

School of Computing and Information Systems
The University of Melbourne
COMP30027 MACHINE LEARNING (Semester 1, 2019)

Tutorial exercises: Week 4

1. Consider the following 10 instances, given so-called “gold standard” labels (assuming a 3-class problem), and the output of four supervised machine learning models:

Instance	Gold	①	②	③	④
1	A	A	A or B	A	A
2	B	A	B or C	A	?
3	A	A	A	A	A
4	C	C	B or C	A	?
5	B	B	A or B or C	A	?
6	C	A	A or C	A	?
7	C	A	A or B or C	A	?
8	A	C	A or B	A	A
9	A	A	A	A	?
10	A	A	A or C	A	A

- (a) Where possible, calculate the **accuracy** and **error rate** of the four models.
- (b) Where possible, calculate the **precision** and **recall**, treating class A as the “positive” class. Do the same for the B and C classes, in turn, and then calculate the **macro-averaged precision and recall**.
2. What is the difference between evaluating using a **holdout** strategy and evaluating using a **cross-validation strategy**?
- (a) What are some reasons we would prefer one strategy over the other?
3. For the following dataset:

ID	Outl	Temp	Humi	Wind	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.
- (b) Classify the test instances using the method of 1-R.
- (c) Classify the test instances using the ID3 **Decision Tree** method:
- Using the **Information Gain** as a splitting criterion
 - Using the **Gain Ratio** as a splitting criterion

1. Consider the following 10 instances, given so-called "gold standard" labels (assuming a ~~3-class~~ problem), and the output of four supervised machine learning models:

Instance	Gold	①	②	③	④
1	A	A	A or B	A	A
2	B	A	B or C	A	?
3	A	A	A	A	A
4	C	C	B or C	A	?
5	B	B	A or B or C	A	?
6	C	A	A or C	A	?
7	C	A	A or B or C	A	?
8	A	C	A or B	A	A
9	A	A	A	A	?
10	A	A	A or C	A	A

- (a) Where possible, calculate the **accuracy** and **error rate** of the four models.
 (b) Where possible, calculate the **precision** and **recall**, treating class A as the "positive" class. Do the same for the B and C classes, in turn, and then calculate the **macro-averaged precision** and **recall**.

For binary class.

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

$$1 - \text{Accuracy} = \text{error rate}$$

Confusion matrix.

		actual	
		I	U
predict	I	TP	FN
	U	FP	TN

(a) System 1, 3 normal :

$$\text{Acc (Sys 1)} = 60\%$$

$$\text{Acc (Sys 3)} = 50\% \quad \text{Acc (Sys 4)} = 40\%$$

For multi-class: use macro-averaging

System 2 : $TP = 10$, $FP = 10$

$$TN = 0$$

$$FN = 0$$

$$\text{Acc (Sys 2)} = \frac{10}{10+10} = 50\%$$

Instance	Gold	①	②	③	④
1	A	A	A or B	A	A
2	B	A	B or C	A	?
3	A	A	A	A	A
4	C	C	B or C	A	?
5	B	B	A or B or C	A	?
6	C	A	A or C	A	?
7	C	A	A or B or C	A	?
8	A	C	A or B	A	A
9	A	A	A	A	?
10	A	A	A or C	A	A

(b) $P = \frac{TP}{TP+FP} = \frac{\text{选对 A}}{\text{选所有 A}} - \text{Precision (查准)}$

$R = \frac{TP}{TP+FN} = \frac{\text{选对 A}}{\text{实际 A}} - \text{Recall (查全)}$

Precision:

$P(\text{sys 1}; A) = \frac{4}{7} \quad P(\text{sys 3}; A) = \frac{5}{10}$

$P(\text{sys 2}; A) = \frac{5}{8} \quad P(\text{sys 4}; A) = \frac{4}{4}$

$P(\text{sys 1}; B) = 1 \quad P(\text{sys 3}; B) = \frac{0}{6}$

$P(\text{sys 2}; B) = \frac{3}{6} \quad P(\text{sys 4}; B) = 0$

Recall:

$R(\text{sys 1}; A) = \frac{4}{5} \quad R(\text{sys 3}; A) = 1$

$R(\text{sys 2}; A) = \frac{5}{8} \quad R(\text{sys 4}; A) = \frac{4}{5}$

Macro-Average:

$\text{macro-P} = \frac{1}{N} \sum_{k \in K} P_k$

$\text{macro-R} = \frac{1}{N} \sum_{k \in K} R_k$

Macro-P of system 3, 4 cannot calculated.

2. What is the difference between evaluating using a **holdout** strategy and evaluating using a **cross-validation** strategy?

(a) What are some reasons we would prefer one strategy over the other?

In "hold-out" eval:

We partition the data into training set and testing set

build the model on the former

then evaluate on the later.

In "cross-validation" eval:

do the same, but a number of times, where each iteration uses one partition of the data as the test data and rest as training set.


(a) hold-hold is subject to some random variation. dependending on which instances are assigned to training set, which are testing set. This could mean our estimate result way-off.

Cross-validation solve this problem, averaging over a bunch of values -
So one weird-partition won't affect result. But takes a long time.
Since we need train a model for each partition.

3. For the following dataset:

ID	Outl	Temp	Humi	Wind	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	?

- Classify the test instances using the method of 0-R.
- Classify the test instances using the method of 1-R.
- Classify the test instances using the ID3 **Decision Tree** method:
 - Using the **Information Gain** as a splitting criterion
 - Using the **Gain Ratio** as a splitting criterion

(a) 

(b) How to choose a attribute?

Counting the errors made on the training instances

Let's choose "Outl" first.

Outl = s : ID A, B, Both label as N, no error

Outl = o : ID C. label as Y, no error

Outl = r : ID D, E, F. label Y/N. make one error (2-1)

total: 1 error for outl

↑
"0-R"

Let's choose "Temp" then.

Temp = h : One error

Temp = m : one error

Temp = c : one error

total : 2 errors for Temp

Assume "Outl" best. ID - G \rightarrow Y . ID - H \rightarrow N

What is 0-R ?

- Baseline classifier

- throw all attributes,

- predict each instance according to which label is most common in training set ("majority class")

What is 1-R ?

- Better Baseline classifier

- choose a single attribute which preferred.

- we will predict according to that attribute which label is most common in training set

(c) Classify the test instances using the ID3 Decision Tree method:

- Using the **Information Gain** as a splitting criterion
- Using the **Gain Ratio** as a splitting criterion

For Information Gain (IG), at each level of decision tree.

we choose the attribute that has the largest IG difference between the HCR entropy of the class distribution at the parent node and the average entropy across the child nodes, which is MI (Mean Information).

$$IG(\underset{\text{attribute}}{A} | \underset{\text{class}}{R}) = \underset{\text{MI}}{H(R)} - \sum_{i \in A} P(A=i) H(A=i)$$

i at top level (root) of the tree.

$$H(R) = - \sum_{k \in R} P(R_k) \log(R_k) = - \left[\frac{3}{6} \log \frac{3}{6} + \frac{3}{6} \log \frac{3}{6} \right] = 1$$

unpredictable.
even distribution

we want daughters to have an uneven distribution.

- means we can select a class with more confidence.
- which means entropy will go down.

For example choose "Out"

$$H(A=s) = - [0 \log 0 + 1 \log 1] = 0$$

$$H(A=o) = - [1 \log 1 + 0 \log 0] = 0$$

$$H(A=r) = - \left[\frac{3}{6} \log \frac{3}{6} + \frac{3}{6} \log \frac{3}{6} \right] = 0.9183$$

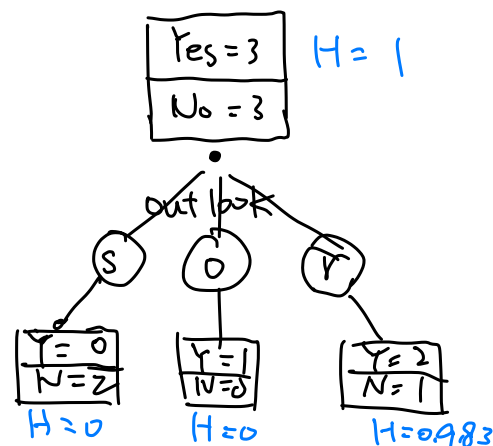
$$MI(o) = \frac{3}{6} \cdot 0 + \frac{1}{6} \cdot 0 + \frac{3}{6} \cdot 0.9183 = 0.4592$$

$$IG(o|R) = 1 - 0.4592 = 0.5408$$

$$MI(X_1, X_2, \dots, X_n) = \sum_{i=1}^n P(x_i) H(x_i)$$

ID	Outl	Temp	Humi	Wind	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	?

this means at this branch we will choose N with confidence



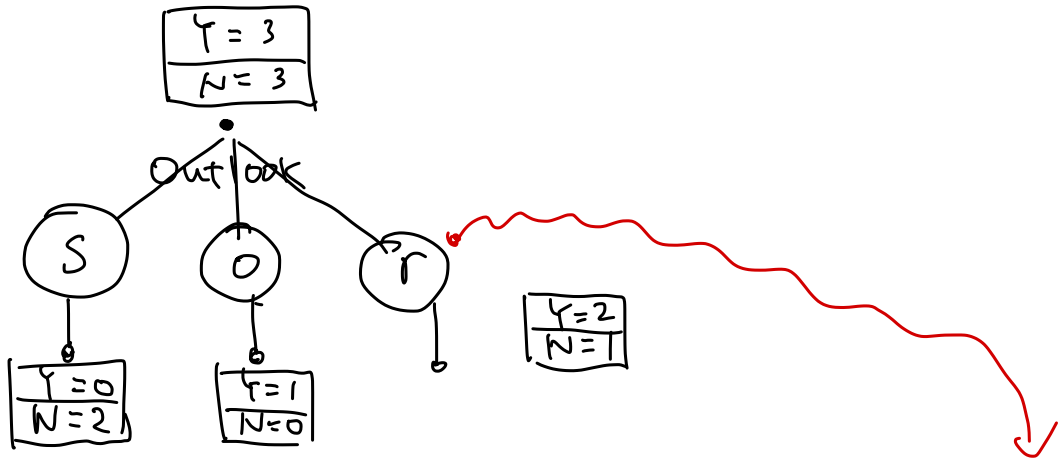
$$MI = 0.4592$$

注: ID 的 IG 最高, 按理说选做 root, 因为 each daughter is purely of a single class, [但我(i)不用]

把所有的IG的算出来, 找出Max IG, 作为 root

	R	Outl			Temp			H		Wind		ID					
		s	o	r	h	m	c	h	n	T	F	A	B	C	D	E	F
Y	3	0	1	2	1	1	1	2	1	0	3	0	0	1	1	1	0
N	3	2	0	1	2	0	1	2	1	2	1	1	1	0	0	0	1
Total	6	2	1	3	3	1	2	4	2	2	4	1	1	1	1	1	1
$P(Y)$	$\frac{1}{2}$	0	1	$\frac{2}{3}$	$\frac{1}{3}$	1	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	0	$\frac{3}{4}$	0	0	1	1	1	0
$P(N)$	$\frac{1}{2}$	1	0	$\frac{1}{3}$	$\frac{2}{3}$	0	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	1	$\frac{1}{4}$	1	1	0	0	0	1
H	1	0	0	0.9183	0.9183	0	1	1	1	0	0.8112	0	0	0	0	0	0
MI				0.4592				1			0.5408						
IG				0.5408				0			0.4592						
SI				1.459				0.9183			0.9183					2.585	
GR				0.3707				0			0.5001					0.3868	

Base on above. choose Outl as root



Now Root is Outlook. $H(R) \neq 1 \Rightarrow H(R) = 0.9183$

If now split by Temp: $MI = \frac{1}{3} \times 0 + \frac{2}{3} \times 1 = 0.6667$

$$IG = H(R) - MI = 0.9183 - 0.6667 = 0.2516$$

... calculate each.

Final all daughters of r are pure.

\Rightarrow
 final tree
Outl = O \cup (Outl = r \wedge wind = F) \Rightarrow Y
Outl = S \cup (Outl = r \wedge wind = T) \Rightarrow N.

第一层

其余层

ii Gain Ratio.

Gain Ratio is similar!

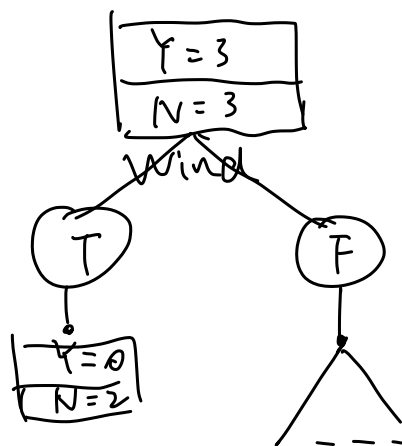
$$\frac{\# \text{ of } i}{\# \text{ of instance}}$$

$$SI(A) = - \sum_{i \in A} P(A=i) \log P(A=i)$$

$$GR(A) = \frac{SI(Y)}{SI(A)}$$

	R	Outl			Temp			H		Wind		ID					
		s	o	r	h	m	c	h	n	T	F	A	B	C	D	E	F
Y	3	0	1	2	1	1	1	2	1	0	3	0	0	1	1	1	0
N	3	2	0	1	2	0	1	2	1	2	1	1	1	0	0	0	1
Total	6	2	1	3	3	1	2	4	2	2	4	1	1	1	1	1	1
P(Y)	$\frac{1}{2}$	0	1	$\frac{2}{3}$	$\frac{1}{3}$	1	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	0	$\frac{3}{4}$	0	0	1	1	1	0
P(N)	$\frac{1}{2}$	1	0	$\frac{1}{3}$	$\frac{2}{3}$	0	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	1	$\frac{1}{4}$	1	1	0	0	0	1
H	1	0	0	0.9183	0.9183	0	1	1	1	0	0.8112	0	0	0	0	0	0
MI				0.4592				1			0.5408						
IG				0.5408				0			0.4592						
SI				1.459				0.9183			0.9183						
GR				0.3707				0			0.5001						

GR of Wind is largest \rightarrow Wind as Root



$$\frac{Y=3}{N=2}$$

$$H(\text{Wind} = F) = 0.8112$$

— — — Calculated

tree final : $\text{Wind} = F \cap (\text{Outl} = u \cup \text{Outl} = r) \rightarrow Y$

$\text{Wind} = T \cup (\text{Wind} = F \cap \text{Outl} = s) \rightarrow N$