

School of Computing and Information Systems
The University of Melbourne
COMP90049 Knowledge Technologies (Semester 2, 2018)
Workshop sample solutions: Week 10

1. Revise Support Vector Machines, paying particular attention to the terms “linear separability” and “maximum margin”.
 - Support vector machines attempt to partition the training data based on the best line (hyperplane) that divide the positive instances of the class that we’re looking for from the negative instances.
 - We say that the data is *linearly separable*, when the points in our k -dimensional vector space corresponding to positive instances can indeed be separated from the negative instances by a line (or, more accurately, a $(k - 1)$ -dimensional hyperplane). This means that all of the positive instances are on one side of the line, and all of the negative instances are on the other side.
 - What is the definition of “best” line? Well, in fact, we choose a pair of parallel lines, one for the positive instances, and one for the negative instances. The pair that we choose is the one that has the greatest perpendicular distance between these parallel lines (calculated using the normal; because the lines are parallel, they have the same normal) — this is called the “margin”.
 - These two lines are called the “support vectors” — we can classify test instances by calculating (again, using the normal) which of the support vectors is closer to the point defined by the test instance. Alternatively, we can just use a single line, halfway between the two support vectors, and use the normal to calculate which “side” of the line the test instance is on.
- (a) What are “soft margins”, and when are they desirable?
 - Soft margins are when we relax the notion of linear separability. A small number of points are allowed to be on the “wrong” side of the line, if we get a (much) better set of support vectors that way (i.e. with a (much) larger margin).
 - Sometimes the data is *almost* linearly separable — if we accept that we probably aren’t going to classify every single test instance correctly anyway, we can produce a classifier that hopefully generalises better to unseen data.
- (b) Why are we interested in “kernel functions” here?
 - The so-called “kernel trick” is often used to transform data.
 - For Support Vector Machines, sometimes the data isn’t linearly separable, but after applying some kind of function — where some useful properties of the data remain (for example, monotonicity of the inner product) — the data becomes linearly separable. If we do this, we can apply the algorithm as usual.
- (c) Why are SVMs “binary classifiers”, and how can we extend them to “multi-class classifiers”?
 - We’re trying to find a line that partitions the data into true and false, suitably interpreted for our data.
 - The point is that for two classes, we only need a single line. (Or more, accurately, a pair of parallel support vectors.)
 - We need (at least) two lines to partition data into three classes. For example, $x < 0$, $0 < x < 1$, $x > 1$ might be a useful partition. But what happens when these lines aren’t parallel?
 - In general, we might want to find three “half-lines”, radiating out from some central point between the (points corresponding to) the three different classes. The problem? This is numerically much more difficult (and usually intractable) to solve.
 - For problems with more than three classes, the entire notion of separating the data becomes poorly-defined — unless we wish to fall back to some clustering approach, in which case, we might as well use Nearest Prototype.

- In practice, then, we don't try to do the above, but rather build *multiple* (binary) SVMs, either by comparing the instances of one class against the instances of all of the other classes (one-vs-many) and then choosing the most confident prediction, or by considering each pair of classes (one-vs-one) and then choosing the class that comes up the most frequently.

2. For the following dataset:

<i>apple</i>	<i>ibm</i>	<i>lemon</i>	<i>sun</i>	CLASS
TRAINING INSTANCES				
Y	N	Y	Y	FRUIT
Y	N	Y	Y	FRUIT
Y	Y	N	N	COMPUTER
Y	Y	Y	Y	COMPUTER

Build a contingency table for each of the four attributes on the data collection above.

- We want four contingency tables: one for each attribute. Notice that the attributes are binary (i.e. Y or N) — we could do this with numerical data, but our methods of interpreting the probabilities would be different¹.
- Now, the *apple* contingency table will be built around the number of instances where *apple* is Y for each of the classes. For example, *apple* is Y for 2 (both) of the *fruit* instances, and N for 0 (neither), and so on.

<i>apple</i>	FRUIT	COMPUTER
Y	2	2
N	0	0

<i>ibm</i>	FRUIT	COMPUTER
Y	0	2
N	2	0

<i>lemon</i>	FRUIT	COMPUTER
Y	2	1
N	0	1

<i>sun</i>	FRUIT	COMPUTER
Y	2	1
N	0	1

Table 1: Contingency tables for the above dataset

- (a) According to “Pointwise Mutual Information”, which attribute has the best correlation with the class COMPUTER?
- Pointwise Mutual Information is calculated according to the following formula (which you can compare with Mutual Information below):

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

- Here, we read these values as: A means $A = Y$ and C , $C = Y$.

¹In the Project 2 “best” datasets, the term frequencies were ignored for Feature Selection — effectively, attributes were treated as Y is the term frequency was greater than 0.

- To assess the PMI of *apple* with COMPUTER, we first need to calculate each of the prior probabilities, and the joint probability (according to the instance-based model):

$$\begin{aligned} P(a) &= \frac{4}{4} \\ P(c) &= \frac{2}{4} \\ P(a, c) &= \frac{2}{4} \end{aligned}$$

- Now we can substitute:

$$\begin{aligned} PMI(a, c) &= \log_2 \frac{P(a, c)}{P(a)P(c)} \\ &= \log_2 \frac{\frac{2}{4}}{\frac{4}{4} \frac{2}{4}} \\ &= \log_2 \frac{0.5}{0.5} = \log_2(1) = 0 \end{aligned}$$

- The PMI of *apple* with respect to the class COMPUTER is 0, which means that there is no correlation, and (probably) no predictive capacity.
- What about the other attributes?

$$\begin{aligned} P(i) &= \frac{2}{4} \\ P(i, c) &= \frac{2}{4} \\ PMI(i, c) &= \log_2 \frac{P(i, c)}{P(i)P(c)} \\ &= \log_2 \frac{\frac{2}{4}}{\frac{2}{4} \frac{2}{4}} \\ &= \log_2 \frac{0.5}{0.25} = \log_2(2) = 1 \end{aligned}$$

$$\begin{aligned} P(l) &= \frac{3}{4} \\ P(l, c) &= \frac{1}{4} \\ PMI(l, c) &= \log_2 \frac{P(l, c)}{P(l)P(c)} \\ &= \log_2 \frac{\frac{1}{4}}{\frac{3}{4} \frac{2}{4}} \\ &= \log_2 \frac{0.25}{0.375} = \log_2\left(\frac{2}{3}\right) \approx -0.58 \\ P(s) &= \frac{3}{4} \\ P(s, c) &= \frac{1}{4} \\ PMI(s, c) &= \log_2 \frac{P(s, c)}{P(s)P(c)} \\ &= \log_2 \frac{\frac{1}{4}}{\frac{3}{4} \frac{2}{4}} \\ &= \log_2 \frac{0.25}{0.375} = \log_2\left(\frac{2}{3}\right) \approx -0.58 \end{aligned}$$

- All in all, *ibm* has the greatest Pointwise Mutual Information with the class COMPUTER, as we might expect, because they are perfectly correlated.
 - What if we had found the PMI of FRUIT instead? Well, in that case, we would have thought that *lemon* and *sun* were the best attributes (with a PMI of 0.42); we wouldn't have found *ibm* as a good attribute because it is reverse-correlated with FRUIT.
- (b) Use the method of “Mutual Information” to rank the “goodness” of the four features in predicting this two-class problem, according to the following formula:

$$MI(A, C) = \sum_{i \in \{a, \bar{a}\}} \sum_{j \in \{c, \bar{c}\}} P(i, j) \log_2 \frac{P(i, j)}{P(i)P(j)}$$

- This is very similar to PMI, but we are going to combine the PMIs of the attribute occurring together with each class, as well as not occurring (having a frequency of 0) with each class, weighted by the number of times this has happened.
- All of this information can be trivially read off the contingency tables above. For example, for *apple*:

$$\begin{aligned} P(a) &= \frac{4}{4} \\ P(\bar{a}) &= 0 \\ P(F) &= \frac{2}{4} \\ P(C) &= \frac{2}{4} \\ P(a, F) &= \frac{2}{4} \\ P(a, C) &= \frac{2}{4} \\ P(\bar{a}, F) &= 0 \\ P(\bar{a}, C) &= 0 \end{aligned}$$

- Substituting in (and taking $0 \log x \equiv 0$), we find:

$$\begin{aligned} MI(a) &= P(a, F) \log_2 \frac{P(a, F)}{P(a)P(F)} + P(\bar{a}, F) \log_2 \frac{P(\bar{a}, F)}{P(\bar{a})P(F)} + \\ &\quad P(a, C) \log_2 \frac{P(a, C)}{P(a)P(C)} + P(\bar{a}, C) \log_2 \frac{P(\bar{a}, C)}{P(\bar{a})P(C)} \\ &= \frac{2}{4} \log_2 \frac{\frac{2}{4}}{\frac{4}{4} \frac{2}{4}} + 0 \log_2 \frac{0}{(0) \frac{2}{4}} + \frac{2}{4} \log_2 \frac{\frac{2}{4}}{\frac{4}{4} \frac{2}{4}} + 0 \log_2 \frac{0}{(0) \frac{2}{4}} \\ &= \frac{2}{4}(0) + 0 + \frac{2}{4}(0) + 0 = 0 \end{aligned}$$

- The value of 0 here means that, regardless of whether *apple* is present or absent from the instance, it has no predictive power for either of these classes.
- What about the other attributes? *ibm*:

$$\begin{aligned} P(i) &= \frac{2}{4} \\ P(\bar{i}) &= \frac{2}{4} \\ P(i, F) &= 0 \\ P(i, C) &= \frac{2}{4} \\ P(\bar{i}, F) &= \frac{2}{4} \\ P(\bar{i}, C) &= 0 \end{aligned}$$

$$\begin{aligned}
MI(i) &= P(i, f) \log_2 \frac{P(i, f)}{P(i)P(f)} + P(\bar{i}, f) \log_2 \frac{P(\bar{i}, f)}{P(\bar{i})P(f)} + \\
&\quad P(i, c) \log_2 \frac{P(i, c)}{P(i)P(c)} + P(\bar{i}, c) \log_2 \frac{P(\bar{i}, c)}{P(\bar{i})P(c)} \\
&= 0 \log_2 \frac{0}{\frac{2}{4}\frac{2}{4}} + \frac{2}{4} \log_2 \frac{\frac{2}{4}}{\frac{2}{4}\frac{2}{4}} + \frac{2}{4} \log_2 \frac{\frac{2}{4}}{\frac{2}{4}\frac{2}{4}} + 0 \log_2 \frac{0}{\frac{2}{4}\frac{2}{4}} \\
&= 0 + \frac{2}{4}(1) + \frac{2}{4}(1) + 0 = 1
\end{aligned}$$

- *lemon*:

$$\begin{aligned}
P(l) &= \frac{3}{4} \\
P(\bar{l}) &= \frac{1}{4} \\
P(l, f) &= \frac{2}{4} \\
P(l, c) &= \frac{1}{4} \\
P(\bar{l}, f) &= 0 \\
P(\bar{l}, c) &= \frac{1}{4}
\end{aligned}$$

$$\begin{aligned}
MI(l) &= P(l, f) \log_2 \frac{P(l, f)}{P(l)P(f)} + P(\bar{l}, f) \log_2 \frac{P(\bar{l}, f)}{P(\bar{l})P(f)} + \\
&\quad P(l, c) \log_2 \frac{P(l, c)}{P(l)P(c)} + P(\bar{l}, c) \log_2 \frac{P(\bar{l}, c)}{P(\bar{l})P(c)} \\
&= \frac{2}{4} \log_2 \frac{\frac{2}{4}}{\frac{3}{4}\frac{2}{4}} + 0 \log_2 \frac{0}{\frac{1}{4}\frac{2}{4}} + \frac{1}{4} \log_2 \frac{\frac{1}{4}}{\frac{3}{4}\frac{2}{4}} + \frac{1}{4} \log_2 \frac{\frac{1}{4}}{\frac{1}{4}\frac{2}{4}} \\
&\approx \frac{2}{4}(0.42) + 0 + \frac{1}{4}(-0.58) + \frac{1}{4}(1) \approx 0.32
\end{aligned}$$

- *sun* is exactly the same as *lemon* (check the contingency tables above), so I won't repeat the calculations here.
- All in all, *ibm* appears to be the best attribute for this dataset, as we might expect. *lemon* and *sun* are both marginally useful, and *apple* is completely useless (so maybe we should get rid of it! :)).

3. For the following dataset:

<i>ID</i>	<i>Outl</i>	<i>Temp</i>	<i>Humi</i>	<i>Wind</i>	PLAY
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	?

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

- 0-R classifies a test instance based on the majority class of the training data.
 - In this case, the training data has 3 N instances and 3 Y instances: there is no majority class.
 - Consequently, we believe that both classes are equally adequate at describing this data, and we would just choose one arbitrarily. Let's say we choose the lexicographically smaller one: N. Both test instances would then be classified as N.
- (b) Classify the test instances using the method of 1-R.
- 1-R chooses a single attribute on which to make the decision for classifying all of the test instances. The attribute that it chooses is the one which makes the fewest errors, when used to classify the training data.
 - The class chosen for a given attribute value is the equivalent of applying the 0-R method, for the instances with that attribute value. Consequently, each error it makes on the training data corresponds to an instance **not** of the majority class.
 - Let's look at it in action:
 - First, *Outl*:
 - *Outl* takes 3 values: **s**, **o** and **r**.
 - There are 2 **s** instances, which both have the class N. If we apply 0-R here, we would choose N for both instances, and make 0 errors.
 - There is just 1 **o** instance, which has the class Y. 0-R would (trivially) choose Y, and make 0 errors.
 - There are 3 **r** instances, 2 Y and 1 N. If we apply 0-R here, we would choose Y for all 3 instances, and make 1 error (for the **r** instance with the class Y).
 - Altogether, this is $0 + 0 + 1 = 1$ error.
 - *Temp*:
 - *Temp* takes 3 values: **h**, **m** and **c**.
 - There are 3 **h** instances, 2 Y and 1 N. If we apply 0-R here, we would choose Y for all 3 instances, and make 1 error.
 - There is just 1 **m** instance, which has the class Y. 0-R makes 0 errors.
 - There are 2 **c** instances, 1 Y and 1 N. If we apply 0-R here, we would make 1 error, regardless of whether we chose Y or N.
 - Altogether, this is $1 + 0 + 1 = 2$ errors.
 - *Humi*:
 - *Humi* takes 2 values: **h** and **n**.
 - There are 4 **h** instances, 2 Y and 2 N. 0-R makes 2 errors.
 - There are 2 **n** instances, 1 Y and 1 N. 0-R makes 1 error.
 - Altogether, this is $2 + 1 = 3$ errors.
 - *Wind*:
 - *Wind* takes 2 values: **F** and **T**.
 - There are 4 **F** instances, 3 Y and 1 N. 0-R makes 1 error.
 - There are 2 **T** instances, both N. 0-R makes 0 errors.
 - Altogether, this is $1 + 0 = 1$ error.
 - So, we choose the attribute with the fewest errors; in this case, there is a tie between *Outl* and *Wind* (1 error in total). We would need to break the tie arbitrarily.
 - Let's say we chose *Outl*, how would we classify the test instances? Well, we would read off the value of *Outl* for each test instance, and classify it using 0-R (for the training attributes with the same value of *Outl*), similar to the above. For the test instance G, *Outl* is **o**, and we would choose Y, based on the single training instance with **o**. For H, *Outl* is **s**, and we would classify it as N.
 - If instead we had chosen *Wind*, G takes the value T, so we choose N, and H takes the value F, so we choose Y. Note that these are the opposite classifications as the ones we would have made based on *Outl*!

- I can't help but notice that we skipped one of the attributes: *ID*. What happens for that attribute:
 - *ID* takes 6 different values (in the training data): **A**, **B**, **C**, **D**, **E**, and **F**.
 - When *ID* is **A**, there is just a single instance, of class **N**: 0-R makes 0 errors.
 - When *ID* is **B**, there is just a single instance, of class **N**: 0-R makes 0 errors.
 - When *ID* is **C**, there is just a single instance, of class **Y**: 0-R makes 0 errors.
 - And so on. You can see that we obviously aren't going to make any errors with this attribute, because each attribute value only occurs with a single training instance.
- So, 1-R thinks that, actually *ID* is the best attribute — unfortunately, though, it's completely useless for classifying the test data, because we haven't seen those values (**G** and **H**) in the training data!