# Evaluating Classifier Accuracy: Holdout & Cross-Validation Methods

- ❑ **Holdout method**
  - ❑ Given data is randomly partitioned into two independent sets
    - ❑ Training set (e.g., 2/3) for model construction
    - ❑ Test set (e.g., 1/3) for accuracy estimation
  - ❑ <u>Random sampling</u>: a variation of holdout
    - ❑ Repeat holdout k times, accuracy = avg. of the accuracies obtained
- ❑ **Cross-validation** (*k*-fold, where k = 10 is most popular)
  - ❑ Randomly partition the data into *k mutually exclusive* subsets, each approximately equal size
  - ❑ At *i*-th iteration, use $D_i$ as test set and others as training set
  - ❑ <u>Leave-one-out</u>: *k* folds where *k* = # of tuples, for small sized data
  - ❑ **<u>*Stratified cross-validation*</u>**: folds are stratified so that class dist. in each fold is approx. the same as that in the initial data

# Evaluating Classifier Accuracy: Bootstrap

❑ **Bootstrap**

  ❑ Works well with small data sets

  ❑ Samples the given training tuples uniformly *with replacement*

    ❑ Each time a tuple is selected, it is equally likely to be selected again and re-added to the training set
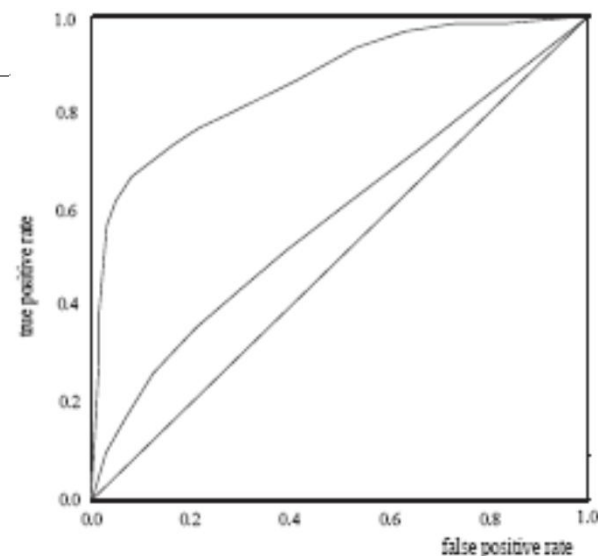
❑ Several bootstrap methods, and a common one is **.632 bootstrap**

  ❑ A data set with *d* tuples is sampled *d* times, with replacement, resulting in a training set of *d* samples.  The data tuples that did not make it into the training set end up forming the test set.  About 63.2% of the original data end up in the bootstrap, and the remaining 36.8% form the test set (since $(1 - 1/d)^d \approx e^{-1} = 0.368$)

  ❑ Repeat the sampling procedure *k* times, overall accuracy of the model:

$$Acc(M) = \frac{1}{k} \sum_{i=1}^{k} (0.632 \times Acc(M_i)_{test\_set} + 0.368 \times Acc(M_i)_{train\_set})$$

# Model Selection: ROC Curves



- ❑ **ROC** (Receiver Operating Characteristics) curves: for visual comparison of classification models

- ❑ Originated from signal detection theory

- ❑ Shows the trade-off between the true positive rate and the false positive rate

- ❑ The area under the ROC curve is a measure of the accuracy of the model

- ❑ Rank the test tuples in decreasing order: the one that is most likely to belong to the positive class appears at the top of the list

- ❑ The closer to the diagonal line (i.e., the closer the area is to 0.5), the less accurate is the model

- ❑ Vertical axis represents the true positive rate (TP/P=sensitivity)

- ❑ Horizontal axis rep. the false positive rate (FP/N=1-specifity)

- ❑ The plot also shows a diagonal line

- ❑ A model with perfect accuracy will have an area of 1.0

# Issues Affecting Model Selection

❑ **Accuracy**

   ❑ classifier accuracy: predicting class label

❑ **Speed**

   ❑ time to construct the model (training time)

   ❑ time to use the model (classification/prediction time)

❑ **Robustness**: handling noise and missing values

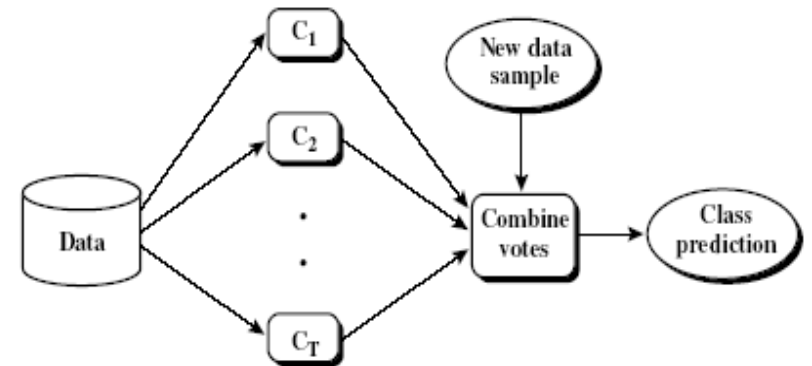❑ **Scalability**: efficiency in disk-resident databases

❑ **Interpretability**

   ❑ understanding and insight provided by the model

❑ Other measures, e.g., goodness of rules, such as decision tree size or compactness of classification rules

# Chapter 8. Classification: Basic Concepts

❑ Classification: Basic Concepts

❑ Decision Tree Induction

❑ Bayes Classification Methods

❑ Model Evaluation and Selection

❑ Techniques to Improve Classification Accuracy: Ensemble Methods

❑ Summary

49

# Ensemble Methods: Increasing the Accuracy



❑ Ensemble methods

   ❑ Use a combination of models to increase accuracy

   ❑ Combine a series of k learned models, $M_1$, $M_2$, …, $M_k$, with the aim of creating an improved model M*

❑ Popular ensemble methods

   ❑ Bagging: averaging the prediction over a collection of classifiers

   ❑ Boosting: weighted vote with a collection of classifiers

   ❑ Ensemble: combining a set of heterogeneous classifiers

# Bagging: Boostrap Aggregation

- ❑ Analogy: Diagnosis based on multiple doctors' majority vote
- ❑ Training
  - ❑ Given a set D of $d$ tuples, at each iteration $i$, a training set $D_i$ of $d$ tuples is sampled with replacement from D (i.e., bootstrap)
  - ❑ A classifier model $M_i$ is learned for each training set $D_i$
- ❑ Classification: classify an unknown sample **X**
  - ❑ Each classifier $M_i$ returns its class prediction
  - ❑ The bagged classifier M* counts the votes and assigns the class with the most votes to **X**
- ❑ Prediction: can be applied to the prediction of continuous values by taking the average value of each prediction for a given test tuple
- ❑ Accuracy: Proved improved accuracy in prediction
  - ❑ Often significantly better than a single classifier derived from D
  - ❑ For noise data: not considerably worse, more robust

# Boosting

❑ Analogy: Consult several doctors, based on a combination of weighted diagnoses—weight assigned based on the previous diagnosis accuracy

❑ How boosting works?

    ❑ **Weights** are assigned to each training tuple

    ❑ A series of k classifiers is iteratively learned

    ❑ After a classifier $M_i$ is learned, the weights are updated to allow the subsequent classifier, $M_{i+1}$, to **pay more attention to the training tuples that were misclassified** by $M_i$

    ❑ The final **M\* combines the votes** of each individual classifier, where the weight of each classifier's vote is a function of its accuracy

❑ Boosting algorithm can be extended for numeric prediction

❑ Comparing with bagging: Boosting tends to have greater accuracy, but it also risks overfitting the model to misclassified data

# Adaboost (Freund and Schapire, 1997)

- ❑ Given a set of $d$ class-labeled tuples, $(\mathbf{X_1}, y_1), ..., (\mathbf{X_d}, y_d)$
- ❑ Initially, all the weights of tuples are set the same (1/d)
- ❑ Generate k classifiers in k rounds.  At round i,
  - ❑ Tuples from D are sampled (with replacement) to form a training set $D_i$ of the same size
  - ❑ Each tuple's chance of being selected is based on its weight
  - ❑ A classification model $M_i$ is derived from $D_i$
  - ❑ Its error rate is calculated using $D_i$ as a test set
  - ❑ If a tuple is misclassified, its weight is increased, o.w. it is decreased
- ❑ Error rate: err($\mathbf{X_j}$) is the misclassification error of tuple $\mathbf{X_j}$. Classifier $M_i$ error rate is the sum of the weights of the misclassified tuples:

$$error\ (M_i) = \sum_j^d w_j \times err\ (\mathbf{X_j})$$

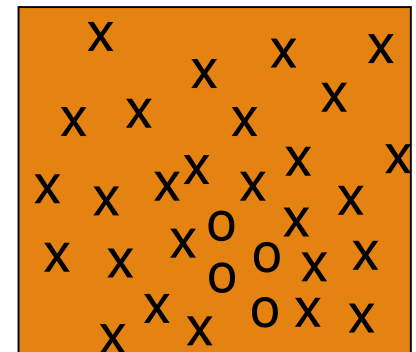- ❑ The weight of classifier $M_i$'s vote is

$$\log \frac{1 - error(M_i)}{error(M_i)}$$

# Random Forest (Breiman 2001)

- ❑ Random Forest:
  - ❑ Each classifier in the ensemble is a *decision tree* classifier and is generated using a random selection of attributes at each node to determine the split
  - ❑ During classification, each tree votes and the most popular class is returned
- ❑ Two Methods to construct Random Forest:
  - ❑ Forest-RI (*random input selection*):  Randomly select, at each node, F attributes as candidates for the split at the node. The CART methodology is used to grow the trees to maximum size
  - ❑ Forest-RC (*random linear combinations*):  Creates new attributes (or features) that are a linear combination of the existing attributes (reduces the correlation between individual classifiers)
- ❑ Comparable in accuracy to Adaboost, but more robust to errors and outliers
- ❑ Insensitive to the number of attributes selected for consideration at each split, and faster than bagging or boosting

# Classification of Class-Imbalanced Data Sets

❑ Class-imbalance problem: Rare positive example but numerous negative ones, e.g., medical diagnosis, fraud, oil-spill, fault, etc.

❑ Traditional methods assume a balanced distribution of classes and equal error costs: not suitable for class-imbalanced data

❑ Typical methods in two-class classification:

  ❑ **Oversampling**: re-sampling of data from positive class

  ❑ **Under-sampling**: randomly eliminate tuples from negative class

  ❑ **Threshold-moving**: move the decision threshold, t, so that the rare class tuples are easier to classify, and hence, less chance of costly false negative errors

  ❑ **Ensemble techniques**: Ensemble multiple classifiers introduced above

❑ Still difficult for class imbalance problem on multiclass tasks

55

# Chapter 8. Classification: Basic Concepts

❑ Classification: Basic Concepts

❑ Decision Tree Induction

❑ Bayes Classification Methods

❑ Model Evaluation and Selection

❑ Techniques to Improve Classification Accuracy: Ensemble Methods

❑ Summary

# Summary

❑ Classification: Extracting models describing important data classes

❑ Effective and scalable methods

  ❑ Decision tree induction, Naive Bayesian classification, rule-based classification, and many other classification methods

❑ Evaluation metrics:

  ❑ Accuracy, sensitivity, specificity, precision, recall, $F$ measure, and $F_{\beta}$ measure

  ❑ Stratified k-fold cross-validation is recommended for accuracy estimation

❑ Enemble: Bagging and boosting can be used to increase overall accuracy by learning and combining a series of individual models

❑ Model selection: Significance tests and ROC curves

❑ No single method has been found to be superior over all others for all data sets

# References (1)

- C. Apte and S. Weiss. **Data mining with decision trees and decision rules**. Future Generation Computer Systems, 13, 1997

- P. K. Chan and S. J. Stolfo. **Learning arbiter and combiner trees from partitioned data for scaling machine learning**. KDD'95

- A. J. Dobson. **An Introduction to Generalized Linear Models**. Chapman & Hall, 1990.

- R. O. Duda, P. E. Hart, and D. G. Stork. **Pattern Classification**, 2ed. John Wiley, 2001

- U. M. Fayyad. **Branching on attribute values in decision tree generation**. AAAI'94.

- Y. Freund and R. E. Schapire. **A decision-theoretic generalization of on-line learning and an application to boosting**. J. Computer and System Sciences, 1997.

- J. Gehrke, R. Ramakrishnan, and V. Ganti. **Rainforest: A framework for fast decision tree construction of large datasets**. VLDB'98.

- J. Gehrke, V. Gant, R. Ramakrishnan, and W.-Y. Loh, **BOAT -- Optimistic Decision Tree Construction**. SIGMOD'99.

- T. Hastie, R. Tibshirani, and J. Friedman. **The Elements of Statistical Learning: Data Mining, Inference, and Prediction.** Springer-Verlag, 2001.

# References (2)

- T.-S. Lim, W.-Y. Loh, and Y.-S. Shih. **A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms.** Machine Learning, 2000

- J. Magidson. **The Chaid approach to segmentation modeling: Chi-squared automatic interaction detection**. In R. P. Bagozzi, editor, Advanced Methods of Marketing Research, Blackwell Business, 1994

- M. Mehta, R. Agrawal, and J. Rissanen. **SLIQ : A fast scalable classifier for data mining**. EDBT'96

- T. M. Mitchell. **Machine Learning**. McGraw Hill, 1997

- S. K. Murthy, **Automatic Construction of Decision Trees from Data: A Multi-Disciplinary Survey**, Data Mining and Knowledge Discovery 2(4): 345-389, 1998

- J. R. Quinlan. **Induction of decision trees**. *Machine Learning*, 1:81-106, 1986.

- J. R. Quinlan. **C4.5: Programs for Machine Learning**. Morgan Kaufmann, 1993.

- J. R. Quinlan. **Bagging, boosting, and c4.5**. AAAI'96.

# References (3)

❑ R. Rastogi and K. Shim. **Public: A decision tree classifier that integrates building and pruning**. VLDB'98

❑ J. Shafer, R. Agrawal, and M. Mehta. **SPRINT : A scalable parallel classifier for data mining**. VLDB'96

❑ J. W. Shavlik and T. G. Dietterich. **Readings in Machine Learning**. Morgan Kaufmann, 1990

❑ P. Tan, M. Steinbach, and V. Kumar. **Introduction to Data Mining**. Addison Wesley, 2005

❑ S. M. Weiss and C. A. Kulikowski. **Computer Systems that Learn: Classification and Prediction Methods from Statistics, Neural Nets, Machine Learning, and Expert Systems**. Morgan Kaufman, 1991

❑ S. M. Weiss and N. Indurkhya. **Predictive Data Mining**. Morgan Kaufmann, 1997

❑ I. H. Witten and E. Frank. **Data Mining: Practical Machine Learning Tools and Techniques**, 2ed. Morgan Kaufmann, 2005