

Adaptive Computation and Machine Learning

4. MODEL TESTING

Recall that before training a model on a given dataset, a portion of the data is set aside to be used as testing; this is called the test dataset. The test dataset must be independent of the training process; usually, it comprises anywhere between 20% and 50% of the original dataset. The test dataset is different to the validation dataset, which is used in the training process to avoid overfitting or for tuning hyperparameters.

4.1. Evaluation metrics.

In this section we discuss methods that use the test dataset to evaluate how good (or how fit) a model is and, thereby, to decide if it should be kept or discarded.

Consider a classification problem in which we have a dataset consisting of datapoints and corresponding target classes and we seek to create a model that will correctly predict the class of each datapoint. Suppose that a model has been trained for this dataset. Evaluating such a model on a test dataset is done by applying the model to every datapoint in the test dataset and comparing the class predicted by the model with the actual class of the datapoint in the target dataset.

Note that the output of the model must be converted into a class prediction.

For example, consider a dataset in which the targets are 0's and 1's, and the model is a neural network with a single output node that uses the sigmoid activation function. If the output obtained from the model for some datapoint is 0.65, then this must be converted to either 0 or 1 (probably 1 in this case).

Next, suppose a dataset has three classes and a one-hot encoding is used, and the model is a neural network with three output nodes, each using sigmoid activation. If the output obtained for some datapoint is $(0.7, 0.2, 0.1)$, then the predicted class is the one corresponding to $(1, 0, 0)$.

The following metric can be used on the test dataset:

$$\mathbf{Accuracy} = \frac{\text{number of correct predictions}}{\text{total number of predictions}}$$

The accuracy of a model on the test dataset is a value between 0 and 1, which is often converted to a percentage. Accuracy is not always the best metric and there are other metrics that can be used to determine how good a model is.

Consider the case of **binary classification** in which there are only two possible classes to predict. Let the two classes be 0 and 1.

The results of applying a model on the test dataset can be represented as in the following example:

		prediction	
		1	0
actual	1	70	30
	0	15	50

The above table is called a **confusion matrix**. The table shows the predicted class versus the actual class for each datapoint. In the table we can read off that 70 of the datapoints in class 1 were correctly predicted, 30 of the datapoints in class 1 were incorrectly predicted as 0, 15 of the datapoints in class 0 were incorrectly predicted as 1, and 50 of the datapoints in class 0 were correctly predicted. There are 165 datapoints in the test dataset. The accuracy of the model on the test dataset is $\frac{70+50}{165} = 0.7273$, or 72.73%.

To see why accuracy is not always a good metric, consider the next example. Suppose the results of applying the model on the test dataset are given by the following confusion matrix:

		prediction	
		1	0
actual	1	65	9
	0	4	2

In this example, there are 80 datapoints in the test dataset; 65 datapoints in class 1 were correctly predicted as 1, 9 datapoints in class 1 were incorrectly predicted as 0, 4 datapoints in class 0 were incorrectly predicted as 1 and 2 datapoints in class 0 were correctly predicted as 0. The accuracy of the model on this test dataset is $\frac{65+2}{80} = 83.75\%$.

Although 83.75% seems like a good accuracy, there is simpler model that does better.

Consider the model that always predicts class 1 regardless of the datapoint.

Since the test dataset has 74 actual 1's out of 80 datapoints, this model will have accuracy of $\frac{74}{80} = 92.5\%$.

However, such a model does not help with the original dataset, since there is (presumably) a need to predict class 0 in some cases. Thus, accuracy is not a good metric in this case. The reason for this is that the (test) dataset is very unbalanced and the large number of 1 classes dominates the accuracy metric. For datasets with an even distribution of classes, i.e., a **balanced** dataset, accuracy is a good metric.

To discuss other metrics, we introduce some notation. We still consider binary classification using classes 1 and 0. A datapoint with classification 1 will be called **positive** and a datapoint with classification 0 will be called **negative**. We define the following:

True Positive (TP): when the predicted class is 1 and the actual class is 1.

True Negative (TN): when the predicted class is 0 and the actual class is 0.

False Positive (FP): when the predicted class is 1 and the actual class is 0.

False Negative (FN): when the predicted class is 0 and the actual class is 1.

Using the above notation we can define the **confusion matrix** as the following table:

		prediction	
		1	0
actual	1	TP	FN
	0	FP	TN

The formula for calculating accuracy from TP , TN , FP and FN is as follows:

$$\mathbf{Accuracy} = \frac{TP+TN}{TP+FP+TN+FN}$$

Other metrics that are based on TP , TN , FP and FN are:

$$\mathbf{Precision} = \frac{TP}{TP+FP}$$

$$\mathbf{Sensitivity} = \mathbf{Recall} = \frac{TP}{TP+FN}.$$

Precision is the positive predictive value.

Sensitivity, or Recall, is the true positive rate.

In the above example,

$$\text{Precision} = \frac{65}{65+4} = 0.94$$

$$\text{Sensitivity} = \frac{65}{65+9} = 0.88.$$

Deciding whether to use precision or sensitivity depends on the dataset. Precision is used in datasets where you want few false positives, i.e., you want to be sure your positive predictions are correct (e.g., in medical tests, it's important that there are few false positives, or incorrect diagnoses of illnesses). Sensitivity is used if the number of false negatives is more important, that is, you don't want to miss many 1's (e.g., in fraud detection).

A metric that provides a balance between precision and sensitivity is the **F1 score**, defined as follows:

$$F1 = \frac{2 * \text{precision} * \text{sensitivity}}{\text{precision} + \text{sensitivity}}$$

The $F1$ score is an average of precision and sensitivity, in fact, it is the **harmonic mean** of precision and sensitivity. The harmonic mean of numbers a and b is defined as the reciprocal of the average of the reciprocals:

$$\text{harmonic mean of } a \text{ and } b = \frac{1}{\frac{1}{2}(\frac{1}{a} + \frac{1}{b})} = \frac{2ab}{a+b}.$$

The $F1$ score in the above example is $\frac{2(0.94)(0.88)}{0.94+0.88} = 0.91$.

In the first example above, precision = $\frac{70}{85} = 0.82$, sensitivity = $\frac{70}{100} = 0.7$ and $F1 = \frac{2(0.82)(0.7)}{0.82+0.7} = 0.76$.

Next, consider the case of multi-class classification. Suppose that a dataset has three classes, say A , B and C . If a model is trained on the training dataset and tested on the test dataset, then the confusion matrix would look something like:

		prediction		
		A	B	C
actual	A	42	6	2
	B	3	50	7
	C	0	6	34

In the above confusion matrix, the first row indicates that there are 50 datapoints with target A ; of these, 42 are correctly classified as A , 6 are misclassified as B and 2 are misclassified

as C . The second row indicates that there are 60 datapoints with target B ; of these, 50 are correctly classified as B , 3 are misclassified as A and 7 are misclassified as C . The third row indicates that there are 40 datapoints with target C ; of these, 34 are correctly classified as C , 0 are misclassified as A and 6 are misclassified as B .

The accuracy of the model is $\frac{42+50+34}{42+50+34+6+2+3+7+0+6} = 0.84$, that is, 84%.

The notions of sensitivity and precision can be defined for each of the classes.

For example, for class A , if we consider A as positive, then

$$\begin{aligned}\text{precision} &= \frac{42}{42+3+0} = 0.93 \\ \text{sensitivity} &= \frac{42}{42+6+2} = 0.84 \\ F1 &= \frac{2*(0.93)(0.84)}{0.93+0.84} = 0.88.\end{aligned}$$

For class B , if we consider B as positive, then

$$\begin{aligned}\text{precision} &= \frac{50}{6+50+6} = 0.81 \\ \text{sensitivity} &= \frac{50}{3+50+7} = 0.83 \\ F1 &= \frac{2*(0.81)(0.83)}{0.81+0.83} = 0.82.\end{aligned}$$

For class C , if we consider C as positive, then

$$\begin{aligned}\text{precision} &= \frac{34}{2+7+34} = 0.79 \\ \text{sensitivity} &= \frac{34}{0+6+34} = 0.85 \\ F1 &= \frac{2*(0.79)(0.85)}{0.79+0.85} = 0.82.\end{aligned}$$

Lastly, for each of sensitivity, precision and $F1$, an average value over the three classes can be obtained, which is denoted as the **macro** score. In the above example,

$$\begin{aligned}\text{macro-precision} &= \frac{0.93+0.81+0.79}{3} = 0.84 \\ \text{macro-sensitivity} &= \frac{0.84+0.83+0.85}{3} = 0.84 \\ \text{macro-}F1 &= \frac{0.88+0.82+0.82}{3} = 0.84.\end{aligned}$$

4.2. The ROC and the AUC.

The **receiver operating characteristic curve**, or the **ROC curve** for short, is the plot of the true positive rate against the false positive rate for each threshold setting in a binary classification model.

The **true positive rate** (TPR) is the proportion of true positives that are *correctly* classified as positives, which is just sensitivity:

$$\text{TPR} = \frac{TP}{TP + FN}.$$

The **false positive rate** (FPR) is the proportion of true negatives that are *incorrectly* classified as positives:

$$\text{FPR} = \frac{FP}{FP + TN}.$$

As an example, consider a dataset in which the targets are 0's and 1's, for which a neural network with a single output node that uses the sigmoid activation function is trained. Then for every input from the test dataset, the model produces an output value between 0 and 1, which must be converted into a prediction, i.e., either 0 or 1. A natural way to do this is to use a threshold value, say t . Then, if y is the output from the model and $y > t$ the prediction is 1, whereas if the $y \leq t$, the prediction is 0. In a sense, the output y can be considered the probability that the given input has target 1. If the probability is greater than the given threshold t then 1 is predicted, otherwise 0 is predicted.

A common choice for the threshold value is $t = 0.5$ since this is the middle point of the range of possible outputs. However, there may be good reasons for using other choices of t . For example, if predicting a rare disease, a higher threshold value of t could be chosen, say $t = 0.9$, to prevent making false diagnoses. But if predicting a contagious disease such as monkeypox, a low threshold value such as $t = 0.1$ could be chosen so as not to miss any possible cases.

The **ROC curve** is a plot that considers all possible choices of threshold value t in the range $[0, 1]$. For each such value of t , using the test dataset, the TPR is computed, and also the FPR. These two values are plotted as a single point on a system of axes with FPR on the horizontal axis and TPR on the vertical axis.

Once the points on the axes have been plotted for each possible value of t in $[0, 1]$, the points can be connected by straight lines to form a continuous curve, which is the ROC curve.

For example, consider the table below.

The first two columns show the input values and targets of the datapoints in the test dataset.

The ‘output’ column gives the output values from the model for each input x_i .

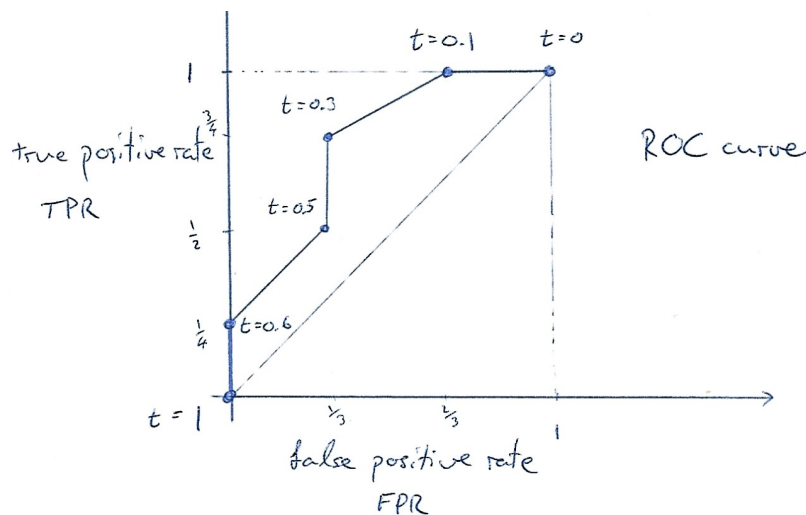
The column labeled ‘ $t = 0$ ’ gives the predicted values for each input if the threshold value of $t = 0$ used. Since every output is greater than 0, the prediction is 1 in every case.

Similarly, the column labeled ‘ $t = 0.1$ ’ gives the predicted values for each input if the threshold value of $t = 0.1$ used. In this case, the input x_7 gets predicted as 0 since its output value from the model is not greater than 0.1. The remaining columns are similarly obtained.

input	target	output	$t = 0$	$t = 0.1$	$t = 0.3$	$t = 0.5$	$t = 0.6$	$t = 1$
x_1	1	0.9	1	1	1	1	1	0
x_2	1	0.6	1	1	1	1	0	0
x_3	0	0.3	1	1	0	0	0	0
x_4	1	0.3	1	1	0	0	0	0
x_5	1	0.5	1	1	1	0	0	0
x_6	0	0.6	1	1	1	1	0	0
x_7	0	0.1	1	0	0	0	0	0
TPR			1	1	$\frac{3}{4}$	$\frac{2}{4}$	$\frac{1}{4}$	0
FPR			1	$\frac{2}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	0	0

The last two rows in the table give the true positive rate and false positive rate for each choice of t . These points are plotted on the graph below and the points are connected by straight lines to give the ROC curve.

Note that all other choices of t will result in one of these points (please check a few).



An ideal situation would be a t value that gives TPR of 1, meaning that all 1's are correctly predicted, and FPR of 0, meaning that no 0's are incorrectly predicted as 1's (so all 0's are correctly predicted as 0). This would give a point at $(0, 1)$ on the axes. Although this is unlikely, we can nevertheless try find a t value that gives a point closest to the point $(0, 1)$. In the example above, the value $t = 0.3$ is the best choice; it gives the point $(\frac{1}{3}, \frac{3}{4})$.

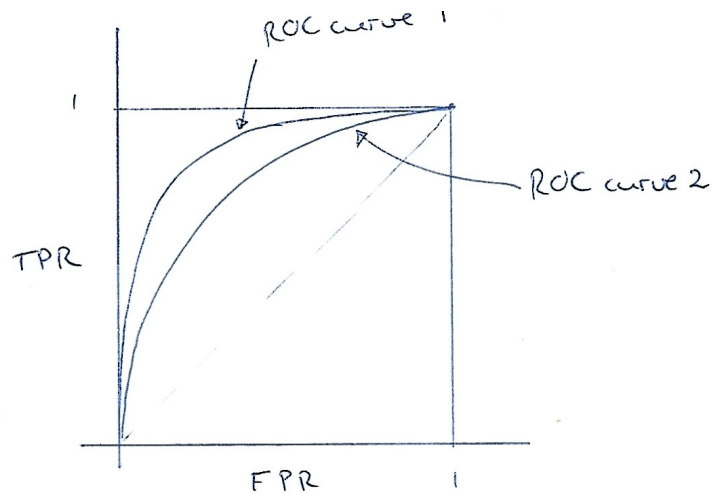
A ROC curve can be constructed for any binary classifier model with respect to a test dataset as long as some threshold value is used in the model to obtain predictions.

For another example, consider a random forest of decision trees and suppose the random forest is a binary classifier, meaning that there are only two target classes, say 0 and 1. Recall that for a given input, a prediction is obtained from each tree in the forest and the final prediction is the majority class obtained from the individual trees. Instead of simply choosing the majority class, a threshold t could be set such that a 1 is predicted only if the proportion of trees predicting 1 is greater than t . Then a ROC curve could be constructed for the random forest using a range of t values.

The **area under the curve**, or **AUC** for short, is a measure used to compare various models for which we have a ROC curve. AUC simply refers to the area under the ROC curve.

Suppose that for some dataset with target classes 0 and 1, both a neural network model and a random forest model are constructed from the training dataset. Then the ROC curve can be constructed for each of the models using the test dataset. To decide which model is better, the AUC is computed for each ROC curve — the model with the greater AUC is the better one.

Recall that the ideal ROC curve goes through the point $(1, 0)$; this ROC curve would have AUC of 1. Thus, the model with the greater AUC is the one closest to the ideal and is the preferred one. The diagram below shows two ROC curves; we can see that the AUC for ROC curve 1 is greater than the AUC for ROC curve 2.



EXERCISE

Draw the ROC curve for the following table, which gives the input values and targets for datapoints in the test dataset and also the corresponding output values obtained from the model.

input	target	output
x_1	1	0.85
x_2	0	0.25
x_3	0	0.3
x_4	1	0.5
x_5	1	0.25
x_6	0	0.1
x_7	1	0.9
x_8	1	0.6
x_9	0	0.5
x_{10}	0	0.2

4.3. Cross-validation.

Another method for testing a model on a dataset is **k-fold cross-validation**.

(Confusingly, ‘validation’ in this context actually means testing, so ‘k-fold cross testing’ would be a better name.)

This method is often used if the dataset is not sufficiently large to accommodate a large enough training dataset and test dataset. An advantage of this method is that every datapoint is used both in training and in testing in some way.

The method is as follows.

Choose the type of model to be applied.

Choose a number k (usually in the range 5 to 10, but any value will do).

Split the dataset into k approximately equal parts, which are called the **folds**.

Choose $k - 1$ of the folds as the training dataset. The remaining fold will be the test dataset.

Train the model on the training dataset. (If necessary, use a portion of the dataset to prevent overtraining).

Use the remaining fold as the test dataset to obtain the accuracy, or any other metric.

Repeat the above process k times, each time with a different fold selected for testing.

Note that for each split of the data, a new model is trained from scratch and is independent of the models trained on other splits.

The final metric is obtained by averaging the metrics obtained on each fold.

A good final metric obtained by the cross validation method provides confidence in the type of model being used. Note that the above method results in k different models of the selected type. The best model can be chosen as the final model, or a fresh model can be trained.

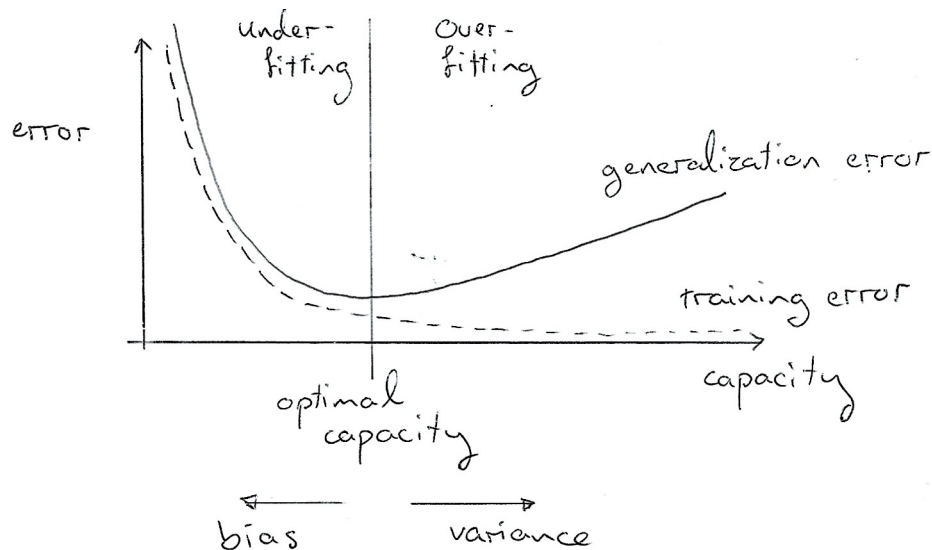
If the final metric obtained from cross validation is not good enough, then it can mean that the type of model is not suitable and should be adjusted. For example, in the case of a neural network, it can mean that the number of nodes or layers in the network should be changed.

4.4. Underfitting vs Overfitting, Bias vs Variance.

The following diagram summarises the issues of overfitting and underfitting when choosing a model type. Loosely speaking, the **capacity** of a model type correlates to the complexity of the datasets that it is able to accurately model.

For example, consider a regression problem in one variable. A model type with low capacity is a straight line model since such models can only accurately fit data that is linear. On more complex datasets, a straight line model may do very poorly; we call this **underfitting**. Quadratic models have greater capacity to model datasets so underfitting is less of a problem, but now the risk of overfitting increases. As the model type becomes more complex the capacity increases, but so does the risk of overfitting. If we use neural networks as an example, the capacity is increased by adding more layers and more nodes to the model. However, it is difficult to know what a good capacity is. The optimal capacity is obtained by minimising the generalisation error, i.e., the error on the test set.

The issue of underfitting versus overfitting is often phrased as a trade-off between **bias** and **variance**. Reducing the capacity of the model means increasing the bias – the trained model is then more biased towards the type of model chosen and, hence, more susceptible to underfitting. Increasing the capacity of the model means there is greater variance in what the trained model can be, but then it is more susceptible to overfitting.



5. Probabilistic Methods

Consider the following dataset that has only one attribute, the COMS2 mark, and one target value which is either *Pass* or *Fail* in COMS3:

COMS2 COMS3

A	$Pass$
C	$Fail$
C	$Pass$
B	$Pass$
B	$Fail$
C	$Pass$
A	$Fail$
B	$Pass$

Suppose a student got a B in COMS2 and is currently doing COMS3. Based on the above dataset, what prediction could we make for this student: *Pass* or *Fail*?

We can calculate the probability of *Pass* by looking at all the rows in which $COMS2 = B$.

There are 3 such rows, and of the 3, there are 2 *Pass*'s, so the probability of *Pass* given that the student got a B in COMS2 is $\frac{2}{3}$. We write this using conditional probabilities as:

$$P(Pass|COMS2 = B) = \frac{2}{3}, \quad \text{or} \quad P(Pass|B) = \frac{2}{3}.$$

Alternative, we can use **Bayes' Rule** which states: for events X and Y ,

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}.$$

Bayes' rule gives us a way of calculating the probability of X given Y , i.e., $P(X|Y)$, from the probability of Y given X , i.e., $P(Y|X)$. In Bayes' rule,

$P(X)$ is called the **prior** probability,

$P(X|Y)$ is called the **posterior** probability,

$P(Y|X)$ is called the **likelihood** (of Y given X).

In the above example, using Bayes' rule, we calculate:

$$P(Pass|B) = \frac{P(B|Pass)P(Pass)}{P(B)}.$$

Now, $P(B|Pass) = \frac{2}{5}$. This comes from looking at all the rows which have $COMS3 = Pass$, of which there are 5, and counting the number of these rows in which $COMS2 = B$, which is 2. Also, $P(Pass) = \frac{5}{8}$, since we have 8 datapoints and in 5 of these $COMS3 = Pass$. Lastly, $P(B) = \frac{3}{8}$, since we have 8 datapoints and in 3 of these $COMS2 = B$. Putting this into the above equation, we get:

$$P(Pass|B) = \frac{P(B|Pass)P(Pass)}{P(B)} = \frac{(\frac{2}{5})(\frac{5}{8})}{\frac{3}{8}} = \frac{2}{3}.$$

We can similarly work out that $P(Fail|B) = \frac{1}{3}$. Based on these probabilities, we hypothesise that the most likely outcome for the student who got B in $COMS2$ is $Pass$.

5.1. MAP hypothesis.

Consider a dataset S in which the targets are the classes C_1, \dots, C_n . Given a new datapoint $\mathbf{x} = (x_1, \dots, x_m)$ we want to decide which C_i to classify it as. Suppose we can obtain the following conditional probabilities:

$$P(C_1|\mathbf{x}), \dots, P(C_n|\mathbf{x}).$$

Here, $P(C_i|\mathbf{x})$ is the probability that the given datapoint \mathbf{x} is classified as C_i .

Based on these probabilities, we would classify \mathbf{x} as the C_i for which $P(C_i|\mathbf{x})$ is maximum. This is the **maximum a posteriori** hypothesis, or **MAP** hypothesis.

Observe that, by Bayes' rule, for each class C_i , we have:

$$P(C_i|\mathbf{x}) = \frac{P(\mathbf{x}|C_i)P(C_i)}{P(\mathbf{x})}$$

Thus, to find the MAP hypothesis, we need to find the maximum of:

$$\frac{P(\mathbf{x}|C_1)P(C_1)}{P(\mathbf{x})}, \dots, \frac{P(\mathbf{x}|C_n)P(C_n)}{P(\mathbf{x})}.$$

Since each of these terms has the same denominator, it is sufficient to find the maximum of:

$$P(\mathbf{x}|C_1)P(C_1), \dots, P(\mathbf{x}|C_n)P(C_n).$$

Example: Consider the following dataset on students doing COMS3:

$$S = \begin{bmatrix} \text{COMS2} & \text{doing labs?} & \text{doing tuts?} & \text{COMS3} \\ A & N & Y & Pass \\ C & Y & N & Fail \\ C & N & Y & Pass \\ B & Y & Y & Pass \\ B & N & Y & Fail \\ C & Y & N & Pass \\ A & N & N & Fail \\ B & Y & N & Pass \end{bmatrix}$$

In the above dataset, the attributes of the data are: COMS2, doing labs?, and doing tuts?. The target value is in the column COMS3.

Suppose a current student had a B for COMS2, is not doing labs and is not doing tuts, i.e., $\mathbf{x} = (B, N, N)$. We want to predict the outcome of COMS3, i.e., $Pass$ or $Fail$, by using the MAP hypothesis. That is, we want to determine which of the following is larger: $P(Pass|B, N, N)$ or $P(Fail|B, N, N)$.

First, apply Bayes' rule:

$$P(Pass|B, N, N) = \frac{P(B, N, N|Pass)P(Pass)}{P(B, N, N)},$$

$$P(Fail|B, N, N) = \frac{P(B, N, N|Fail)P(Fail)}{P(B, N, N)}$$

To determine which of the above is larger, we need only find the larger of:

$$P(B, N, N|Pass)P(Pass) \text{ and } P(B, N, N|Fail)P(Fail).$$

From the target column of the dataset we get $P(Pass) = \frac{5}{8}$ and $P(Fail) = \frac{3}{8}$.

To find $P(B, N, N|Pass)$ and $P(B, N, N|Fail)$ is more difficult since we do not even have a situation like this in our dataset. We use a simplifying assumption here. Assume that the attributes COMS2, doing labs? and doing tuts? are conditionally independent with respect to $Pass$ and $Fail$. This means that

$$P(B, N, N|Pass) = P(\text{COMS2} = B|Pass)P(\text{doing labs?} = N|Pass)P(\text{doing tuts?} = N|Pass).$$

Thus, we can calculate as follows:

$$P(B, N, N|Pass) = \left(\frac{2}{5}\right) \left(\frac{2}{5}\right) \left(\frac{2}{5}\right) = \frac{8}{125}$$

hence

$$P(B, N, N|Pass)P(Pass) = \left(\frac{8}{125}\right) \left(\frac{5}{8}\right) = \frac{1}{25}.$$

Similarly, we assume that

$$P(B, N, N|Fail) = P(\text{COMS2} = B|Fail)P(\text{doing labs?} = N|Fail)P(\text{doing tuts?} = N|Fail).$$

so we calculate

$$P(B, N, N|Fail) = \left(\frac{1}{3}\right) \left(\frac{2}{3}\right) \left(\frac{2}{3}\right) = \frac{4}{27}$$

hence

$$P(B, N, N|Fail)P(Fail) = \left(\frac{4}{27}\right) \left(\frac{3}{8}\right) = \frac{1}{18}.$$

Since $\frac{1}{18} > \frac{1}{25}$, we hypothesise that the student is more likely to *Fail*, i.e., we classify the student as *Fail*.

The method used in the above example is the Naïve Bayes Classifier.

5.2. Naïve Bayes Classifier.

Suppose we have a dataset S in which the targets are the classes C_1, \dots, C_n , and, for a new datapoint $\mathbf{x} = (x_1, \dots, x_m)$, we want to decide which C_i to classify it as. In order to use the MAP hypothesis, we must find the maximum of:

$$P(\mathbf{x}|C_1)P(C_1), \dots, P(\mathbf{x}|C_n)P(C_n).$$

If we make the simplifying assumption that the attributes x_i are all conditionally independent with respect to the C_i s, then:

$$P(\mathbf{x}|C_i) = P(x_1|C_i)P(x_2|C_i) \cdots P(x_m|C_i).$$

With this simplification, we classify \mathbf{x} as the class C_i which gives the maximum of

$$P(\mathbf{x}|C_1)P(C_1), \dots, P(\mathbf{x}|C_n)P(C_n).$$

This is known as the **Naïve Bayes Classifier**.

The word ‘naïve’ refers to the naïve assumption that the attributes are conditionally independent. This is usually not the case.

Example: Consider the dataset in the previous example, and now suppose that a student got C for COMS2, is doing labs and is not doing tuts, i.e., $\mathbf{x} = (C, Y, N)$. In the dataset, there are two previous cases of such students - one of which *Passed* and one which *Failed* COMS3.

So we cannot make a prediction based on these two cases. Let's use the Naïve Bayes Classifier to make a classification.

We want to find the larger of:

$$P(C, Y, N|Pass)P(Pass) \text{ and } P(C, Y, N|Fail)P(Fail).$$

Using the conditional independence assumption we calculate:

$$\begin{aligned} & P(C, Y, N|Pass)P(Pass) \\ &= P(\text{COMS2} = C|Pass)P(\text{doing labs?} = Y|Pass)P(\text{doing tuts?} = N|Pass)P(Pass) \\ &= \left(\frac{2}{5}\right) \left(\frac{3}{5}\right) \left(\frac{2}{5}\right) \left(\frac{5}{8}\right) \\ &= \frac{3}{50} \end{aligned}$$

$$\begin{aligned} & P(C, Y, N|Fail)P(Fail) \\ &= P(\text{COMS2} = C|Fail)P(\text{doing labs?} = Y|Fail)P(\text{doing tuts?} = N|Fail)P(Fail) \\ &= \left(\frac{1}{3}\right) \left(\frac{1}{3}\right) \left(\frac{2}{3}\right) \left(\frac{3}{8}\right) \\ &= \frac{1}{36} \end{aligned}$$

Thus, the Naïve Bayes Classifier would classify this student as *Pass*, since $\frac{3}{50} > \frac{1}{36}$.

5.3. Applying Bayes' Rule.

Bayes' rule can be used directly to make predictions about a dataset. The method of applying Bayes' rule is usually described using urns filled with balls of different colours. We give such an example below. In the exercises below there is an example that shows how this method may be applied to a data problem.

Suppose we have two urns containing some *Blue* balls and some *Red* balls, as follows:

Urn 1 contains 5 balls, 3 of which are *Red* and 2 of which are *Blue*.

Urn 2 contains 8 balls, 3 of which are *Red* and 5 of which are *Blue*.

A friend will select one of the urns by (secretly) flipping a coin.

If the coin lands heads up, he will choose Urn 1 and if it lands tails up, he will choose Urn 2.

You must guess which of the urns he chose.

You know that the urn was chosen by a coin-flip, so before you receive any new information, your best guess is that there is a $\frac{1}{2}$ probability that Urn 1 was chosen and a $\frac{1}{2}$ probability that Urn 2 was chosen. This is the **prior** probability distribution: $P(\text{Urn 1}) = \frac{1}{2}$ and $P(\text{Urn 2}) = \frac{1}{2}$.

Now your friend (secretly) takes a random ball from the chosen urn, looks at it, and returns it to the urn. He tells you that it was a *Red* ball. You must now estimate the probability that Urn 1 was chosen and the probability that Urn 2 was chosen. Naïvely, you can guess, correctly, that Urn 1 is more likely since it has a higher proportion of *Red* balls than Urn 2 has. We use Bayes' Rule to calculate the exact probabilities as follows:

$$P(\text{Urn 1}|\text{Red}) = \frac{P(\text{Red}|\text{Urn 1})P(\text{Urn 1})}{P(\text{Red})}$$

We can work out that $P(\text{Red}|\text{Urn 1}) = \frac{3}{5}$ since 3 of the 5 balls in Urn 1 are *Red*. Also $P(\text{Urn 1}) = \frac{1}{2}$ as discussed above. To find $P(\text{Red})$ we use marginalisation:

$$\begin{aligned} P(\text{Red}) &= P(\text{Red}|\text{Urn 1})P(\text{Urn 1}) + P(\text{Red}|\text{Urn 2})P(\text{Urn 2}) \\ &= \left(\frac{3}{5}\right)\left(\frac{1}{2}\right) + \left(\frac{3}{8}\right)\left(\frac{1}{2}\right) \\ &= \frac{39}{80}. \end{aligned}$$

Thus,

$$\begin{aligned} P(\text{Urn 1}|\text{Red}) &= \frac{P(\text{Red}|\text{Urn 1})P(\text{Urn 1})}{P(\text{Red})} \\ &= \frac{\left(\frac{3}{5}\right)\left(\frac{1}{2}\right)}{\frac{39}{80}} \\ &= \frac{24}{39}. \end{aligned}$$

It follows that $P(\text{Urn 2}|\text{Red}) = 1 - \frac{24}{39} = \frac{15}{39}$.

Thus, our posterior probability distribution is: $P(\text{Urn 1}) = \frac{24}{39}$ and $P(\text{Urn 2}) = \frac{15}{39}$.

Let us suppose now that the friend (secretly) takes a second random ball from the chosen urn, and again, the ball is *Red*. This is more evidence that the chosen urn is Urn 1, but we can calculate the exact probability using Bayes' Rule again, i.e.:

$$P(\text{Urn 1}|\text{Red}) = \frac{P(\text{Red}|\text{Urn 1})P(\text{Urn 1})}{P(\text{Red})}$$

We still have that $P(\text{Red}|\text{Urn 1}) = \frac{3}{5}$. Our prior probability distribution now has $P(\text{Urn 1}) = \frac{15}{39}$. To get $P(\text{Red})$ we marginalise again:

$$\begin{aligned} P(\text{Red}) &= P(\text{Red}|\text{Urn 1})P(\text{Urn 1}) + P(\text{Red}|\text{Urn 2})P(\text{Urn 2}) \\ &= \left(\frac{3}{5}\right)\left(\frac{24}{39}\right) + \left(\frac{3}{8}\right)\left(\frac{15}{39}\right) \\ &= 0.513 \end{aligned}$$

Thus,

$$\begin{aligned} P(\text{Urn 1}|\text{Red}) &= \frac{P(\text{Red}|\text{Urn 1})P(\text{Urn 1})}{P(\text{Red})} \\ &= \frac{\left(\frac{3}{5}\right)\left(\frac{24}{39}\right)}{0.513} \\ &= 0.720. \end{aligned}$$

It follows that $P(\text{Urn 2}|\text{Red}) = 1 - 0.720 = 0.280$.

Thus, our posterior probability distribution is $P(\text{Urn 1}) = 0.720$ and $P(\text{Urn 2}) = 0.280$.

This makes sense since it much more likely that Urn 1 is the chosen urn.

5.4. Marginalisation.

The formula we used for marginalisation to get $P(\text{Red})$ above is given more generally as follows.

For an event X that can be classified as any one of the classes C_1, \dots, C_n , we have that

$$P(X) = \sum_{i=1}^n P(X|C_i)P(C_i).$$

We say that we are **marginalising** X over the C_i 's.

The above formula can be obtained as follows:

Event X must be classified as one of the classes C_1, \dots, C_n , and it cannot be classified as more than one, so

$$P(C_1|X) + P(C_2|X) + \dots + P(C_n|X) = 1$$

hence, by Bayes' Rule:

$$\frac{P(X|C_1)P(C_1)}{P(X)} + \frac{P(X|C_2)P(C_2)}{P(X)} + \dots + \frac{P(X|C_n)P(C_n)}{P(X)} = 1$$

and the marginalisation formula follows if we multiply through by $P(X)$.

EXERCISES

(1) Consider the following dataset from Lecture 5:

$$S = \left[\begin{array}{c|c|c|c} F_1 & F_2 & F_3 & \text{target} \\ \hline T & 0 & B & Yes \\ F & 2 & A & No \\ F & 1 & C & Yes \\ T & 2 & B & Yes \\ F & 0 & A & No \\ F & 2 & C & No \\ F & 0 & A & Yes \\ T & 1 & B & No \\ T & 0 & B & Yes \\ F & 1 & C & Yes \\ T & 2 & A & No \\ T & 2 & C & No \end{array} \right]$$

Use the Naïve Bayes Classifier to classify input $(F, 1, A)$ as either *Yes* or *No*.

Do the same for inputs $(F, 0, A)$ and for $(T, 2, C)$. Compare your answers with those obtained using decision trees in Lecture 5.

(2) Consider the following dataset from Lecture 5:

$$S = \left[\begin{array}{c|c|c|c} F_1 & F_2 & F_3 & \text{target} \\ \hline a & 0 & Y & A \\ b & 0 & N & C \\ c & 1 & Y & B \\ b & 1 & Y & B \\ c & 0 & N & A \\ a & 0 & N & C \\ c & 1 & N & B \\ a & 1 & Y & A \\ b & 0 & Y & B \\ a & 1 & N & C \end{array} \right]$$

Use the Naïve Bayes Classifier to classify input $(b, 1, N)$ as either A , B or C .

Do the same with input $(c, 1, Y)$. Compare your answers with those obtained by decision trees in Lecture 5.

- (3) Consider an urn problem similar to the example in the notes, but where Urn 1 has 1 *Red* ball and 3 *Blue* balls, and Urn 2 has 7 *Red* balls and 3 *Blue* balls.
 - (a) A friend uses a fair coin to choose an urn. Give the prior probability distribution on the urns.
 - (b) Your friend secretly takes a ball out of the chosen urn (so you can't see which urn) and tells you the ball is *Blue* and then returns it to the urn. Use Bayes' rule to work out $P(\text{Urn 1}|\text{Blue})$ and $P(\text{Urn 2}|\text{Blue})$.
 - (c) Your friend then draws another ball out of the chosen urn, and tells you it is *Red*. Now work out $P(\text{Urn 1}|\text{Red})$ and $P(\text{Urn 2}|\text{Red})$.
- (4) Suppose you have a dataset compiled from 200 randomly chosen people. For each person in the dataset you have some personal information and you also know if the person prefers cricket or soccer. There are 80 people in the dataset that prefer cricket and 120 people that prefer soccer.
 - (a) Suppose you want to know if a new random person prefers cricket or soccer. Based on the sample dataset, what is your prior probability that the person prefers cricket and the prior probability that the person prefers soccer?
 - (b) You are now informed that the new person is female. In your dataset, you know that 45 of the 80 people that prefer cricket are female, and 50 of the 120 people that prefer soccer are female. Use Bayes' rule to update your probability that the person prefers cricket and the probability that the person prefers soccer.
 - (c) You now receive some more information - the person is under 25 years old. In your dataset, of the 80 that prefer cricket, 30 are under 25 and 15 of them are female, and of the 120 that prefer soccer, 50 are under 25 and 30 of them are female. Use Bayes' rule to determine the probability that the person prefers cricket and the probability that the person prefers soccer.
- (5) Suppose you are given a dataset S with N datapoints. Suppose S has attributes x_1, x_2, x_3 and target t , where attribute x_1 can take values in $\{T, F\}$, attribute x_2 can take values in $\{a, b, c\}$ and attribute x_3 can take values in $\{0, 1, 2, 3\}$. The target t can take values in $\{C_1, C_2, C_3\}$. Describe how to implement the Naïve Bayes Classifier for any input $\mathbf{x} = (x_1, x_2, x_3)$, where $x_1 \in \{T, F\}$, $x_2 \in \{a, b, c\}$ and $x_3 \in \{0, 1, 2, 3\}$.