



---

# CONTENTS

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>The Training/Validation/Testing Process</b> | <b>3</b>  |
| 1.0.1    | Training Set                                   | 3         |
| 1.0.2    | Validation Set                                 | 4         |
| 1.0.3    | Testing Set                                    | 4         |
| <b>2</b> | <b>Evaluating Classification Systems</b>       | <b>5</b>  |
| 2.1      | Definition of a Confusion Matrix               | 5         |
| 2.2      | Performance Metrics                            | 6         |
| <b>3</b> | <b>Evaluation Binary Classifiers</b>           | <b>7</b>  |
| 3.0.1    | Binary Classification                          | 7         |
| 3.0.2    | Accuracy                                       | 8         |
| 3.0.3    | Precision and Recall                           | 8         |
| 3.0.4    | F1 Score                                       | 9         |
| <b>4</b> | <b>The k-Nearest Neighbor Classifier</b>       | <b>11</b> |
| 4.1      | The k-NN Algorithm                             | 11        |
| 4.2      | Choosing the Right 'k'                         | 12        |
| 4.3      | Considerations                                 | 12        |
| <b>5</b> | <b>Naive Bayes Classifier</b>                  | <b>13</b> |
| 5.1      | Bayes' Theorem                                 | 13        |
| 5.2      | The Naivety of Naive Bayes                     | 14        |
| 5.3      | Working of Naive Bayes Classifier              | 14        |

---

|              |                             |           |
|--------------|-----------------------------|-----------|
| <b>A</b>     | <b>Answers to Exercises</b> | <b>15</b> |
| <b>Index</b> |                             | <b>17</b> |



## CHAPTER 1

---

# The Training/Validation/Testing Process

In machine learning, it's essential to assess the performance of a model accurately. This assessment helps us choose the best model and tune its parameters. The data used to develop machine learning models is typically divided into three sets: training, validation, and testing.

### 1.0.1 Training Set

The training set is used to train the model, i.e., to adjust the model's weights and biases in the case of neural networks, or to determine the best split in decision trees, among other things. The model learns from this data, which is why it's called the "training" set.

### 1.0.2 Validation Set

The validation set is used to tune model parameters (hyperparameters), to choose the model architecture (for example, the number of hidden layers in a neural network), or to determine the degree of the polynomial in polynomial regression, among other uses. This set provides an unbiased evaluation of a model fit on the training dataset while tuning model hyperparameters.

### 1.0.3 Testing Set

The testing set is used to provide an unbiased evaluation of the final model fit on the training dataset. The test set serves as a proxy for real-world data that the model has not seen before. It's important to only use the test set once, after all training and validation is complete, to avoid "leaking" information from the test set into the model.

The key to this process is that each set of data is separate and independent. Mixing data between the sets can lead to overly optimistic or pessimistic assessments of a model's performance.

In practice, the division of data into these three sets can be done randomly (often with 70% for training, 15% for validation, and 15% for testing), or using more structured methods like cross-validation, depending on the amount and nature of the available data.



## CHAPTER 2

---

# Evaluating Classification Systems

The confusion matrix is a tabular method used in machine learning to evaluate the performance of a classification model. It allows for the visualization of the model's performance and to compute various performance metrics.

### 2.1 Definition of a Confusion Matrix

A confusion matrix is a specific table layout that presents the performance of a classification model. For a binary classification problem, it is a 2x2 matrix that compares the actual and the predicted classifications.

|                    | Actual Positive     | Actual Negative     |
|--------------------|---------------------|---------------------|
| Predicted Positive | True Positive (TP)  | False Positive (FP) |
| Predicted Negative | False Negative (FN) | True Negative (TN)  |

## 2.2 Performance Metrics

Using the confusion matrix, we can compute several performance metrics:

- **Accuracy:** The proportion of correct predictions (both true positives and true negatives) among the total number of cases examined. It is calculated as  $(TP + TN) / (TP + TN + FP + FN)$ .
- **Precision:** The proportion of positive identifications that were actually correct. It is calculated as  $TP / (TP + FP)$ .
- **Recall (Sensitivity):** The proportion of actual positives that were identified correctly. It is calculated as  $TP / (TP + FN)$ .
- **Specificity:** The proportion of actual negatives that were identified correctly. It is calculated as  $TN / (TN + FP)$ .
- **F1 Score:** The harmonic mean of precision and recall. It tries to find the balance between precision and recall.  $F1 = 2 * (Precision * Recall) / (Precision + Recall)$ .



## CHAPTER 3

---

# Evaluation Binary Classifiers

Accuracy, recall, precision, and the F1 score are widely used metrics to measure and compare the performance of binary classifiers. This chapter will delve into these evaluation measures, providing insights into their interpretation and practical applications.

### 3.0.1 Binary Classification

Before diving into the evaluation metrics, let's clarify the concept of binary classification. In binary classification, we aim to assign each instance in a dataset to one of two mutually exclusive classes. For example, classifying emails as spam or not spam, identifying whether a patient has a specific medical condition or not, or predicting whether a credit card transaction is fraudulent or legitimate are common binary classification tasks.

To evaluate the performance of a binary classifier, we need metrics that can provide insights into how well the classifier performs in distinguishing between the two classes.

### 3.0.2 Accuracy

Accuracy is a widely used metric for evaluating binary classifiers. It measures the overall correctness of the classifier's predictions by calculating the ratio of correctly classified instances to the total number of instances in the dataset. Mathematically, accuracy can be expressed as:

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

While accuracy provides a general overview of the classifier's performance, it may not be sufficient in certain scenarios. This is especially true when the dataset is imbalanced, meaning that one class significantly outweighs the other in terms of the number of instances.

### 3.0.3 Precision and Recall

Precision and recall are evaluation metrics that provide insights into the classifier's performance on specific classes, allowing us to identify potential trade-offs between false positives and false negatives.

Precision measures the proportion of correctly predicted positive instances (true positives) out of all instances predicted as positive (true positives + false positives). It can be expressed as:

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

Recall, also known as sensitivity or true positive rate, measures the proportion of correctly predicted positive instances (true positives) out of all actual positive instances (true positives + false negatives). Mathematically, recall can be represented as:

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

Precision and recall are complementary metrics. Precision focuses on the quality of positive predictions, while recall emphasizes the classifier's ability to identify positive instances. The choice between precision and recall depends on the specific requirements of the problem at hand.



### 3.0.4 F1 Score

The F1 score combines precision and recall into a single metric, providing a balanced evaluation measure that considers both false positives and false negatives. It is the harmonic mean of precision and recall, and it can be calculated as:

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

The F1 score ranges between 0 and 1, where a value of 1 represents perfect precision and recall. It is particularly useful when we want to strike a balance between precision and recall, considering both the false positives and false negatives in the classifier's predictions.





## CHAPTER 4

---

# The k-Nearest Neighbor Classifier

The k-nearest neighbors (k-NN) algorithm is a type of instance-based learning algorithm used for classification and regression. Given a new, unknown observation, k-NN algorithm searches through the entire dataset to find the 'k' training examples that are closest to the new instance, and predicts the label based on these 'k' nearest neighbors.

### 4.1 The k-NN Algorithm

The algorithm can be summarized as follows:

1. Given a new observation  $x$ , compute the distance between  $x$  and all points in the training set.
2. Identify the 'k' points in the training data that are closest to  $x$ .
3. If k-NN is used for classification, output the most common class among these 'k'

points as the prediction. If k-NN is used for regression, output the average of the values of these 'k' points as the prediction.

The distance between points can be calculated using various metrics, the most common one being the Euclidean distance:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

where  $n$  is the number of features, and  $x_i$  and  $y_i$  are the corresponding features of  $\mathbf{x}$  and  $\mathbf{y}$ .

## 4.2 Choosing the Right 'k'

The choice of 'k' has a significant impact on the k-NN algorithm. A small 'k' (like 1) can capture a lot of noise and lead to overfitting, while a large 'k' can smooth over many details and potentially lead to underfitting. Cross-validation is typically used to select an optimal 'k'.

## 4.3 Considerations

Although the k-NN algorithm is simple to understand and implement, it can be computationally intensive for large datasets, as it requires computing the distance between every pair of points. Additionally, it's sensitive to the choice of the distance metric and the scale of the features.



## CHAPTER 5

---

# Naive Bayes Classifier

The Naive Bayes classifier is a simple yet effective algorithm for classification tasks. It is a probabilistic classifier based on Bayes' theorem and some additional simplifying assumptions, which make it particularly suitable for high-dimensional datasets.

### 5.1 Bayes' Theorem

Bayes' theorem describes the relationship of conditional probabilities of statistical quantities. In the context of classification, it can be written as:

$$P(C|X) = \frac{P(X|C) \cdot P(C)}{P(X)}$$

where:

- $P(C|X)$  is the posterior probability of class  $C$  given predictor (features)  $X$ .

- $P(C)$  is the prior probability of class.
- $P(X|C)$  is the likelihood which is the probability of predictor given class.
- $P(X)$  is the prior probability of predictor.

## 5.2 The Naivety of Naive Bayes

The "Naive" in Naive Bayes comes from the assumption that each feature in the dataset is independent of all other features, given the class. This is a strong (and often unrealistic) assumption, hence the name "naive". Despite this unrealistic assumption, the Naive Bayes classifier often performs well in practice.

In the context of text classification, this naivety translates into assuming that every word in a document is independent of all other words, given the document's class.

## 5.3 Working of Naive Bayes Classifier

When given an instance to classify, the Naive Bayes classifier calculates the posterior probability of that instance belonging to each possible class. The classifier then outputs the class with the highest posterior probability.

For computational reasons, and because the denominator  $P(X)$  is constant given the input, we typically use the following simplification in practice:

$$P(C|X) \propto P(X|C) \cdot P(C)$$

which means that we can focus on maximizing  $P(X|C) \cdot P(C)$ .

Naive Bayes classifiers are highly scalable and are known for their simplicity, speed, and suitability for high-dimensional datasets.



## APPENDIX A

---

# Answers to Exercises







---

# INDEX

accuracy, [8](#)

confusion matrix, [5](#)

f1 score, [9](#)

k-nearest neighbor, [11](#)

naive Bayes classifier, [13](#)

precision, [8](#)

recall, [8](#)