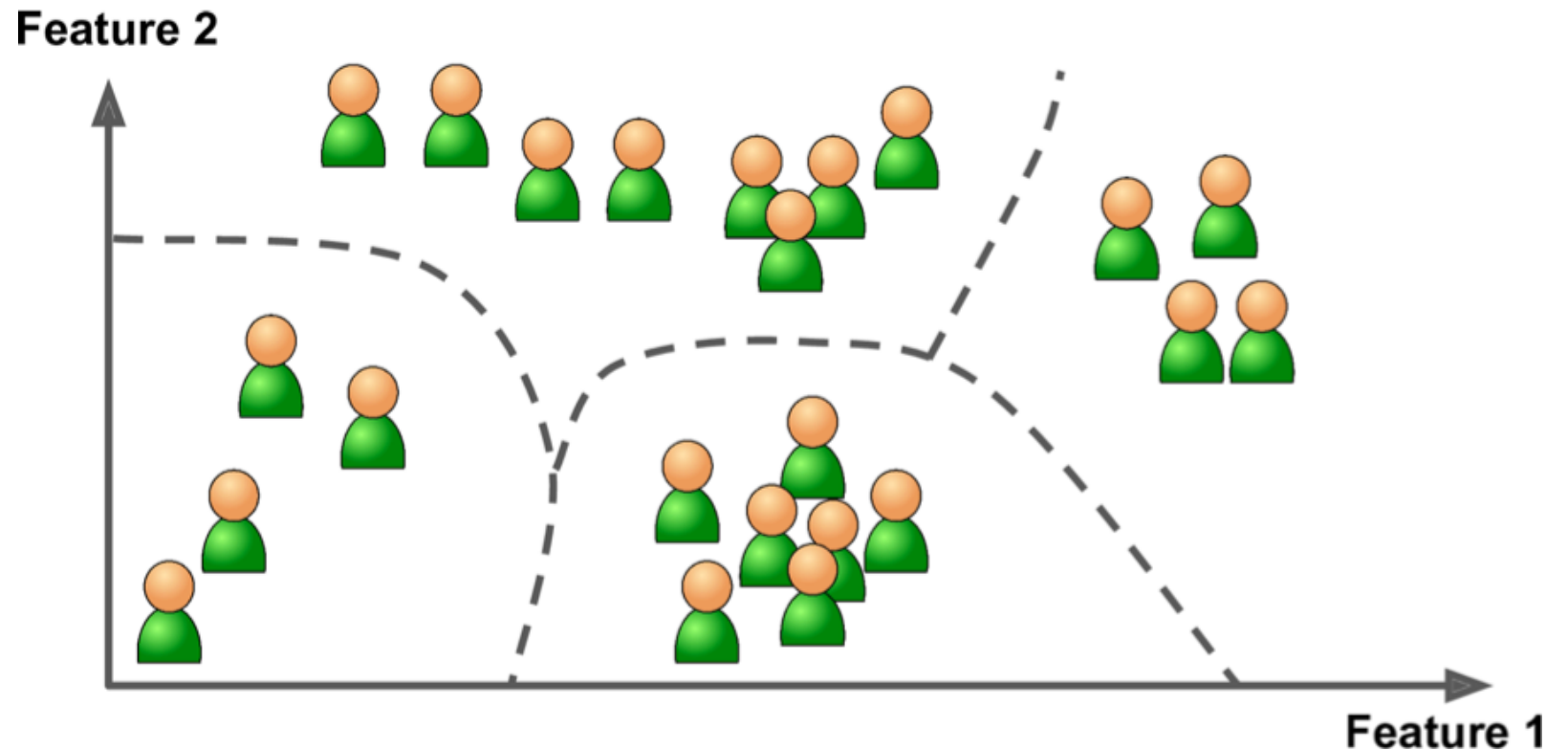
# Clustering

Blog's visitors dataset:

Clustering is to group similar visitors without labeling.

For example, the visitors can be clustered based on their gender and the visited sections.

This may help you target your posts for each group.



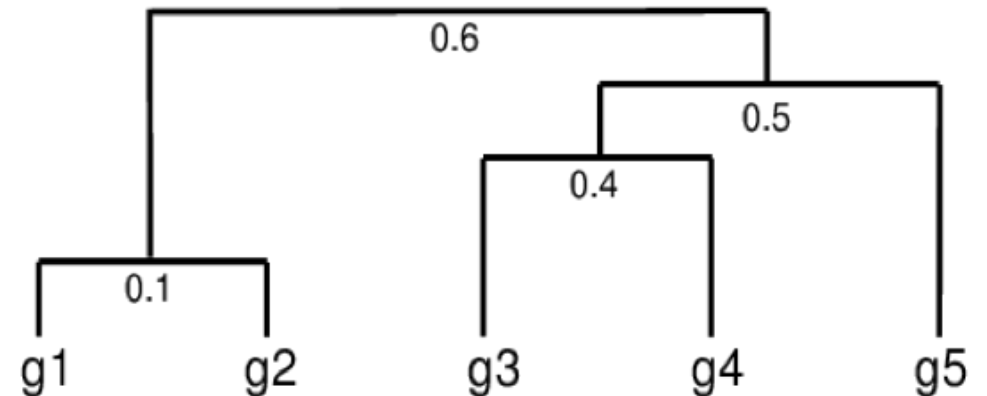
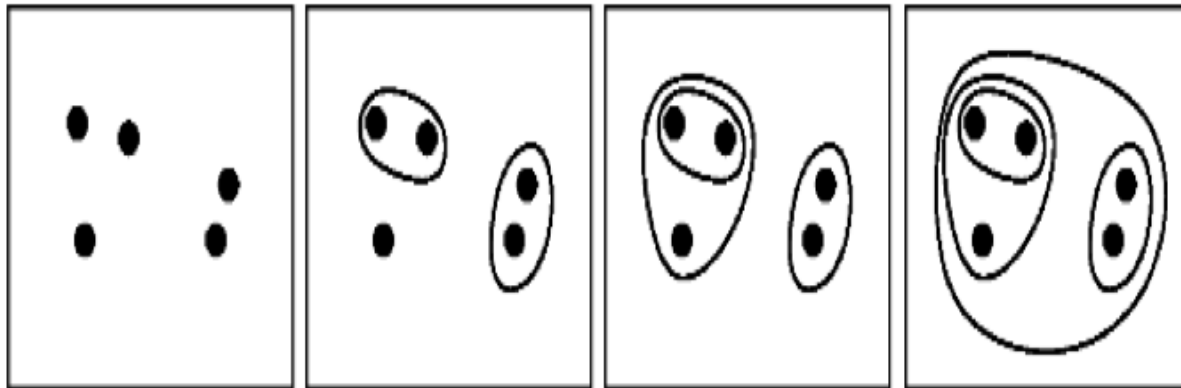
# Hierarchical Clustering

- **Agglomerative** or "bottom up" approach.

Each observation starts as a cluster, and pairs of clusters are grouped in one cluster, moves up the hierarchy.

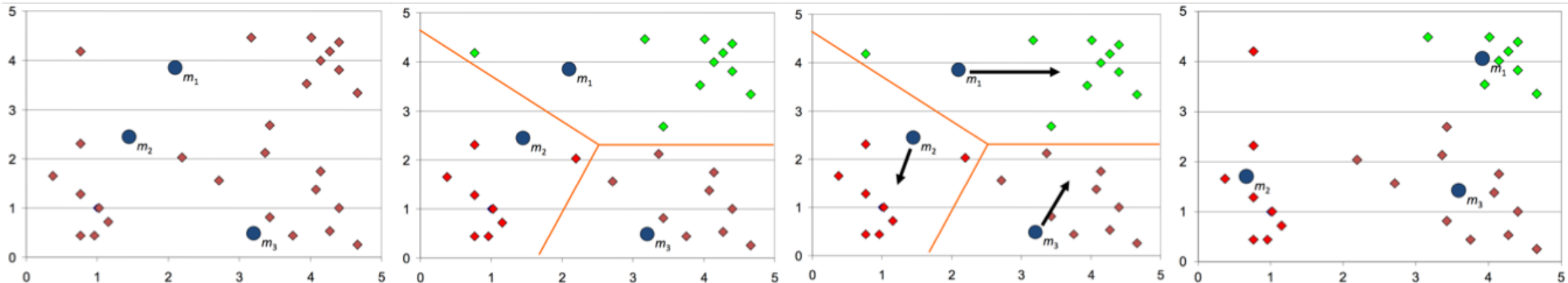
- **Divisive** or "top down" approach.

All observations start in one cluster, and splits are performed recursively as one moves down the hierarchy.



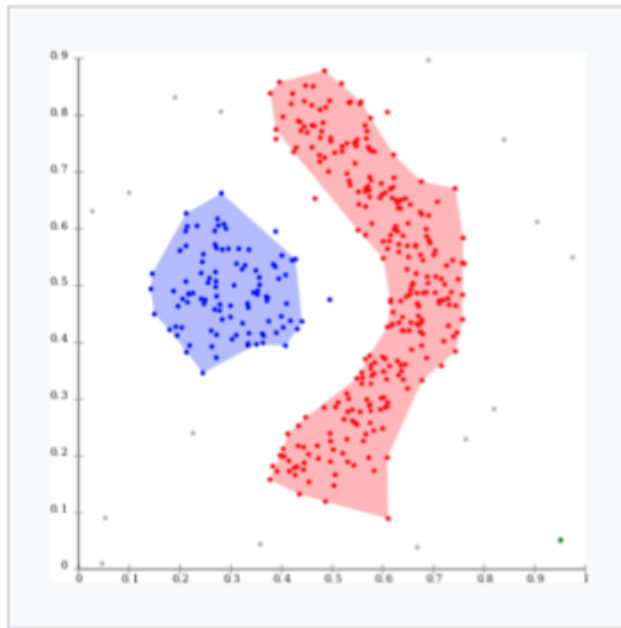
# Partitional / Centroid-based Clustering

- K-means:
  - Define K number of clusters.
  - For each object  $x_i$ 
    - Calculate the distances between  $x_i$  and the K centroids
    - (Re)assign  $x_i$  to the cluster whose centroid is the closest to  $x_i$
  - Update the cluster centroids based on current assignment

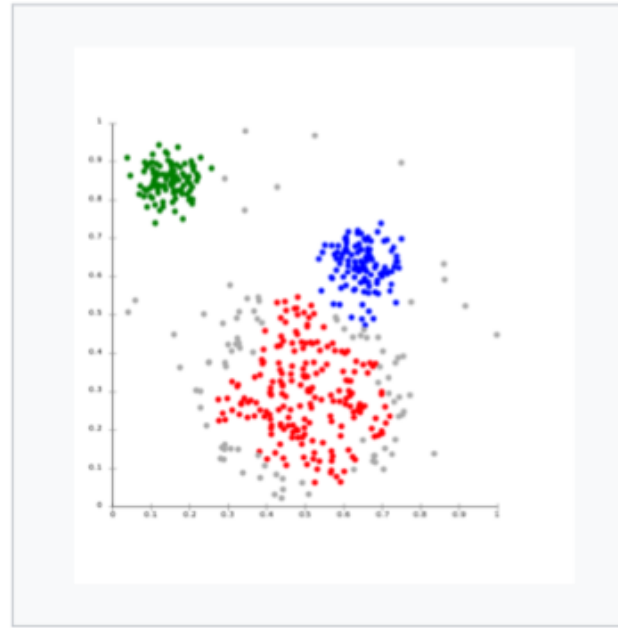


# Density-based clustering

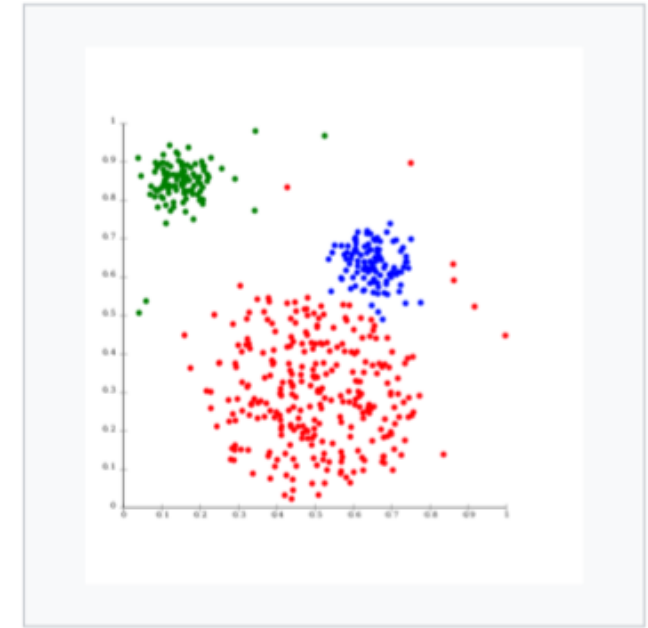
- The objective is to cluster the objects in areas of higher density.
- Objects that are not close to any cluster considered as noise and border points between clusters.



Density-based clustering with  
**DBSCAN**.



**DBSCAN** assumes clusters of similar  
density, and may have problems  
separating nearby clusters

























**OPTICS** is a DBSCAN variant that  
handles different densities much better



# Association Rules / Apriori Algorithm

Association rules analysis is a technique to uncover how items are associated to each other. There are three common ways to measure association.

Transaction 1	   
Transaction 2	  
Transaction 3	 
Transaction 4	 
Transaction 5	   
Transaction 6	  
Transaction 7	 
Transaction 8	 

$$\text{Support} \{ \text{apple} \} = \frac{4}{8}$$

how popular an itemset is?

$$\text{Confidence} \{ \text{apple} \rightarrow \text{Coca-Cola} \} = \frac{\text{Support} \{ \text{apple}, \text{Coca-Cola} \}}{\text{Support} \{ \text{apple} \}}$$

how likely item Y is purchased when item X?

$$\text{Lift} \{ \text{apple} \rightarrow \text{Coca-Cola} \} = \frac{\text{Support} \{ \text{apple}, \text{Coca-Cola} \}}{\text{Support} \{ \text{apple} \} \times \text{Support} \{ \text{Coca-Cola} \}}$$

how likely item Y is purchased when item X is purchased, while controlling for how popular item Y is.

# Thank You !

Video tutorials: <https://goo.gl/JqevFk>

Source code: <https://goo.gl/FAYr51>