

A Classification Problem

Is a loan to a person who is 45 years old, divorced, renting an apartment, with two kids and annual income of 100K high risk or low risk?





(Classification vs. Regression							
	 Classification predicts categorical attribute values Regression predicts <i>continuous</i> numerical attribute values 							
	SID	HW1	HW2	HW3	Final	Pass/Fail	_	
	1	40	60	70	95	Passed		
	2	10	15	11	65	Failed		
	3	30	45	40	75	Passed		
	4	35	50	35	?	?		

TID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes









Spli ◆A	Splitting Attribute Selection After a split, ideally each subset would "pure", i.e. contains only one class of					
ſ	ECOLOS Gender	Age	Preferred color			
-	female	20	pink			
	male	20	black			
	female	15	pink			
	male	15	black			

Attribute Selection Measures

Entropy (Information Gain)Gini indexGain Ratio

Entropy

$$Entropy(D) = -\sum_{i=1}^{m} p_i \log_2(p_i)$$

 $\label{eq:pi_i} \ensuremath{\overset{\circ}{\mathcal{P}}}_i \ensuremath{\text{is the fraction of records in D that} \\ \ensuremath{\text{belongs to class C_i}} \ensuremath{\overset{\circ}{\mathcal{P}}}_i \ensurema$



Preferred color

- 2 black and 2 pink??
- 3 black and 1 pink??
- 4 black??











- Training error
 - Misclassification of training records
- Testing (Generalization) error
 - Misclassification of testing records









Tree Pruning – Postpruning

- Buttom-up pruning of a fully constructed tree
 - Replace a subtree with a leaf node if it reduces testing error
 - How do we know whether it reduces testing error or not??
 - Pruning based on Minimum Description Length (MDL)

Estimate Testing Errors

- Use a *pruning set* in addition to the training set
- Optimistic error estimation
 - The training set is a good representation of the overall data (optimistic!), so the training error is the testing error
- Pessimistic error estimation
 - Training error + penalty term for model complexity











About Decision Tree Classification ...

- Inexpensive to construct
- Extremely fast at classifying unknown records
- Easy to interpret for small-sized trees
- Accuracy is comparable to other classification techniques for many simple data sets











 $P(\mathbf{x}_k | \mathbf{C}_i)$ is the fraction of number of records in \mathbf{C}_i with value \mathbf{x}_k for attribute \mathbf{A}_k







P(No|HO=No,MS=M,AI=120K) vs. P(Yes|HO=No,MS=M,AI=120K)



Avoid Zero $P(x_k|C_i)$

- A zero $P(\mathbf{x}_k | C_i)$ would make the whole $P(\mathbf{x} | C_i)$ zero
- To avoid this problem, add 1 to to each count, assuming the training set is sufficiently large that the effect of adding one is negligible
- Example
 - Low income:0
 - Medium income: 990
 - High income: 10

About Naive Bayesian Classification

- The most accurate classification if the conditional independence assumption holds
- In practice, some attributes may be correlated
 - E.g. education level and annual income

Bayesian Belief Network (BBN)

- A directed acyclic graph (dag) encoding the dependencies among a set of variables
- A conditional probability table (CPT) for each node given its immediate parent nodes



BBN Terminology

- If there is a directed arc from x to y
 - x is a parent of y
 - $\hfill\blacksquare$ \hfill \hfill hfill \hfill hfill hfill hfill hfill hfill \hfill hfill h
- If there is a directed path from x to y
 - x is an ancestor of Y
 - Y is a descendent of X

Conditional Independence in BBN

A node in a Bayesian network is conditionally independent of its nondescendants if its parents are known











Bayesian Classification Examples – 3

$$\begin{split} P(HD = Yes \mid BP = High, D = Healthy, E = Yes) \\ &= \frac{P(BP = High \mid HD = Yes, D = Healthy, E = Yes)P(HD = Yes \mid D = Healthy, E = Yes)}{P(BP = High \mid HD = Yes)P(HD = Healthy, E = Yes)} \\ &= \frac{P(BP = High \mid HD = Yes)P(HD = Yes \mid D = Healthy, E = Yes)}{\sum_{i=1}^{n} P(BP = High \mid HD = a_i)P(HD = a_i \mid D = Healthy, E = Yes)} \\ &= 0.59 \end{split}$$

Other Classification Methods

- Rule-based
- Artificial Neural Network (ANN)
- Support Vector Machine (SVM)
- Association rule analysis
- Nearest neighbor
- Genetic algorithms
- $\ensuremath{\circledast}$ Rough Set and Fuzzy Set theory
- **.**

Ensemble Methods

- Use a number of *base* classifiers, and make a predication by combining the predications of all the classifiers
- Example
 - Binary classification
 - 25 classifiers, each with error rate 35%
 - Predict by majority vote
 - Error rate of the ensemble classifier??

Construct an Ensemble Classifier

- Train k classifiers with one dataset
 - Use the same dataset for each classifier??
 - Divide the dataset into k subsets??
 - Bagging and Boosting

Bagging (Bootstrap Aggregation)

- ♦ A *bootstrap* sampling of |D|
 - Uniform sampling with replacement (vs. without replacement)
 - Allow duplicates in the sample
 - |D| samples
 - Roughly contains 63.2% of the original records. Why??
- Bagging
 - Use a bootstrap sample as the training set for each classifier

Bagging Example

- Record (x,y)
 - x: attribute
 - y: class label
- Ensemble classifier: 10 classifiers, majority vote

```
©Tan, Steinbach, Kumar Introduction to Data Mining 2004
```

Baggin	g Exar	nple	– Dataset	t
	x	Y	_	
	0.1	1		
	0.2	1		
	0.3	1		
	0.4	-1		
	0.5	-1		
	0.6	-1		
	0.7	-1		
	0.8	1		
	0.9	1		
	1.0	1		





About Bagging

- Reduces the errors associated with random fluctuations in the training data for *unstable classifiers*, e.g. decision trees, rule-based classifiers, ANN
- May degrade the performance of stable classifiers, e.g. Bayesian network, SVM, k-NN

Intuition for Boosting

- Sample with weights
 - hard-to-classify records should be chosen more often
- Combine the prediction of the base classifiers with weights
 - Classifiers with lower error rates get more voting power

Boosting – Training For k classifiers, do k rounds of Assign a weight to each record Sample with replacement according to the weights Train a classifier M_i Calculate error(M_i) Update the weights of the records Increase the weights of the orrectly classified records Decrease the weights of the correctly classified records



Boosting Implementation

- How the record weights are updated
- How the classifier weights are calculated



Example of Accuracy Measures

Example

- Two classes C₁ and C₂
- 100 testing records with 50 $\rm C_1$ records and 50 $\rm C_2$ records
- 20 C₁ records misclassified as C₂, and 10 C₂ records misclassified as C₁

Accuracy measures

- Accuracy and error rates??
- Confusion matrix??
- Precision and Recall??

Evaluate the Accuracy of a Classifier

- The Holdout Method
 - Given a set of records with known class labels, use half of them for training and the other half for testing (or 2/3 for training and 1/3 for testing)

Problems of the Holdout Method

- More records for training means less for testing, and vice versa
- Distribution of the data in the training/testing set may be different from the original dataset
- Some classifiers are sensitive to random fluctuations in the training data

Random Subsampling

- $\ensuremath{\textcircled{\sc Repeat}}$ the holdout method $\ensuremath{\Bbbk}$ times
- $Take the average accuracy over the <math display="inline">{\bf k}$ iterations
- Random subsampling methods
 - Cross-validation
 - Bootstrap

K-fold Cross-validation

- Divide the original dataset into k nonoverlapping subsets
- Each iteration uses (k-1) subsets for training, and the remaining subset for testing
- Total errors are the sum of the errors in each iteration

Bootstrap (.632 Bootstrap)

- Each iteration uses a bootstrap sample to train the classifier, and the remaining records for testing
- Calculate the overall accuracy:

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

Predicating Continuous Values

- Regression methods
 - Linear regression
 - Non-linear regression
- Other methods
 - Some classification methods can be adapted to predict continuous values



Linear Regression Using Least-Squares Method $w_{1} = \frac{\sum_{i=1}^{|D|} (x_{i} - \overline{x})(y_{i} - \overline{y})}{\sum_{i=1}^{|D|} (x_{i} - \overline{x})^{2}}$ $w_{0} = \overline{y} - w_{1}\overline{x}$



