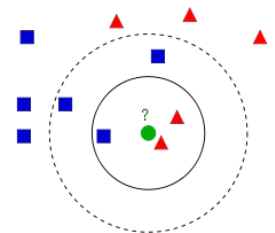


K-Nearest Neighbour

Tushar B. Kute,
<http://tusharkute.com>



What sort of Machine Learning?

- An idea that can be used for machine learning—as does another maxim involving poultry: "birds of a feather flock together."
- In other words, things that are alike are likely to have properties that are alike.
- We can use this principle to classify data by placing it in the category with the most similar, or "nearest" neighbors.

Nearest Neighbor Classification

- In a single sentence, nearest neighbor classifiers are defined by their characteristic of classifying unlabeled examples by assigning them the class of the most similar labeled examples. Despite the simplicity of this idea, nearest neighbor methods are extremely powerful. They have been used successfully for:
 - Computer vision applications, including optical character recognition and facial recognition in both still images and video
 - Predicting whether a person enjoys a movie which he/she has been recommended (as in the Netflix challenge)
 - Identifying patterns in genetic data, for use in detecting specific protein or diseases

The kNN Algorithm

- The kNN algorithm begins with a training dataset made up of examples that are classified into several categories, as labeled by a nominal variable.
- Assume that we have a test dataset containing unlabeled examples that otherwise have the same features as the training data.
- For each record in the test dataset, kNN identifies k records in the training data that are the "nearest" in similarity, where k is an integer specified in advance.
- The unlabeled test instance is assigned the class of the majority of the k nearest neighbors.

The kNN Algorithm

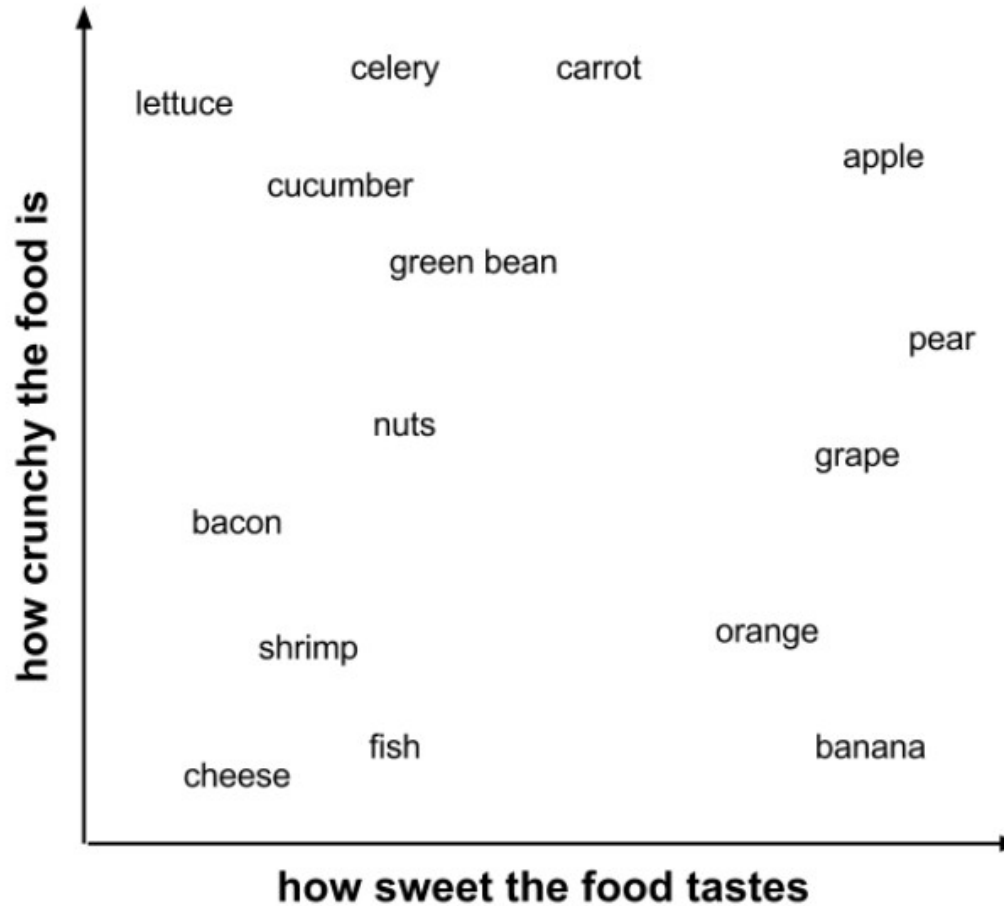


Example:

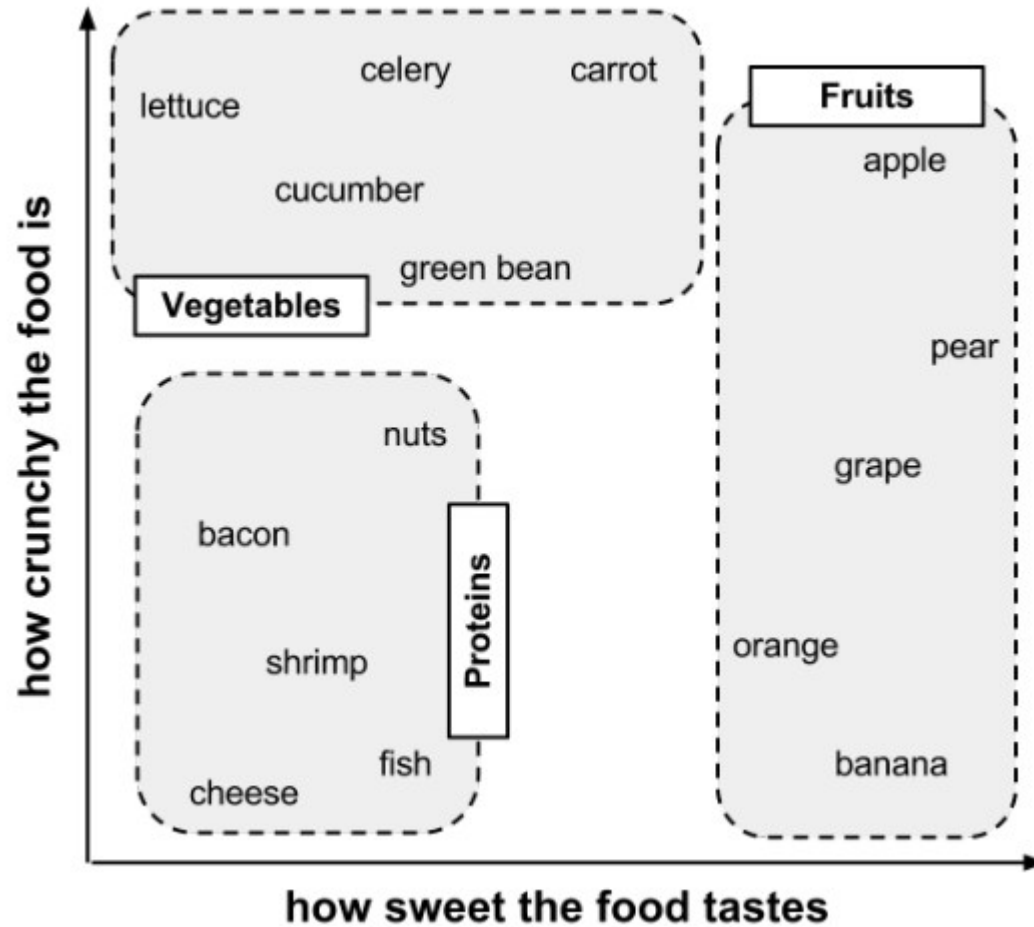
ingredient	sweetness	crunchiness	food type
apple	10	9	fruit
bacon	1	4	protein
banana	10	1	fruit
carrot	7	10	vegetable
celery	3	10	vegetable
cheese	1	1	protein

Reference: Machine Learning with R, Brett Lantz, Packt Publishing

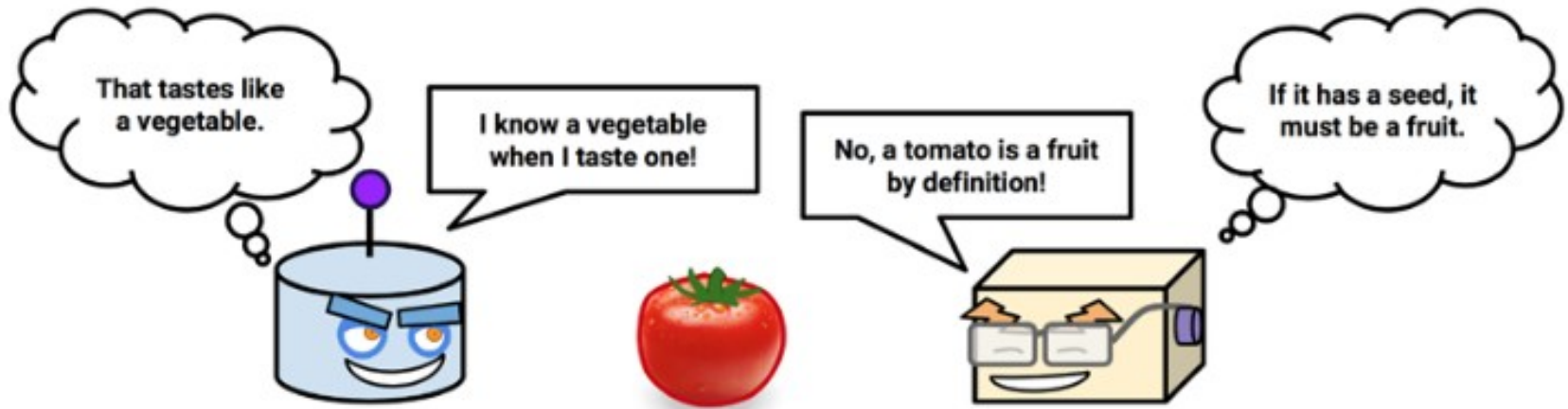
Example:



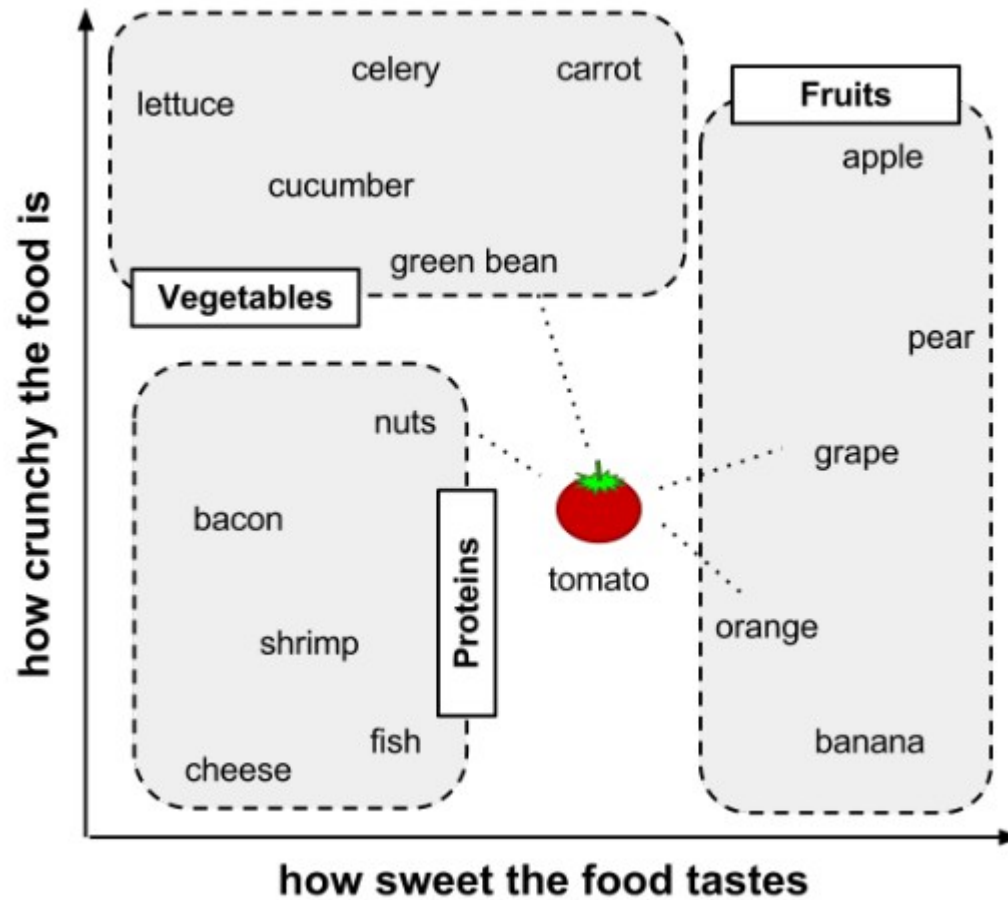
Example:



Classify me now!



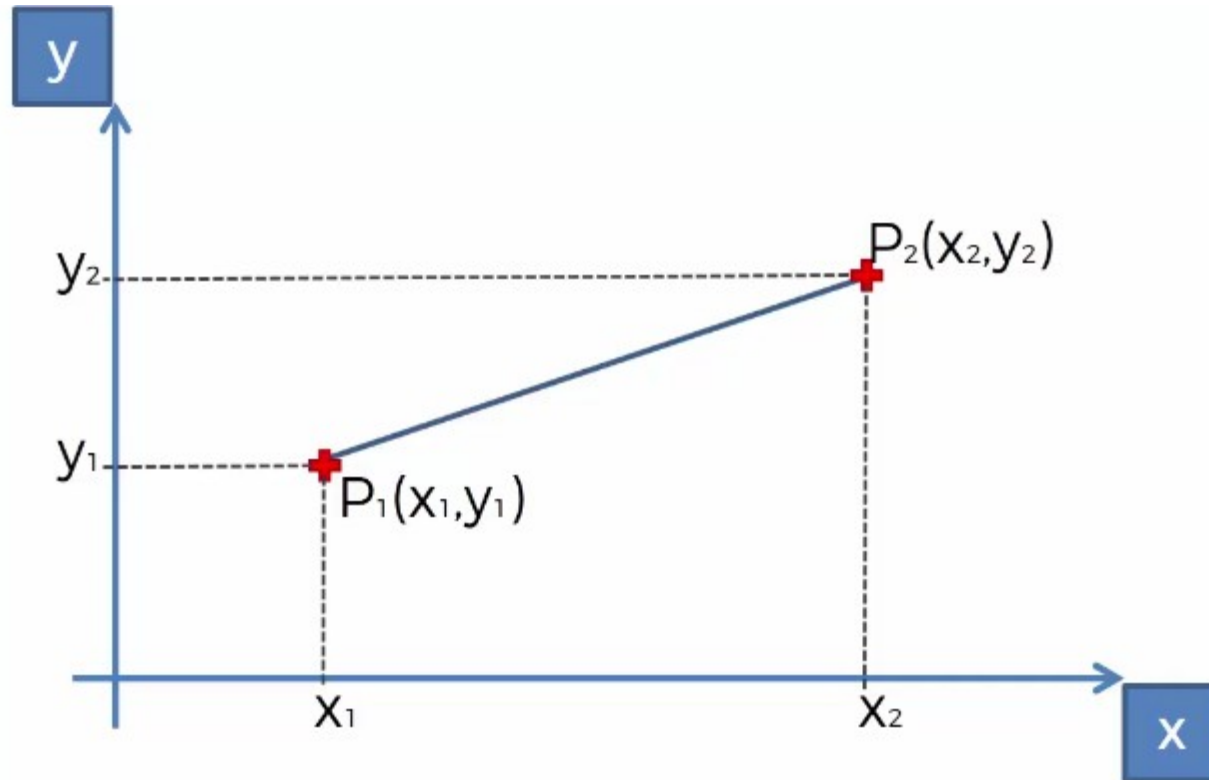
Example:



Calculating Distance

- Locating the tomato's nearest neighbors requires a distance function, or a formula that measures the similarity between two instances.
- There are many different ways to calculate distance.
- Traditionally, the kNN algorithm uses Euclidean distance, which is the distance one would measure if you could use a ruler to connect two points, illustrated in the previous figure by the dotted lines connecting the tomato to its neighbors.

Calculating Distance



$$\text{Euclidean Distance between } P_1 \text{ and } P_2 = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Reference: Super Data Science

Calculating Distance

- Euclidean distance is specified by the following formula, where p and q are the examples to be compared, each having n features. The term p_1 refers to the value of the first feature of example p , while q_1 refers to the value of the first feature of example q :

$$\text{dist}(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$$

- The distance formula involves comparing the values of each feature. For example, to calculate the distance between the tomato (sweetness = 6, crunchiness = 4), and the green bean (sweetness = 3, crunchiness = 7), we can use the formula as follows:

$$\text{dist}(\text{tomato}, \text{green bean}) = \sqrt{(6 - 3)^2 + (4 - 7)^2} = 4.2$$

Distance Functions

Euclidean

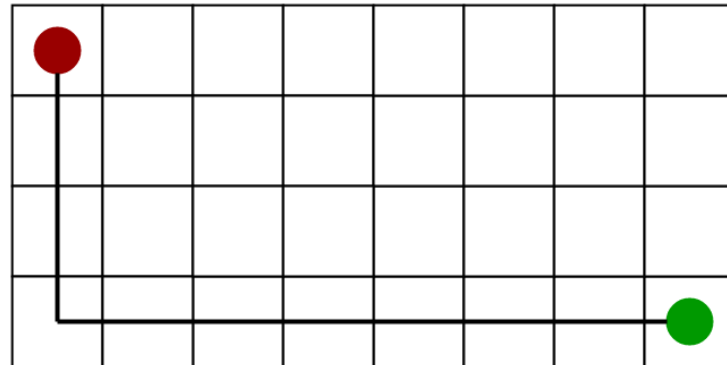
$$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

Manhattan

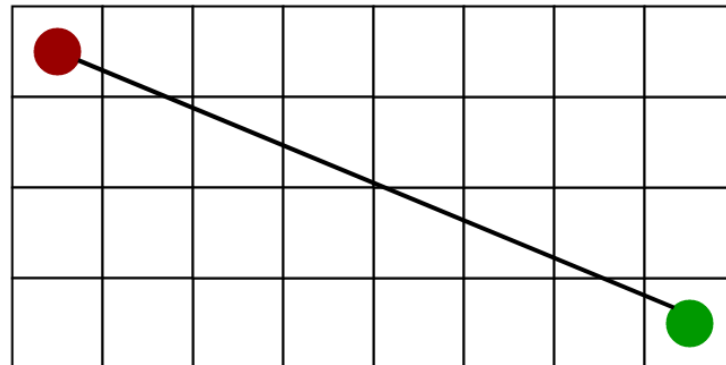
$$\sum_{i=1}^k |x_i - y_i|$$

Distance

Manhattan Distance



Euclidean Distance



Closest Neighbors

ingredient	sweetness	crunchiness	food type	distance to the tomato
grape	8	5	fruit	$\sqrt{(6 - 8)^2 + (4 - 5)^2} = 2.2$
green bean	3	7	vegetable	$\sqrt{(6 - 3)^2 + (4 - 7)^2} = 4.2$
nuts	3	6	protein	$\sqrt{(6 - 3)^2 + (4 - 6)^2} = 3.6$
orange	7	3	fruit	$\sqrt{(6 - 7)^2 + (4 - 3)^2} = 1.4$

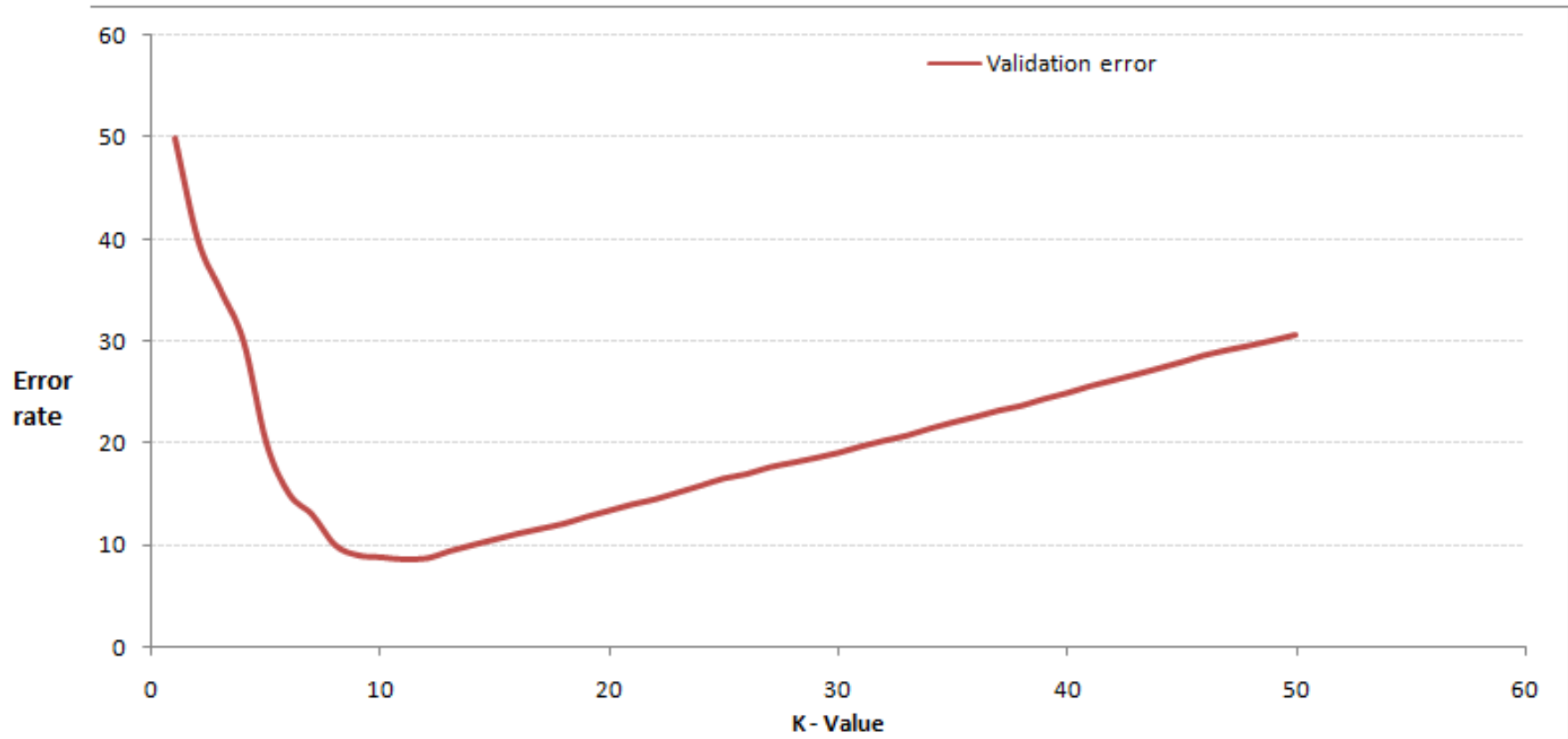
Summary Steps

- Step-1: Select the number K of the neighbors
- Step-2: Calculate the Euclidean distance of K number of neighbors
- Step-3: Take the K nearest neighbors as per the calculated Euclidean distance.
- Step-4: Among these k neighbors, count the number of the data points in each category.
- Step-5: Assign the new data points to that category for which the number of the neighbor is maximum.
- Step-6: Our model is ready.

Choosing Appropriate K

- Deciding how many neighbors to use for kNN determines how well the model will generalize to future data.
- The balance between overfitting and underfitting the training data is a problem known as the bias-variance tradeoff.
- Choosing a large k reduces the impact or variance caused by noisy data, but can bias the learner such that it runs the risk of ignoring small, but important patterns.

Choosing Appropriate K



Choosing Appropriate K

- In practice, choosing k depends on the difficulty of the concept to be learned and the number of records in the training data.
- Typically, k is set somewhere between 3 and 10. One common practice is to set k equal to the square root of the number of training examples.
- In the classifier, we might set $k = 4$, because there were 15 example ingredients in the training data and the square root of 15 is 3.87.

Min-Max normalization

- The traditional method of rescaling features for kNN is min-max normalization.
- This process transforms a feature such that all of its values fall in a range between 0 and 1. The formula for normalizing a feature is as follows. Essentially, the formula subtracts the minimum of feature X from each value and divides by the range of X :

$$X_{new} = \frac{X - \min(X)}{\max(X) - \min(X)}$$

- Normalized feature values can be interpreted as indicating how far, from 0 percent to 100 percent, the original value fell along the range between the original minimum and maximum.

The Lazy Learning

- Using the strict definition of learning, a lazy learner is **not** really learning anything.
- Instead, it merely stores the training data in it. This allows the training phase to occur very rapidly, with a potential downside being that the process of making predictions tends to be relatively slow.
- Due to the heavy reliance on the training instances, lazy learning is also known as **instance-based learning** or **rote learning**.

Few Lazy Learning Algorithms

- K Nearest Neighbors
- Local Regression
- Lazy Naive Bayes

The kNN Algorithm

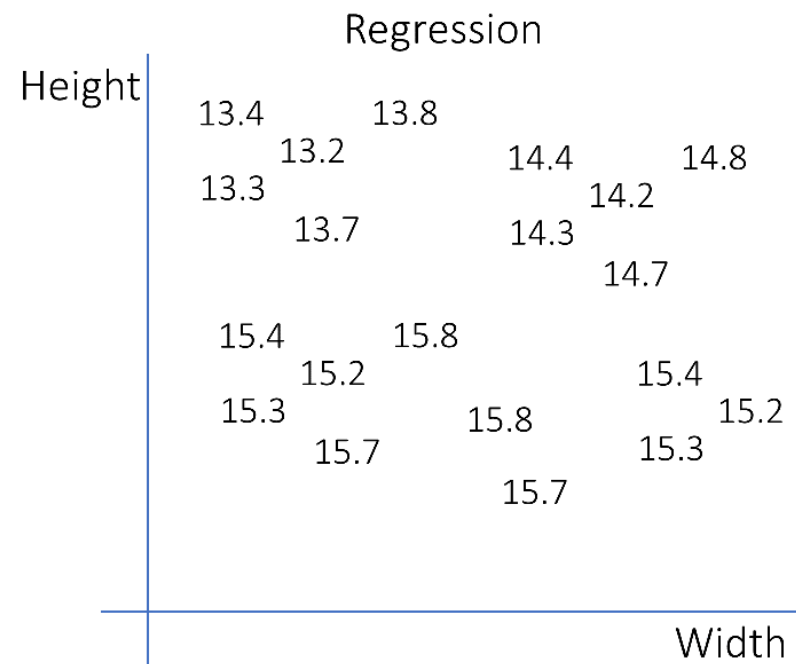
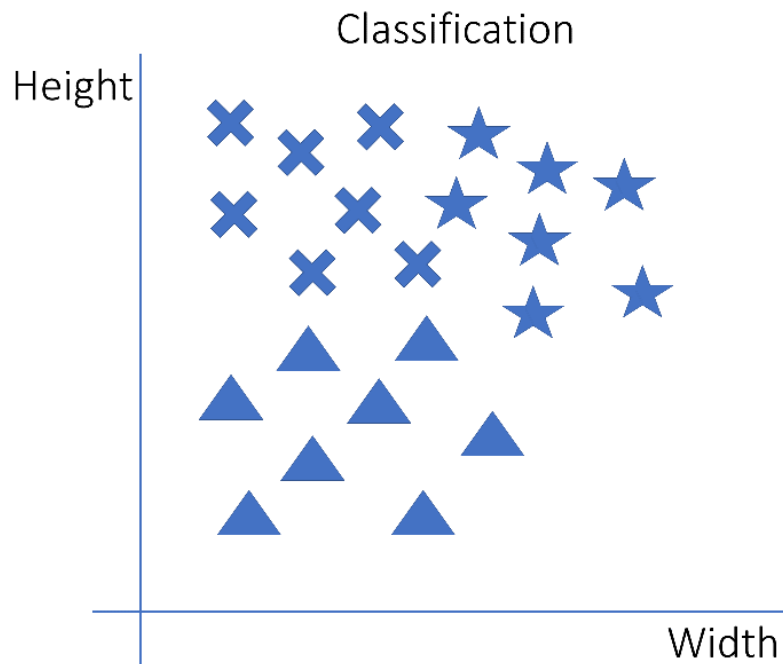
Strengths

- Simple and effective
- Makes no assumptions about the underlying data distribution
- Fast training phase

Weaknesses

- Does not produce a model, which limits the ability to find novel insights in relationships among features
 - Slow classification phase
 - Requires a large amount of memory
 - Nominal features and missing data require additional processing
-

Classification vs. Regression



Python Packages needed

- pandas
 - Data Analytics
- numpy
 - Numerical Computing
- matplotlib.pyplot
 - Plotting graphs
- sklearn
 - Classification and Regression Classes

Useful resources

- www.mitu.co.in
- www.pythonprogramminglanguage.com
- www.scikit-learn.org
- www.towardsdatascience.com
- www.medium.com
- www.analyticsvidhya.com
- www.kaggle.com
- www.stephacking.com
- www.github.com

Thank you

This presentation is created using LibreOffice Impress 5.1.6.2, can be used freely as per GNU General Public License



@mitu_skillologies



/MITuSkillologies



@mitu_group



/company/mitu-
skillologies



MITUSkillologies

Web Resources

<https://mitu.co.in>

<http://tusharkute.com>

contact@mitu.co.in

tushar@tusharkute.com