

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, random subsampling
 - Cross-validation
 - Bootstrap
- Comparing classifiers:
 - Cost-benefit analysis and ROC Curves

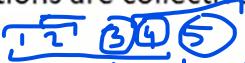
① Holdout Method and Random Subsampling

- the given data are randomly partitioned into two independent sets, a training set and a test set. Typically,
 - 2/3 of data = training set
 - 1/3 of data = test set.
- The training set is used to derive the model. The model's accuracy is then estimated with the test set
 - Random subsampling -- variation of the holdout method, holdout method is repeated k times.
 - The overall accuracy estimate is taken as the average of the accuracies obtained from each iteration



Cross-Validation

Take data ~~some~~ iter.

- In k-fold cross-validation, the initial data are randomly partitioned into k mutually exclusive subsets or "folds," D_1, D_2, \dots, D_k , each of approximately equal size.
- Training and testing is performed k times.
- In iteration i , partition D_i is reserved as the test set, and the remaining partitions are collectively used to train the model.

- each sample is used the same number of times for training and once for testing.
- Accuracy is the overall number of correct classifications from the k iterations, divided by the total number of tuples in the initial data.

[all data used
100 → 100 are]

20

-
- Leave-one-out is a special case of k -fold cross-validation where k is set to the number of initial tuples. That is, only one sample is "left out" at a time for the test set.
 - In stratified cross-validation, the folds are stratified so that the class distribution of the tuples in each fold is approximately the same as that in the initial data



correct accuracy

21

(3)

Bootstrap

- the bootstrap method samples the given training tuples uniformly with replacement.

✓ $Acc(M) = \frac{1}{k} \sum_{i=1}^k (0.632 \times Acc(M_i)_{test_set} + 0.368 \times Acc(M_i)_{train_set})$,

Some tuple can repeat

22

Classifier Evaluation Metrics: Confusion

Matrix

Confusion Matrix:

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

Example of Confusion Matrix:

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

- Given m classes, an entry, $CM_{i,j}$ in a **confusion matrix** indicates # of tuples in class i that were labeled by the classifier as class j
- May have extra rows/columns to provide totals

23

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

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

- Classifier Accuracy, or recognition rate: percentage of test set tuples that are correctly classified

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

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

- 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
- Sensitivity: True Positive recognition rate
 - Sensitivity = TP/P
- Specificity: True Negative recognition rate
 - Specificity = TN/N

24

Limitation of Accuracy

- Consider a 2-class problem
 - Number of Class 0 examples = 9990
 - Number of Class 1 examples = 10
- If model predicts everything to be class 0, accuracy is $9990/10000 = 99.9\%$
 - Accuracy is misleading because model does not detect any class 1 example

25

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

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

$$\text{precision} = \frac{TP}{TP + FP}$$

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

$$\text{recall} = \frac{TP}{TP + FN}$$

- Perfect score is 1.0

- Inverse relationship between precision & recall

- **F measure (F_1 or F-score):** harmonic mean of precision and recall,

$$F = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

- **F_β :** weighted measure of precision and recall

- assigns β times as much weight to recall as to precision

$$F_\beta = \frac{(1 + \beta^2) \times \text{precision} \times \text{recall}}{\beta^2 \times \text{precision} + \text{recall}}$$