# K-Nearest Neighbour

## MACHINE LEARNING

# K-Nearest Neighbour

- K-NN is a simple algorithms that store all available cases.

- The k-nearest neighbors (KNN) algorithm is a simple, easy-to-implement supervised machine learning algorithm that can be used to solve both classification and regression problems.

- Classification is done using similarity.

- Also, It is known as:
  - ✓ Memory-Based Reasoning
  - ✓ Example-Based Reasoning
  - ✓ Instance-Based Learning
  - ✓ Case-Based Reasoning
  - ✓ Lazy Learning

# Different Learning Methods

- **Eager Learning**
  - Explicit description of target function on the whole training set

- **Instance-based Learning**
  - Learning=storing all training instances
  - Classification=assigning target function to a new instance
  - Referred to as "Lazy" learning
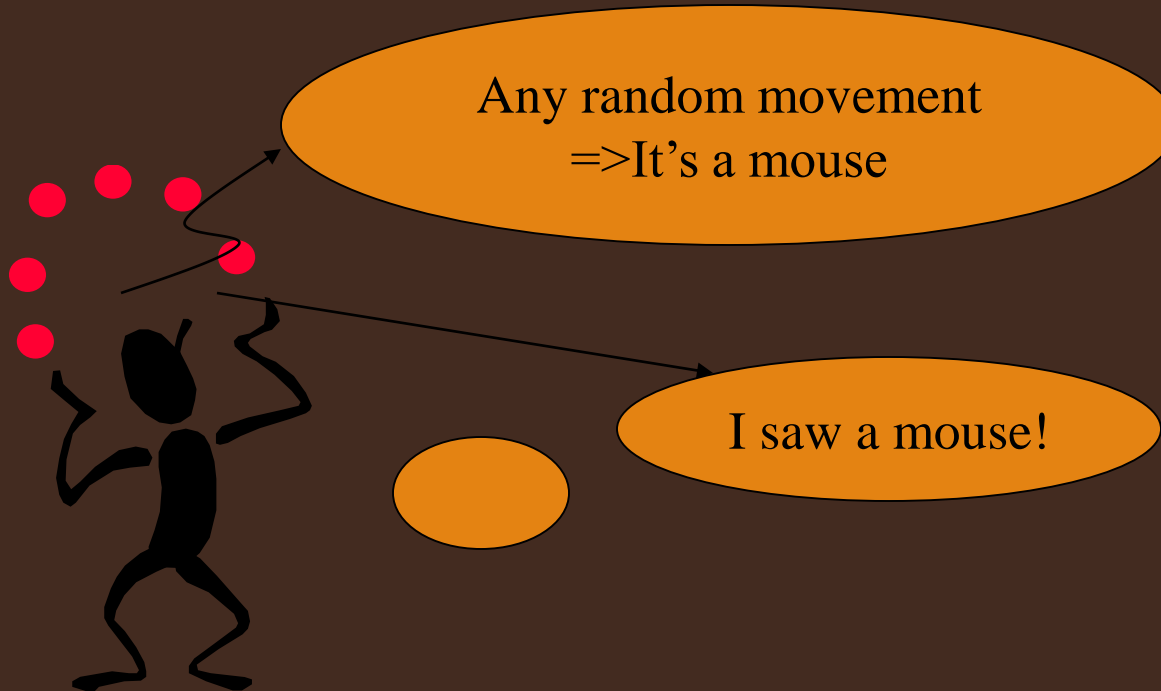
# K-NN Fundamentals

- **Requires three things**
  - The set of stored records.
  - Distance Metric to compute distance between records.
  - The value of k, the number of nearest neighbors to retrieve.
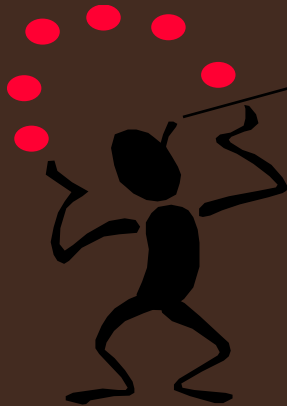
- **To classify an unknown record:**
  - Compute distance to other training records.
  - Identify K- nearest neighbors.
  - Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote).
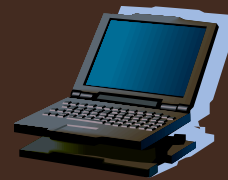
# Different Learning Methods

- Eager Learning

# Instance-based Learning
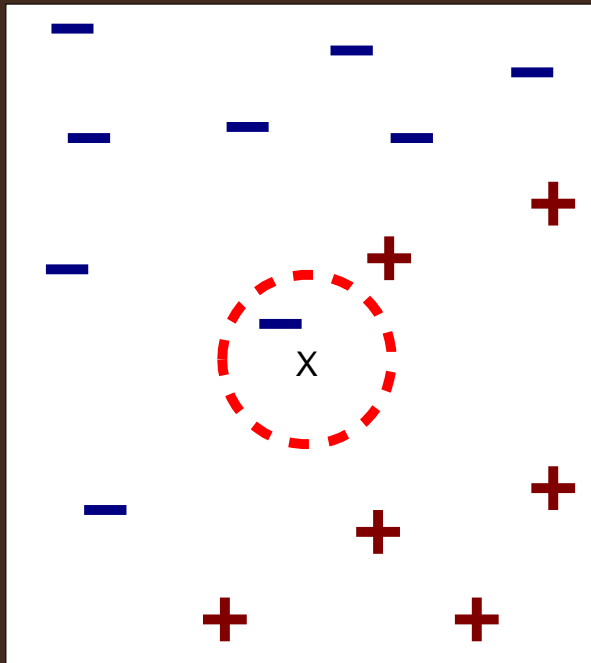
Its very similar to a Desktop!!

# Instance-based Learning

- K-Nearest Neighbor Algorithm
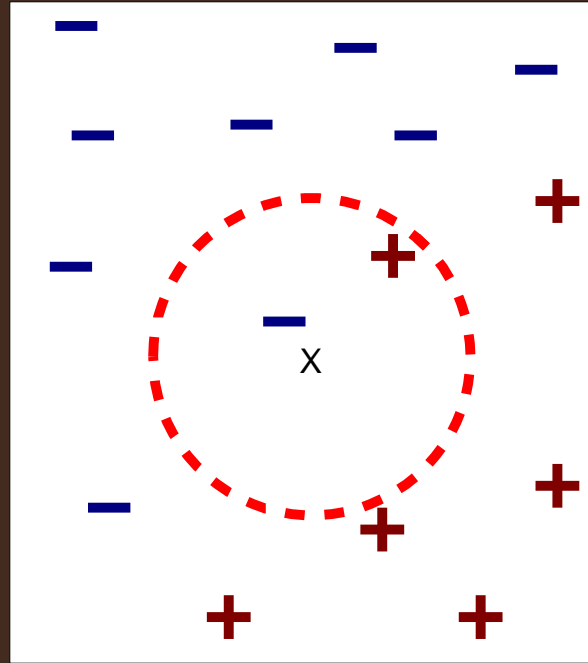- Weighted Regression
- Case-based reasoning

# K-Nearest Neighbor

- Features
  - All instances correspond to points in an n-dimensional Euclidean space
  - Classification is delayed till a new instance arrives
  - Classification done by comparing feature vectors of the different points
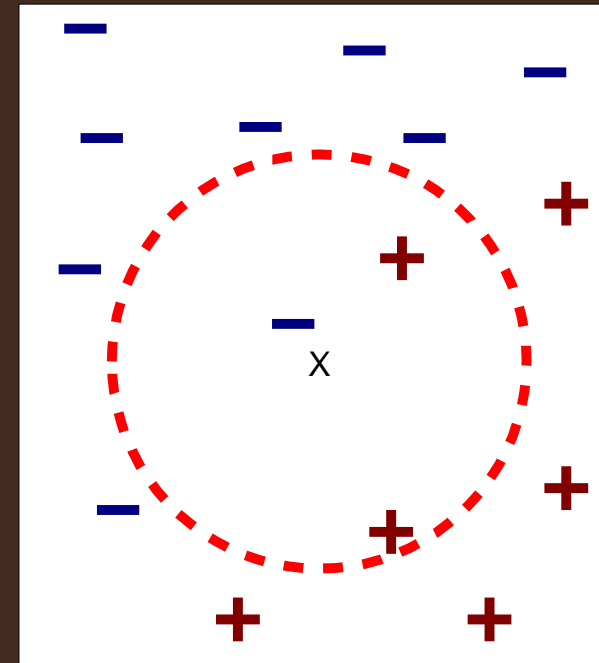  - Target function may be discrete or real-valued
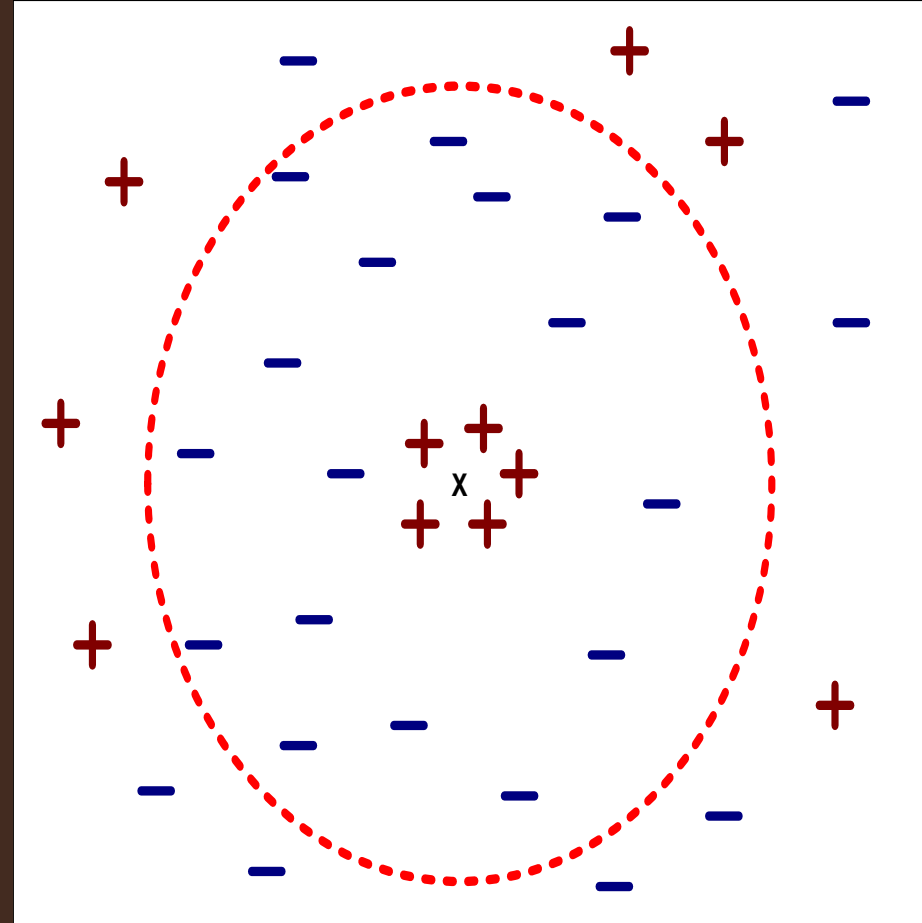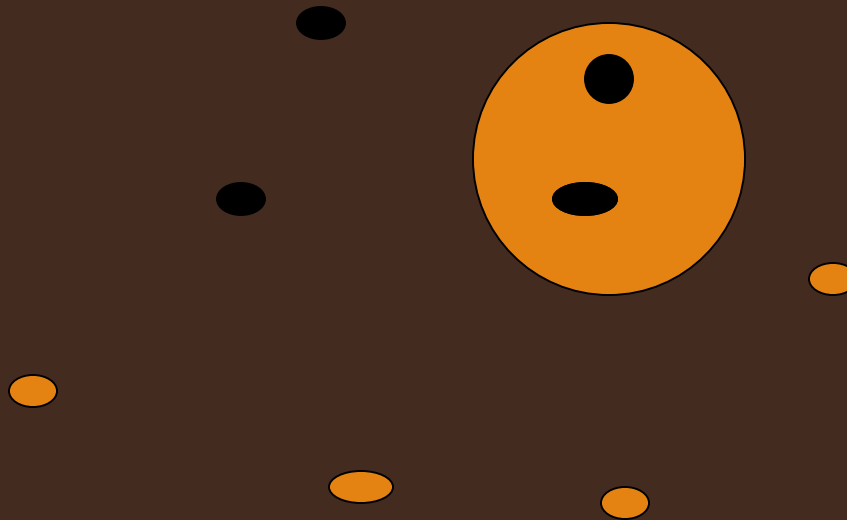
# K-? Nearest Neighbor
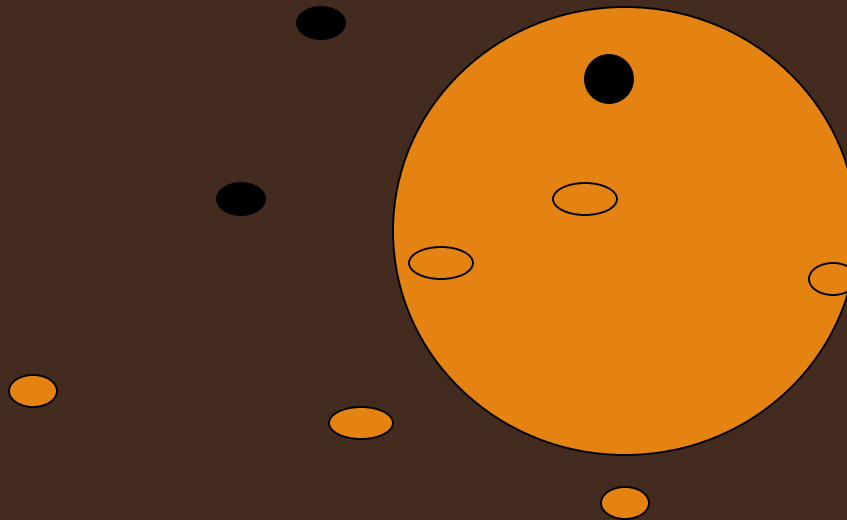


1-NN

2-NN

3-NN

# Select K

- If K is too small, sensitive to noise points.

- If K is too large, neighborhood may include points from other classes.

- Always taken a odd value.

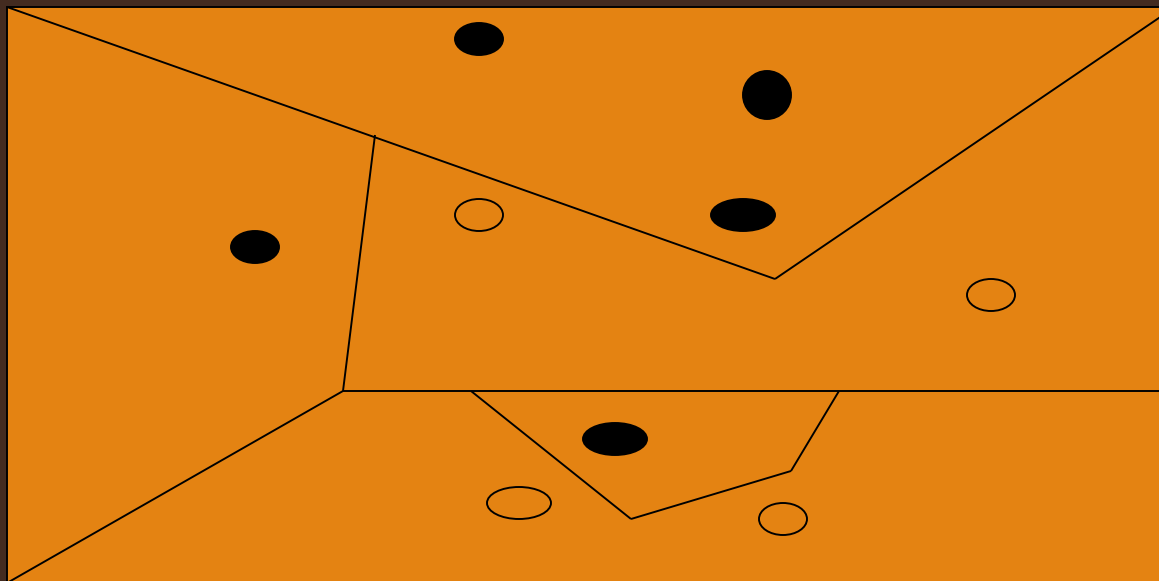# K-Nearest Neighbor

- An arbitrary instance is represented by $(a_1(x), a_2(x), a_3(x),.., a_n(x))$
  - $a_i(x)$ denotes features

- Euclidean distance between two instances

  $d(x_i, x_j) = sqrt$ (sum for r=1 to n $(a_r(x_i) - a_r(x_j))^2$)

- Continuous valued target function
  - mean value of the k nearest training examples

# Voronoi Diagram

- Decision surface formed by the training examples

# Distance-Weighted Nearest Neighbor Algorithm

- Assign weights to the neighbors based on their 'distance' from the query point
  - Weight 'may' be inverse square of the distances

→ All training points may influence a particular instance
  - Shepard's method

# Remarks

- Curse of Dimensionality

# Remarks

- Curse of Dimensionality

# Remarks

- Curse of Dimensionality

# Remarks

- Efficient memory indexing
  - To retrieve the stored training examples (kd-tree)

# K-NN Model



**Compute Distance**

**Test Record**

**Training Records**

**Choose k of the "nearest" records**

# K-NN Algorithm

1. Load the data
2. Initialize K to your chosen number of neighbors
3. For each example in the data
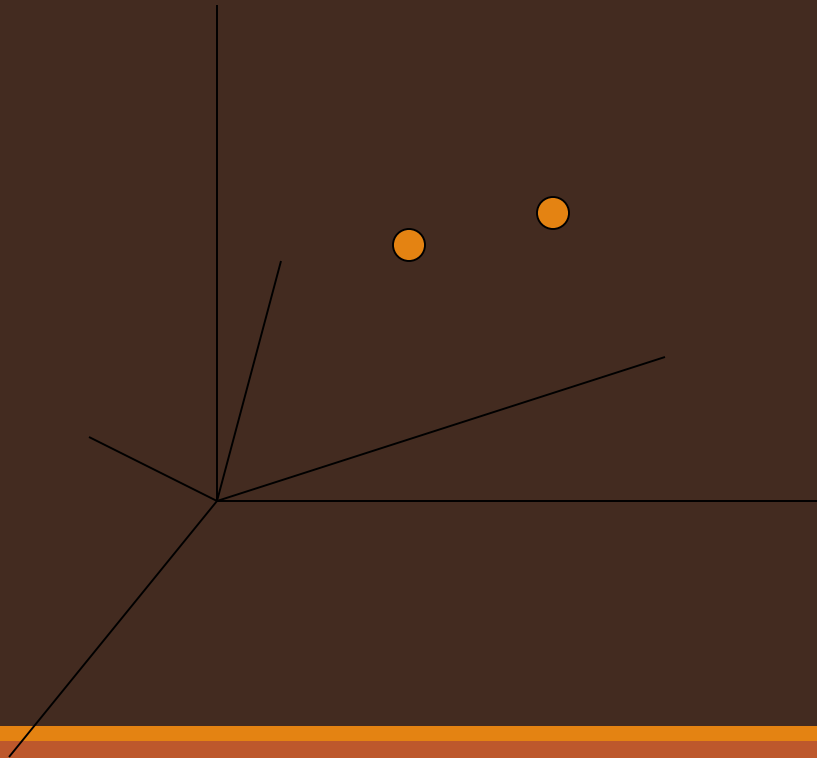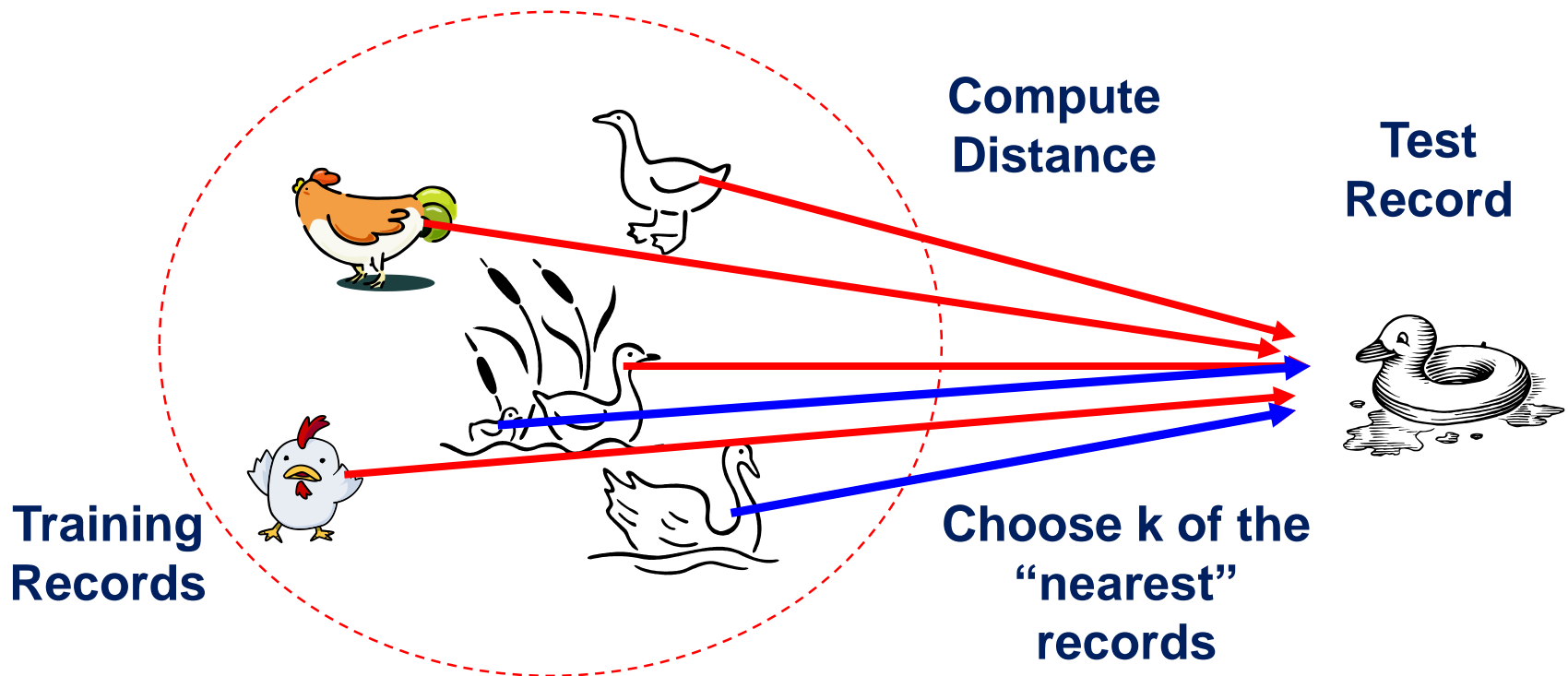   1. Calculate the distance between the query example and the current example from the data.
   2. Add the distance and the index of the example to an ordered collection
4. Sort the ordered collection of distances and indices from smallest to largest (in ascending order) by the distances.
5. Pick the first K entries from the sorted collection.
6. Get the labels of the selected K entries
7. If regression, return the mean of the K labels
8. If classification, return the mode of the K labels

# **Advantage**

- The algorithm is simple and easy to implement.

- There's no need to build a model, tune several parameters, or make additional assumptions.

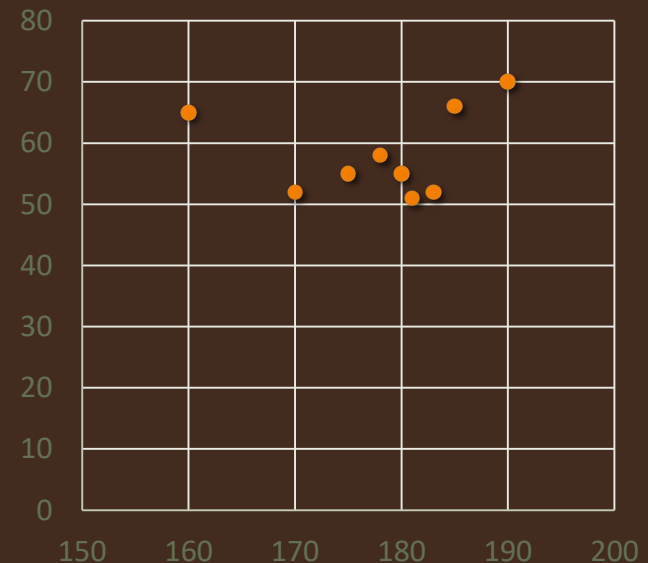- The algorithm is versatile. It can be used for classification, regression, and search.

# Disadvantage

- The algorithm gets significantly slower as the number of examples and/or predictors/independent variables increase.

# Example1

- Consider a dataset of health status as given in Table. What will be the status of a sample *(180, 65).*

| S.No. | Height | Weight | Status | Distance |
|-------|--------|--------|-----------|----------|
| 1 | 160 | 65 | Unhealthy | 20 |
| 2 | 170 | 52 | Healthy | 16.40122 |
| 3 | 175 | 55 | Healthy | 11.18034 |
| 4 | 178 | 58 | Healthy | 7.28011 |
| 5 | 180 | 55 | Unhealthy | 10 |
| 6 | 181 | 51 | Unhealthy | 14.03567 |
| 7 | 183 | 52 | Unhealthy | 13.34166 |
| 8 | 185 | 66 | Healthy | 5.09902 |
| 9 | 190 | 70 | Healthy | 11.18034 |



$$\text{Euclidean Distance } (D) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

# Example 1...

*Test Sample :(180, 65).*

| S.No. | Height | Weight | Status | Distance | |
|-------|--------|--------|--------|----------|---|
| 8 | 185 | 66 | Healthy | 5.09902 | K=1 |
| 4 | 178 | 58 | Healthy | 7.28011 | K=2 |
| 5 | 180 | 55 | Unhealthy | 10 | K=3 |
| 3 | 175 | 55 | Healthy | 11.18034 | |
| 9 | 190 | 70 | Healthy | 11.18034 | |
| 7 | 183 | 52 | Unhealthy | 13.34166 | |
| 6 | 181 | 51 | Unhealthy | 14.03567 | |
| 2 | 170 | 52 | Healthy | 16.40122 | |
| 1 | 160 | 65 | Unhealthy | 20 | |

K=1, Healthy; K=2, Healthy; K=3, Healthy

10/14/2020

# Example 2

- **Sanction of loan amount**

# K-NN Classification

| Age | Loan | Default | Distance |
|-----|------|---------|----------|
| 25 | $40,000 | N | 102000 |
| 35 | $60,000 | N | 82000 |
| 45 | $80,000 | N | 62000 |
| 20 | $20,000 | N | 122000 |
| 35 | $120,000 | N | 22000 |
| 52 | $18,000 | N | 124000 |
| 23 | $95,000 | Y | 47000 |
| 40 | $62,000 | Y | 80000 |
| 60 | $100,000 | Y | 42000 |
| 48 | $220,000 | Y | 78000 |
| 33 | $150,000 | Y | 8000 |
| | | | |
| 48 | $142,000 | ? | |

$$\text{Euclidean Distance } (D) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$