



Lecture Notes on
Pattern Recognition

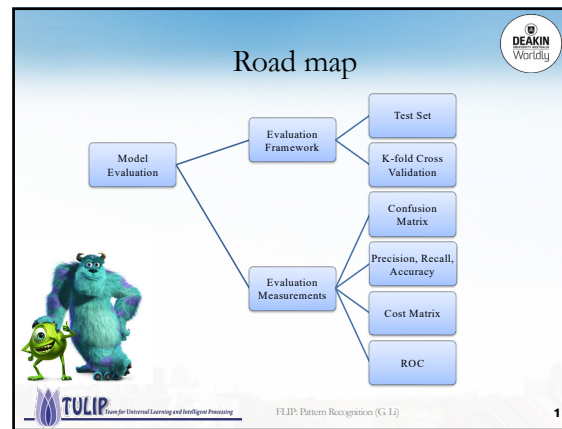
Session 08D: Model Evaluation

Gang Li
School of Information Technology
Deakin University, VIC 3125, Australia

FLIP: Pattern Recognition (G. Li)


0



1

Model Evaluation

- What do we care?
- The Test Set method
- The Leave One Out method
- The K-fold Cross-Validation



TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

2

What do we care?

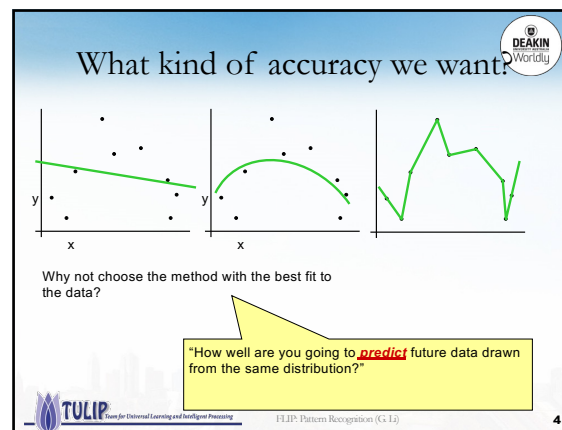
- Regular Performance of Algorithms
 - Accuracy
 - Scope
 - Efficiency
- Which ones we care?

1. Accuracy ? Fitting Accuracy, or
 2. Efficiency ? Prediction Accuracy
 3. Scope

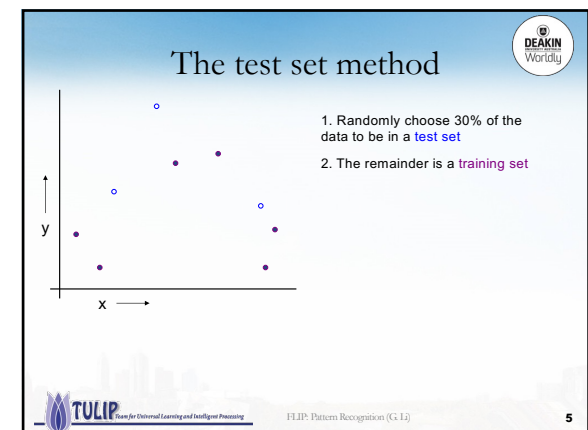
TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

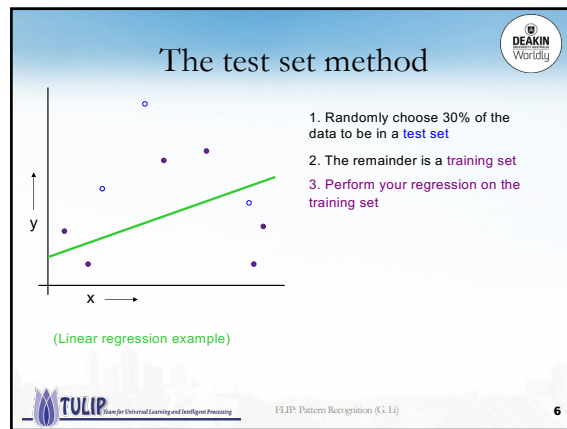
3



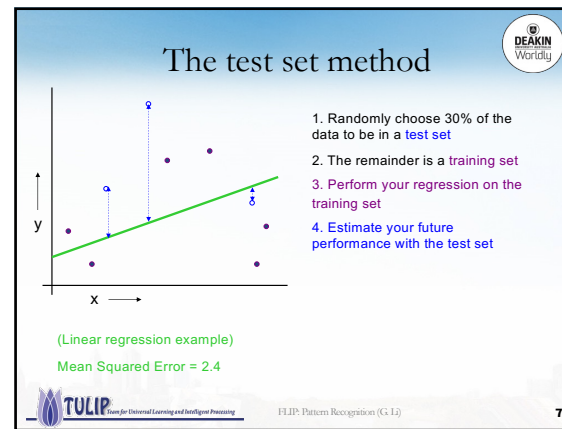
4



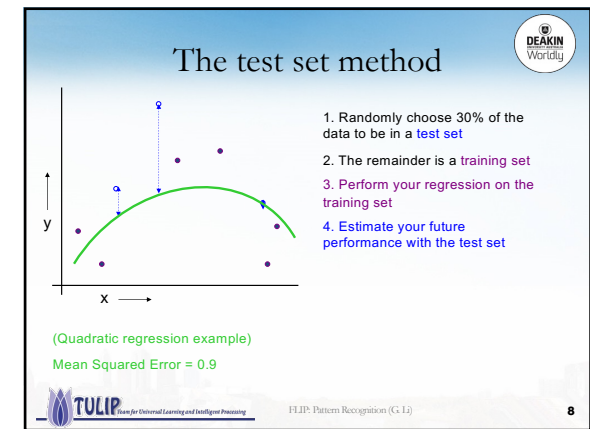
5



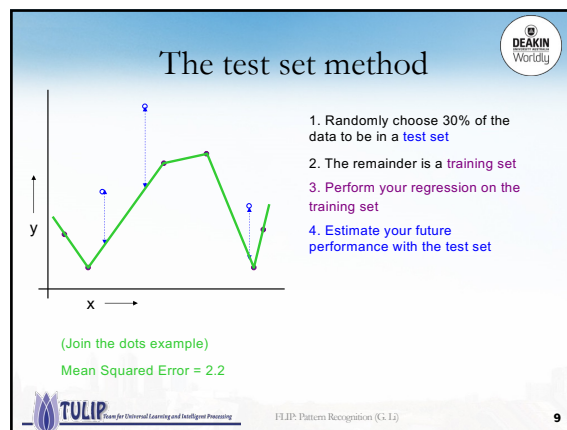
6



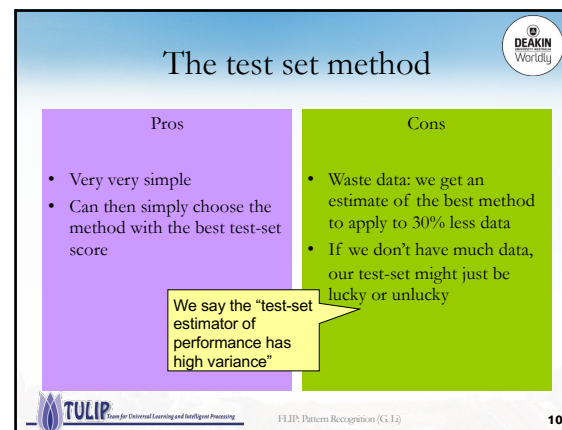
7



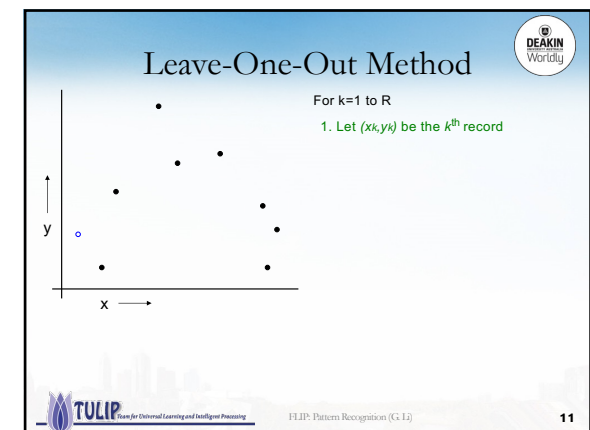
8



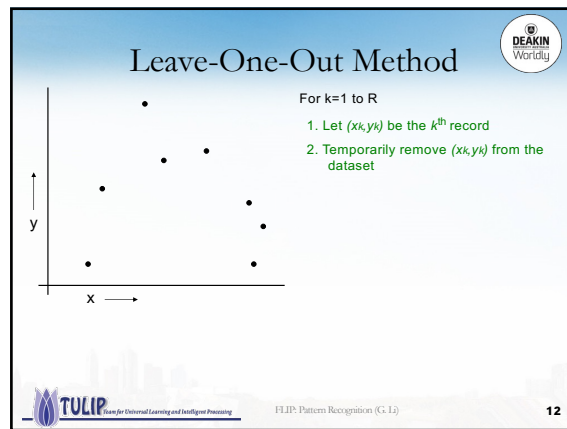
9



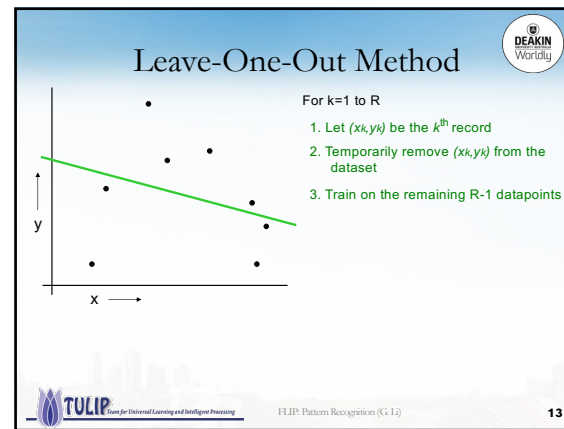
10



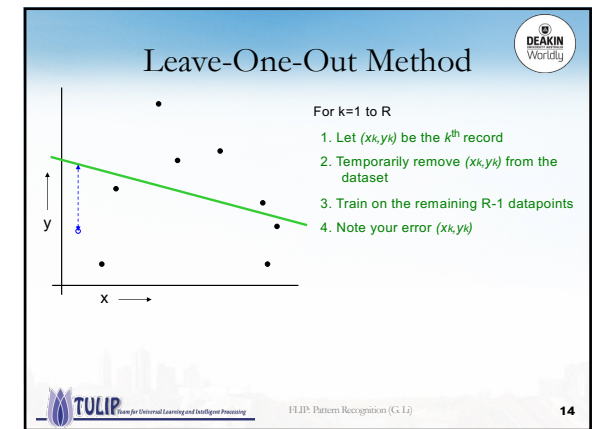
11



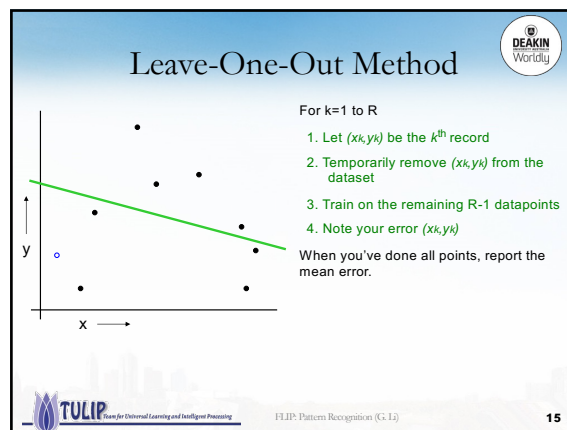
12



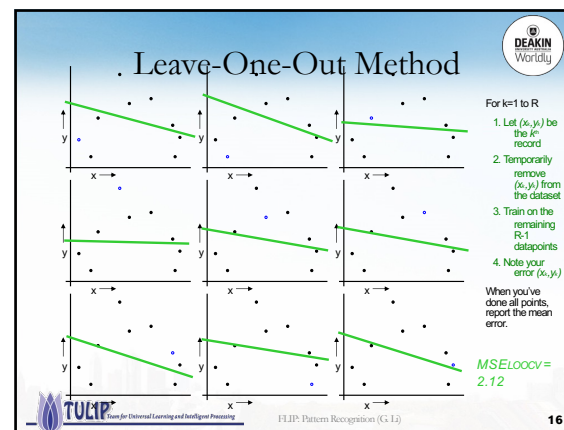
13



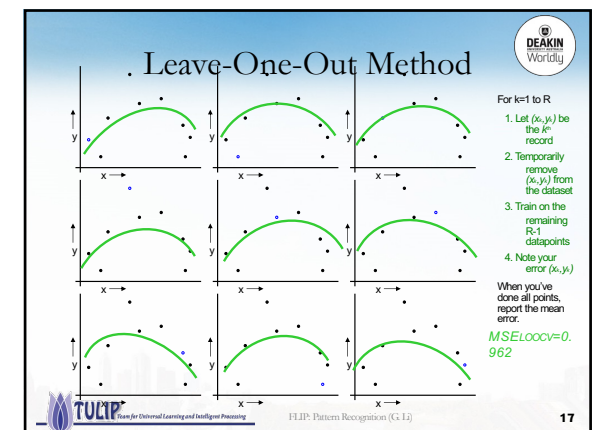
14



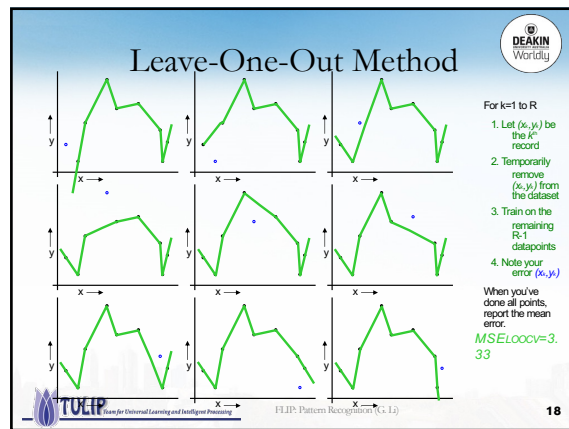
15



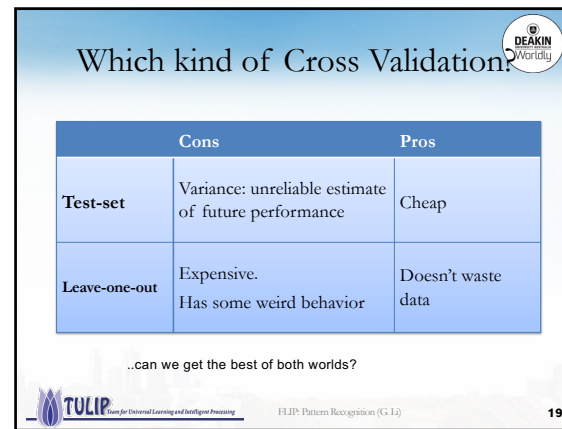
16



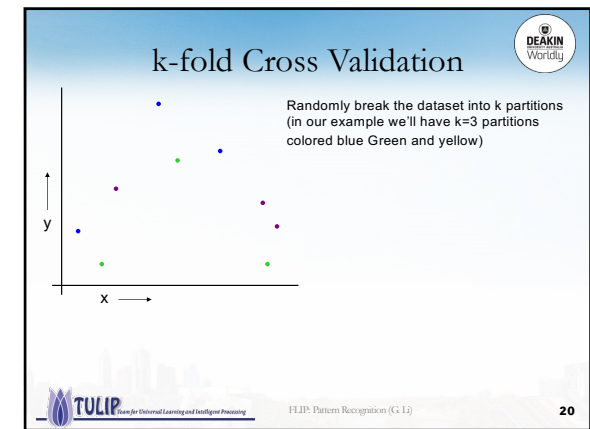
17



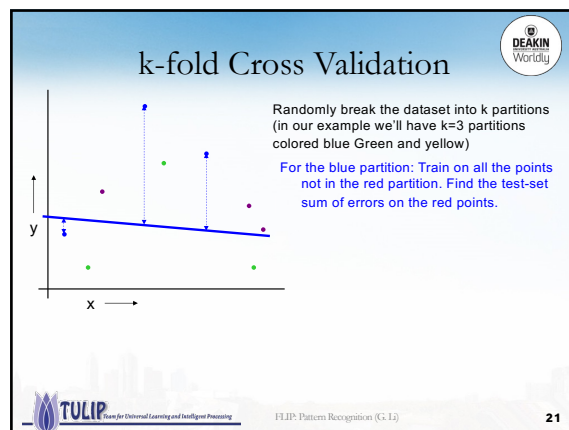
18



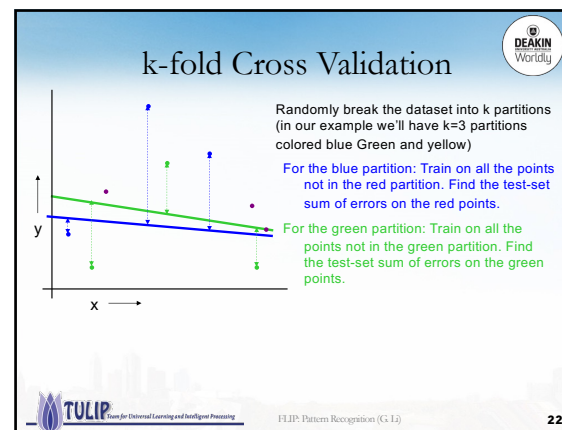
19



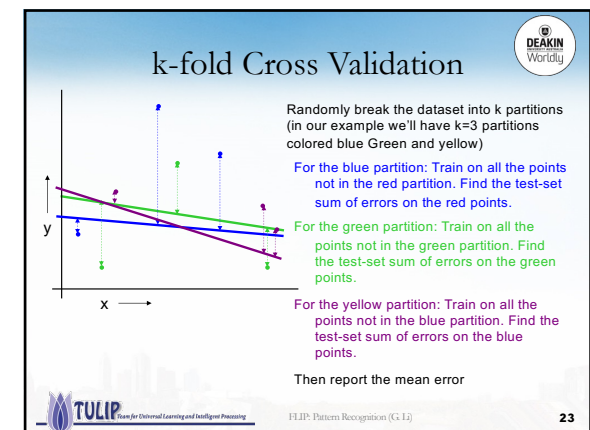
20



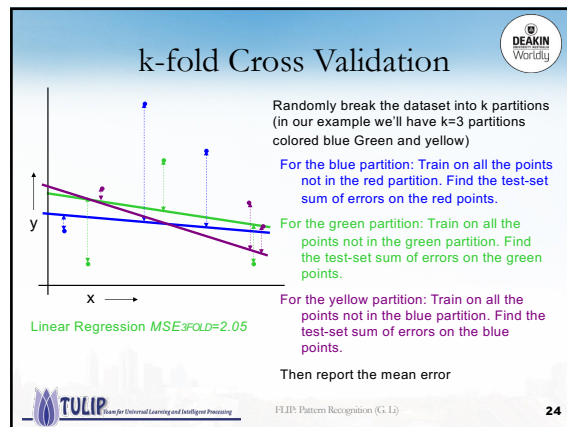
21



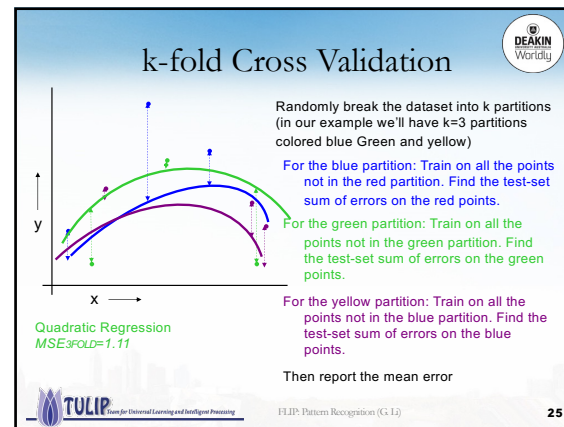
22



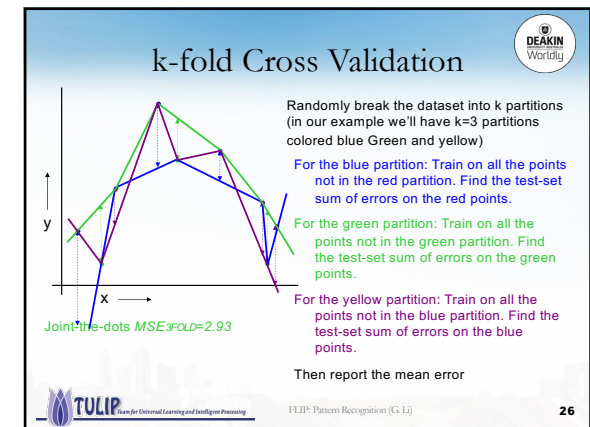
23



24



25



26

Which kind of Cross Validation?

	Downside	Upside
Test-set	Variance: unreliable estimate of future performance	Cheap
Leave-one-out	Expensive. Has some weird behavior	Doesn't waste data
10-fold	Wastes 10% of the data. 10 times more expensive than test set	Only wastes 10%. Only 10 times more expensive instead of R times.
3-fold	Wastier than 10-fold. Expensivier than test set	Slightly better than test-set
R-fold	Identical to Leave-one-out	

FLIP: Pattern Recognition (G. Li) 27

27

CV-based Model Selection

- We're trying to decide which algorithm to use.
- We train each model and make a table...

i	f_i	TRAIN-ERR	10-FOLD-CV-ERR	Choice
1	f_1			
2	f_2			
3	f_3			
4	f_4			
5	f_5			
6	f_6			

FLIP: Pattern Recognition (G. Li) 28

28

Other Model Selection Methods

- Model selection methods:
 - Cross-validation**
 - Occam's Razor Type: The smaller, the better**
 - Minimum Description Length (MDL)
 - Minimum Messaging Length (MML)
 - AIC (Akaike Information Criterion)
 - BIC (Bayesian Information Criterion)
 - VC-dimension (Vapnik-Chervonenkis Dimension)

Only directly applicable to choosing classifiers

FLIP: Pattern Recognition (G. Li) 29

29

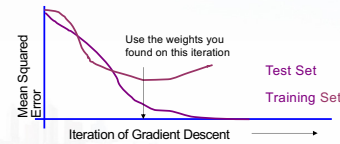
Cross-Validation for regression

- Choosing the number of hidden units in a neural net
- Feature selection
- Choosing a polynomial degree
- Choosing which regressor to use

30

Supervising Gradient Descent

- This is a weird but common use of Test-set validation
 - Suppose you have a neural net with too many hidden units. It will overfit.
 - As gradient descent progresses, maintain a graph of **MSE-testset-error** vs. **Iteration**



31

Cross-validation for classification

- Instead of computing the mean squared errors (MSE) on a test set, you should compute various measurements ...
 - Error rate (or its dual part Accuracy):
 - The total number of misclassifications on a test-set

32

Evaluation and performance



- Confusion Matrix
- Accuracy
- Model Comparison

33

Evaluation and performance

- How well is your classification algorithm?
 - Focus on the predictive capability of a model, rather than how fast it takes to classify or build models, scalability, etc.
- Confusion matrix
 - Detail classification result for each class.
- Accuracy
 - How well we can predict for each class

34

Evaluation

- Take 'Yes' as the positive class (class of interest)

True Class	Algorithm returns
Yes	Yes
Yes	Yes
Yes	No
Yes	Yes
No	No
No	No
No	Yes
No	No
No	Yes
Yes	Yes
Yes	Yes
No	No
No	No
Yes	Yes

Classified As →	A=Yes	B=No
A=Yes	6	1
B=No	2	5

Actual Class	Predicted class		
	+	-	
+	a	b	a: TP (true positive) b: FN (false negative) c: FP (false positive) d: TN (true negative)
-	c	d	

35

Evaluation

- Most widely-used metric:
Accuracy = $(a+d) / (a+b+c+d)$
- Consider a binary classification problem
 - # of Class "Positive" instances = 9990
 - # of Class "Negative" instances = 10
- A naïve model that predicts everything positive with accuracy 99.9%

Actual Class	Predicted class		
	+	-	
+	a	b	a: TP (true positive) b: FN (false negative) c: FP (false positive) d: TN (true negative)
-	c	d	

Evaluation

- Alternative Measures
Precision = $a / (a+c)$
Recall = $a / (a+b)$
F-measure or F1 = $2a / (2a+b+c)$
- Other measures
 - ROC
 - AUC

Actual Class	Predicted class		
	+	-	
+	a	b	a: TP (true positive) b: FN (false negative) c: FP (false positive) d: TN (true negative)
-	c	d	

Cost Matrix

- $C(i, j)$
 - Cost of misclassifying class i instance as class j

Actual Class	Predicted class		
	+	-	
+	a	b	a: TP (true positive) b: FN (false negative) c: FP (false positive) d: TN (true negative)
-	c	d	

Cost Matrix: example

Actual Class	Predicted class		
	+	-	
+	-1	100	
-	1	0	

Actual Class	Predicted class		
	+	-	
+	180	20	
-	60	240	

accuracy = 82%
Cost = 1880

Actual Class	Predicted class		
	+	-	
+	140	60	
-	10	290	

accuracy = 86%
Cost = 5870

Cost-Sensitive Classification

- Traditional Classification**
 - Given a query t , we will classify it as class i if

$$i = \arg \max_j p(j | t)$$
 - For Binary Classification, classify it as "positive" if

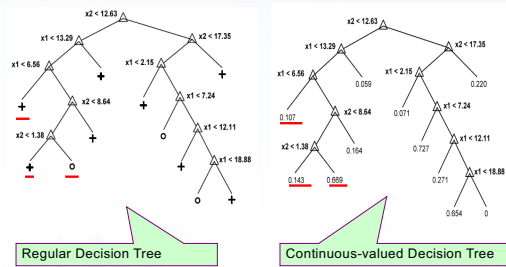
$$p(+|t) > p(-|t)$$
- Cost-Sensitive Classification**
 - Classify t as j

$$i = \arg \min_j \sum_k p(k | t) \times C(k, j)$$

ROC: Receiver Operating Characteristic

- A graphical approach for displaying trade-off between detection rate and false alarm rate
 - Developed in 1950s for signal detection theory to analyze noisy signals
- To draw ROC curve, classifier must produce continuous-valued output
 - Outputs are used to rank test records, from the most likely record to be classified as a positive class to the most likely record to be classified as a negative class
 - Performance of each classifier is represented as a point on the ROC curve

Continuous-valued outputs from a decision tree

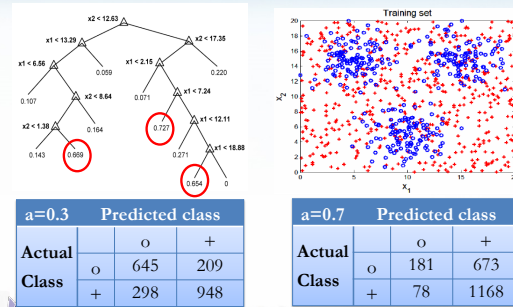


Team for Universal Learning and Intelligent Processing

FLJP: Pattern Recognition (G. Li)

42

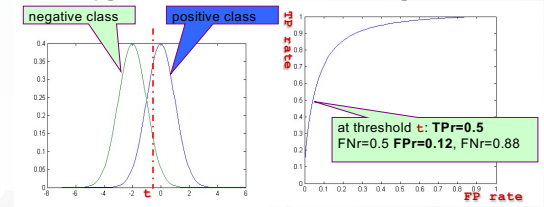
ROC Curve Example



43

ROC Curve Example

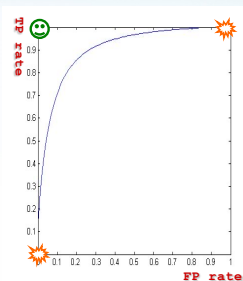
- 1-dimensional data set containing 2 classes
 - **Positive** vs. **negative**
 - Any points located at $x > t$ is classified as positive



44

ROC Curve

- (TPr, FPr)**
 - $TPr = TP / (TP + FN) = \text{RECALL}$
 - $FPr = FP / (FP + TN)$
- (0, 0)
 - Declare everything to be negative
- (1, 1)
 - Declare everything to be positive class
- (1, 0)
 - Ideal
- Diagonal line:
 - Random guessing
- Below diagonal line
 - Prediction is opposite of the true class

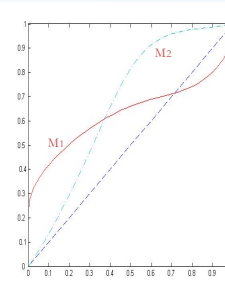


Team for Universal Learning and Intelligent Processing

FLJP: Pattern Recognition (G. Li)

45

Using ROC for Model Comparison



- No Model consistently outperforms the other
 - M_1 is better for small FP Rate
 - M_2 is better for large FP Rate
- Area under the ROC curve
 - Ideal:
 - **Area = 1**
 - Random guess
 - **Area = 0.5**

46

How to Construct an ROC Curve

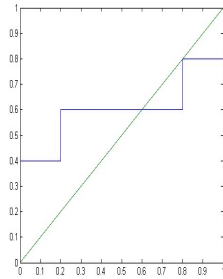
Instance	P(+ A)	True
1	0.95	+
2	0.93	+
3	0.87	-
4	0.85	-
5	0.85	-
6	0.85	+
7	0.76	-
8	0.53	+
9	0.43	-
10	0.25	+

- Use classifier that produces continuous valued output for each test instance $P(+|A)$
- Sort the instances according to $P(+|A)$ in decreasing order
- Apply threshold at each unique value of $P(+|A)$
- Count the number of **TP**, **FP**, **TN**, **FN** at each threshold
- $TPr = TP / (TP + FN)$
- $FPr = FP / (FP + TN)$

47

How to Construct an ROC Curve

Class	Thrd	TP	FP	TN	FN	TPR	FPR
+	0.25	5	5	0	0	1	1
-	0.43	4	5	0	1	0.8	1
+	0.53	4	4	1	1	0.8	0.8
-	0.76	3	4	1	2	0.6	0.8
-	0.85	3	3	2	2	0.6	0.6
-	0.85	3	2	3	2	0.6	0.4
+	0.85	3	1	4	2	0.6	0.2
-	0.87	2	1	4	3	0.4	0.2
+	0.93	2	0	5	3	0.4	0
+	0.95	1	0	5	4	0.2	0
-	1.00	0	0	5	5	0	0



TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

48

48

Model Comparison

- Task 1
 - If we have two models M_1 and M_2 tested on different data sets, how can we tell the relative performance of these two models?
- Task 2
 - If we have two models M_1 and M_2 tested on same data sets, how can we tell the relative performance of these two models?

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

49

49

Model Comparison

- Given two models:
 - M_1 : accuracy = 85%, tested on 30 instances
 - M_2 : accuracy = 75%, tested on 5000 instances
- Can we say M_1 is better than M_2 ?
 - How much confidence can we place on accuracy of M_1 and M_2 ?
 - Can the difference in performance measure be explained as a result of random fluctuations in the test set?

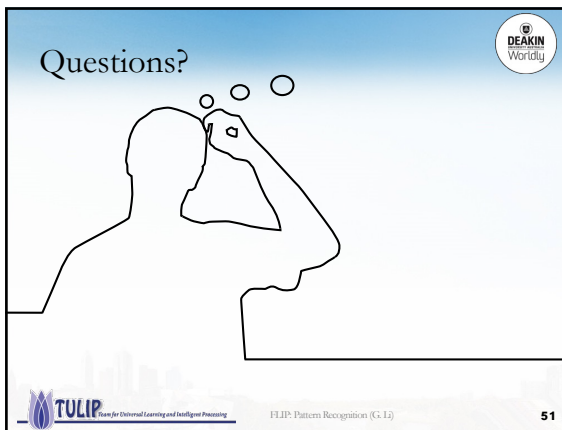
TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

50

50

Questions?



TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

51

51