

Selection of the appropriate machine learning technique

The selection of an appropriate machine learning technique depends on several factors, including the nature of the problem, the characteristics of the data, and the specific goals of the task.

It's important to note that model selection may involve experimentation and comparison of multiple algorithms to identify the one that performs best for a specific task. Additionally, the choice of a machine learning technique is not fixed, and it may evolve as the understanding of the problem and the data deepens.

Some key considerations to help guide the choice of a machine learning technique are the following:

Nature of the Problem:

- **Supervised vs. Unsupervised:** Determine whether the problem is supervised (labeled data for training) or unsupervised (unlabeled data).
- **Classification vs. Regression:** Decide whether the task involves predicting categories (classification) or numerical values (regression).

Type of Data:

- **Structured vs. Unstructured:** For structured data (tabular), traditional machine learning methods like decision trees, random forests, and support vector machines may be effective. For unstructured data (text, images, audio), deep learning methods such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs) might be suitable.
- **Dimensionality:** If the dataset has a high dimensionality, techniques for dimensionality reduction (e.g., PCA for unsupervised learning) may be beneficial.
- **Feature Characteristics:** Depending on the characteristics of the features, different algorithms may be more suitable. For example, decision trees can handle categorical features well, while neural networks can automatically learn representations from raw data.

Data Size:

- **Small vs. Large Datasets:** Deep learning models often require large amounts of data for effective training. If the dataset is small, traditional machine learning algorithms may be more appropriate.

Interpretability:

- **Interpretability Requirements:** Consider the need for interpretable models. If interpretability is crucial, simpler models like decision trees or linear regression may be preferred over complex models like deep neural networks.

Computational Resources:

- **Computational Constraints:** Evaluate the available computational resources. Deep learning models can be computationally intensive, especially for training large neural networks.

Model Performance Metrics:

For classification tasks, metrics like accuracy, precision, recall, and F1 score are common. For regression, metrics like mean squared error (MSE) or mean absolute error (MAE) may be used.

Training, Test and Validation

In machine learning (ML), the use of training and test data is fundamental to the model development process.

The **training data** is the portion of the dataset used to train the machine learning model. It consists of input data points (features) and corresponding labels (target variables). During the training phase, the model learns patterns and relationships in the data, adjusting its parameters to minimize the difference between predicted and actual values.

The **test data**, on the other hand, is a separate subset of the dataset that is not used during the training phase. It is reserved for evaluating the performance of the trained model. The test data also consists of input features and corresponding labels, but the model uses this data solely to assess its generalization ability — how well it performs on new, unseen data.

In addition to training and test data, a third subset called **validation data** is often used during the model development process. The validation data is used to fine-tune model hyperparameters and to further evaluate model performance during training.

Cross-validation is an important technique used to assess model performance, especially when the dataset is limited. It involves splitting the dataset into multiple subsets (folds), with each fold used both as a training set and as a test set in separate iterations. This helps provide a more reliable estimate of the model's performance across different subsets of the data.

Regression metrics

1. Mean Square Error (MSE)

It measures how closely the expected response value for a certain observation matches the actual response value for that observation.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

2. Mean Absolute Error (MAE)

It measures how closely the expected response value for a certain observation matches the actual response value for that observation in an absolute value.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Classification metrics

Error rate, accuracy, precision, recall, and F1 score

1. Error rate

The most popular method for measuring the precision of our estimate in classification contexts is the error rate, which is the percentage of errors that occur when we apply our estimate to the observations:

$$\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$$

2. Confusion matrix, precision, recall, specificity, sensitivity, accuracy and F1-score

A confusion matrix is a performance evaluation tool used in machine learning to assess the performance of a classification model. It provides a summary of the predictions made by the model compared to the actual ground truth labels in a tabular format, as shown in Table 1.

		Predicted class		
		– or Null	+ or Non-null	Total
True class	– or Null	True Neg. (TN)	False Pos. (FP)	N
	+ or Non-null	False Neg. (FN)	True Pos. (TP)	P
Total		N*	P*	

Table 1: Confusion matrix

- True Positive (TP): The model correctly predicts a positive class instance as positive.
- False Positive (FP): The model incorrectly predicts a negative class instance as positive (Type I error).
- True Negative (TN): The model correctly predicts a negative class instance as negative.
- False Negative (FN): The model incorrectly predicts a positive class instance as negative (Type II error).

From the confusion matrix, various performance metrics can be derived to evaluate the model, such as:

- *Accuracy*: The proportion of correct predictions out of the total predictions.

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}}$$

- *Precision*: The proportion of true positive predictions out of all positive predictions.

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

- *Recall (Sensitivity)*: The proportion of true positive predictions out of all actual positive instances.

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

- *F1 Score*: The harmonic mean of precision and recall, providing a balanced measure of both metrics.

$$F1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- *Specificity*: The proportion of true negative predictions out of all actual negative instances.
- *False Positive Rate (FPR)*: The proportion of false positive predictions out of all actual negative instances.