

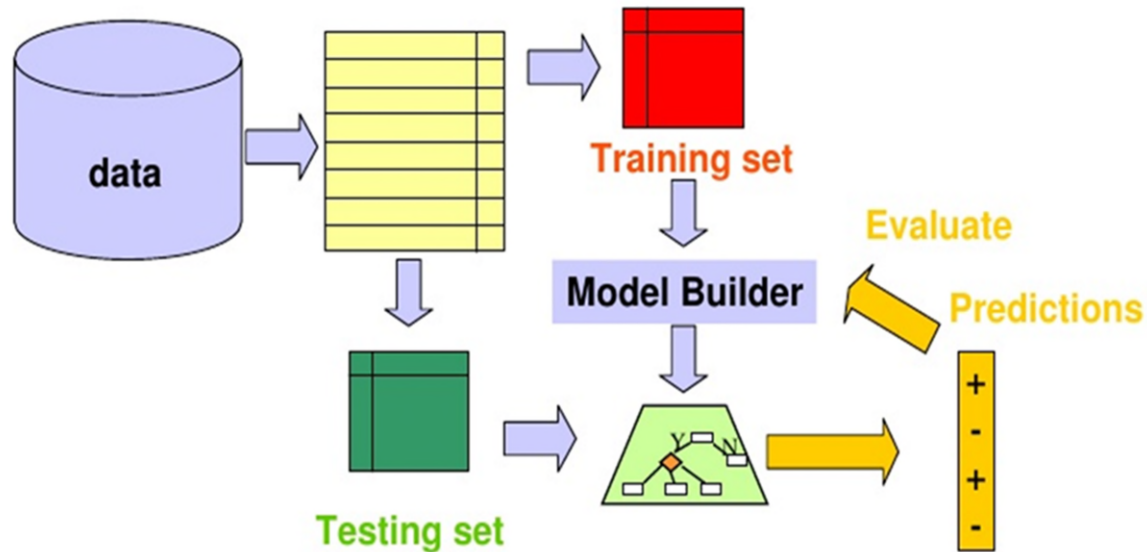
7. Performance Evaluation for Classifiers

In supervised learning, a classification model is trained using the labelled data. How can we understand the performance of the model?

7.1 Evaluation Procedures

7.1.1 Holdout method

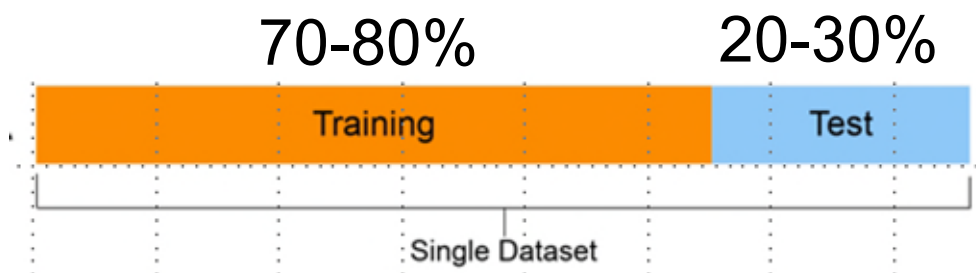
In this method, a **subset of the available data is held back**. **This subset** of the data is used as the **test data** for evaluating the performance of a trained classifier.



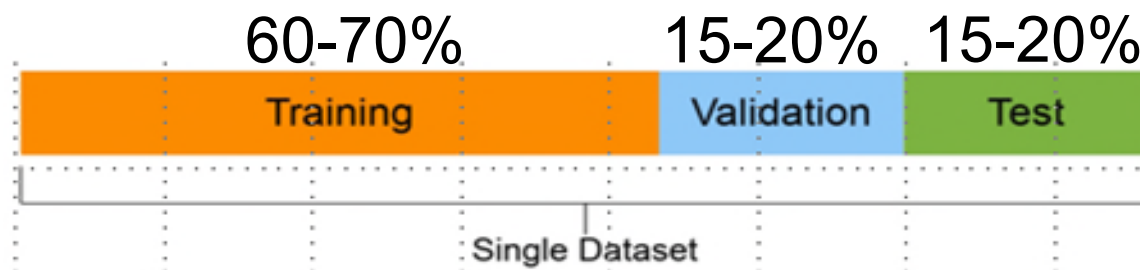
The labels of the test data are predicted using the trained model.

Then the predicted value is compared with the actual value of the label.

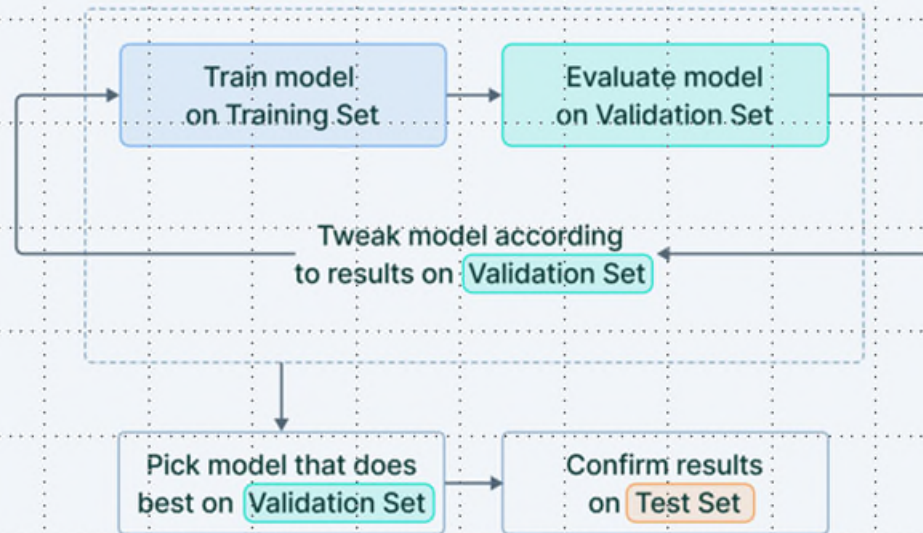
Usually, 70%–80% of the available data (randomly sampled) is used for classifier training, and the remaining 20%–30% is used as test data.



Quite often, the available data is partitioned into **3 portions** – **a training and a test data, and a validation data.**



Training data/validation/test



- ❑ The validation data is used for measuring the model performance **during training stage**.
- ❑ The test data is used **after model training**, to measure and report the final performance of the model.

An obvious problem in this method is that the division of data of different classes into the **training and test data may not be proportional**.

This situation is worse if the overall percentage of data related to certain classes is much less than other classes (i.e. **class imbalance** problem).

This problem can be addressed to some extent by applying **stratified random sampling**: random sampling is conducted on each individual class.

For example, 70% and 30% of the samples in each class are randomly selected as training and test data.

7.1.2 Repeated holdout method

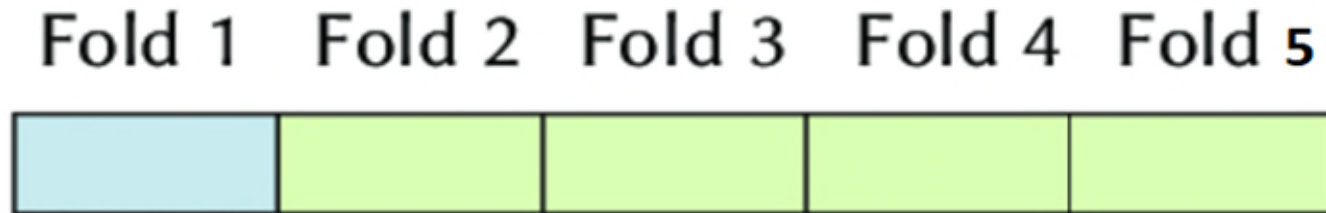
A special variant of the holdout method, called **repeated holdout**, is often employed to ensure the randomness of the composed training and test data sets.

In repeated holdout, **several random holdouts** are used to measure the model performance. In the end, **the average of all performance** is taken.

As multiple holdouts have been drawn, the training and test data (and also validation data, in case it is drawn) are more likely to contain representative data from all classes and resemble the original input data closely.

7.1.3 K-fold cross-validation method

In k -fold cross-validation, the data set is divided into k -completely distinct or non-overlapping random partitions called folds:



Each of the k folds (i.e. 5 in the above example) is given an opportunity to be used as a held back test set, whilst all other folds collectively are used as a training dataset. A total of k models are fit and evaluated on the k hold-out test sets and the **mean performance is reported**.



5-fold cross validation

7.1.4 Leave-one-out cross-validation (LOOCV)

An extreme case of k -fold cross-validation is to set $k = N$, N is the total number of samples.

In each fold, only 1 sample is used as test data, the remaining $N - 1$ samples are used as training data. This extreme case is called **Leave-one-out cross-validation**.

It is obvious that the number of iterations is equal to the total number of samples. Hence, it is **computationally very expensive** and is used only when the number of samples is small.

7.1.5 Repeated k -fold cross-validation

The estimate of model performance via k – fold cross-validation can be noisy: each time the procedure is run, a **different split of the dataset** into k folds is used, and in turn, a **different mean estimate of model performance** is obtained.

One solution to the problem is to repeat the k -fold cross-validation process multiple times and report the mean performance across all folds and all repeats. This approach is generally referred to as **repeated k – fold cross-validation**.

Procedure of Repeated k -Fold Cross-Validation

- (1) Set $repeat = 1$
- (2) Perform k -fold cross validation and obtain the evaluation performance (averaged over k folds) of the current repeat.
- (3) Let $repeat = repeat + 1$ and reshuffle the data
- (4) Go to Step (2) until the total number of repeats reach to a pre-set value such as 50.
- (5) Calculate the average performance over all the repeats

7.2 Performance Evaluation Metrics

For a two-class problem, we assume one class is positive, and another is negative. Usually, **the class of interest is named as positive**, for example, “cancer” is named positive, while “normal” is named as negative. The classification result of a sample may have the 4 scenarios:

- (i) **True positive (TP)**: a sample of positive class is correctly classified into positive class
- (ii) **False positive (FP)**: a sample of negative class is wrongly classified into positive class
- (iii) **True negative (TN)**: a sample of negative class is correctly classified into negative class
- (iv) **False negative (FN)**: a sample of positive class is wrongly classified into negative class

7.2.1 Confusion matrix

A confusion matrix is a summary of classification results as shown below:

		Actual Values	
		Positive	Negative
Predicted Values	Positive	True Positive	False Positive
	Negative	False Negative	True Negative

An example:

	Actual Positive	Actual Negative
Predicted Positive	85	4
Predicted Negative	2	9

7.2.2 Accuracy

Accuracy is the proportion of correctly classified samples:

$$Acc = \frac{TP + TN}{TP + FP + TN + FN} = \frac{85 + 9}{100} = 0.94 = 94\%$$

7.2.3 Error rate

Error rate is the proportion of wrongly classified samples:

$$ER = \frac{FP + FN}{TP + FP + TN + FN} = \frac{4 + 2}{100} = 0.06 = 6\%$$

7.2.4 Sensitivity and Specificity

Sensitivity of a classifier measures the proportion of positive samples that have been correctly classified. It is measured as:

$$Sensitivity = \frac{TP}{TP + FN} = \frac{85}{85 + 2} = 97.7\%$$

Specificity of a classifier measures the proportion of negative samples that have been correctly classified. It is measured as:

$$\textit{Specificity} = \frac{TN}{TN + FP} = \frac{9}{9 + 4} = 69.2\%$$

7.2.5 Precision and Recall

Precision gives the proportion of positive predictions which are truly positive:

$$\textit{Precision} = \frac{TP}{TP + FP} = \frac{85}{85 + 4} = 95.5\%$$

Recall measures the proportion of correct prediction of positives to the total number of positives:

$$Recall = \frac{TP}{TP + FN} = \frac{85}{85 + 2} = 97.7\%$$

Example:

Pick all correct answers from the following

- (a) EE6407 class has male and female students
- (b) EE6407 class comprise students from multiple countries
- (c) All students in EE6407 class are PhD students
- (d) EE6407 exam is closed-book.

There are 3 correct answers: (a), (b) and (d).

- (i) If you pick (a) only, you get 1 correct answer (a), but miss 2 correct answers (b) and (d), then the recall is $1/3=33.33\%$, while your precision is $1/1=100\%$.
- (ii) If you pick (a)(b)(c)(d), you do not miss any correct answers but pick 1 wrong answer, then the recall is $3/3=100\%$, while your precision is $3/4=75\%$.
- (iii) If you pick (a)(b)(d), you do not miss any correct answers, then the recall is $3/3=100\%$, while your precision is $3/3=100\%$.

7.2.6 F-score

F-score is measure of classifier performance that combines the precision and recall:

$$F = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

In the above example:

Case (i):

$$F = \frac{2 \times 1 \times 0.3333}{1 + 0.3333} = 0.5$$

Case (ii):

$$F = \frac{2 \times 0.75 \times 1}{0.75 + 1} = 0.8571$$

Case (iii):

$$F = \frac{2 \times 1 \times 1}{1 + 1} = 1$$

7.2.7 Receiver operating characteristic (ROC) curve

We use discriminant function for classification. For example, in the Bayes decision rule, this discriminant function is the posterior probability. For a two-class problem:

Decide ω_1 if $P(\omega_1|x) > P(\omega_2|x)$;

Decide ω_2 if $P(\omega_2|x) > P(\omega_1|x)$.

In this decision rule, **we implicitly assume the threshold is 0.5.**

But in some applications, for example, the diagnosis of a patient into “cancer” or “normal”, a different threshold instead of **0.5 might be a more prudent choice.**

Setting different thresholds for classifying positive class will inevitably change the 4 metrics in the confusion matrix: TP , FP , TN , FN , and hence the other evaluation metrics of the classification model. One of these thresholds will probably give a better result than the others, depending on whether we are aiming to lower the number of False Negatives or False Positives.

We can generate different confusion matrices and compare the various metrics. But that would not be a prudent thing to do. Instead, what we can do is generate a plot between some of these metrics so that we can easily visualize which threshold gives us a better result.

The ROC curve solves just this problem!

Receiver Operating Characteristic (ROC) visualizes the efficiency of a classification model in the detection of true positives while avoiding the occurrence of false positives.

We define the following measures:

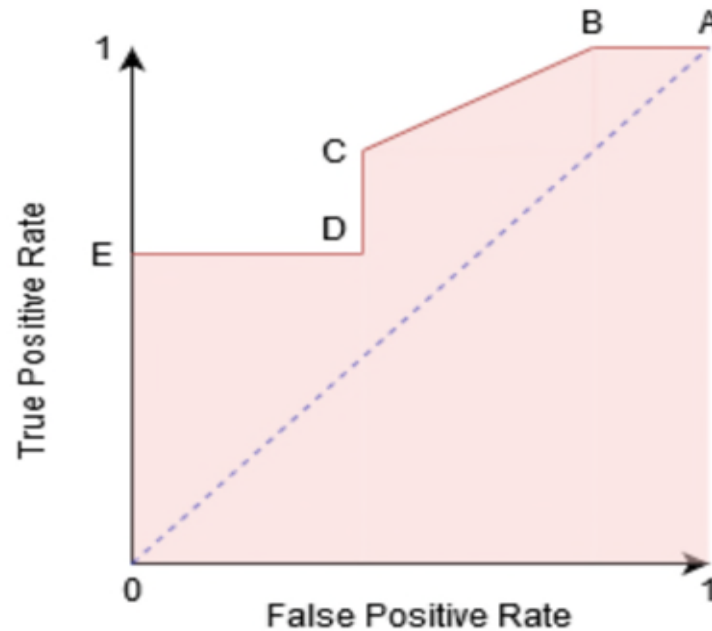
True Positive Rate (TPR)

$$TPR = \frac{TP}{TP + FN} = \textit{Sensitivity} = \textit{Recall}$$

False Positive Rate (FPR)

$$FPR = \frac{FP}{FP + TN} = 1 - \textit{specificity}$$

In the ROC curve, the false positive rate is plotted (in the horizontal axis) against true positive rate (in the vertical axis) at different classification thresholds.



The **Area Under the Curve (AUC)**, i.e. the shaded part in the above figure, is the measure of the **ability of a classifier to distinguish between classes.**

- 1) When $AUC=1$, the classifier is able to perfectly distinguish between all the Positive and the Negative class points correctly;
- 2) When $AUC=0$, the classifier would be predicting all Negatives as Positives, and all Positives as Negatives.
- 3) When $AUC=0.5$, then the classifier is not able to distinguish between Positive and Negative class points, meaning the classifier is predicting random class.
- 4) When $0.5 < AUC < 1$, there is a high chance that the classifier will be able to distinguish the positive class from the negative class, meaning the classifier is able to detect more numbers of True positives and True negatives than False negatives and False positives.

