

Lecture 13:  
Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

## Lecture 13: Evaluation

# Assignment Project Exam Help

<https://eduassistpro.github.io>

Sarah Erfani and Karin Verspoor, CIS

Add WeChat Search edu\_assist\_pr



THE UNIVERSITY OF  
MELBOURNE

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

### Evaluation

Measures

Model Validation

Results

comparison

Random Baseline

Zero-R

One-R

# Assignment Project Exam Help

■ <https://eduassistpro.github.io/>

- **Consistency:** is the classifier able to  
classify all training instances?

# Add WeChat edu\_assist\_pro

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

- **Under-fitting:** model not expressive enough to capture patterns in the data.
- **Over-fitting:** model too complicated; capture noise in the data.
- **Appropriate-fitting** model captures essential patterns in the data.

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

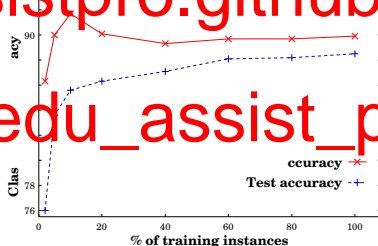
- Usually, the given data set is partitioned into two disjoint sets. The *training set* is used to build the model; the *test set* is used to validate it.

## Inductive Learning Hypothesis:

Any hypothesis found to approximate the target function well over a sufficiently large training data set will also approximate the Id-out

■

- Learning curves represent the performance of a fixed learning strategy over different sizes of training data, relative to a fixed evaluation metric.



## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

For a two class problem:

There are Positive and Negative cases

A classifier may classify

- a Negative instance as Positive (False Positive, FP)
- a Negative instance as Negative (True Negative, TN)

		Y	N
Actual	Y	true positive (TP)	false negative (FN)
	N	false positive (FP)	true negative (TN)

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results

comparison

Random Baseline

Zero-R

One-R

Outlook	Temperature	Humidity	Windy	Cluster	Play
sunny	hot	high	FALSE	0	no
sunny	hot	high	TRUE	0	no
overcast	hot	high	FALSE	0	yes
		high	FAL	1	yes
		n			yes
		n			no
		hi	FA	0	no
sunny	cool	normal	FALSE	1	
rainy	mild	normal	FALSE	1	
sunny	mild	normal	TRUE	1	
overcast	mild	high	TRUE	1	
overcast	hot	normal	FALSE	0	yes
rainy	mild	high	TRUE	0	no

Cluster 0 = "no", Cluster 1 = "yes"

# Clustering accuracy

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

Outlook	Temperature	Humidity	Windy	Cluster	Play
sunny	hot	high	FALSE	0	no
sunny	hot	high	TRUE	0	no
overcast	hot	high	FALSE	0	yes
rain	mild	high	FALSE	1	yes
rainy	cool	normal	FALSE	1	yes
rainy	cool	normal	TRUE	1	no
		hi	FA	0	
overcast	mild	high	TRUE	1	
overcast	hot	normal	FALSE	0	
rainy	mild	high	TRUE	0	

Cluster 0 = "no", Cluster 1 = "yes"

		<i>Predicted</i>	
		<i>Y</i>	<i>N</i>
<i>Actual</i>	<i>Y</i>	TP (7)	FN (2)
	<i>N</i>	FP (1)	TN (4)

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

### Evaluation

#### Measures

Model Validation

#### Results comparison

Random Baseline

Zero-R

One-R

# Assignment Project Exam Help

<https://eduassistpro.github.io>

$$ACC = \frac{TP + we}{TP + FN + FP + TN}$$

# Add WeChat edu\_assist\_pr



## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation  
Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

Assignment Project Exam Help

- Alternatively, we sometimes talk about the *error rate*.

<https://eduassistpro.github.io>

- comparing the

error rate ER for a given method with  
method  $ER_0$

Add WeChat edu\_assist\_pro

$$ERR = \frac{ER - ER_0}{ER_0}$$

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation  
Measures  
Model Validation

Results  
comparison  
Random Baseline  
Zero-R  
One-R

- If we wish to know what we have positively identified **not** what we have correctly ignored (or equivalently performance relative to a single class of interest), we use *precision* and *recall*

$$\frac{TP}{FP}$$

$$\frac{P}{FN}$$

<https://eduassistpro.github.io>

- Precision: Proportion of positive predictions that are correct
- Recall: Accuracy with respect to positive c  
also called true positive rate

Add WeChat edu\_assist\_pr

- *Specificity* is the accuracy with resp

$$\text{Specificity} = \frac{TN}{TN + FP}$$

(sensitivity/specificity is often used in scientific applications)

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

- To compute an overall P/R value over multiple categories:

■ *micro-averaging*

$$\text{Precision}_\mu = \frac{\sum_{i=1}^c TP_i}{\sum_{i=1}^c (TP_i + FP_i)}$$

■ *macro-averaging*

$\text{Precision}_M$

$\text{Recall}_M$

$c$

- In what situations are they the same/different?

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results

comparison

Random Baseline

Zero-R

One-R

# Assignment Project Exam Help

<https://eduassistpro.github.io>

# Add WeChat edu\_assist\_pr

[https://en.wikipedia.org/wiki/Sensitivity\\_and\\_specificity](https://en.wikipedia.org/wiki/Sensitivity_and_specificity)

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation  
Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

- In applications where we make individual decisions for each data point rather than generating a monolithic ranking, *F-score* gives us an overall picture of system performance:

<https://eduassistpro.github.io/>  
c mean]

- Set  $\beta$  depending on how much we care about false positives

- Conventionally  $\beta = 1$ , called the

$$\text{F1-score} = 2 \frac{PR}{P + R}$$

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation  
Measures  
Model Validation

Results  
comparison  
Random Baseline  
Zero-R  
One-R

You may see people refer to AUC and ROC.

The **ROC** = Receiver Operating  
Characteristic

- A plot illustrating the performance of a classifier as

**AUC** = Area Under the Curve

- sometimes called **AUROC**
- equal to the probability that a classifier will rank a randomly  
tance higher  
sen

<https://eduassistpro.github.io>

(Recall/Sensitivity) vs. False  
Positive Rate ( $1 - \text{Specificity}$ )

- The best possible prediction method would yield a point in the upper left corner or coordinate (0,1) of the ROC space, representing 100% sensitivity (no false negatives) and 100% specificity (no false positives).

Add WeChat edu\_assist\_pro

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

# Assignment Project Exam Help

<https://eduassistpro.github.io>

- If we use all of our data to train a model, how haven't *overfit* our model to our data

# Add WeChat edu\_assist\_pr

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

- The (training) *bias* of a classifier is the average distance between the expected value and the estimated value
  - Bias is large if the learning method produces classifiers that are consistently wrong.
  - Bias is small if (i) the classifiers are consistently right or (ii) different

<https://eduassistpro.github.io>

Add WeChat: edu\_assist\_pro

- Variance is large if different training sets classifiers.
  - It is small if the training set has a minor effect on the decisions made, be they correct or incorrect.
  - Variance measures how inconsistent they are correct or incorrect.
- The *noise* in a dataset is the inherent variability of the training data
- In evaluation, we aim to minimise classifier bias and variance (but there's not a lot we can do about noise!)



## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation  
Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

# Assignment Project Exam Help

- Train a classifier over a fixed training dataset, and evaluate it over a

<https://eduassistpro.github.io>

- Disadvantages:

- trade off between more training and
- representativeness of the training

Add WeChat edu\_assist\_pro

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

# Assignment Project Exam Help

- Perform holdout over multiple iterations, randomly selecting the

<https://eduassistpro.github.io>

- reduction in variance and bias over “

- Disadvantages:

- reproducibility

# Add WeChat edu\_assist\_pr

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

Let us assume we have  $N$  data points for which we know the labels.

We choose each data point as test case and the rest as training data.

This means we have to train the system  $N$  times and the average performance is computed across the  $N$  predictions.

<https://eduassistpro.github.io>

- will be unique and repeatable.

- The method also generally gives high accuracy as all  $(N - 1)$  points are used in training.  
(It is typically the case that having more training data leads to a more accurate classifier can be built.)

### Bad point:

- It is infeasible if we have large data set and the training is itself very expensive.

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

# Assignment Project Exam Help

Let us assume we have  $N$  data points for which we know the labels.

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

This means we have to train the system  $M$  time  
performance is computed across the  $M$  run

Typical values for  $M$ : 5 or 10 (i.e. 5-fold cross-  
cross-validation)

# Add WeChat edu\_assist\_pro



## Cross Validation: Partitioning

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

## Model Validation

- Split up into  $N$  equal-sized partitions  $A_i$ :

<https://eduassistpro.github.io>

# Add WeChat edu\_assist\_pr

ld WeChat edu\_assist\_p

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

- For each  $i = 1 \dots N$ , take  $P_i$  as the test data and  $\{P_j : j \neq i\}$  as the training data

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

Add WeChat edu\_assist\_pr

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

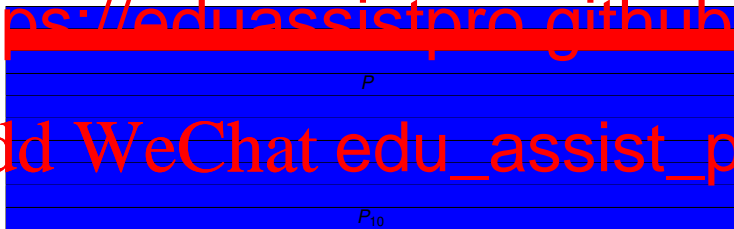
Zero-R

One-R

- For each  $i = 1 \dots N$ , take  $P_i$  as the test data and  $\{P_j : j \neq i\}$  as the training data

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr



## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

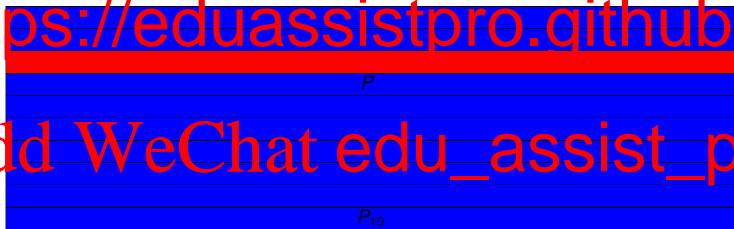
Zero-R

One-R

- For each  $i = 1 \dots N$ , take  $P_i$  as the test data and  $\{P_j : j \neq i\}$  as the training data

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr





## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

# Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pr

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation  
Measures  
Model Validation

Results  
comparison  
Random Baseline  
Zero-R  
One-R

### Good points:

- We need to train the system only  $M$  times unlike Leave-One-Out which requires training  $N$  times.
- We can measure the stability of the system across different

<https://eduassistpro.github.io>

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io)

- There can be a bias in evaluating the system due to sampling, how data is distributed among the  $M$  partitions
- The results will not be unique unless we repeat the process identically. One solution is to repeat the process randomly shuffling the data  $M/2$  times
- The results will give slightly lower accuracy values as only  $\frac{M-1}{M}$  of the data is used for training.
- For small data sets it is not always possible to partition the data properly such that each partition represents the data IID (Independently Distributed).

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

## Results comparison

Random Baseline

Zero-R

One-R

Assignment Project Exam Help

- *Baseline* = naive method which we would expect any reasonably

*walk 42km*  
<https://eduassistpro.github.io>  
ainst

*e.g. for a marathon runner, t  
run/the world record time/*  
Add WeChat edu\_assist\_pr

- “Baseline” often used as umbrella ter

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

## Results comparison

Random Baseline

Zero-R

One-R

- Baselines are important in establishing whether any proposed method is doing better than “dumb and simple”

<https://eduassistpro.github.io>

- In formulating a baseline, we need to be to the importance of positives and neg

Add WeChat [edu\\_assist\\_pr](https://eduassistpro.github.io)  
*limited utility of a baseline on  
classification task aimed  
new diamond mines (as nearly all sites are unsuitable)*

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

**Method 1:** randomly assign a class to each test instance

- Often the only option in unsupervised/semi-supervised contexts

<https://eduassistpro.github.io>

- Assumes we know the prior probability
- Alleviate effects of variance by:
  - running method  $N$  times and calculate
  - OR
  - arriving at a deterministic estimate of the accuracy of random assignment =  $\sum_i P(C_i)^2$

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

# Assignment Project Exam Help

- **Method:** classify all instances according to the most common class

<https://eduassistpro.github.io>

- Inappropriate if the majority class is is  
to identify needles in the haystack

Add WeChat [edu\\_assist\\_pr](https://eduassistpro.github.io)

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

Creates one rule for each attribute in the training data, then selects the rule with the smallest error rate as its one rule

- **Method:** create a “decision stump” for each attribute, with branches

<https://eduassistpro.github.io>

■

For each attribute,

For each value of the attribute, make a rule

1 count how often each class appears

2 find the most frequent class

3 make the rule assign that class to this attribute-value

Calculate the error rate of the rules

Choose the rules with the smallest error rate

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

Outlook	Temperature	Humidity	Windy	Play
sunny	hot	high	FALSE	no
sunny	hot	high	TRUE	no
overcast	hot	high	FALSE	yes
		high	FAL	

sunny	mild	high	FALSE	no
sunny	cool	normal	FALSE	yes
rainy	mild	normal	FALSE	yes
sunny	mild	normal	TRUE	yes
overcast	mild	high	TRUE	yes
overcast	hot	normal	FALSE	yes
rainy	mild	high	TRUE	no



## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

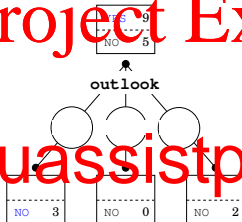
Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R



<https://eduassistpro.github.io>

Add WeChat  edu\_assist\_pro

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

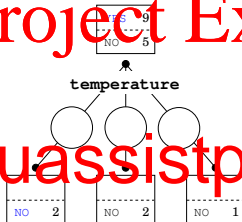
Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat  edu\_assist\_pr

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

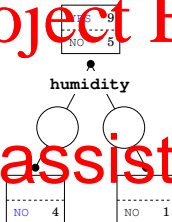
Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat  edu\_assist\_pr

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

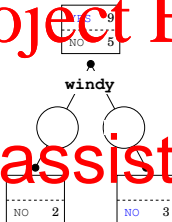
Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R



Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat  edu\_assist\_pr

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

# Assignment Project Exam Help

<https://eduassistpro.github.io>

- unable to capture attribute interaction
- bias towards high arity attributes (

# Add WeChat edu\_assist\_pr

## Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation  
Measures  
Model Validation

Results  
comparison  
Random Baseline  
Zero-R  
One-R

# Assignment Project Exam Help

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

- What is a baseline? What are some examples of reasonable baselines to compare with?

# Add WeChat edu\_assist\_pr

### Lecture 13: Evaluation

COMP90049  
Knowledge  
Technologies

Evaluation

Measures

Model Validation

Results  
comparison

Random Baseline

Zero-R

One-R

Evaluation in IR (unranked retrieval): Manning, Raghavan and Schütze  
Introduction to Information Retrieval, Cambridge University Press. 2008.

**Section 8.** <http://nlp.stanford.edu/IR-book/html/htmledition/>

<https://eduassistpro.github.io/4.6>

[http://nlp.stanford.edu/IR-book/html/htmledition/  
the-bias-variance-tradeoff-1.html](http://nlp.stanford.edu/IR-book/html/htmledition/the-bias-variance-tradeoff-1.html)

ROC: Tom Fawcett, "An introduction to ROC a  
Letters 27 (2006) [https:](https://ccrma.stanford.edu/workshops/mir2009/references/ROCintro.pdf)

[//ccrma.stanford.edu/workshops/mir2009/references/ROCintro.pdf](https://ccrma.stanford.edu/workshops/mir2009/references/ROCintro.pdf)