

Step 07

# MODEL SELECTION

- Model selection
- Key considerations
- Classification of machine learning
- Different algorithms with their use cases and examples



## Model Selection

- Choosing the right machine learning model is essential to obtain accurate predictions.
- The selection of the appropriate model depends on the type of problem to be solved and the nature of data.

So how do we know which one to choose for our specific problem?

Here are a few things to consider when selecting a model:

**Problem type:** The type of problem we are trying to solve (e.g. classification, regression, clustering) can help narrow down the possible algorithms that are suitable for our task.

- **Size of data:** Some algorithms work better with large datasets, while others may perform well even with small amounts of data.
- **Interpretability:** Some models, like decision trees and linear regression, are easy to interpret and understand, while others like neural networks may be more difficult to interpret.
- **Performance:** The performance of a model is crucial, and this can be evaluated through metrics such as accuracy, precision, recall, and F1-score.
- **Time and resource constraints:** Some models may require more time and computational resources to train and run, so it's important to consider the available resources when selecting a model.

Once you understand the nature of the data, the next step is to choose the type of algorithm that best fits the data. There are three types of algorithms,

## **01 - Supervised Learning Algorithms**

These algorithms are used when the data is labeled, and the goal is to predict an output variable based on input variables. Some of the most popular algorithms are,

- Linear Regression
- Logistic Regression
- Decision Trees
- Random Forest
- Naive Bayes
- K-Nearest Neighbors (KNN)
- Support Vector Machines (SVM)

## **02 - Unsupervised Learning Algorithms**

These algorithms are used when the data is unlabeled, and the goal is to find hidden patterns or groups in the data. Some of the common algorithms like,

- K-Means Clustering
- Hierarchical Clustering
- Principal Component Analysis (PCA)
- Apriori Algorithm

## **03 - Semi-Supervised Learning Algorithms**

These algorithms are used when only some of the data is labeled, and the goal is to predict the labels of the remaining data. Examples include,

- Label Propagation
- Label Spreading,
- Self-Training.

Let's discuss some popular supervised machine learning algorithms and their use cases based on the nature of data and type of problem:

### **Linear Regression**

- **Use case:** Used for predicting a continuous variable based on one or more independent variables and the relation between input and output is linear.
- **Example:** It can be used to predict the price of a house based on its size, location, and other relevant factors.

### **Logistic Regression**

- **Use case:** Logistic regression is used for binary classification problems, where the output variable takes only two values (e.g., yes or no, true or false).

- **Example:** Predicting whether a customer will churn or not based on their past purchase behavior.

## **Decision trees**

- **Use case:** Used for both classification and regression problems where the data has a hierarchical structure.
- **Example:** Predicting whether a customer will buy a particular product based on their demographic details, browsing history, etc.

## **Random Forest**

- It is an ensemble learning method that uses multiple decision trees to improve the accuracy of predictions.

- **Use case:** Random forest is useful when dealing with high-dimensional datasets with many input variables, where decision trees may suffer from overfitting.
- **Example:** Predicting whether a customer will default on a loan based on their financial history, credit score, etc.

## Naive Bayes

- **Use case :** Used for classification problems based on Bayes' theorem, which assumes that the features are independent of each other.
- **Example:** Spam filtering based on the text content of an email.



## **K-Nearest Neighbors (KNN)**

- **Use case:** Used for both classification and regression problems. It determines the class or value of a data point by looking at its k-nearest neighbors in the training dataset.
- **Example:** Predicting the price of a house based on the prices of similar houses in the neighborhood.

## **Support Vector Machines (SVM)**

- **Use case:** Used for classification and regression problems. It separates the data into classes by finding the optimal hyperplane with the maximum margin between the classes.
- **Example:** Predicting whether a customer will default on a loan based on their financial history, credit score, etc.

Let's, discuss some common unsupervised learning algorithms with their use case and examples,

### **K-Means**

- **Use case:** Used for clustering problems. It divides the data into  $k$  clusters based on their similarity.
- **Example:** Segmenting customers into different groups based on their purchase behavior.

### **Hierarchical Clustering**

**Use case:** Used for clustering problems. It creates a hierarchy of clusters by recursively merging or splitting them based on their similarity.

**Example:** Grouping documents based on their similarity in natural language processing.

## Principal Component Analysis (PCA)

- **Use case:** Used for dimensionality reduction. It transforms high-dimensional data into a lower-dimensional space while retaining most of the variance in the data.
- **Example:** Reducing the dimensionality of image data for faster processing in computer vision applications.

Here are the some popular semi-supervised learning algorithms with examples,

## Self-training

- **Use case:** In this algorithm, a model is trained on a small labeled dataset and then used to label a larger unlabeled dataset. The newly labeled data is added to the labeled dataset, and the process is repeated.

- **Example:** A model may be trained on a small set of labeled images of dogs and cats, and then used to label a larger set of unlabeled images. The newly labeled data is added to the labeled set and the model is retrained on the larger labeled set.

### Co-training

- **Use case:** This algorithm uses two classifiers, each of which is trained on a different subset of the features. The classifiers then label the unlabeled data, and the newly labeled data is used to retrain the classifiers.
- **Example:** In a text classification problem, one classifier may be trained on the text itself, while the other is trained on the metadata associated with the text, such as the author or publication date.

**Note: There are many other machine learning algorithms out there, but these are some of the most popular ones.**

It's important to keep in mind that there is no one-size-fits-all solution, and the best algorithm for a particular problem may vary depending on the data and the objectives.

In summary, selecting the right algorithm for our data is crucial for building an accurate and reliable machine learning model.

By understanding the nature of the data and the type of problem, we can narrow down our choices and choose the most suitable algorithm.

Don't settle for good enough, take the time to explore and experiment with different algorithms to find the best one for our data.

Happy learning!😊

# THANK YOU

@ Gangadhar Neelam