


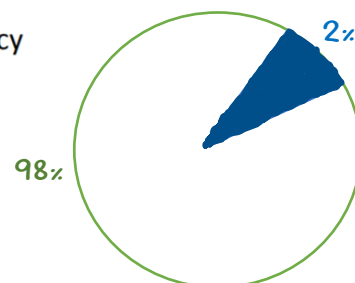
Chapter 8. Classification: Basic Concepts

- ❑ Classification: Basic Concepts
- ❑ Decision Tree Induction
- ❑ Bayes Classification Methods
- ❑ Linear Classifier
- ❑ Model Evaluation and Selection 
- ❑ Techniques to Improve Classification Accuracy: Ensemble Methods
- ❑ Additional Concepts on Classification
- ❑ Summary

47

Model Evaluation and Selection

- ❑ Evaluation metrics
 - ❑ How can we measure accuracy?
 - ❑ Other metrics to consider?
- ❑ Use **validation test set** of class-labeled tuples instead of training set when assessing accuracy
- ❑ Methods for estimating a classifier's accuracy
 - ❑ Holdout method
 - ❑ Cross-validation
 - ❑ Bootstrap
- ❑ Comparing classifiers:
 - ❑ ROC Curves



Classifier ไม่เป็น ความแม่นยำ 98%

48

Classifier Evaluation Metrics: Confusion Matrix

Confusion Matrix: คำตอบที่แท้จริง

Actual class \ Predicted class	Precision	
	C_1	$\neg C_1$
C_1	True Positives (TP)	False Negatives (FN) Recall
$\neg C_1$	False Positives (FP)	True Negatives (TN)

- In a confusion matrix w. m classes, $CM_{i,j}$ indicates # of tuples in class i that were labeled by the classifier as class j

- May have extra rows/columns to provide totals

Example of Confusion Matrix:

Actual class \ Predicted class	Positive		Negative	Total
	buy_computer = yes	buy_computer = no		
buy_computer = yes	Positive 6954	46		7000
buy_computer = no	Negative 412	2588		3000
Total	7366	2634		10000

49

Classifier Evaluation Metrics: Accuracy, Error Rate, Sensitivity and Specificity

A \ P	C	$\neg C$	
C	TP	FN	P
$\neg C$	FP	TN	N
	P'	N'	All

Classifier accuracy, or recognition rate

- Percentage of test set tuples that are correctly classified

$$\text{Accuracy} = (TP + TN) / \text{All}$$

- Error rate: $1 - \text{accuracy}$, or
Error rate = $(FP + FN) / \text{All}$

Class imbalance problem

- One class may be rare
 - E.g., fraud, or HIV-positive
- Significant *majority of the negative class* and minority of the positive class
- Measures handle the class imbalance problem
 - Sensitivity** (recall): True positive recognition rate
 - $\text{Sensitivity} = TP / P$
 - Specificity**: True negative recognition rate
 - $\text{Specificity} = TN / N$

50

Classifier Evaluation Metrics: Precision and Recall, and F-measures

- **Precision:** Exactness: what % of tuples that the classifier labeled as positive are actually positive?

$$P = \text{Precision} = \frac{TP}{TP + FP}$$

ตัวที่ Modal ทายว่าเป็น Pos ถูกต้องมากน้อยแค่ไหน

- **Recall:** Completeness: what % of positive tuples did the classifier label as positive?

$$R = \text{Recall} = \frac{TP}{TP + FN}$$

ตัวที่เป็น Pos จริงๆ เจอมากน้อยแค่ไหน

- Range: [0, 1]
- The “inverse” relationship between precision & recall
- **F measure (or F-score):** harmonic mean of precision and recall
- In general, it is the weighted measure of precision & recall

$$F_\beta = \frac{1}{\alpha \cdot \frac{1}{P} + (1 - \alpha) \cdot \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

Assigning β times as much weight to recall as to precision)

- **F1-measure (balanced F-measure)**

- That is, when $\beta = 1$, $F_1 = \frac{2PR}{P + R}$ R ยิ่งสูง ยิ่งดี , R ยิ่งต่ำ ยิ่งไม่ดี ต้องให้บาลานซ์กัน