

Assignment Project Exam Help

COMP9318 Tutorial 2: Classification

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Q1 I

Consider the following training dataset and the original decision tree induction algorithm (ID3).

Risk is the class label attribute. The *Height* values have been already discretized into disjoint ranges.

1. Calculate the information gain if *Gender* is chosen as the test attribute.
2. Calculate the information gain if *Height* is chosen as the test attribute.
3. Dr
4. Ge

F	(1.5, 1.6]	Low
M	(1.9, 2.0]	High
F	(1.8, 1.9]	Medium
F	(1.8, 1.9]	Medium
M	(1.6, 1.7]	Low
M	(1.8, 1.9]	Medium
F	(1.5, 1.6]	Low
M	(1.6, 1.7]	Low
M	(2.0, ∞]	High
M	(2.0, ∞]	High
F	(1.7, 1.8]	Medium
M	(1.9, 2.0]	Medium
F	(1.8, 1.9]	Medium
F	(1.7, 1.8]	Medium
F	(1.7, 1.8]	Medium

Solution to Q1 I

1. The original entropy is $I_{Risk} = I(Low, Medium, High) = I(4, 8, 3) = 1.4566$. Consider *Gender*.

Gender	entropy
F	$I(3, 6, 0)$
M	$I(1, 2, 3)$
<u>9</u>	<u>6</u>

2. <https://eduassistpro.github.io>

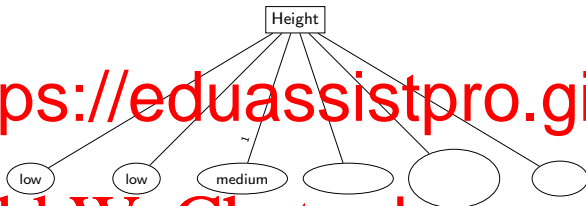
Height	entropy
(1.5, 1.6]	$I($
(1.6, 1.7]	$I($
(1.7, 1.8]	$I($
(1.8, 1.9]	$I($
(1.9, 2.0]	$I(0, 1, 1)$
(2.0, ∞]	$I(0, 0, 2)$

The expected entropy is $\frac{2}{15} \cdot I(2, 0, 0) + \frac{2}{15} \cdot I(2, 0, 0) + \frac{3}{15} \cdot I(0, 3, 0) + \frac{4}{15} \cdot I(0, 4, 0) + \frac{2}{15} \cdot I(0, 1, 1) + \frac{2}{15} \cdot I(0, 0, 2) = 0.1333$. The information gain is $1.4566 - 0.1333 = 1.3233$

Solution to Q1 II

3. ID3 decision tree:

- ▶ According to the computation above, we should first choose *Height* to split
 - ▶ After split, the only problematic partition is the (1.9, 2.0] one. However, the only remaining attribute *Gender* cannot divide them. As there is a draw, we can use any label.
- The final tree is shown in the figure below.



4. The rules are

- ▶ IF $height \in (1.5, 1.6]$, THEN $Rish = Lo$
- ▶ IF $height \in (1.6, 1.7]$, THEN $Rish = Low$.
- ▶ IF $height \in (1.7, 1.8]$, THEN $Rish = Medium$.
- ▶ IF $height \in (1.8, 1.9]$, THEN $Rish = Medium$.
- ▶ IF $height \in (1.9, 2.0]$, THEN $Rish = Medium$ (or High).
- ▶ IF $height \in (2.0, \infty]$, THEN $Rish = High$.

Q2 I

Consider applying the SPRINT algorithm on the following training dataset

<i>Age</i>	<i>CarType</i>	<i>Risk</i>
23	family	High
17	sports	High
43	sports	High
68	family	Low

Answ

1. Wr actively.
2. Assume the first split criterion is $Age <$
lists for the left child node (i.e. corresponding to the pa
 $Age < 27.5$)
3. Assume that the two attribute lists for the root node are sto
relational tables name AL_Age and $AL_CarType$, respectively. We can in
fact generate the attribute lists for the child nodes using standard SQL
statements. Write down the SQL statements which will generate the
attribute lists for the left child node for the split criterion $Age < 27.5$.
4. Write down the final decision tree constructed by the SPRINT algorithm.

Solution to Q2 I

- Attribute list of Age is:

Age	class	Index
17	High	2
20	High	6
23	High	1
32	Low	5
43	High	3

<https://eduassistpro.github.io>

Add WeChat: edu_assist_pro

family	High	1
sports	High	
sports	High	
family	Low	
truck	Low	
family	High	6

- Attribute list of Age is:

Age	class	Index
17	High	2
20	High	6
23	High	1

Solution to Q2 II

Attribute list of *CarType* is:

<i>CarType</i>	class	Index
family	High	1
sports	High	2
family	High	6

SQL for the attribute list of *Age*:

```
SELECT Age, Class, Index
```

```
FROM AL_Age A, AL_CarType C
WHERE A.Age < 27.5
AND A.index = C.index
```

- Consider the attribute list of *Age*: there are 5 p each of them have gini index value as:

<i>Age</i>	above	below	<i>gini_{split}</i>
17 – 20	(1, 0)	(3, 2)	0.40
20 – 23	(2, 0)	(2, 2)	0.33
23 – 32	(3, 0)	(1, 2)	0.22
32 – 43	(3, 1)	(1, 1)	0.42
43 – 68	(4, 1)	(0, 1)	0.27

Solution to Q2 III

therefore, the best split should be $Age > 27.5$.

Consider the attribute list of *CarType*:

<i>CarType</i>	High	Low
f	2	1
s	2	0
t	0	1

Assignment Project Exam Help

<https://eduassistpro.github.io>

<i>CarType</i>	Hig	
t	0	
f, s	4	

Add WeChat edu_assist_pr

Each of them have gini index value as: 0.44, 0.33, 0.27, re

Therefore, the best split is *CarType* in ('truck').

Obviously, splitting on *Age* is better. Therefore, we shall split by $Age > 27.5$.

The attribute lists for each of the child node have already been computed.

Since the tuples in the partition for $Age < 27.5$ are all "high", we only need to look at the partition for $Age \geq 27.5$.

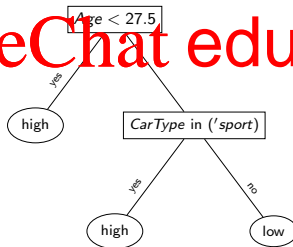
Age	class	Index
32	Low	5
43	High	3
68	Low	4

Car Type	class	Index
sports	High	3
family	Low	4

Assignment Project Exam Help

<https://eduassistpro.github.io>

The final tree is:



Consider a (simplified) email classification example. Assume the training dataset contains 1000 emails in total, 100 of which are spams.

1. Calculate the class prior probability distribution. How would you classify a new incoming email?

2. A friend of you suggests that whether the email contains a \$ char is a good feature to detect spam emails. You look into the training dataset

a $\bar{\$}$
a $\$$
<https://eduassistpro.github.io>

SPAM	91	9
NO SPAM	63	837

Describe the (naive) Bayes Classifier you can build to “evidence”. How would this classifier predict the class of an incoming email that contains a \$ character?

3. Another friend of you suggest looking into the feature of whether the email's length is longer than a fixed threshold (e.g., 500 bytes). You obtain the following results (this feature denoted as L (\bar{L})).

Assignment Project Exam Help

Class	L	\bar{L}
-------	-----	-----------

<https://eduassistpro.github.io>

incoming email that contains a \$ character and is shorter than the threshold?

Add WeChat edu_assist_pr

Assignment Project Exam Help

100

<https://eduassistpro.github.io>

2. In order to build a (naïve) bayes classifier, we need to calculate (and store) the likelihood of the feature for each class.

Add WeChat $\frac{P(\$ | \text{SPAM})}{P(\$ | \text{NOSPAM})}$ edu_assist_pro

Solution to Q3 II

To classify the new object, we calculate the posterior probability for both classes as:

$$\begin{aligned} P(\text{SPAM} | X) &= \frac{1}{P(X)} \cdot P(X | \text{SPAM}) \cdot P(\text{SPAM}) \\ &= \frac{1}{P(X)} \cdot P(\$ | \text{SPAM}) \cdot P(\text{SPAM}) \end{aligned}$$

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](#)

So the prediction will be SPAM.

3. The likelihood of the new feature for each class is:

$P(L \text{SPAM})$	$\frac{40}{100} = 0.40$
$P(L \text{NOSPAM})$	$\frac{400}{900} = 0.44$

Solution to Q3 III

(Note: we can easily obtain probabilities, e.g.,

$$P(\bar{L} \mid \text{SPAM}) = 1 - P(L \mid \text{SPAM}) = 0.60$$

To classify the new object, we calculate the posterior probability for both classes as:

$$P(\text{SPAM} \mid X) = \frac{P(X \mid \text{SPAM}) \cdot P(\text{SPAM})}{1}$$

<https://eduassistpro.github.io>

$$= \frac{1}{P(X)} \cdot 0.60 \cdot 0.91 \cdot \frac{1}{1}$$

Add WeChat [edu_assist_pro](https://eduassistpro.github.io)

$$= \frac{1}{P(X)} \cdot P(\$, \bar{L} \mid \text{NOSPAM}) \cdot P(\text{NOSPAM})$$

$$= \frac{1}{P(X)} \cdot P(\$ \mid \text{NOSPAM}) \cdot P(\bar{L} \mid \text{NOSPAM}) \cdot P(\text{NOSPAM})$$

$$= \frac{1}{P(X)} \cdot 0.56 \cdot 0.07 \cdot 0.90 = \frac{1}{P(X)} \cdot 0.035$$

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Based on the data in the following table,

1. estimate a Bernoulli Naive Bayes classifier (using the add-one smoothing)
2. apply the classifier to the test document
3. estimate a multinomial Naive Bayes classifier (using the add-one smoothing)
4. ap

You do
test do

	docID	words in document	ina?
training set	1	Taipei Taiwan	
	2	Macao Taiwan Shanghai	
	3	Japan Sapporo	No
	4	Sapporo Osaka Taiwan	No
test set	5	Taiwan Taiwan Taiwan Sapporo Bangkok	?

Solution to Q3 I

We use the following abbreviations to denote the words, i.e., TP = Taipei, TW = Taiwan, MC = Macao, SH = Shanghai, JP = Japan, SP = Sapporo, OS = Osaka. The size of the vocabulary is 7.

1. (Bernoulli NB) We take each word in the vocabulary as a feature/attribute, and hence can obtain the following “rational” training s

	TP	TW	MC	SH	JP	SP	OS	class
2	0	1	1	1	0	0	0	Y
3	0	0	0	0	1	1	0	N
4	0	1	0	0	0	1	1	N

The testing document is (ignoring the unknown tok

docID	TP	TW	MC	SH	JP	SP	OS	class
5	0	1	0	0	0	1	0	?

Solution to Q3 II

By looking at the test data, we calculate the *necessary* probabilities for the 'Y' class as (note that there are 2 possible values for each variable)

Assignment Project Exam Help

$$P(TP = 0 | Y) = \frac{1 + 1}{2 + 2}$$

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

$$\frac{P(SH = 0 | Y)}{P(JP = 0 | Y)} = \frac{1}{1}$$

$$P(SP = 1 | Y) = \frac{0 + 1}{2 + 2}$$

$$P(OS = 0 | Y) = \frac{2 + 1}{2 + 2}$$

We calculate the *necessary* probabilities for the 'N' class as

$$P(N) = \frac{2}{2+1}$$
$$P(TP = 0|N) = \frac{2+1}{2+2}$$

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](#)

$$P(SP = 1|N) = \frac{2+1}{2+2}$$

$$P(OS = 0|N) = \frac{1+1}{2+2}$$

Solution to Q3 V

Finally,

$$P(N|X) \propto P(N) \cdot P(TP = 0|N) \cdot P(TW = 1|N) \cdot P(MC = 0|N) \cdot P(SH = 0|N) \\ \cdot P(JP = 0|N) \cdot P(SP = 1|N) \cdot P(OS = 0|N) \\ = \frac{1}{2} \cdot \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{2} \cdot \frac{1}{3} \cdot \frac{1}{2} = \frac{1}{81} \cdot 0.020$$

2. (<https://eduassistpro.github.io>)

Doc	class						
TP TW MC TW SH	Y						
JP SP TP OS FW	N						

The testing document is (ignoring the out-of-vocabulary words):

Doc	class				
TW TW TW SP	?				

By looking at the test data, we calculate the *necessary* probabilities for the 'Y' class as (note that there are 7 possible values for the variable w_i)

Assignment Project Exam Help

$$P(Y) = \frac{2}{2+1}$$

<https://eduassistpro.github.io>

Finally,

Add WeChat edu_assist_pro

$$\begin{aligned} P(Y|X) &\propto P(Y) \cdot P(w_i = T) \\ &\cdot P(w_i = TW|Y) \cdot P(w_i = SP|Y) \\ &= \frac{1}{2} \frac{1}{4} \frac{1}{4} \frac{1}{4} \frac{1}{12} = \frac{1}{1536} \approx 0.000651 \end{aligned}$$

We calculate the *necessary* probabilities for the 'Y' class as

$$P(N) = \frac{2}{4}$$
$$P(w_i = T|N) = \frac{1+1}{5+7}$$

<https://eduassistpro.github.io>

Finally,

$$P(N|X_5) \propto P(N) \cdot P(w_1 = T) \cdot P(w_2 = T) \cdot P(w_3 = T) \cdot P(w_4 = T) \cdot P(w_5 = T)$$
$$= \frac{1}{2} \frac{1}{6} \frac{1}{6} \frac{1}{6} \frac{1}{4} = \frac{1}{1728} \approx 0.000579$$

Therefore, doc 5 should belong to the 'Yes' class.

Assignment Project Exam Help

Consider a binary classification problem.

1. Fir

2

a

<https://eduassistpro.github.io>

to the positive class?

2. We then identify a feature x , and rearrange t based on their x values. The result is shown in the t

Add WeChat edu_assist_pr

x	y	count
1	-	6
1	+	2
2	-	5
2	+	2
3	-	7

Assignment Project Exam Help

<https://eduassistpro.github.io>

Table: Training Data

Add WeChat edu_assist_pro

For each of the group of training examples with the same x value, compute its probability p_i and $\text{logit}(p) := \log \frac{p}{1-p}$.

- What is your estimate of the probability that a novel test instance belongs to the positive class if its x value is 1?
- We can run a linear regression on the (x, logit) pairs from each group. Will this be the same as what Logistic Regression does?

1. $\Pr(+)=\frac{25}{47}$.

2. See table below.

x	$\text{cnt}(y=0)$	$\text{cnt}(y=1)$	p	$\text{logit}(p)$
		2		
		2		
		6		
		7		
5	1	8		

3. $\Pr(+|x=1)=\frac{3}{8}$

4. Not the same. The main reason is that Logistic regression the likelihood of the data, and this is in generally different from minimizing the SSE as in Linear Regression.

Assignment Project Exam Help

Consider two-dimensional vectors $\mathbf{A} = \begin{pmatrix} 1 \\ 3 \end{pmatrix}$ and $\mathbf{B} = \begin{pmatrix} -1 \\ 0 \end{pmatrix}$ and $\mathbf{C} = \mathbf{A} + \mathbf{B}$.



<https://eduassistpro.github.io> hat if
me as

- ▶ Can you construct a matrix \mathbf{M} such that its i
in polar coordinates exhibit "linearity"? i.e.,

Add WeChat edu_assist_pr

Assignment Project Exam Help

$$\mathbf{C} = \begin{pmatrix} 1 \\ 3 \end{pmatrix}$$

$$\begin{pmatrix} 2 \\ 3 \end{pmatrix} = \mathbf{A}' \quad \mathbf{A}' = \begin{pmatrix} 5 \\ 1.5 \end{pmatrix}$$

<https://eduassistpro.github.io>

Obviously, we still have $\mathbf{C} = \mathbf{A} + \mathbf{B}$.

► (Obviously) No.

► One possible $\mathbf{M} = \begin{pmatrix} 5 & 2 \\ 0 & -\frac{1}{2} \end{pmatrix}$. Then $\mathbf{M} \begin{pmatrix} 1 \\ \rho \end{pmatrix} =$
"linearity" holds.

Add WeChat [edu_assist_pro](https://eduassistpro.github.io)

Assignment Project Exam Help

Consider a set of d -dimensional points arranged in a *data matrix*

\mathbf{o}_1
 \mathbf{o}_2

$\mathbf{X}_{n \times d}$

s to

a m -d

$\pi(\mathbf{o}_i)$

<https://eduassistpro.github.io>

- Computer $r := \frac{\|\pi(\mathbf{o}_i)\|^2}{\|\mathbf{o}_i\|^2}$. Can you guess what will be the minimum values of r ?

Add WeChat edu_assist_pro

► Since

$$\begin{aligned}\| \pi(\mathbf{o}) \|^2 &= \pi(\mathbf{o})^\top \pi(\mathbf{o}) = \mathbf{A}^\top(\mathbf{o})^\top \mathbf{A}(\mathbf{o}) \\ &= \mathbf{o}^\top \mathbf{A} \mathbf{A}^\top \mathbf{o} = \mathbf{o}^\top \mathbf{A} \mathbf{A}^\top \mathbf{o}\end{aligned}$$

<https://eduassistpro.github.io>

Comment: The above is the Rayleigh Quotient (c.f.

where $\mathbf{M} = \mathbf{A} \mathbf{A}^\top$). The maximum and minimum are attained by the maximum and minimum eigenvalues of \mathbf{M} . This property is also used in the technical proof of the spectral clustering too (not required).