

School of Computing and Information Systems
The University of Melbourne
COMP90049 INTRODUCTION to MACHINE LEARNING
(Semester 1, 2021)

Week 5: Sample Solutions

1. How is **holdout** evaluation different to **cross-validation** evaluation? What are some reasons we would prefer one strategy over the other?

In a holdout evaluation strategy, we partition the data into a training set and a test set: we build the model on the former and evaluate on the latter.

In a cross-validation evaluation strategy, we do the same as above, but a number of times, where each iteration uses one partition of the data as a test set and the rest as a training set (and the partition is different each time).

Why we prefer cross-validation to holdout evaluation strategy? Because, holdout is subject to some random variation, depending on which instances are assigned to the training data, and which are assigned to the test data. Any instance that forms part of the model is excluded from testing, and vice versa. This could mean that our estimate of the performance of the model is way off or changes a lot from data set to data set.

While Cross-validation mostly solves this problem: we're averaging over a bunch of values, so that one weird partition of the data won't throw our estimate of performance completely off; also, each instance is used for testing, but also appears in the training data for the models built on the other partitions. It usually takes much longer to cross-validate, however, because we need to train a model for every test partition.

2. A **confusion matrix** is a summary of the performance of a (supervised) classifier over a set of development ("test") data, by counting the various instances:

		Actual			
		<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>
Classified	<i>a</i>	10	2	3	1
	<i>b</i>	2	5	3	1
	<i>c</i>	1	3	7	1
	<i>d</i>	3	0	3	5

- (i). Calculate the classification **accuracy** of the system. Find the **error rate** for the system.

In this context, Accuracy is defined as the fraction of correctly identified instances, out of all of the instances. In the case of a confusion matrix, the correct instances are the ones enumerated along the main diagonal (classified as *a* and actually *a* etc.):

$$\begin{aligned} \text{Accuracy} &= \frac{\text{\# of correctly identified instance}}{\text{total \# of instances}} \\ &= \frac{10 + 5 + 7 + 5}{10 + 2 + 3 + 1 + 2 + 5 + 3 + 1 + 1 + 3 + 7 + 1 + 3 + 0 + 3 + 5} \\ &= \frac{27}{50} = 54\% \end{aligned}$$

Error rate is just the complement of accuracy:

$$\begin{aligned} \text{Error Rate} &= \frac{\text{\# of incorrectly identified instance}}{\text{total \# of instances}} = 1 - \text{Accuracy} = 1 - 54\% \\ &= 46\% \end{aligned}$$

- (ii). Calculate the **precision**, **recall** and **F-score** (where $\beta = 1$), for class d .

Precision for a given class is defined as the fraction of correctly identified instances of that class, from the times that class was attempted to be classified. We are interested in the true positives (TP) where we attempted to classify an item as an instance of said class (in this case, d) and it was actually of that class (d): in this case, there are 5 such instances. The false positives (FP) are those items that we attempted to classify as being of class d , but they were actually of some other class: there are $3 + 0 + 3 = 6$ of those.

$$Precision = \frac{TP}{TP + FP} = \frac{5}{5 + 3 + 0 + 3} = \frac{5}{11} \approx 45\%$$

Recall for a given class is defined as the fraction of correctly identified instance of that class, from the times that class actually occurred. This time, we are interested in the true positives, and the false negatives (FN): those items that were actually of class d , but we classified as being of some other class; there are $1 + 1 + 1 = 3$ of those.

$$Recall = \frac{TP}{TP + FN} = \frac{5}{5 + 1 + 1 + 1} = \frac{5}{8} \approx 62\%$$

F-score is a measure which attempts to combine Precision (P) and Recall (R) into a single score. In general, it is calculated as:

$$F_{\beta} = \frac{(1 + \beta^2) P.R}{(\beta^2.P) + R}$$

By far, the most typical formulation is where the parameter β is set to 1: this means that Precision and Recall are equally important to the score, and that the score is a harmonic mean.

In this case, we have calculated the Precision of class d to be 0.45 and the Recall to be 0.62. The F-score where ($\beta = 1$) of class d is then:

$$F_1 = \frac{2 P.R}{P + R} = \frac{2 \times 0.45 \times 0.62}{0.45 + 0.62} \approx 53\%$$

- (iii). Why can't we do this for the whole system? How can we consider the whole system?

The concept of precision and recall is defined per-class on the bases of one (interesting) class vs the rest of (not interesting) classes.

Since this system, similar to all other multiclass classifiers, considers all classes (a , b , c and d) as interesting, we need to calculate the precision and recall per each class (vs the rest) and then find the average for the whole system.

As covered in the lectures, there are multiple methods for calculating the average precision and recall for the whole system. We can use methods like *Macro* Averaging, *Micro* Averaging or *Weighted* Averaging.

Our choice of method depends on our goal and the domain of the system. In cases where we want to emphasis on identifying the behaviour of the system for small classes *Macro* averaging can be a better choice. While in situations that we want to evaluate the system mostly based on its behaviour in detecting the large classes *Micro* averaging would be a better option.

3. For the following dataset:

<i>ID</i>	<i>Outl</i>	<i>Temp</i>	<i>Humi</i>	<i>Wind</i>	<i>PLAY</i>
TRAINING INSTANCES					
A	s	h	h	F	N
B	s	h	h	T	N
C	o	h	h	F	Y
D	r	m	h	F	Y
E	r	c	n	F	Y
F	r	c	n	T	N
TEST INSTANCES					
G	o	c	n	T	?
H	s	m	h	F	?

(i). Classify the test instances using the method of **0-R**.

0-R is the quintessentially baseline classifier: we throw away all of the attributes, other than the class labels, and just predict each test instance according to whichever label is most common in the training data. (Hence, it is also common called the “majority class classifier”.)

In this case, the two labels Y and N are equally common in the training data — so we are required to apply a tiebreaker. Remember that we’re simply choosing one label to be representative of the entire collection, and both labels seem equally good for that here: so, let’s say N.

Consequently, both test instances are classified as N.

(ii). Classify the test instances using the method of **1-R**.

1-R is a slightly better baseline, which requires us to choose a single attribute to represent the entire decision-making process. For example, if *Outl* is our preferred attribute, then we base the classification of each test instance solely based on its value of *Outl*. (This is sometimes called a “Decision Stump”.)

Given our preferred attribute, we will label a test instance according to whichever label in the training data was most common, for training instances with the corresponding attribute value.

How do we decide which attribute to choose? Well, the most common method is simply by counting the errors made on the training instances.

Let’s say we chose *Outl*: first, we need to observe the predictions for the 3 values (s, o, and r), and then we will count up the errors that those predictions will make on the training data (it turns out that we can do this simultaneously):

- When *Outl* is s, there are two such instances in the training data. Both of these instances are labelled as N: our label for this attribute value will be N. We will therefore predict the label of both of these training instances correctly (we will predict N, and both of these instances are actually N).
- When *Outl* is o, there is just a single instance in the training data, labelled as Y: our label for this attribute value will be Y. We will therefore predict the label of this training instance correctly.
- When *Outl* is r, there are three such instances in the training data. Two of these instances are labelled as Y, the other is N: our label for this attribute value will be Y (as there are more Y instances than N instances — you can

see that we're applying the method of 0-R here). We will therefore make one error: the instance which was actually n.

In total, that is 1 error for `Outl`. Hopefully, you can see that this is a very wordy explanation of a very simple idea. (We will go through the other attributes more quickly.)

For `Temp`:

- When `Temp` is h, there are two N instances and one Y: we will make one mistake.
- When `Temp` is m, there is one Y instance: we won't make any mistakes.
- When `Temp` is c, there is one N instance and one Y instance: we will make one mistake.

This is 2 errors in total for `Temp`, which means that it's less good than `Outl`.

We'll leave the other attributes as an exercise, and assume that `Outl` was the best attribute (it's actually tied with `Wind`): how do the test instances get classified?

- For test instance G, `Outl` is o — the 0-R classifier over the o instances from the training data gives Y, so we predict Y for this instance.
- For test instance H, we don't have a value for `Outl` so we cannot find the label. But if we use the assumption of `Outl = s` (as suggested by the question) — the 0-R classifier over the s instances from the training data gives N, so we predict N for instance H.

4. Given the above dataset, we wished to perform feature selection on this dataset, where the class is `PLAY`:

<i>ID</i>	<i>Outl</i>	<i>Temp</i>	<i>Humi</i>	<i>Wind</i>	<i>PLAY</i>
A	s	h	h	F	N
B	s	h	h	T	N
C	o	h	h	F	Y
D	r	m	h	F	Y
E	r	c	n	F	Y
F	r	c	n	T	N

- (i). Which of `Humi` and `Wind` has the greatest *Pointwise Mutual Information* for the class Y? What about N?

To determine Pointwise Mutual Information (PMI), we compare the joint probability to the product of the prior probabilities as follows:

$$PMI(A, C) = \log_2 \frac{P(A \cap C)}{P(A)P(C)}$$

Note that this formulation only really makes sense for binary attributes and binary classes (which we have here.)

Suitably interpreting $P(X)$ as $P(X = Y)$ (or T or h here), we find:

$$PMI(Humi, Play) = \log_2 \frac{P(Humi \cap Play)}{P(Humi)P(Play)}$$

$$\begin{aligned}
&= \log_2 \frac{\frac{2}{6}}{\frac{4}{6} \frac{3}{6}} = \log_2(1) = 0 \\
PMI(Wind, Play) &= \log_2 \frac{P(Wind \cap Play)}{P(Wind)P(Play)} \\
&= \log_2 \frac{\frac{0}{6}}{\frac{2}{6} \frac{3}{6}} = \log_2(0) = -\infty
\end{aligned}$$

So, we find that `Wind` is (perfectly) *negatively correlated* with `PLAY`; whereas `Humi` is (perfectly) *uncorrelated*.

You should compare these results with the negative class `PLAY=N`, where `Wind` is positively correlated, but `Humi` is still uncorrelated.

(ii). Which of the attributes has the greatest *Mutual Information* for the class, as a whole?

A general form of the Mutual Information (MI) is as follows:

$$MI(X, C) = \sum_{x \in X} \sum_{c \in \{Y, N\}} P(x, c) PMI(x, c)$$

Effectively, we're going to consider the PMI of every possible attribute value–class combination, weighted by the proportion of instances that actually had that combination.

To handle cases like $PMI(Wind)$ above, we are going to equate $0 \log 0 \equiv 0$ (which is true in the limit anyway).

For `Outl`, this is going to look like:

$$\begin{aligned}
MI(Outl) &= P(s, Y)PMI(s, Y) + P(o, Y)PMI(o, Y) + P(r, Y)PMI(r, Y) + \\
&\quad P(s, N)PMI(s, N) + P(o, N)PMI(o, N) + P(r, N)PMI(r, N) \\
&= \frac{0}{6} \log_2 \frac{\frac{0}{6}}{\frac{2}{6} \frac{3}{6}} + \frac{1}{6} \log_2 \frac{\frac{1}{6}}{\frac{1}{6} \frac{3}{6}} + \frac{2}{6} \log_2 \frac{\frac{2}{6}}{\frac{3}{6} \frac{3}{6}} + \\
&\quad \frac{2}{6} \log_2 \frac{\frac{2}{6}}{\frac{2}{6} \frac{3}{6}} + \frac{0}{6} \log_2 \frac{\frac{0}{6}}{\frac{1}{6} \frac{3}{6}} + \frac{1}{6} \log_2 \frac{\frac{1}{6}}{\frac{3}{6} \frac{3}{6}} \\
&\approx 0 \log_2 0 + (0.1667)(1) + (0.3333)(0.4150) + \\
&\quad (0.3333)(1) + 0 \log_2 0 + (0.1667)(-0.5850) \\
&\approx 0.541
\end{aligned}$$

It's worth noting that while some individual terms can be negative, the sum must be at least zero.

For `Temp`, this is going to look like:

$$\begin{aligned}
MI(Temp) &= P(h, Y)PMI(h, Y) + P(m, Y)PMI(m, Y) + P(c, Y)PMI(c, Y) + \\
&\quad P(h, N)PMI(h, N) + P(m, N)PMI(m, N) + P(c, N)PMI(c, N) \\
&= \frac{1}{6} \log_2 \frac{\frac{1}{6}}{\frac{3}{6} \frac{3}{6}} + \frac{1}{6} \log_2 \frac{\frac{1}{6}}{\frac{1}{6} \frac{3}{6}} + \frac{1}{6} \log_2 \frac{\frac{1}{6}}{\frac{2}{6} \frac{3}{6}} +
\end{aligned}$$

$$\begin{aligned}
& \frac{2}{6} \log_2 \frac{\frac{2}{6}}{\frac{2}{6} \frac{3}{6}} + \frac{0}{6} \log_2 \frac{\frac{0}{6}}{\frac{1}{6} \frac{3}{6}} + \frac{1}{6} \log_2 \frac{\frac{1}{6}}{\frac{2}{6} \frac{3}{6}} \\
& \approx (0.1667)(-0.5850) + (0.1667)(1) + (0.1667)(0) + \\
& \quad (0.3333)(0.4150) + 0 \log_2 0 + (0.1667)(-0.5850) \\
& \approx 0.110
\end{aligned}$$

We will leave the workings as an exercise, but the Mutual Information for `Humi` is 0, and for `Wind` is 0.459.

Consequently, `Outl` appears to be the best attribute (perhaps as we might expect), and `Wind` also seems quite good; whereas `Temp` is not very good, and `Humi` is completely unhelpful.