# DATA SCIENCE 102: CLUSTERING

# AGENDA

- Unsupervised Learning
- Clustering
  - Use Cases
  - Types of Clustering Algorithms
- K-Means Clustering
  - K-Means Algorithm
  - Optimal K
  - Coded Example

# UNSUPERVISED LEARNING

# UNSUPERVISED LEARNING

- Unsupervised learning is used when there is **no outcome variable** (**y**) to predict or classify
- Attempts to learn patterns in the data other than predicting **y**
- Unsupervised learning methods include:
  - Clustering Techniques
  - Association Rules
  - Dimension Reduction Methods

# CLUSTERING TECHNIQUES

- Use Cases

- Types of Clustering Algorithm

- Distance Scoring

# USAGE

- Segmentation of data into sets of homogenous clusters of records to generate insight
- Clustering can help improve the performance of supervised methods by modelling each cluster rather than the entire heterogeneous dataset
- Cluster analysis helps to form groups (clusters) of similar observations based on several measurements made on those observations
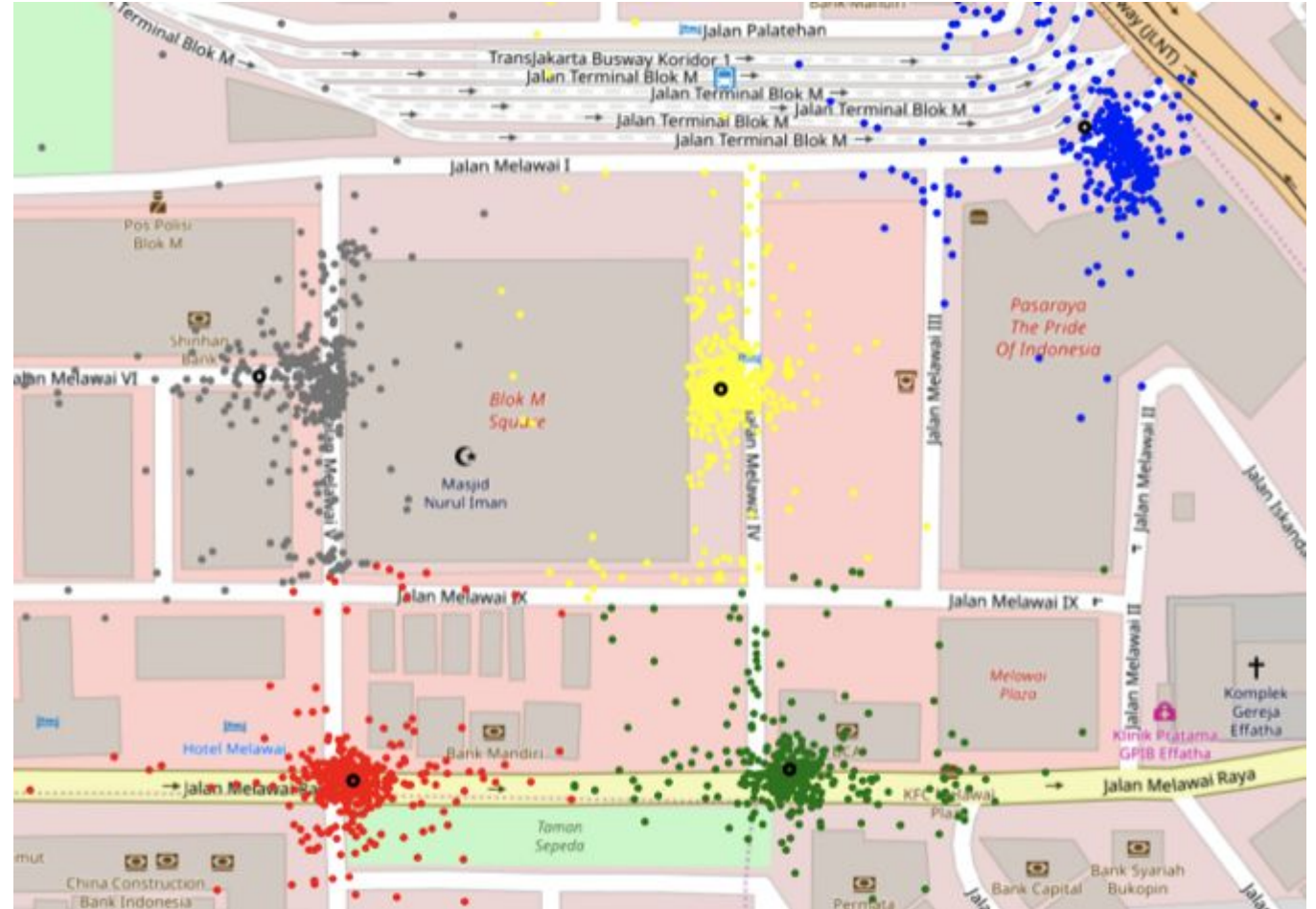
# USE CASES

- Marketing
  - **Market segmentation** of customers based on demographic and transaction history and tailored marketing strategy
  - **Market structure analysis** identifying groups of similar products according to competitive measures of similarity
- Finance
  - **Balanced portfolios** by choosing stocks from different clusters
  - **Industry Analysis** finding similar firms through "market measures"
- Accounting
  - **Group transactions** by type
  - **Anomaly detection**

*Data Mining for Business Analytics: Concepts, Techniques, and Applications in R* by
GalitShmueli, Peter C. Bruce, InbalYahav, Nitin R. Patel, Kenneth C. LichtendahlJr. (2018)

# USE CASE - GO-JEK FANTASTIC DRIVERS

- Go-Jek used K-Means algorithm to identify their better drivers
- It also helped them "pin" pick up points at popular locations



Find out more here: **https://blog.gojekengineering.com/fantastic-drivers-and-how-to-find-them-a88239ef3b29**

# TYPES OF CLUSTERING ALGORITHMS

- There are two general types of clustering algorithms:
  a. Hierarchical
    - Agglomerative - begins with *n* clusters and sequentially merge similar clusters until a single cluster is obtained
    - Divisive - starts with a single cluster including all records and does the opposite
  b. Non-hierarchical *(Focused for this class; k-means clustering)*
    - Using predetermined number of clusters to assign observations to each cluster
    - Less computationally intensive and preferred for larger datasets

*Data Mining for Business Analytics: Concepts, Techniques, and Applications in R* by
GalitShmueli, Peter C. Bruce, InbalYahav, Nitin R. Patel, Kenneth C. LichtendahlJr. (2018)

# K-MEANS CLUSTERING

- K-Means Algorithm
- Distance Scoring
- Optimal K
- Limitations of K-Means

# K-MEANS ALGORITHM

- Start with *k* initial clusters (*k* needs to be pre-defined)
- At every step, each record is reassigned to the cluster with the "closest" centroid
- Recompute the centroids of clusters that lost or gained a record, and repeat Step 2
- Stop when moving any more records between clusters **increases cluster dispersion**

- Suppose two different observations are i and j, the distance metric for them is $d_{ij}$
- The formula to calculate the distance between two observed points is the **Euclidean Distance**:
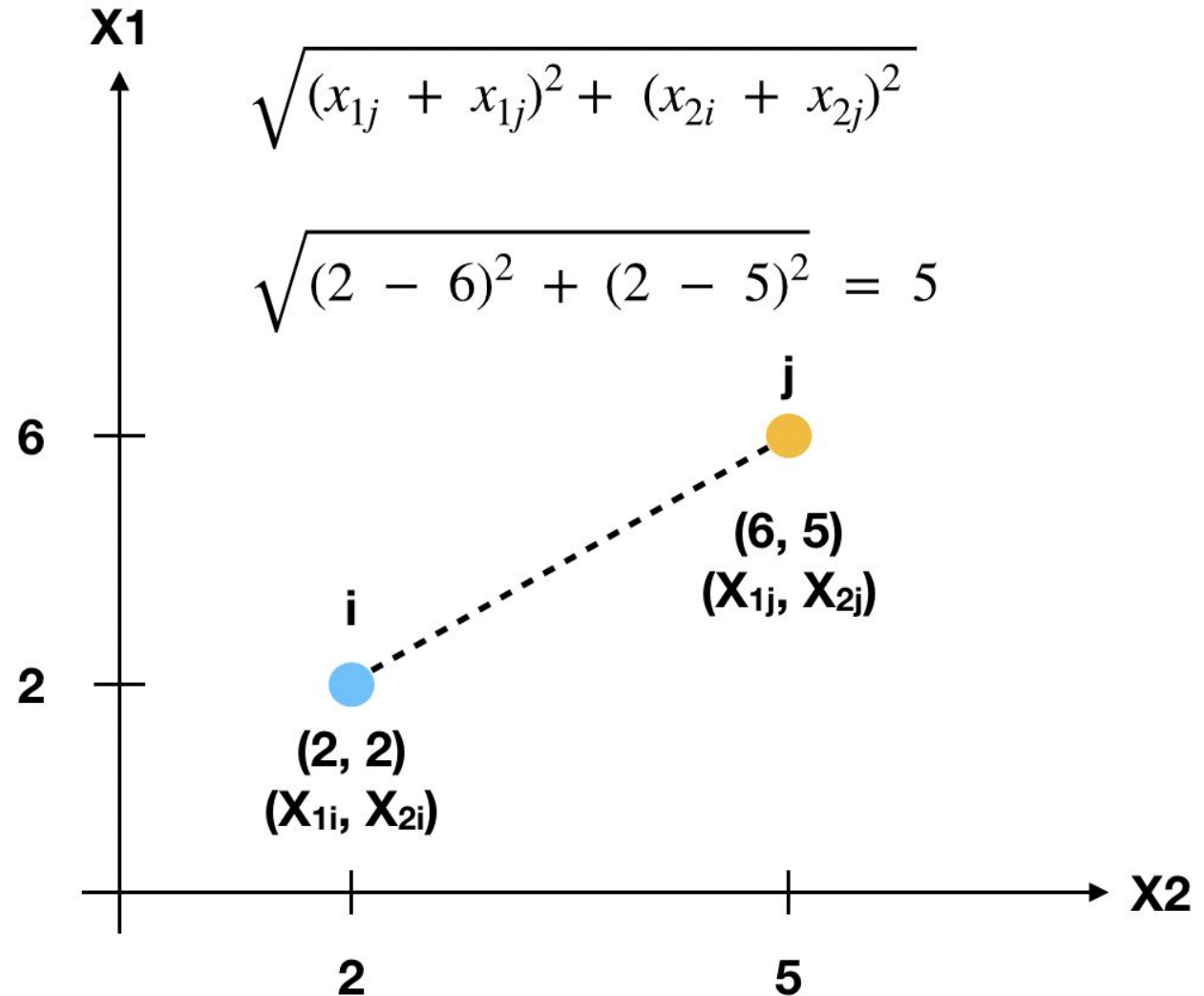
$$d_{ij} = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \cdots + (x_{ip} - x_{jp})^2}$$

*Data Mining for Business Analytics: Concepts, Techniques, and Applications in R* by
GalitShmueli, Peter C. Bruce, InbalYahav, Nitin R. Patel, Kenneth C. LichtendahlJr. (2018)

| ID | X1 | X2 |
|----|----|----|
| i  | 2  | 5  |
| j  | 2  | 6  |

$$\sqrt{(x_{1j} + x_{1j})^2 + (x_{2i} + x_{2j})^2}$$

$$\sqrt{(2 - 6)^2 + (2 - 5)^2} = 5$$

X1

j

6

(6, 5)
($X_{1j}$, $X_{2j}$)

i

(2, 2)
($X_{1i}$, $X_{2i}$)

2

X2

2

5

# K-MEANS ALGORITHM - VISUALIZED



Click here for an interactive k-means algorithm animation: https://www.naftaliharris.com/blog/visualizing-k-means-clustering/
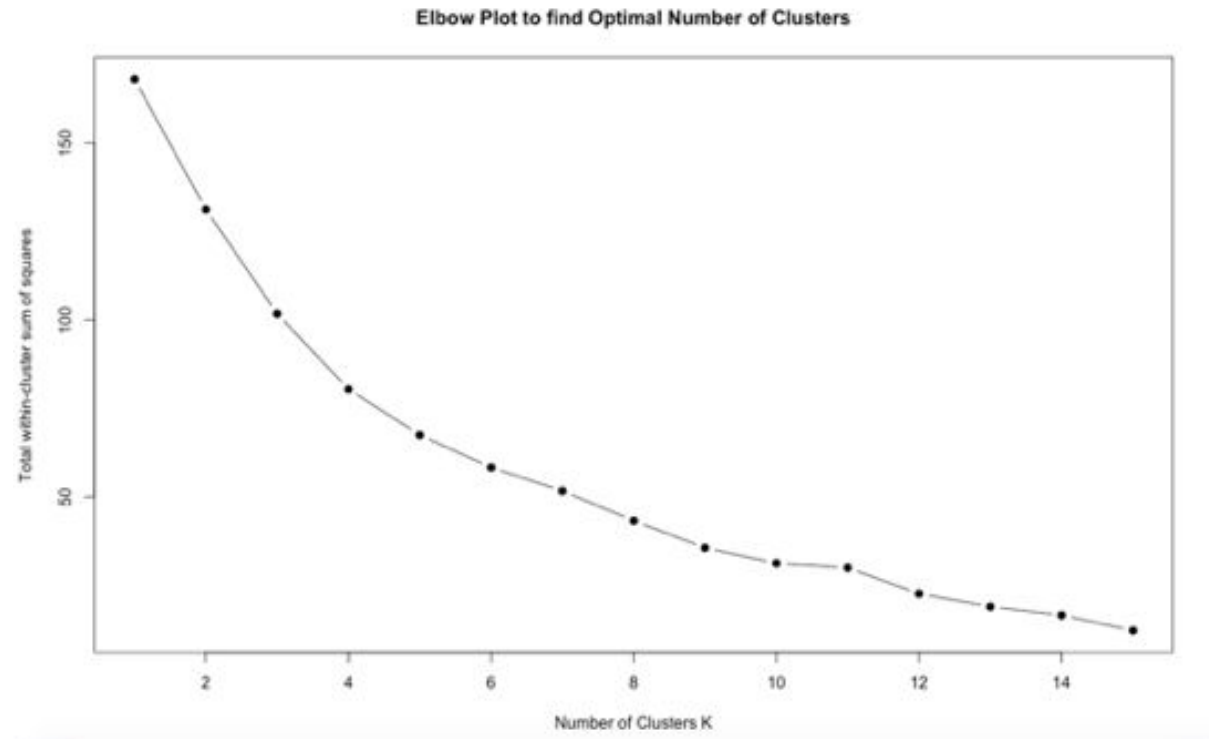
# OPTIMAL K

- In many cases, there is a lack of information to be used for the initial number of *k*.
- K-means algorithm, like other clustering algorithms, work towards compressing the *k* number of data-points by summarising them by their "means" (hence, k-means)
- <u>There is no right number of *k*'s.</u> However, information on within-cluster dispersion can assist in determining the optimal number of *k*
- With an "elbow chart", we can graphically evaluate whether there is a decline in cluster heterogeneity when more clusters are added (i.e *k* is increased)

*Data Mining for Business Analytics: Concepts, Techniques, and Applications in R* by
GalitShmueli, Peter C. Bruce, InbalYahav, Nitin R. Patel, Kenneth C. LichtendahlJr. (2018)

- Select a K where the next increase in K has **little to no improvement** in the total within cluster sum of squares (WSS)
- Higher the K, the more computation required

  choose the k with the sharpest drop

**Elbow Plot to find Optimal Number of Clusters**

Total within-cluster sum of squares vs Number of Clusters K

# SPOTIFY DATASET - SELECTING K*

```python
4  X = df_spotify_cluster # <<< Numerical DataFrame here
5  distorsions = []
6  for k in range(2, 20):
7      kmeans = KMeans(n_clusters=k)
8      kmeans.fit(X)
9      distorsions.append(kmeans.inertia_)
10
11 fig = plt.figure(figsize=(15, 5))
12 plt.plot(range(2, 20), distorsions)
13 plt.grid(True)
14 plt.title('Elbow curve')
```

Loops for clusters from 2 to 19

Initialises based on K clusters

Calculates within cluster sum of squared

Plots elbow plot

# SPOTIFY DATASET - K-MEANS ALGORITHM*

```python
1  k = 5
2  model = KMeans(n_clusters=k,   # < Initialise Number Of Clusters here
3                 random_state=0)
4
5  spotify_kmeans = model.fit(df_spotify_cluster) # < DataFrame of All Variables
6  print(spotify_kmeans)
```

```
KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
    n_clusters=5, n_init=10, n_jobs=1, precompute_distances='auto',
    random_state=0, tol=0.0001, verbose=0)
```

Initialises the KMeans model based on K clusters

Trains the model by fitting in all variables into the model and **returns a kmeans result set**

# IN-CLASS PRACTICE - CLUSTERING*

- Try out the in-class practice with the credit card spending behaviour
- Do a summary statistics on the different clusters of credit card spending behaviours