

Binary Classification

Binary Classification

Defining true/false positives

In general, Positive = identified and negative = rejected. Therefore:

TN True negative = correctly rejected

FP False positive = incorrectly identified

FN False negative = incorrectly rejected

TP True positive = correctly identified

Medical testing example

- True positive = Sick people correctly diagnosed as sick
- False positive= Healthy people incorrectly identified as sick
- True negative = Healthy people correctly identified as healthy
- False negative = Sick people incorrectly identified as healthy.

Definitions

Accuracy Rate

The accuracy rate calculates the proportion of observations being allocated to the **correct** group by the predictive model. It is calculated as follows:

$$\frac{\text{Number of Correct Classifications}}{\text{Total Number of Classifications}} \\ = \frac{TP + TN}{TP + FP + TN + FN}$$

Misclassification Rate

The misclassification rate calculates the proportion of observations being allocated to the **incorrect** group by the predictive model. It is calculated as follows:

$$\frac{\text{Number of Incorrect Classifications}}{\text{Total Number of Classifications}} \\ = \frac{FP + FN}{TP + FP + TN + FN}$$

0.1 Binary Classification

Binary or binomial classification is the task of classifying the elements of a given set into two groups on the basis of a Classification rule. Some typical binary classification tasks are

- medical testing to determine if a patient has certain disease or not (the classification property is the presence of the disease)
- quality control in factories; i.e. deciding if a new product is good enough to be sold, or if it should be discarded (the classification property is being good enough)
- deciding whether a page or an article should be in the result set of a search or not (the classification property is the relevance of the article, or the usefulness to the user)

Statistical classification in general is one of the problems studied in computer science, in order to automatically learn classification systems; some methods suitable for learning binary classifiers include the decision trees, Bayesian networks, support vector machines, neural networks, probit regression, and logit regression.

Sometimes, classification tasks are trivial. Given 100 balls, some of them red and some blue, a human with normal color vision can easily separate them into red ones and blue ones. However, some tasks, like those in practical medicine, and those interesting from the computer science point of view, are far from trivial, and may produce faulty results if executed imprecisely.

Binary Classification is the task of classifying the members of a given set of objects into two groups on the basis if them having a particular set of characteristics.

Sensitivity and specificity are statistical measures of the performance of a binary classification test, also known in statistics as classification function. Sensitivity (also called the true positive rate, or the recall rate in some fields) measures the proportion of actual positives which are correctly identified as such (e.g., the percentage of sick people who are correctly identified as having the condition), and is complementary to the false negative rate. Specificity (sometimes called the true negative rate) measures the proportion of negatives which are correctly identified as such (e.g., the percentage of healthy people who are correctly identified as not having the condition), and is complementary to the false positive rate. A perfect predictor would be described as 100% sensitive (e.g., all sick are identified as sick) and 100% specific (e.g., all healthy are identified as healthy); however, theoretically any predictor will possess a minimum error bound known as the Bayes error rate. For any test, there is usually a trade-off between the measures. For instance, in an airport security setting in which one is testing for potential threats to safety, scanners may be set to trigger on low-risk items like belt buckles and keys (low specificity), in order to reduce the risk of missing objects that do pose a threat to the aircraft and those aboard (high sensitivity). This trade-off can be represented graphically as a receiver operating characteristic curve.

Binary Classification

What Is Classification

Classification is the problem of identifying to which of a set of categories (sub-populations) a new observation belongs, on the basis of a training set of data containing observations (or instances) whose category membership is known.

- To train (create) a classifier, the fitting function estimates the parameters of a Gaussian distribution for each class.
- To predict the classes of new data, the trained classifier finds the class with the smallest misclassification cost.

Types I and II Error

A type I error is the incorrect rejection of a true null hypothesis. A type II error is the failure to reject a false null hypothesis. A type I error is a false positive. Usually a type I error leads one to conclude that a thing or relationship exists when really it doesn't. A type II error is a false negative. Null hypothesis (H_0) is true Null hypothesis (H_0) is false Reject Type I error Correct outcome null hypothesis False positive True positive Fail to reject Correct outcome Type II error null hypothesis True negative False negative

False Positive and False Negative Error

- A false positive error, commonly called a false alarm is a result that indicates a given condition has been fulfilled, when it actually has not been fulfilled. A false positive error is a Type I error where the test is checking a single condition, and results in an affirmative or negative decision usually designated as true or false.
- A false negative error is where a test result indicates that a condition failed, while it actually was successful. A false negative error is a Type II error occurring in test steps where a single condition is checked for and the result can either be positive or negative.

Confusion Matrix

- A confusion matrix, is a table with two rows and two columns that reports the number of false positives, false negatives, true positives, and true negatives.
- This allows more detailed analysis than mere proportion of correct guesses (accuracy). Accuracy is not a reliable metric for the real performance of a classifier, because it will yield misleading results if the data set is unbalanced (that is, when the number of samples in different classes vary greatly).
- For example, if there were 95 cats and only 5 dogs in the data set, the classifier could easily be biased into classifying all the samples as cats. The overall accuracy would be 95%, but in practice the classifier would have a 100% recognition rate for the cat class but a 0% recognition rate for the dog class.

Sensitivity and Specificity

Sensitivity and specificity are measures of the performance of a binary classification test.

- Sensitivity (also called the true positive rate, or the recall rate) measures the proportion of actual positives which are correctly identified as such (e.g. the percentage of sick people who are correctly identified as having the condition).
sensitivity (Recall) = $\frac{\text{number of true positives}}{\text{number of true positives} + \text{number of false negatives}}$
 - *(Remark: We will use the terms Sensitivity and Recall interchangeably. Sensitivity is more commonly used in a medical context, while recall is more commonly used in data science.)*
- Specificity measures the proportion of negatives which are correctly identified as such (e.g. the percentage of healthy people who are correctly identified as not having the condition, sometimes called the true negative rate).

$$\text{Specificity} = \frac{TN}{TP + FN}$$

Receiver Operating Characteristic (ROC) curve

In a Receiver Operating Characteristic (ROC) curve the true positive rate (Sensitivity) is plotted in function of the false positive rate (100-Specificity) for different cut-off points. Each point on the ROC curve represents a sensitivity/specificity pair corresponding to a particular decision threshold. A test with perfect discrimination (no overlap in the two distributions) has a ROC curve that passes through the upper left corner (100specificity). Therefore the closer the ROC curve is to the upper left corner, the higher the overall accuracy of the test (Zweig and Campbell, 1993).

Accuracy, Recall and Precision: An Example

Calculating precision and recall is actually quite easy. Imagine there are 135 positive cases among 10,000 cases. You want to predict which ones are positive, and you pick 265 to have a better chance of catching many of the 135 positive cases. You record the IDs of your predictions, and when you get the actual results you sum up how many times you were right or wrong. There are four ways of being right or wrong:

- TN / True Negative: case was negative and predicted negative
- TP / True Positive: case was positive and predicted positive
- FN / False Negative: case was positive but predicted negative
- FP / False Positive: case was negative but predicted positive

Now count how many of the 10,000 cases fall in each category:

	Predicted Negative	Predicted Positive
Negative Cases	TN: 9,700	FP: 165
Positive Cases	FN: 35	TP: 100

What percent of your predictions were correct? The accuracy was $(9,760+60)$ out of $10,000 = 98.00\%$ What percent of the positive cases did you catch? The recall was 100 out of $135 = 74.07\%$ What percent of positive predictions were correct? The precision was 100 out of $265 = 37.74\%$ What percent of negative predictions were correct? The specificity was 9700 out of $9735 = 99.64$

The F Score

The F-score or F-measure is a measure of a classification procedures accuracy. It considers both the precision and the recall to compute the score.

$$F = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

1 Class Imbalance

A data set said to be highly skewed if sample from one class is in higher number than other. In an imbalanced data set the class having more number of instances is called as major class while the one having relatively less number of instances are called as minor class . Applications such as medical diagnosis prediction of rare but important disease is very important than regular treatment. Similar situations are observed in other areas, such as detecting fraud in banking operations, detecting network intrusions, managing risk and predicting failures of technical equipment. In such situation most of the binary classifier are biased towards the major classes and hence show very poor classification rates on minor classes. It is also possible that classifier predicts everything as major class and ignores the minor class completely. The Accuracy measure is an example of an metric that is affected by this bias. As the F-measure is not computed using the True Negatives, it is less Biased.

2 Model Accuracy

Prediction error refers to the discrepancy or difference between a predicted value (based on a model) and the actual value. In the standard regression situation, prediction error refers to how well our regression equation predicts the outcome variable scores of new cases based on applying the model (coefficients) to the new cases predictor variable scores.

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN} \quad (1)$$

$$\text{Precision} = \frac{TP}{TP + FP} \quad (2)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3)$$

Cross Validation

The confusion table is a table in which the rows are the observed categories of the dependent and the columns are the predicted categories. When prediction is perfect all cases will lie on the diagonal. The percentage of cases on the diagonal is the percentage of correct classifications. The cross validated set of data is a more honest presentation of the power of the discriminant function than that provided by the original classifications and often produces a poorer outcome. The cross validation is often termed a jack-knife classification, in that it successively classifies all cases but one to develop a discriminant function and then categorizes the case that was left out. This process is repeated with each case left out in turn. This is known as leave-1-out cross validation. This cross validation produces a more reliable function. The argument behind it is that one should not use the case you are trying to predict as part of the categorization process.

Error Rates

- We can evaluate error rates by means of a training sample (to construct the discrimination rule) and a test sample.

- An optimistic error rate is obtained by reclassifying the training data. (In the training data sets, how many cases were misclassified). This is known as the apparent error rate.
- The apparent error rate is obtained by using in the training set to estimate the error rates. It can be severely optimistically biased, particularly for complex classifiers, and in the presence of over-fitted models.
- If an independent test sample is used for classifying, we arrive at the true error rate.
- The true error rate (or conditional error rate) of a classifier is the expected probability of misclassifying a randomly selected pattern. It is the error rate of an infinitely large test set drawn from the same distribution as the training data.

Misclassification Cost

As in all statistical procedures it is helpful to use diagnostic procedures to assess the efficacy of the discriminant analysis. We use cross-validation to assess the classification probability. Typically you are going to have some prior rule as to what is an acceptable misclassification rate. Those rules might involve things like, what is the cost of misclassification? Consider a medical study where you might be able to diagnose cancer. There are really two alternative costs. The cost of misclassifying someone as having cancer when they don't. This could cause a certain amount of emotional grief. Additionally there would be the substantial cost of unnecessary treatment. There is also the alternative cost of misclassifying someone as not having cancer when in fact they do have it. A good classification procedure should

- result in few misclassifications
- take prior probabilities of occurrence into account
- consider the cost of misclassification

For example, suppose there tend to be more financially sound firms than bankrupt firm. If we really believe that the prior probability of a financially distressed and ultimately bankrupted firm is very small, then one should classify a randomly selected firm as non-bankrupt unless the data overwhelmingly

favor bankruptcy. There are two costs associated with discriminant analysis classification: The true misclassification cost per class, and the expected misclassification cost (ECM) per observation.

Suppose there we have a binary classification system, with two classes: class 1 and class 2. Suppose that classifying a class 1 object as belonging to class 2 represents a more serious error than classifying a class 2 object as belonging to class 1. There would an assignable cost to each error. $c(i-j)$ is the cost of classifying an observation into class j if its true class is i . The costs of misclassification can be defined by a cost matrix. Predicted

	Class 1	Class 2
Class 1	0	$c(2-1)$
Class 2	$c(1-2)$	0

Expected cost of misclassification (ECM)

Let p_1 and p_2 be the prior probability of class 1 and class 2 respectively. Necessarily $p_1 + p_2 = 1$.

The conditional probability of classifying an object as class 1 when it is in fact from class 2 is denoted $p(1-2)$. Similarly the conditional probability of classifying an object as class 2 when it is in fact from class 1 is denoted $p(2-1)$.

$$ECM = c(2|1)p(2|1)p_1 + c(1|2)p(1|2)p_2$$

(In other words: the sum of the cost of misclassification times the (joint) probability of that misclassification. A reasonable classification rule should have ECM as small as possible.