*classification → classifying observation into labels*

*Regression → predicting a number*

# K-Nearest Neighbors (K-NN)
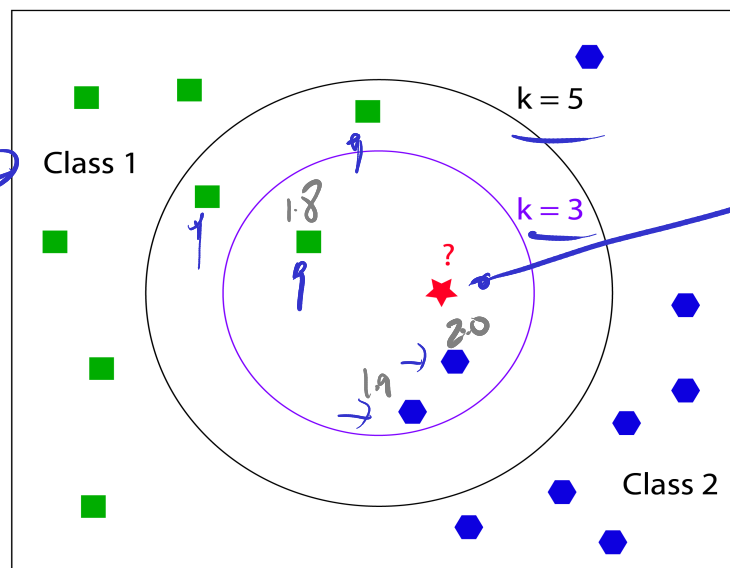
## 1. Introduction to K-Nearest Neighbors (KNN)

K-Nearest Neighbors (KNN) is a simple, yet powerful supervised machine learning algorithm used for both classification and regression tasks. It makes predictions based on the 'k' closest data points in the feature space.

## 2. How KNN Works

1. Choose the number of neighbors: $k$ ✓

2. Compute the distance between the new data point and all points in the training set (commonly using Euclidean distance).

3. Select the $k$ nearest neighbors.

4. For classification, assign the majority label among the $k$ neighbors.

5. For regression, compute the average of the neighbors' values.

*Classification task →*

*if k=3, then the predicted label is blue hexagon*



Class 1

k = 5

1.8

1.7

1.9

2.0

1.9

? 

k = 3

*new data point*

Class 2

*if k=5, Then The predicted label is green square*

# 3. Common Applications of KNN in Data Science

- **Customer segmentation:** Classifying users based on purchasing behavior.

- **Medical diagnosis:** Predicting disease based on patient features.

# 4. Advantages of KNN

- **Simple to understand and implement.**

- **Non-parametric:** No assumptions about data distribution.

- **Versatile:** Works for both classification and regression.

# 5. Disadvantages of KNN

- **Computationally expensive:** Slow with large datasets since it stores the entire training data.

- **Sensitive to noise and outliers.**

- **Feature scaling is required:** Performance is affected if features are on different scales.

- **Not interpretable:** Doesn't provide insights into feature importance.

# 6. When to Use KNN

- When your dataset is small to medium in size.

- When the data is labeled and not too noisy.

- When interpretability is not a primary concern.

# 7. Summary

K-Nearest Neighbors is a valuable algorithm for data scientists to understand due to its simplicity and versatility. While not ideal for large or high-dimensional datasets, it offers a practical introduction to the concepts of distance-based learning and instance-based algorithms.