

Data Mining Techniques: Classification and Prediction

Mirek Riedewald

Some slides based on presentations by
Han/Kamber/Pei, Tan/Steinbach/Kumar, and Andrew
Moore

Classification and Prediction Overview

- Introduction
- Decision Trees
- Statistical Decision Theory
- Nearest Neighbor
- Bayesian Classification
- Artificial Neural Networks
- Support Vector Machines (SVMs)
- Prediction
- Accuracy and Error Measures
- Ensemble Methods

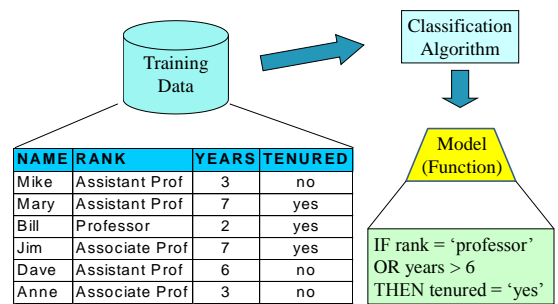
2

Classification vs. Prediction

- Assumption: after data preparation, we have a data set where each record has attributes X_1, \dots, X_n , and Y .
- Goal: learn a function $f: (X_1, \dots, X_n) \rightarrow Y$, then use this function to predict y for a given input record (x_1, \dots, x_n) .
 - **Classification:** Y is a discrete attribute, called the **class label**
 - Usually a categorical attribute with small domain
 - **Prediction:** Y is a continuous attribute
- Called **supervised learning**, because true labels (Y -values) are known for the initially provided data
- Typical applications: credit approval, target marketing, medical diagnosis, fraud detection

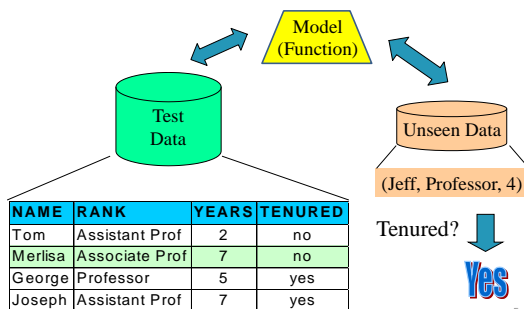
3

Induction: Model Construction



4

Deduction: Using the Model



5

Classification and Prediction Overview

- Introduction
- **Decision Trees**
- Statistical Decision Theory
- Bayesian Classification
- Artificial Neural Networks
- Support Vector Machines (SVMs)
- Nearest Neighbor
- Prediction
- Accuracy and Error Measures
- Ensemble Methods

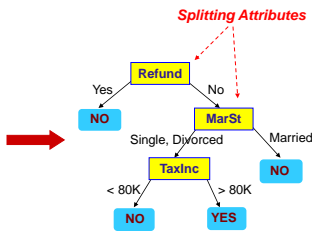
6

Example of a Decision Tree

categorical
categorical
continuous
class

Tid	Refund	Marital Status	Taxable Income	Cheat
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

Training Data

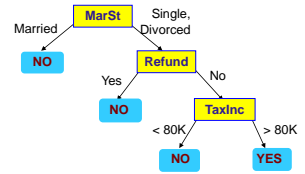


Model: Decision Tree

Another Example of Decision Tree

categorical
categorical
continuous
class

Tid	Refund	Marital Status	Taxable Income	Cheat
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

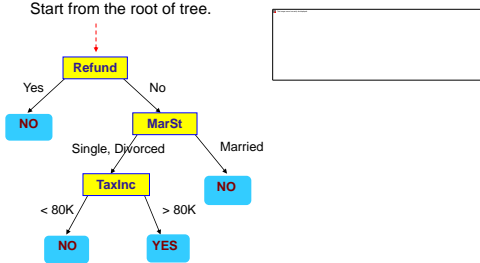


There could be more than one tree that fits the same data!

Apply Model to Test Data

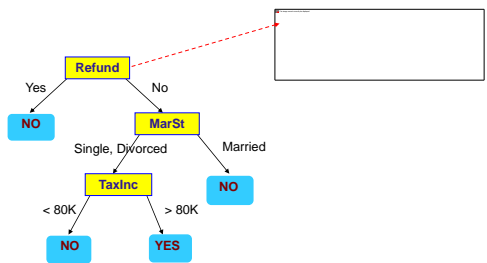
Start from the root of tree.

Test Data



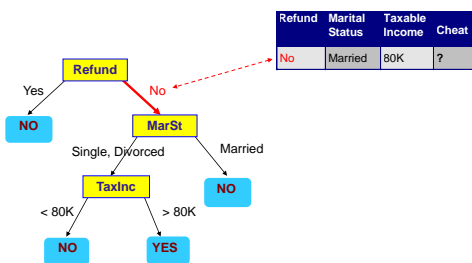
Apply Model to Test Data

Test Data



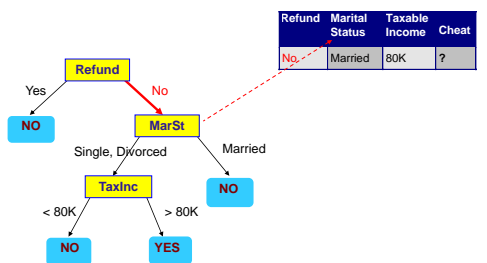
Apply Model to Test Data

Test Data

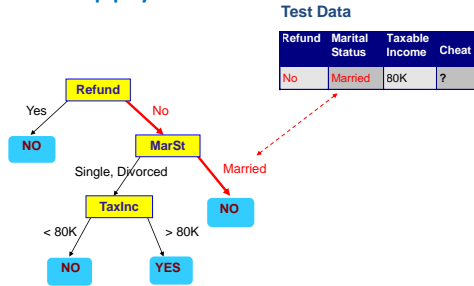


Apply Model to Test Data

Test Data

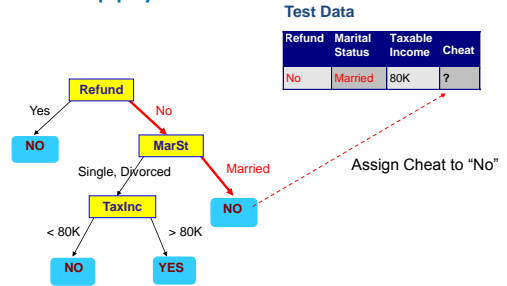


Apply Model to Test Data



13

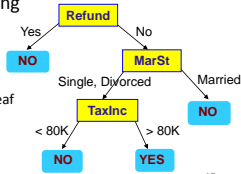
Apply Model to Test Data



14

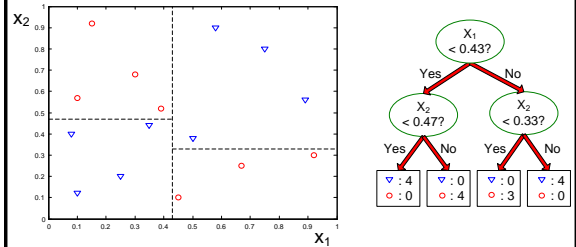
Decision Tree Induction

- Basic greedy algorithm
 - Top-down, recursive divide-and-conquer
 - At start, all the training records are at the root
 - Training records partitioned recursively based on split attributes
 - Split attributes selected based on a heuristic or statistical measure (e.g., information gain)
- Conditions for stopping partitioning
 - Pure node (all records belong to same class)
 - No remaining attributes for further partitioning
 - Majority voting for classifying the leaf
 - No cases left



15

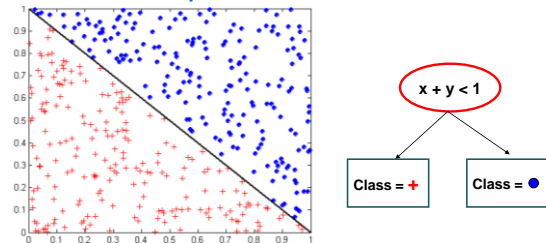
Decision Boundary



Decision boundary = border between two neighboring regions of different classes.
For trees that split on a single attribute at a time, the decision boundary is parallel to the axes.

16

Oblique Decision Trees



- Test condition may involve multiple attributes
- More expressive representation
- Finding optimal test condition is computationally expensive

17

How to Specify Split Condition?

- Depends on attribute types
 - Nominal
 - Ordinal
 - Numeric (continuous)
- Depends on number of ways to split
 - 2-way split
 - Multi-way split

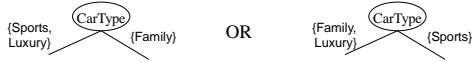
18

Splitting Nominal Attributes

- **Multi-way split:** use as many partitions as distinct values.



- **Binary split:** divides values into two subsets; need to find optimal partitioning.



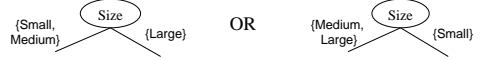
19

Splitting Ordinal Attributes

- **Multi-way split:**



- **Binary split:**



- **What about this split?**



20

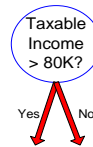
Splitting Continuous Attributes

- Different options

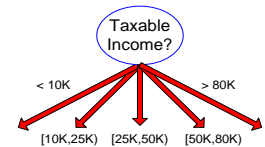
- **Discretization** to form an ordinal categorical attribute
 - Static – discretize once at the beginning
 - Dynamic – ranges found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.
- **Binary Decision:** ($A < v$) or ($A \geq v$)
 - Consider all possible splits, choose best one

21

Splitting Continuous Attributes



(i) Binary split

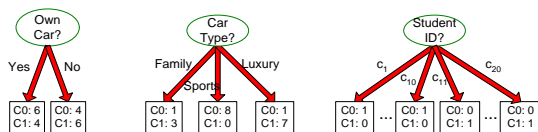


(ii) Multi-way split

22

How to Determine Best Split

Before Splitting: 10 records of class 0,
10 records of class 1



Which test condition is the best?

23

How to Determine Best Split

- Greedy approach:
 - Nodes with **homogeneous** class distribution are preferred
- Need a measure of node impurity:

C0: 5
C1: 5

Non-homogeneous,
High degree of impurity

C0: 9
C1: 1

Homogeneous,
Low degree of impurity

24

Attribute Selection Measure: Information Gain

- Select attribute with highest information gain
- p_i = probability that an arbitrary record in D belongs to class C_i , $i=1, \dots, m$
- Expected information (entropy) needed to classify a record in D :

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

- Information needed after using attribute A to split D into v partitions D_1, \dots, D_v :

$$\text{Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \text{Info}(D_j)$$

- Information gained by splitting on attribute A :

$$\text{Gain}_A(D) = \text{Info}(D) - \text{Info}_A(D)$$

25

Example

- Predict if somebody will buy a computer
- Given data set:

Age	Income	Student	Credit rating	Buys computer
≤ 30	High	No	Bad	No
≤ 30	High	No	Good	No
31...40	High	No	Bad	Yes
> 40	Medium	No	Bad	Yes
> 40	Low	Yes	Bad	Yes
> 40	Low	Yes	Good	No
31...40	Low	Yes	Good	Yes
≤ 30	Medium	No	Bad	No
≤ 30	Low	Yes	Bad	Yes
> 40	Medium	Yes	Bad	Yes
≤ 30	Medium	Yes	Good	Yes
31...40	Medium	No	Good	Yes
31...40	High	Yes	Bad	Yes
> 40	Medium	No	Good	No

26

Information Gain Example

- Class P: buys_computer = "yes"
- Class N: buys_computer = "no"

$$\text{Info}(D) = I(9,5) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.940$$

Age	Yes	No	I(Yes, No)
≤ 30	2	3	0.971
31...40	4	0	0
> 40	3	2	0.971

Age	Income	Student	Credit rating	buys computer
≤ 30	High	No	Bad	No
≤ 30	High	No	Good	No
31...40	High	No	Bad	Yes
> 40	Medium	No	Bad	Yes
> 40	Low	Yes	Bad	Yes
> 40	Low	Yes	Good	No
31...40	Low	Yes	Good	Yes
≤ 30	Medium	No	Bad	No
≤ 30	Low	Yes	Bad	Yes
> 40	Medium	Yes	Bad	Yes
≤ 30	Medium	Yes	Good	Yes
31...40	Medium	No	Good	Yes
31...40	High	Yes	Bad	Yes
> 40	Medium	No	Good	No

$$\text{Info}_{\text{age}}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

- $\frac{5}{14} I(2,3)$ means "age ≤ 30" has 5 out of 14 samples, with 2 yes'es and 3 no's.
- Similar for the other terms

- Hence $\text{Gain}_{\text{age}}(D) = \text{Info}(D) - \text{Info}_{\text{age}}(D) = 0.246$

- Similarly, $\text{Gain}_{\text{income}}(D) = 0.029$
 $\text{Gain}_{\text{student}}(D) = 0.151$
 $\text{Gain}_{\text{credit_rating}}(D) = 0.048$

- Therefore we choose **age** as the splitting attribute

27

Gain Ratio for Attribute Selection

- Information gain is biased towards attributes with a large number of values
- Use gain **ratio** to normalize information gain:

$$\text{GainRatio}_A(D) = \text{Gain}_A(D) / \text{SplitInfo}_A(D)$$

$$\text{SplitInfo}_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \log_2 \left(\frac{|D_j|}{|D|} \right)$$

- E.g., $\text{SplitInfo}_{\text{income}}(D) = -\frac{4}{14} \log_2 \frac{4}{14} - \frac{6}{14} \log_2 \frac{6}{14} - \frac{4}{14} \log_2 \frac{4}{14} = 0.926$

- $\text{GainRatio}_{\text{income}}(D) = 0.029/0.926 = 0.031$
- Attribute with maximum gain ratio is selected as splitting attribute

28

Gini Index

- Gini index, $\text{gini}(D)$, is defined as $\text{gini}(D) = 1 - \sum_{i=1}^m p_i^2$
- If data set D is split on A into v subsets D_1, \dots, D_v the gini index $\text{gini}_A(D)$ is defined as

$$\text{gini}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \text{gini}(D_j)$$

- Reduction in Impurity:

$$\Delta \text{gini}_A(D) = \text{gini}(D) - \text{gini}_A(D)$$

- Attribute that provides smallest $\text{gini}_{\text{split}}(D)$ (= largest reduction in impurity) is chosen to split the node

29

Comparing Attribute Selection Measures

- No clear winner (and there are many more)

- Information gain:

- Biased towards multivalued attributes

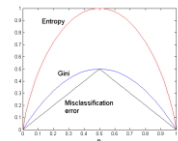
- Gain ratio:

- Tends to prefer unbalanced splits where one partition is much smaller than the others

- Gini index:

- Biased towards multivalued attributes

- Tends to favor tests that result in equal-sized partitions and purity in both partitions



30

Practical Issues of Classification

- Underfitting and overfitting
- Missing values
- Computational cost
- Expressiveness

31

How Good is the Model?

- **Training set error:** compare prediction of training record with true value
 - Not a good measure for the error on unseen data. (Discussed soon.)
- **Test set error:** for records that were **not** used for training, compare model prediction and true value
 - Use holdout data from available data set

32

Training versus Test Set Error

- We'll create a training dataset

Five inputs, all bits, are generated in all 32 possible combinations

Output y = copy of e , except a random 25% of the records have y set to the opposite of e

32 records

a	b	c	d	e	y
0	0	0	0	0	0
0	0	0	0	1	0
0	0	0	1	0	0
0	0	0	1	1	1
0	0	1	0	0	1
:	:	:	:	:	:
1	1	1	1	1	1

33

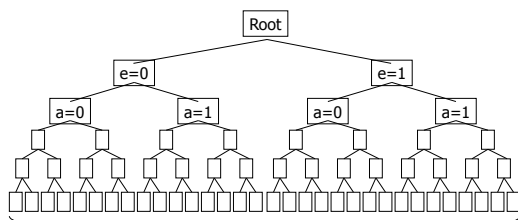
Test Data

- Generate test data using the same method: copy of e , 25% inverted; done independently from previous noise process
- Some y 's that were corrupted in the training set will be uncorrupted in the testing set.
- Some y 's that were uncorrupted in the training set will be corrupted in the test set.

a	b	c	d	e	y (training data)	y (test data)
0	0	0	0	0	0	0
0	0	0	0	1	0	1
0	0	0	1	0	0	1
0	0	0	1	1	1	1
0	0	1	0	0	1	1
:	:	:	:	:	:	:
1	1	1	1	1	1	1

34

Full Tree for The Training Data



25% of these leaf node labels will be corrupted

Each leaf contains exactly one record, hence **no error** in predicting the training data!

35

Testing The Tree with The Test Set

	1/4 of the tree nodes are corrupted	3/4 are fine
1/4 of the test set records are corrupted	1/16 of the test set will be correctly predicted for the wrong reasons	3/16 of the test set will be wrongly predicted because the test record is corrupted
3/4 are fine	3/16 of the test predictions will be wrong because the tree node is corrupted	9/16 of the test predictions will be fine

In total, we expect to be wrong on 3/8 of the test set predictions

36

What's This Example Shown Us?

- Discrepancy between training and test set error
- But more importantly
 - ...it indicates that there is something we should do about it if we want to predict well on future data.

37

Suppose We Had Less Data

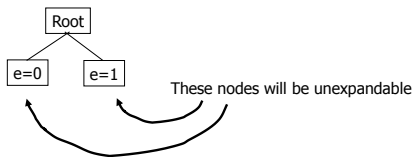
These bits are hidden

Output $y =$ copy of e , except a random 25% of the records have y set to the opposite of e

	a	b	c	d	e	y
32 records	0	0	0	0	0	0
	0	0	0	0	1	0
	0	0	0	1	0	0
	0	0	0	1	1	1
	0	0	1	0	0	1
	:	:	:	:	:	:
	1	1	1	1	1	1

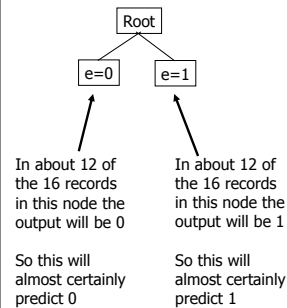
38

Tree Learned Without Access to The Irrelevant Bits



39

Tree Learned Without Access to The Irrelevant Bits



40

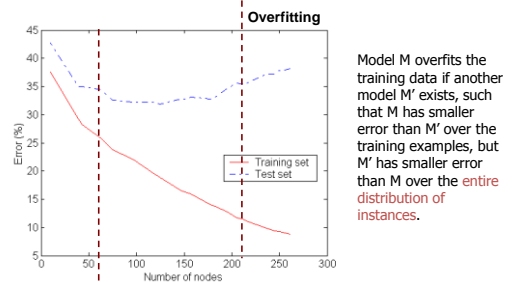
Tree Learned Without Access to The Irrelevant Bits

	almost certainly none of the tree nodes are corrupted	almost certainly all are fine
1/4 of the test set records are corrupted	n/a	1/4 of the test set will be wrongly predicted because the test record is corrupted
3/4 are fine	n/a	3/4 of the test predictions will be fine

In total, we expect to be wrong on only 1/4 of the test set predictions

41

Typical Observation



Underfitting: when model is too simple, both training and test errors are large

42

Reasons for Overfitting

- Noise
 - Too closely fitting the training data means the model's predictions reflect the noise as well
- Insufficient training data
 - Not enough data to enable the model to generalize beyond idiosyncrasies of the training records
- Data fragmentation (special problem for trees)
 - Number of instances gets smaller as you traverse down the tree
 - Number of instances at a leaf node could be too small to make any confident decision about class

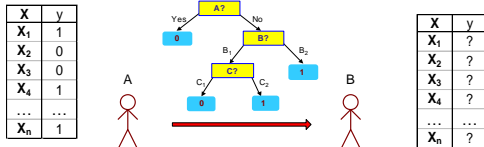
43

Avoiding Overfitting

- General idea: make the tree smaller
 - Addresses all three reasons for overfitting
- *Prepruning*: Halt tree construction early
 - Do not split a node if this would result in the goodness measure falling below a threshold
 - Difficult to choose an appropriate threshold, e.g., tree for XOR
- *Postpruning*: Remove branches from a "fully grown" tree
 - Use a set of data different from the training data to decide when to stop pruning
 - **Validation data**: train tree on training data, prune on validation data, then test on test data

44

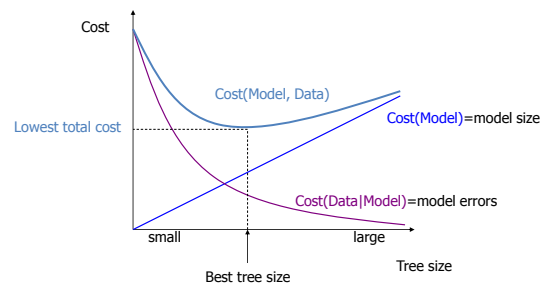
Minimum Description Length (MDL)



- Alternative to using validation data
 - Motivation: data mining is about finding regular patterns in data; regularity can be used to compress the data; method that achieves greatest compression found most regularity and hence is best
- Minimize $\text{Cost}(\text{Model}, \text{Data}) = \text{Cost}(\text{Model}) + \text{Cost}(\text{Data} | \text{Model})$
 - Cost is the number of bits needed for encoding.
 - $\text{Cost}(\text{Data} | \text{Model})$ encodes the misclassification errors.
 - $\text{Cost}(\text{Model})$ uses node encoding plus splitting condition encoding.

45

MDL-Based Pruning Intuition



46

Handling Missing Attribute Values

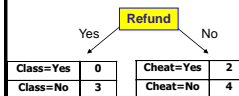
- Missing values affect decision tree construction in three different ways:
 - How impurity measures are computed
 - How to distribute instance with missing value to child nodes
 - How a test instance with missing value is classified

47

Distribute Instances

Tid	Refund	Marital Status	Taxable Income	Class
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

Tid	Refund	Marital Status	Taxable Income	Class
10	?	Single	90K	Yes



Class	Count
Class=Yes	0 + 3/9
Class=No	3

Class	Count
Class=Yes	2 + 6/9
Class=No	4

Probability that Refund=Yes is 3/9
 Probability that Refund=No is 6/9
 Assign record to the left child with weight = 3/9 and to the right child with weight = 6/9

48

Computing Impurity Measure

Tid	Refund	Marital Status	Taxable Income	Class
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	?	Single	90K	Yes

Before Splitting: Entropy(Parent)
= $-0.3 \log(0.3) - (0.7) \log(0.7) = 0.881$

Split on Refund: assume records with missing values are distributed as discussed before

3/9 of record 10 go to Refund=Yes

6/9 of record 10 go to Refund=No

Entropy(Refund=Yes)

$$= -(1/3 / 10/3) \log(1/3 / 10/3)$$

$$= -(3 / 10/3) \log(3 / 10/3) = 0.469$$

Entropy(Refund=No)

$$= -(8/3 / 20/3) \log(8/3 / 20/3)$$

$$= -(4 / 20/3) \log(4 / 20/3) = 0.971$$

Entropy(Children)

$$= 1/3 * 0.469 + 2/3 * 0.971 = 0.804$$

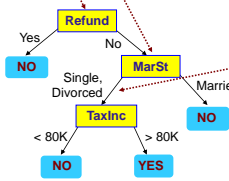
$$\text{Gain} = 0.881 - 0.804 = 0.077$$

49

Classify Instances

New record:

Tid	Refund	Marital Status	Taxable Income	Class
11	No	?	85K	?



	Married	Single	Divorced	Total
Class=No	3	1	0	4
Class=Yes	6/9	1	1	2.67
Total	3.67	2	1	6.67

Probability that Marital Status = Married is 3.67/6.67

Probability that Marital Status = {Single, Divorced} is 3/6.67

50

Tree Cost Analysis

- Finding an optimal decision tree is NP-complete
 - Optimization goal: minimize expected number of binary tests to uniquely identify any record from a given finite set
- Greedy algorithm
 - $O(\#attributes * \#training_instances * \log(\#training_instances))$
 - At each tree depth, all instances considered
 - Assume tree depth is logarithmic (fairly balanced splits)
 - Need to test each attribute at each node
 - What about binary splits?
 - Sort data once on each attribute, use to avoid re-sorting subsets
 - Incrementally maintain counts for class distribution as different split points are explored
- In practice, trees are considered to be fast both for training (when using the greedy algorithm) and making predictions

51

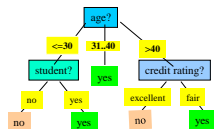
Tree Expressiveness

- Can represent any finite discrete-valued function
 - But it might not do it very efficiently
 - Example: parity function
 - Class = 1 if there is an even number of Boolean attributes with truth value = True
 - Class = 0 if there is an odd number of Boolean attributes with truth value = True
 - For accurate modeling, must have a complete tree
- Not expressive enough for modeling continuous attributes
 - But we can still use a tree for them in practice; it just cannot accurately represent the true function

54

Rule Extraction from a Decision Tree

- One rule is created for each path from the root to a leaf
 - Precondition: conjunction of all split predicates of nodes on path
 - Consequent: class prediction from leaf
- Rules are mutually exclusive and exhaustive
- Example: Rule extraction from buys_computer decision-tree
 - IF age = young AND student = no THEN buys_computer = no
 - IF age = young AND student = yes THEN buys_computer = yes
 - IF age = mid-age THEN buys_computer = yes
 - IF age = old AND credit_rating = excellent THEN buys_computer = yes
 - IF age = young AND credit_rating = fair THEN buys_computer = no



55

Classification in Large Databases

- Scalability: Classify data sets with millions of examples and hundreds of attributes with reasonable speed
- Why use decision trees for data mining?
 - Relatively fast learning speed
 - Can handle all attribute types
 - Convertible to intelligible classification rules
 - Good classification accuracy, but not as good as newer methods (but tree ensembles are top!)

56

Scalable Tree Induction

- High cost when the training data at a node does not fit in memory
- Solution 1: special I/O-aware algorithm
 - Keep only class list in memory, access attribute values on disk
 - Maintain separate list for each attribute
 - Use count matrix for each attribute
- Solution 2: Sampling
 - Common solution: train tree on a sample that fits in memory
 - More sophisticated versions of this idea exist, e.g., *Rainforest*
 - Build tree on sample, but do this for many bootstrap samples
 - Combine all into a single new tree that is guaranteed to be almost identical to the one trained from entire data set
 - Can be computed with two data scans

57

Tree Conclusions

- Very popular data mining tool
 - Easy to understand
 - Easy to implement
 - Easy to use: little tuning, handles all attribute types and missing values
 - Computationally relatively cheap
- Overfitting problem
- Focused on classification, but easy to extend to prediction (future lecture)

58

Classification and Prediction Overview

- Introduction
- Decision Trees
- **Statistical Decision Theory**
- Nearest Neighbor
- Bayesian Classification
- Artificial Neural Networks
- Support Vector Machines (SVMs)
- Prediction
- Accuracy and Error Measures
- Ensemble Methods

60

Theoretical Results

- Trees make sense intuitively, but can we get some hard evidence and deeper understanding about their properties?
- Statistical decision theory can give some answers
- Need some probability concepts first

61

Random Variables

- Intuitive version of the definition:
 - Can take on one of possibly many values, each with a certain probability
 - These probabilities define the probability distribution of the random variable
 - E.g., let X be the outcome of a coin toss, then $\Pr(X=\text{'heads'})=0.5$ and $\Pr(X=\text{'tails'})=0.5$; distribution is uniform
- Consider a discrete random variable X with numeric values x_1, \dots, x_k
 - Expectation: $E[X] = \sum x_i \cdot \Pr(X=x_i)$
 - Variance: $\text{Var}(X) = E[(X - E[X])^2] = E[X^2] - (E[X])^2$

62

Working with Random Variables

- $E[X + Y] = E[X] + E[Y]$
- $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2 \text{Cov}(X, Y)$
- For constants a, b
 - $E[aX + b] = a E[X] + b$
 - $\text{Var}(aX + b) = \text{Var}(aX) = a^2 \text{Var}(X)$
- Iterated expectation:
 - $E[X] = E_x[E_y[Y | X]]$, where $E_y[Y | X] = \sum y_i \cdot \Pr(Y=y_i | X=x)$ is the expectation of Y for a given value x of X , i.e., is a function of X
 - In general for any function $f(X, Y)$:
 $E_{x,y}[f(X, Y)] = E_x[E_y[f(X, Y) | X]]$

63

What is the Optimal Model f(X)?

Let X denote a real-valued random input variable and Y a real-valued random output variable

The squared error of trained model $f(X)$ is $E_{x,y}[(Y - f(X))^2]$

Which function $f(X)$ will minimize the squared error?

Consider the error for a specific value of X and let $\bar{Y} = E_y[Y | X]$:

$$\begin{aligned} E_y[(Y - f(X))^2 | X] &= E_y[(Y - \bar{Y} + \bar{Y} - f(X))^2 | X] \\ &= E_y[(Y - \bar{Y})^2 | X] + E_y[(\bar{Y} - f(X))^2 | X] + 2E_y[(Y - \bar{Y})(\bar{Y} - f(X)) | X] \\ &= E_y[(Y - \bar{Y})^2 | X] + (\bar{Y} - f(X))^2 + 2(\bar{Y} - f(X))E_y[(Y - \bar{Y}) | X] \\ &= E_y[(Y - \bar{Y})^2 | X] + (\bar{Y} - f(X))^2 \end{aligned}$$

(Notice: $E_y[(Y - \bar{Y}) | X] = E_y[Y | X] - E_y[\bar{Y} | X] = \bar{Y} - \bar{Y} = 0$)

64

Optimal Model f(X) (cont.)

The choice of $f(X)$ does not affect $E_y[(Y - \bar{Y})^2 | X]$, but $(\bar{Y} - f(X))^2$ is minimized for $f(X) = \bar{Y} = E_y[Y | X]$.

Note that $E_{x,y}[(Y - f(X))^2] = E_x[E_y[(Y - f(X))^2 | X]]$. Hence

$$E_{x,y}[(Y - f(X))^2] = E_x[E_y[(Y - \bar{Y})^2 | X] + (\bar{Y} - f(X))^2]$$

Hence the squared error is minimized by choosing $f(X) = E_y[Y | X]$ for every X .

(Notice that for minimizing absolute error $E_{x,y}[|Y - f(X)|]$, one can show that the best model is $f(X) = \text{median}(X | Y)$.)

65

Interpreting the Result

- To minimize mean squared error, the best prediction for input $X=x$ is the mean of the Y -values of all training records $(x(i), y(i))$ with $x(i)=x$
 - E.g., assume there are training records (5,22), (5,24), (5,26), (5,28). The optimal prediction for input $X=5$ would be estimated as $(22+24+26+28)/4 = 25$.
- Problem: to reliably estimate the mean of Y for a given $X=x$, we need sufficiently many training records with $X=x$. In practice, often there is only one or no training record at all for an $X=x$ of interest.
 - If there were many such records with $X=x$, we would not need a model and could just return the average Y for that $X=x$.
- The benefit of a good data mining technique is its ability to interpolate and extrapolate from known training records to make good predictions even for X -values that do not occur in the training data at all.
- Classification for two classes: encode as 0 and 1, use squared error as before
 - Then $f(X) = E[Y | X=x] = 1 \cdot \Pr(Y=1 | X=x) + 0 \cdot \Pr(Y=0 | X=x) = \Pr(Y=1 | X=x)$
- Classification for k classes: can show that for 0-1 loss (error = 0 if correct class, error = 1 if wrong class predicted) the optimal choice is to return the majority class for a given input $X=x$
 - This is called the **Bayes classifier**.

66

Implications for Trees

- Since there are not enough, or none at all, training records with $X=x$, the output for input $X=x$ has to be based on records "in the neighborhood"
 - A tree leaf corresponds to a multi-dimensional range in the data space
 - Records in the same leaf are neighbors of each other
- Solution: estimate mean Y for input $X=x$ from the training records in the same leaf node that contains input $X=x$
 - Classification: leaf returns majority class or class probabilities (estimated from fraction of training records in the leaf)
 - Prediction: leaf returns average of Y -values or fits a local model
 - Make sure there are enough training records in the leaf to obtain reliable estimates**

67

Bias-Variance Tradeoff

- Let's take this one step further and see if we can understand overfitting through statistical decision theory
- As before, consider two random variables X and Y
- From a training set D with n records, we want to construct a function $f(X)$ that returns good approximations of Y for future inputs X
 - Make dependence of f on D explicit by writing $f(X; D)$
- Goal: minimize mean squared error over all X, Y , and D , i.e., $E_{X,D,Y}[(Y - f(X; D))^2]$

68

Bias-Variance Tradeoff Derivation

$$\begin{aligned} E_{x,D,y}[(Y - f(X; D))^2] &= E_x E_D E_y[(Y - f(X; D))^2 | X, D] \text{ Now consider the inner term:} \\ E_D E_y[(Y - f(X; D))^2 | X, D] &= E_D[E_y[(Y - EY | X | X, D) + (f(X; D) - EY | X | X, D))^2] \\ &\text{(Same derivation as before for optimal function } f(X)) \\ &= E_D[E_y[(Y - EY | X | X) + (f(X; D) - EY | X | X))^2] \\ &\text{(The first term does not depend on } D, \text{ hence } E_D[E_y[(Y - EY | X | X, D)]] = E_y[(Y - EY | X | X)] \\ \text{Consider the second term:} \\ E_D[(f(X; D) - EY | X | X)]^2 &= E_D[(f(X; D) - E_D[f(X; D)]) + (E_D[f(X; D)] - EY | X | X)]^2 \\ &= E_D[(f(X; D) - E_D[f(X; D)])^2 + E_D[(E_D[f(X; D)] - EY | X | X)]^2 \\ &\quad + 2E_D[(f(X; D) - E_D[f(X; D)])(E_D[f(X; D)] - EY | X | X)] \\ &= E_D[(f(X; D) - E_D[f(X; D)])^2] + (E_D[f(X; D)] - EY | X | X)^2 \\ &\quad + 2E_D[(f(X; D) - E_D[f(X; D)])(E_D[f(X; D)] - EY | X | X)] \\ &= E_D[(f(X; D) - E_D[f(X; D)])^2] + (E_D[f(X; D)] - EY | X | X)^2 \\ &\text{(The third term is zero, because } E_D[f(X; D) - E_D[f(X; D)]] = E_D[f(X; D)] - E_D[f(X; D)] = 0.) \end{aligned}$$

Overall we therefore obtain :

$$E_{x,D,y}[(Y - f(X; D))^2] = E_x[E_D[E_y[(Y - EY | X | X, D)]^2 + E_D[(f(X; D) - E_D[f(X; D)])^2] + E_y[(Y - EY | X | X)]^2]$$

69

Bias-Variance Tradeoff and Overfitting

$(E_D[f(X;D)] - E[Y|X])^2$: bias
 $E_D\{[f(X;D) - E_D[f(X;D)]]^2\}$: variance
 $E_D\{[Y - E[Y|X]]^2 | X\}$: irreducible error (does not depend on f and is simply the variance of Y given X).

- Option 1: $f(X;D) = E[Y|X,D]$
 - Bias: since $E_D[E[Y|X,D]] = E[Y|X]$, bias is zero
 - Variance: $(E[Y|X,D] - E_D[E[Y|X,D]])^2 = (E[Y|X,D] - E[Y|X])^2$ can be very large since $E[Y|X,D]$ depends heavily on D
 - Might overfit!
- Option 2: $f(X;D) = X$ (or other function independent of D)
 - Variance: $(X - E_D[X])^2 = (X - X)^2 = 0$
 - Bias: $(E_D[X] - E[Y|X])^2 = (X - E[Y|X])^2$ can be large, because $E[Y|X]$ might be completely different from X
 - Might underfit!
- Find best compromise between fitting training data too closely (option 1) and completely ignoring it (option 2)

70

Implications for Trees

- Bias decreases as tree becomes larger
 - Larger tree can fit training data better
- Variance increases as tree becomes larger
 - Sample variance affects predictions of larger tree more
- Find right tradeoff as discussed earlier
 - Validation data to find best pruned tree
 - MDL principle

71

Classification and Prediction Overview

- Introduction
- Decision Trees
- Statistical Decision Theory
- Nearest Neighbor
- Bayesian Classification
- Artificial Neural Networks
- Support Vector Machines (SVMs)
- Prediction
- Accuracy and Error Measures
- Ensemble Methods

72

Lazy vs. Eager Learning

- Lazy learning: Simply stores training data (or only minor processing) and waits until it is given a test record
- Eager learning: Given a training set, constructs a classification model before receiving new (test) data to classify
- General trend: Lazy = faster training, slower predictions
- Accuracy: not clear which one is better!
 - Lazy method: typically driven by local decisions
 - Eager method: driven by global and local decisions

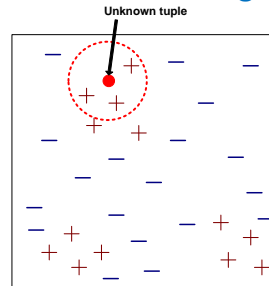
73

Nearest-Neighbor

- Recall our statistical decision theory analysis: Best prediction for input $X=x$ is the mean of the Y -values of all records $(x(i), y(i))$ with $x(i)=x$ (majority class for classification)
- Problem was to estimate $E[Y|X=x]$ or majority class for $X=x$ from the training data
- Solution was to approximate it
 - Use Y -values from training records in neighborhood around $X=x$

74

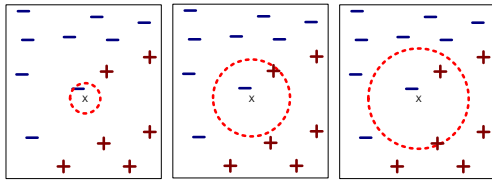
Nearest-Neighbor Classifiers



- Requires:
 - Set of stored records
 - Distance metric for pairs of records
 - Common choice: Euclidean
- $$d(\mathbf{p}, \mathbf{q}) = \sqrt{\sum_i (p_i - q_i)^2}$$
- Parameter k
 - Number of nearest neighbors to retrieve
- To classify a record:
 - Find its k nearest neighbors
 - Determine output based on (distance-weighted) average of neighbors' output

75

Definition of Nearest Neighbor



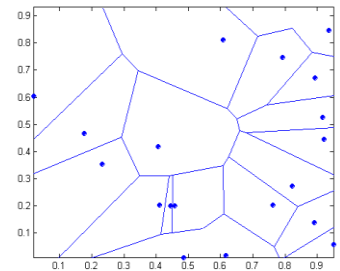
(a) 1-nearest neighbor (b) 2-nearest neighbor (c) 3-nearest neighbor

K-nearest neighbors of a record x are data points that have the k smallest distance to x

76

1-Nearest Neighbor

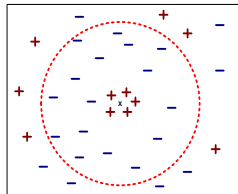
Voronoi Diagram



77

Nearest Neighbor Classification

- Choosing the value of k :
 - k too small: sensitive to noise points
 - k too large: neighborhood may include points from other classes

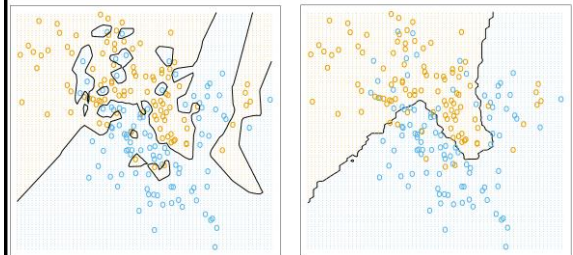


78

Effect of Changing k

1-Nearest Neighbor Classifier

15-Nearest Neighbor Classifier



Source: Hastie, Tibshirani, and Friedman. The Elements of Statistical Learning

79

Explaining the Effect of k

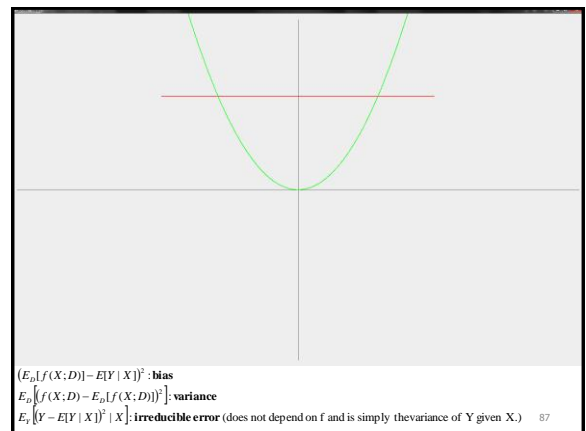
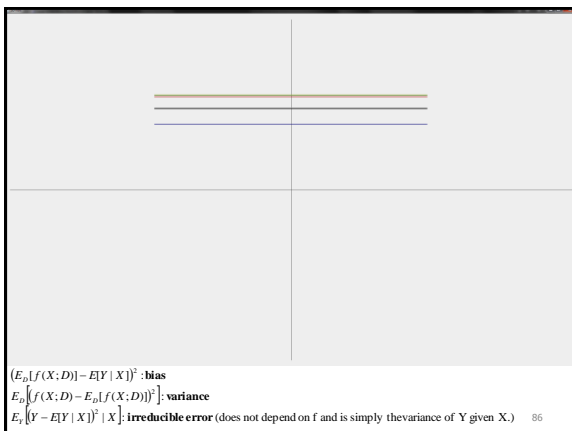
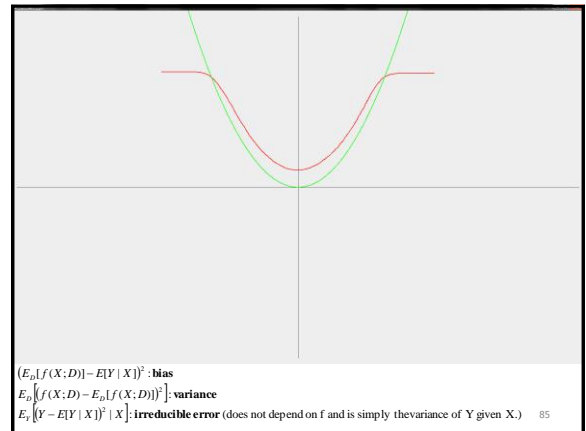
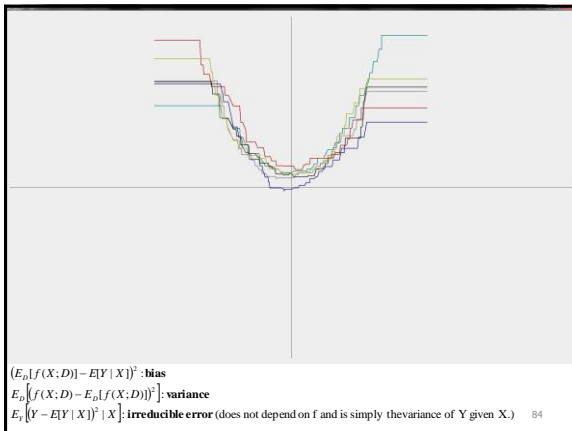
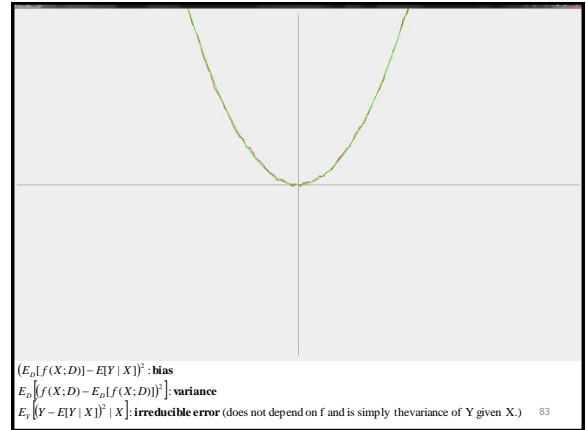
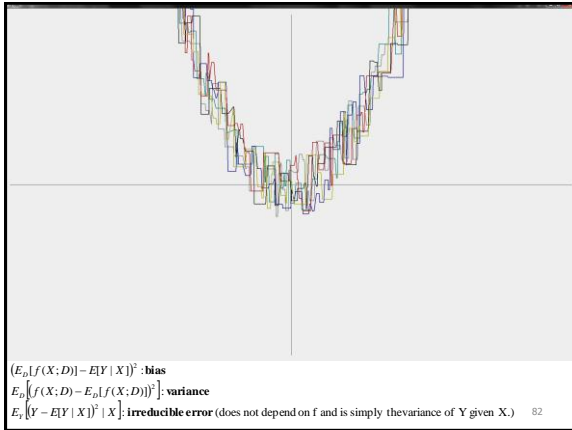
- Recall the bias-variance tradeoff
- Small k , i.e., predictions based on few neighbors
 - High variance, low bias
- Large k , e.g., average over entire data set
 - Low variance, but high bias
- Need to find k that achieves best tradeoff
- Can do that using validation data

80

Experiment

- 50 training points (x, y)
 - $-2 \leq x \leq 2$, selected uniformly at random
 - $y = x^2 + \epsilon$, where ϵ is selected uniformly at random from range $[-0.5, 0.5]$
- Test data sets: 500 points from same distribution as training data, but $\epsilon = 0$
- Plot 1: all $(x, \text{NN1}(x))$ for 5 test sets
- Plot 2: all $(x, \text{AVG}(\text{NN1}(x)))$, averaged over 200 test data set
 - Same for NN20 and NN50

81



Scaling Issues

- Attributes may have to be scaled to prevent distance measures from being dominated by one of the attributes
- Example:
 - Height of a person may vary from 1.5m to 1.8m
 - Weight of a person may vary from 90lb to 300lb
 - Income of a person may vary from \$10K to \$1M
 - Income difference would dominate record distance

88

Other Problems

- Problem with Euclidean measure:
 - High dimensional data: **curse of dimensionality**
 - Can produce counter-intuitive results

1 1 1 1 1 1 1 1 1 1 1 0	vs	1 0 0 0 0 0 0 0 0 0 0 0
0 1 1 1 1 1 1 1 1 1 1 1		0 0 0 0 0 0 0 0 0 0 0 1

$d = 1.4142$

$d = 1.4142$

- Solution: Normalize the vectors to unit length
- Irrelevant attributes might dominate distance
 - Solution: eliminate them

89

Computational Cost

- Brute force: $O(\#trainingRecords)$
 - For each training record, compute distance to test record, keep if among top-k
- Pre-compute Voronoi diagram (expensive), then search spatial index of Voronoi cells: if lucky $O(\log(\#trainingRecords))$
- Store training records in multi-dimensional search tree, e.g., R-tree: if lucky $O(\log(\#trainingRecords))$
- Bulk-compute predictions for many test records using spatial join between training and test set
 - Same worst-case cost as one-by-one predictions, but usually much faster in practice

90

Classification and Prediction Overview

- Introduction
- Decision Trees
- Statistical Decision Theory
- Nearest Neighbor
- **Bayesian Classification**
- Artificial Neural Networks
- Support Vector Machines (SVMs)
- Prediction
- Accuracy and Error Measures
- Ensemble Methods

107

Bayesian Classification

- Performs probabilistic prediction, i.e., predicts class membership probabilities
- Based on Bayes' Theorem
- Incremental training
 - Update probabilities as new training records arrive
 - Can combine prior knowledge with observed data
- Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured

108

Bayesian Theorem: Basics

- \mathbf{X} = random variable for data records ("evidence")
- H = hypothesis that specific record $\mathbf{X}=\mathbf{x}$ belongs to class C
- Goal: determine $P(H | \mathbf{X}=\mathbf{x})$
 - Probability that hypothesis holds given a record \mathbf{x}
- $P(H)$ = **prior** probability
 - The initial probability of the hypothesis
 - E.g., person \mathbf{x} will buy computer, regardless of age, income etc.
- $P(\mathbf{X}=\mathbf{x})$ = probability that data record \mathbf{x} is observed
- $P(\mathbf{X}=\mathbf{x} | H)$ = probability of observing record \mathbf{x} , given that the hypothesis holds
 - E.g., given that \mathbf{x} will buy a computer, what is the probability that \mathbf{x} is in age group 31...40, has medium income, etc.?

109

Bayes' Theorem

- Given data record \mathbf{x} , the **posterior** probability of a hypothesis H , $P(H | \mathbf{X}=\mathbf{x})$, follows from Bayes theorem:

$$P(H | \mathbf{X}=\mathbf{x}) = \frac{P(\mathbf{X}=\mathbf{x} | H)P(H)}{P(\mathbf{X}=\mathbf{x})}$$

- Informally: posterior = likelihood * prior / evidence
- Among all candidate hypotheses H , find the maximally probably one, called **maximum a posteriori (MAP)** hypothesis
- Note: $P(\mathbf{X}=\mathbf{x})$ is the same for all hypotheses
- If all hypotheses are equally probable a priori, we only need to compare $P(\mathbf{X}=\mathbf{x} | H)$
 - Winning hypothesis is called the **maximum likelihood (ML)** hypothesis
- Practical difficulties: requires initial knowledge of many probabilities and has high computational cost

110

Towards Naïve Bayes Classifier

- Suppose there are m classes C_1, C_2, \dots, C_m
- Classification goal: for record \mathbf{x} , find class C_i that has the maximum posterior probability $P(C_i | \mathbf{X}=\mathbf{x})$

- Bayes' theorem:

$$P(C_i | \mathbf{X}=\mathbf{x}) = \frac{P(\mathbf{X}=\mathbf{x} | C_i)P(C_i)}{P(\mathbf{X}=\mathbf{x})}$$

- Since $P(\mathbf{X}=\mathbf{x})$ is the same for all classes, only need to find maximum of $P(\mathbf{X}=\mathbf{x} | C_i)P(C_i)$

111

Computing $P(\mathbf{X}=\mathbf{x} | C_i)$ and $P(C_i)$

- Estimate $P(C_i)$ by counting the frequency of class C_i in the training data
- Can we do the same for $P(\mathbf{X}=\mathbf{x} | C_i)$?
 - Need very large set of training data
 - Have $|X_1| * |X_2| * \dots * |X_d|$ * m different combinations of possible values for X and C_i
 - Need to see every instance \mathbf{x} many times to obtain reliable estimates
- Solution: decompose into lower-dimensional problems

112

Example: Computing $P(\mathbf{X}=\mathbf{x} | C_i)$ and $P(C_i)$

- $P(\text{buys_computer} = \text{yes}) = 9/14$
- $P(\text{buys_computer} = \text{no}) = 5/14$
- $P(\text{age}>40, \text{income}=\text{low}, \text{student}=\text{no}, \text{credit_rating}=\text{bad} | \text{buys_computer}=\text{yes}) = 0 ?$

Age	Income	Student	Credit rating	Buys computer
≤ 30	High	No	Bad	No
≤ 30	High	No	Good	No
31...40	High	No	Bad	Yes
> 40	Medium	No	Bad	Yes
> 40	Low	Yes	Bad	Yes
> 40	Low	Yes	Good	No
31...40	Low	Yes	Good	Yes
≤ 30	Medium	No	Bad	No
≤ 30	Low	Yes	Bad	Yes
> 40	Medium	Yes	Bad	Yes
≤ 30	Medium	Yes	Good	Yes
31...40	Medium	No	Good	Yes
31...40	High	Yes	Bad	Yes
> 40	Medium	No	Good	No

113

Conditional Independence

- X, Y, Z random variables
- X is **conditionally independent** of Y , given Z , if $P(X | Y, Z) = P(X | Z)$
 - Equivalent to: $P(X, Y | Z) = P(X | Z) * P(Y | Z)$
- Example: people with longer arms read better
 - Confounding factor: age
 - Young child has shorter arms and lacks reading skills of adult
 - If age is fixed, observed relationship between arm length and reading skills disappears

114

Derivation of Naïve Bayes Classifier

- Simplifying assumption: all input attributes conditionally independent, given class

$$P(\mathbf{X} = (x_1, \dots, x_d) | C_i) = \prod_{k=1}^d P(X_k = x_k | C_i) = P(X_1 = x_1 | C_i) \cdot P(X_2 = x_2 | C_i) \cdot \dots \cdot P(X_d = x_d | C_i)$$

- Each $P(X_k = x_k | C_i)$ can be estimated robustly

- If X_k is categorical attribute
 - $P(X_k = x_k | C_i) = \frac{\text{\#records in } C_i \text{ that have value } x_k \text{ for } X_k}{\text{\#records of class } C_i \text{ in training data set}}$
- If X_k is continuous, we could discretize it
 - Problem: interval selection
 - Too many intervals: too few training cases per interval
 - Too few intervals: limited choices for decision boundary

115

Estimating $P(X_k=x_k | C_i)$ for Continuous Attributes without Discretization

- $P(X_k=x_k | C_i)$ computed based on Gaussian distribution with mean μ and standard deviation σ :

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

as

$$P(X_k = x_k | C_i) = g(x_k, \mu_{k,C_i}, \sigma_{k,C_i})$$

- Estimate μ_{k,C_i} from sample mean of attribute X_k for all training records of class C_i
- Estimate σ_{k,C_i} similarly from sample

116

Naïve Bayes Example

- Classes:
 - C_1 : buys_computer = yes
 - C_2 : buys_computer = no

Age	Income	Student	Credit_rating	buys_computer
≤ 30	High	No	Bad	No
≤ 30	High	No	Good	No
31...40	High	No	Bad	Yes
> 40	Medium	No	Bad	Yes
> 40	Low	Yes	Bad	Yes
> 40	Low	Yes	Good	No
31...40	Low	Yes	Good	Yes
≤ 30	Medium	No	Bad	No
≤ 30	Low	Yes	Bad	Yes
> 40	Medium	Yes	Bad	Yes
≤ 30	Medium	Yes	Good	Yes
31...40	Medium	No	Good	Yes
31...40	High	Yes	Bad	Yes
> 40	Medium	No	Good	No

- Data sample x
 - age ≤ 30,
 - income = medium,
 - student = yes, and
 - credit_rating = bad

117

Naïve Bayesian Computation

- Compute $P(C)$ for each class:
 - $P(\text{buys_computer} = \text{"yes"}) = 9/14 = 0.643$
 - $P(\text{buys_computer} = \text{"no"}) = 5/14 = 0.357$
- Compute $P(X_k=x_k | C_i)$ for each class
 - $P(\text{age} = \text{"≤ 30"} | \text{buys_computer} = \text{"yes"}) = 2/9 = 0.222$
 - $P(\text{age} = \text{"≤ 30"} | \text{buys_computer} = \text{"no"}) = 3/5 = 0.6$
 - $P(\text{income} = \text{"medium"} | \text{buys_computer} = \text{"yes"}) = 4/9 = 0.444$
 - $P(\text{income} = \text{"medium"} | \text{buys_computer} = \text{"no"}) = 2/5 = 0.4$
 - $P(\text{student} = \text{"yes"} | \text{buys_computer} = \text{"yes"}) = 6/9 = 0.667$
 - $P(\text{student} = \text{"yes"} | \text{buys_computer} = \text{"no"}) = 1/5 = 0.2$
 - $P(\text{credit_rating} = \text{"bad"} | \text{buys_computer} = \text{"yes"}) = 6/9 = 0.667$
 - $P(\text{credit_rating} = \text{"bad"} | \text{buys_computer} = \text{"no"}) = 2/5 = 0.4$
- Compute $P(X=x | C_i)$ using the Naïve Bayes assumption
 - $P(\leq 30, \text{medium, yes, fair} | \text{buys_computer} = \text{"yes"}) = 0.222 * 0.444 * 0.667 * 0.667 = 0.044$
 - $P(\leq 30, \text{medium, yes, fair} | \text{buys_computer} = \text{"no"}) = 0.6 * 0.4 * 0.2 * 0.4 = 0.019$
- Compute final result $P(X=x | C_i) * P(C_i)$
 - $P(X=x | \text{buys_computer} = \text{"yes"}) * P(\text{buys_computer} = \text{"yes"}) = 0.028$
 - $P(X=x | \text{buys_computer} = \text{"no"}) * P(\text{buys_computer} = \text{"no"}) = 0.007$
- Therefore we predict buys_computer = "yes" for input $x = (\text{age} = \text{"≤ 30"}, \text{income} = \text{"medium"}, \text{student} = \text{"yes"}, \text{credit_rating} = \text{"bad"})$

118

Zero-Probability Problem

- Naïve Bayesian prediction requires each conditional probability to be non-zero (why?)

$$P(X = (x_1, \dots, x_d) | C_i) = \prod_{k=1}^d P(X_k = x_k | C_i) = P(X_1 = x_1 | C_i) \cdot P(X_2 = x_2 | C_i) \cdot \dots \cdot P(X_d = x_d | C_i)$$

- Example: 1000 records for buys_computer=yes with income=low (0), income=medium (990), and income = high (10)
 - For input with income=low, conditional probability is zero
- Use Laplacian correction (or Laplace estimator) by adding 1 dummy record to each income level
 - $\text{Prob}(\text{income} = \text{low}) = 1/1003$
 - $\text{Prob}(\text{income} = \text{medium}) = 991/1003$
 - $\text{Prob}(\text{income} = \text{high}) = 11/1003$
- "Corrected" probability estimates close to their "uncorrected" counterparts, but none is zero

119

Naïve Bayesian Classifier: Comments

- Easy to implement
- Good results obtained in many cases
 - Robust to isolated noise points
 - Handles missing values by ignoring the instance during probability estimate calculations
 - Robust to irrelevant attributes
- Disadvantages
 - Assumption: class conditional independence, therefore loss of accuracy
 - Practically, dependencies exist among variables
- How to deal with these dependencies?

120

Probabilities

- Summary of elementary probability facts we have used already and/or will need soon
- Let X be a random variable as usual
- Let A be some predicate over its possible values
 - A is true for some values of X , false for others
 - E.g., X is outcome of throw of a die, A could be "value is greater than 4"
- $P(A)$ is the fraction of possible worlds in which A is true
 - $P(\text{die value is greater than 4}) = 2 / 6 = 1/3$

121

Axioms

- $0 \leq P(A) \leq 1$
- $P(\text{True}) = 1$
- $P(\text{False}) = 0$
- $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$

122

Theorems from the Axioms

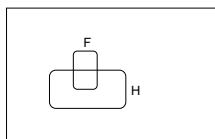
- $0 \leq P(A) \leq 1$, $P(\text{True}) = 1$, $P(\text{False}) = 0$
- $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$
- From these we can prove:
 - $P(\text{not } A) = P(\sim A) = 1 - P(A)$
 - $P(A) = P(A \wedge B) + P(A \wedge \sim B)$

123

Conditional Probability

- $P(A|B)$ = Fraction of worlds in which B is true that also have A true

H = "Have a headache"
F = "Coming down with Flu"



$P(H) = 1/10$
 $P(F) = 1/40$
 $P(H|F) = 1/2$

"Headaches are rare and flu is rarer, but if you're coming down with flu there's a 50-50 chance you'll have a headache."

124

Definition of Conditional Probability

$$P(A|B) = \frac{P(A \wedge B)}{P(B)}$$

Corollary: the **Chain Rule**

$$P(A \wedge B) = P(A|B) P(B)$$

125

Multivalued Random Variables

- Suppose X can take on more than 2 values
- X is a random variable with **arity** k if it can take on exactly one value out of $\{v_1, v_2, \dots, v_k\}$
- Thus

$$P(X = v_i \wedge X = v_j) = 0 \text{ if } i \neq j$$

$$P(X = v_1 \vee X = v_2 \vee \dots \vee X = v_k) = 1$$

126

Easy Fact about Multivalued Random Variables

- Using the axioms of probability
 - $0 \leq P(A) \leq 1$, $P(\text{True}) = 1$, $P(\text{False}) = 0$
 - $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$
- And assuming that X obeys
 - $P(X = v_i \wedge X = v_j) = 0$ if $i \neq j$
 - $P(X = v_1 \vee X = v_2 \vee \dots \vee X = v_k) = 1$
- We can prove that
 - $P(X = v_1 \vee X = v_2 \vee \dots \vee X = v_i) = \sum_{j=1}^i P(X = v_j)$
- And therefore: $\sum_{j=1}^k P(X = v_j) = 1$

127

Useful Easy-to-Prove Facts

$$P(A | B) + P(\sim A | B) = 1$$

$$\sum_{j=1}^k P(X = v_j | B) = 1$$

128

The Joint Distribution Example: Boolean variables A, B, C

Recipe for making a joint distribution of d variables:

129

The Joint Distribution Example: Boolean variables A, B, C

Recipe for making a joint distribution of d variables:

1. Make a truth table listing all combinations of values of your variables (has 2^d rows for d Boolean variables).

A	B	C
0	0	0
0	0	1
0	1	0
0	1	1
1	0	0
1	0	1
1	1	0
1	1	1

130

The Joint Distribution Example: Boolean variables A, B, C

Recipe for making a joint distribution of d variables:

1. Make a truth table listing all combinations of values of your variables (has 2^d rows for d Boolean variables).
2. For each combination of values, say how probable it is.

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

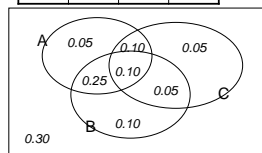
131

The Joint Distribution Example: Boolean variables A, B, C

Recipe for making a joint distribution of d variables:

1. Make a truth table listing all combinations of values of your variables (has 2^d rows for d Boolean variables).
2. For each combination of values, say how probable it is.
3. If you subscribe to the axioms of probability, those numbers must sum to 1.

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10



132

Using the Joint Dist.

gender	hours_worked	wealth	prob
Female	v0.40.5+	poor	0.253122
		rich	0.0245895
v1.40.5+	poor	poor	0.0421768
		rich	0.0116293
Male	v0.40.5-	poor	0.331313
		rich	0.0971295
v1.40.5+	poor	poor	0.134106
		rich	0.105933

Once you have the JD you can ask for the probability of any logical expression involving your attribute

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

133

Using the Joint Dist.

gender	hours_worked	wealth	
Female	v0.40.5-	poor	0.253122
		rich	0.0245895
	v1.40.5+	poor	0.0421768
		rich	0.0116293
Male	v0.40.5-	poor	0.331313
		rich	0.0971295
	v1.40.5+	poor	0.134106
		rich	0.105933

$$P(\text{Poor} \wedge \text{Male}) = 0.4654$$

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

134

Using the Joint Dist.

gender	hours_worked	wealth	
Female	v0.40.5-	poor	0.253122
		rich	0.0245895
	v1.40.5+	poor	0.0421768
		rich	0.0116293
Male	v0.40.5-	poor	0.331313
		rich	0.0971295
	v1.40.5+	poor	0.134106
		rich	0.105933

$$P(\text{Poor}) = 0.7604$$

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

135

Inference with the Joint Dist.

gender	hours_worked	wealth	
Female	v0.40.5-	poor	0.253122
		rich	0.0245895
	v1.40.5+	poor	0.0421768
		rich	0.0116293
Male	v0.40.5-	poor	0.331313
		rich	0.0971295
	v1.40.5+	poor	0.134106
		rich	0.105933

$$P(E_1 | E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\sum_{\text{rows matching } E_1 \text{ and } E_2} P(\text{row})}{\sum_{\text{rows matching } E_2} P(\text{row})}$$

136

Inference with the Joint Dist.

gender	hours_worked	wealth	
Female	v0.40.5-	poor	0.253122
		rich	0.0245895
	v1.40.5+	poor	0.0421768
		rich	0.0116293
Male	v0.40.5-	poor	0.331313
		rich	0.0971295
	v1.40.5+	poor	0.134106
		rich	0.105933

$$P(E_1 | E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\sum_{\text{rows matching } E_1 \text{ and } E_2} P(\text{row})}{\sum_{\text{rows matching } E_2} P(\text{row})}$$

$$P(\text{Male} | \text{Poor}) = 0.4654 / 0.7604 = 0.612$$

137

Joint Distributions

- **Good news:** Once you have a joint distribution, you can answer important questions that involve uncertainty.
- **Bad news:** Impossible to create joint distribution for more than about ten attributes because there are so many numbers needed when you build it.

138

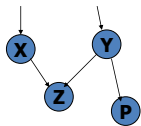
What Would Help?

- Full independence
 - $P(\text{gender}=g \wedge \text{hours_worked}=h \wedge \text{wealth}=w) = P(\text{gender}=g) * P(\text{hours_worked}=h) * P(\text{wealth}=w)$
 - Can reconstruct full joint distribution from a few marginals
- Full conditional independence given class value
 - Naive Bayes
- What about something between Naive Bayes and general joint distribution?

139

Bayesian Belief Networks

- Subset of the variables conditionally independent
- Graphical model of causal relationships
 - Represents dependency among the variables
 - Gives a specification of joint probability distribution

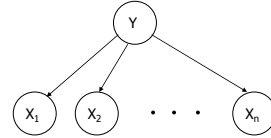


- Nodes: random variables
- Links: dependency
- X and Y are the parents of Z, and Y is the parent of P
- Given Y, Z and P are independent
- Has no loops or cycles

140

Bayesian Network Properties

- Each variable is conditionally independent of its non-descendants in the graph, given its parents
- Naïve Bayes as a Bayesian network:



141

General Properties

- $P(X_1, X_2, X_3) = P(X_1 | X_2, X_3) \cdot P(X_2 | X_3) \cdot P(X_3)$
- $P(X_1, X_2, X_3) = P(X_3 | X_1, X_2) \cdot P(X_2 | X_1) \cdot P(X_1)$
- Network does not necessarily reflect causality



142

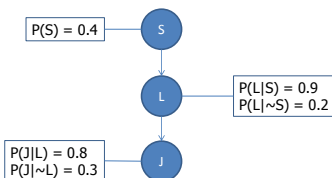
Structural Property

- Missing links simplify computation of $P(X_1, X_2, \dots, X_n)$
- General: $\sum_{i=1}^n P(X_i | X_{i-1}, X_{i-2}, \dots, X_1)$
 - Fully connected: link between every pair of nodes
- Given network: $\sum_{i=1}^n P(X_i | \text{parents}(X_i))$
 - Some links are missing
 - The terms $P(X_i | \text{parents}(X_i))$ are given as **conditional probability tables (CPT)** in the network
- Sparse network allows better estimation of CPT's (fewer combinations of parent values, hence more reliable to estimate from limited data) and faster computation

143

Small Example

- S: Student studies a lot for 6220
- L: Student learns a lot and gets a good grade
- J: Student gets a great job



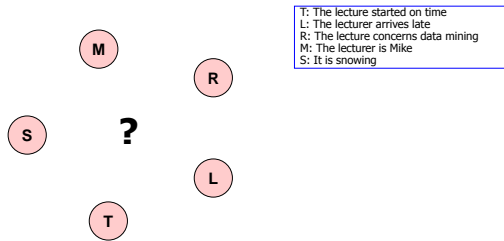
144

Computing $P(S|J)$

- Probability that a student who got a great job was doing her homework
- $P(S|J) = P(S, J) / P(J)$
- $P(S, J) = P(S, J, L) + P(S, J, \sim L)$
- $P(J) = P(J, S, L) + P(J, S, \sim L) + P(J, \sim S, L) + P(J, \sim S, \sim L)$
- $P(J, L, S) = P(J | L, S) \cdot P(L, S) = P(J | L) \cdot P(L | S) \cdot P(S) = 0.8 \cdot 0.9 \cdot 0.4$
- $P(J, \sim L, S) = P(J | \sim L, S) \cdot P(\sim L, S) = P(J | \sim L) \cdot P(\sim L | S) \cdot P(S) = 0.3 \cdot (1-0.9) \cdot 0.4$
- $P(J, L, \sim S) = P(J | L, \sim S) \cdot P(L, \sim S) = P(J | L) \cdot P(L | \sim S) \cdot P(\sim S) = 0.8 \cdot 0.2 \cdot (1-0.4)$
- $P(J, \sim L, \sim S) = P(J | \sim L, \sim S) \cdot P(\sim L, \sim S) = P(J | \sim L) \cdot P(\sim L | \sim S) \cdot P(\sim S) = 0.3 \cdot (1-0.2) \cdot (1-0.4)$
- Putting this all together, we obtain:
- $P(H|J) = (0.8 \cdot 0.9 \cdot 0.4 + 0.3 \cdot 0.1 \cdot 0.4) / (0.8 \cdot 0.9 \cdot 0.4 + 0.3 \cdot 0.1 \cdot 0.4 + 0.8 \cdot 0.2 \cdot 0.6 + 0.3 \cdot 0.8 \cdot 0.6) = 0.3 / 0.54 = 0.56$

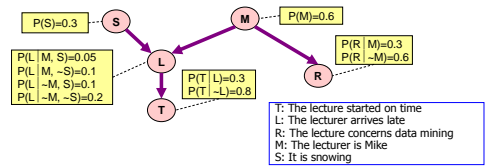
145

More Complex Example



146

Computing with Bayes Net

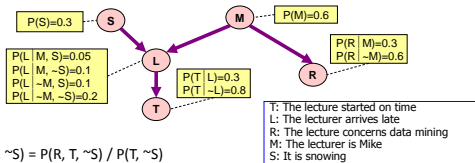


$$P(T, \sim R, L, \sim M, S)$$

$$= P(T | L) \cdot P(\sim R | \sim M) \cdot P(L | \sim M, S) \cdot P(\sim M) \cdot P(S)$$

147

Computing with Bayes Net



$$P(R | T, \sim S) = P(R, T, \sim S) / P(T, \sim S)$$

$$P(R, T, \sim S) = P(L, M, R, T, \sim S) + P(\sim L, M, R, T, \sim S) + P(L, \sim M, R, T, \sim S) + P(\sim L, \sim M, R, T, \sim S)$$

Compute $P(T, \sim S)$ similarly. Problem: There are now 8 such terms to be computed.

148

Inference with Bayesian Networks

- Can predict the probability for any attribute, given any subset of the other attributes
 - $P(M | L, R)$, $P(T | S, \sim M, R)$ and so on
- Easy case: $P(X_i | X_{j_1}, X_{j_2}, \dots, X_{j_k})$ where $\text{parents}(X_i) \subseteq \{X_{j_1}, X_{j_2}, \dots, X_{j_k}\}$
 - Can read answer directly from X_i 's CPT
- What if values are not given for all parents of X_i ?
 - Exact inference of probabilities in general for an arbitrary Bayesian network is **NP-hard**
 - Solutions: probabilistic inference, trade precision for efficiency

149

Training Bayesian Networks

- Several scenarios:
 - Network structure known, all variables observable: learn only the CPTs
 - Network structure known, some hidden variables: gradient descent (greedy hill-climbing) method, analogous to neural network learning
 - Network structure unknown, all variables observable: search through the model space to reconstruct network topology
 - Unknown structure, all hidden variables: No good algorithms known for this purpose
- Ref.: D. Heckerman: Bayesian networks for data mining

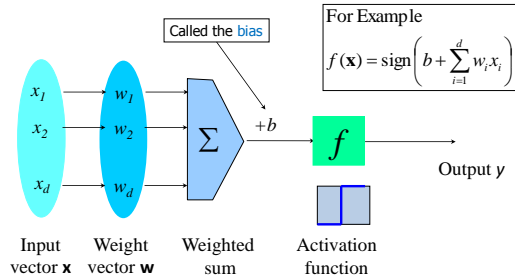
150

Classification and Prediction Overview

- Introduction
- Decision Trees
- Statistical Decision Theory
- Nearest Neighbor
- Bayesian Classification
- **Artificial Neural Networks**
- Support Vector Machines (SVMs)
- Prediction
- Accuracy and Error Measures
- Ensemble Methods

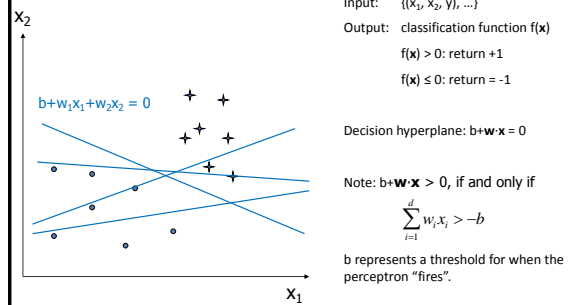
152

Basic Building Block: Perceptron



153

Perceptron Decision Hyperplane



154

Representing Boolean Functions

- AND with two-input perceptron
 - $b = -0.8, w_1 = w_2 = 0.5$
- OR with two-input perceptron
 - $b = -0.3, w_1 = w_2 = 0.5$
- m-of-n function: true if at least m out of n inputs are true
 - All input weights 0.5, threshold weight b is set according to m, n
- Can also represent NAND, NOR
- What about XOR?

155

Perceptron Training Rule

- Goal: correct +1/-1 output for each **training** record
- Start with random weights, constant η (learning rate)
- While some training records are still incorrectly classified do
 - For each training record (\mathbf{x}, y)
 - Let $f_{\text{old}}(\mathbf{x})$ be the output of the current perceptron for \mathbf{x}
 - Set $b := b + \Delta b$, where $\Delta b = \eta(y - f_{\text{old}}(\mathbf{x}))$
 - For all i , set $w_i := w_i + \Delta w_i$, where $\Delta w_i = \eta(y - f_{\text{old}}(\mathbf{x}))x_i$
- Converges to correct decision boundary, if the classes are **linearly separable** and a **small enough η** is used

156

Gradient Descent

- If training records are **not linearly separable**, find best fit approximation
 - Gradient descent to search the space of possible weight vectors
 - Basis for **Backpropagation** algorithm
- Consider **un-thresholded** perceptron (no sign function applied), i.e., $u(\mathbf{x}) = b + \mathbf{w} \cdot \mathbf{x}$
- Measure training error by squared error

$$E(b, \mathbf{w}) = \frac{1}{2} \sum_{(\mathbf{x}, y) \in D} (y - u(\mathbf{x}))^2$$

- D = training data

157

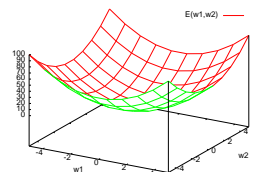
Gradient Descent Rule

- Find weight vector that minimizes $E(b, \mathbf{w})$ by altering it in direction of steepest descent
 - Set $(b, \mathbf{w}) := (b, \mathbf{w}) + \Delta(b, \mathbf{w})$, where $\Delta(b, \mathbf{w}) = -\eta \nabla E(b, \mathbf{w})$
 - $-\nabla E(b, \mathbf{w}) = [\partial E / \partial b, \partial E / \partial w_1, \dots, \partial E / \partial w_n]$ is the **gradient**, hence

$$b := b - \eta \frac{\partial E}{\partial b} = b - \eta \left(- \sum_{(\mathbf{x}, y) \in D} (y - u(\mathbf{x})) \right)$$

$$w_i := w_i - \eta \frac{\partial E}{\partial w_i} = w_i - \eta \sum_{(\mathbf{x}, y) \in D} (y - u(\mathbf{x})) x_i$$

- Start with random weights, iterate until convergence
 - Will converge to global minimum if η is small enough



158

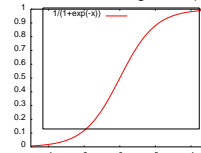
Gradient Descent Summary

- Epoch updating (batch mode)
 - Compute gradient over **entire** training set
 - Changes model once per scan of entire training set
- Case updating (incremental mode, stochastic gradient descent)
 - Compute gradient for a **single** training record
 - Changes model after every single training record immediately
- Case updating can approximate epoch updating arbitrarily close if η is small enough
- What is the difference between perceptron training rule and case updating for gradient descent?
 - Error computation on thresholded vs. unthresholded function

159

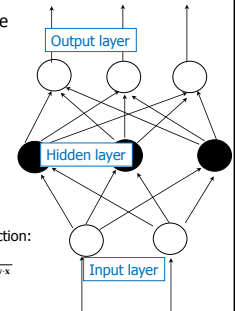
Multilayer Feedforward Networks

- Use another perceptron to combine output of lower layer
 - What about linear units only? Can only construct linear functions!
 - Need nonlinear component
 - sign function: not differentiable (gradient descent!)
 - Use sigmoid: $\sigma(x) = 1/(1+e^{-x})$



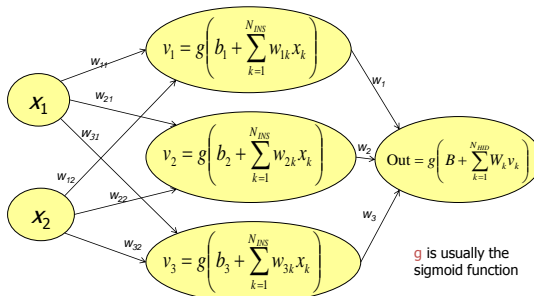
Perceptron function:

$$y = \frac{1}{1 + e^{-b - wx}}$$



160

1-Hidden Layer ANN Example



161

Making Predictions

- Input record fed simultaneously into the units of the input layer
- Then weighted and fed simultaneously to a hidden layer
- Weighted outputs of the last hidden layer are the input to the units in the output layer, which emits the network's prediction
- The network is **feed-forward**
 - None of the weights cycles back to an input unit or to an output unit of a previous layer
- Statistical point of view: neural networks perform nonlinear regression

162

Backpropagation Algorithm

- Earlier discussion: gradient descent for a **single** perceptron using a simple un-thresholded function
- If sigmoid (or other differentiable) function is applied to weighted sum, use **complete function** for gradient descent
- Multiple perceptrons: optimize over all weights of all perceptrons
 - Problems: huge search space, local minima
- **Backpropagation**
 - Initialize all weights with small random values
 - Iterate many times
 - Compute gradient, starting at output and working back
 - Error of hidden unit h : how do we get the true output value? Use weighted sum of errors of each unit influenced by h
 - Update all weights in the network

163

Overfitting

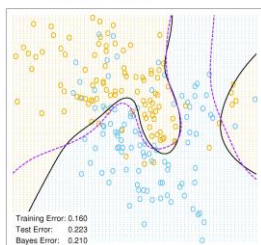
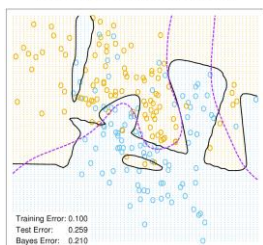
- When do we stop updating the weights?
- Overfitting tends to happen in later iterations
 - Weights initially small random values
 - Weights all similar \Rightarrow smooth decision surface
 - Surface complexity increases as weights diverge
- Preventing overfitting
 - Weight decay: decrease each weight by small factor during each iteration, or
 - Use validation data to decide when to stop iterating

164

Neural Network Decision Boundary

Neural Network - 10 Units, No Weight Decay

Neural Network - 10 Units, Weight Decay=0.02



Source: Hastie, Tibshirani, and Friedman. The Elements of Statistical Learning

165

Backpropagation Remarks

- Computational cost
 - Each iteration costs $O(|D| * |\mathbf{w}|)$, with $|D|$ training records and $|\mathbf{w}|$ weights
 - Number of iterations can be exponential in n , the number of inputs (in practice often tens of thousands)
- Local minima can trap the gradient descent algorithm: convergence guaranteed to *local* minimum, not *global*
- Backpropagation highly effective in practice
 - Many variants to deal with local minima issue, use of case updating

166

Defining a Network

1. Decide network topology
 - #input units, #hidden layers, #units per hidden layer, #output units (one output unit per class for problems with >2 classes)
2. Normalize input values for each attribute to $[0.0, 1.0]$
 - Nominal/ordinal attributes: one input unit *per domain value*
 - For attribute *grade* with values A, B, C, have 3 inputs that are set to 1,0,0 for grade A, to 0,1,0 for grade B, and 0,0,1 for C
 - Why not map it to a single input with domain $[0.0, 1.0]$?
3. Choose learning rate η , e.g., $1 / (\# \text{training iterations})$
 - Too small: takes too long to converge
 - Too large: might never converge (oversteps minimum)
4. Bad results on test data? Change network topology, initial weights, or learning rate; try again.

167

Representational Power

- Boolean functions
 - Each can be represented by a 2-layer network
 - Number of hidden units can grow exponentially with number of inputs
 - Create hidden unit for each input record
 - Set its weights to activate only for that input
 - Implement output unit as OR gate that only activates for desired output patterns
- Continuous functions
 - Every bounded continuous function can be approximated arbitrarily close by a 2-layer network
- Any function can be approximated arbitrarily close by a 3-layer network

168

Neural Network as a Classifier

- Weaknesses
 - Long training time
 - Many non-trivial parameters, e.g., network topology
 - Poor interpretability: What is the meaning behind learned weights and hidden units?
 - Note: hidden units are alternative representation of input values, capturing their relevant features
- Strengths
 - High tolerance to noisy data
 - Well-suited for continuous-valued inputs and outputs
 - Successful on a wide array of real-world data
 - Techniques exist for extraction of rules from neural networks

169

Classification and Prediction Overview

- Introduction
- Decision Trees
- Statistical Decision Theory
- Nearest Neighbor
- Bayesian Classification
- Artificial Neural Networks
- Support Vector Machines (SVMs)
- Prediction
- Accuracy and Error Measures
- Ensemble Methods

171

SVM—Support Vector Machines

- Newer and very popular classification method
- Uses a nonlinear mapping to transform the original training data into a higher dimension
- Searches for the optimal separating hyperplane (i.e., “decision boundary”) in the new dimension
- SVM finds this hyperplane using support vectors (“essential” training records) and margins (defined by the support vectors)

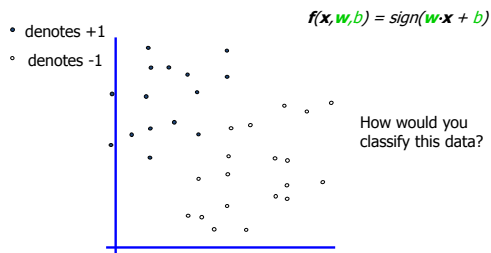
172

SVM—History and Applications

- Vapnik and colleagues (1992)
 - Groundwork from Vapnik & Chervonenkis’ statistical learning theory in 1960s
- Training can be slow but accuracy is high
 - Ability to model complex nonlinear decision boundaries (margin maximization)
- Used both for classification and prediction
- Applications: handwritten digit recognition, object recognition, speaker identification, benchmarking time-series prediction tests

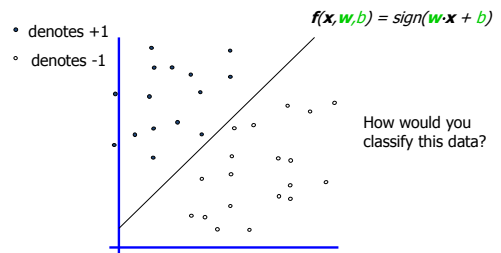
173

Linear Classifiers



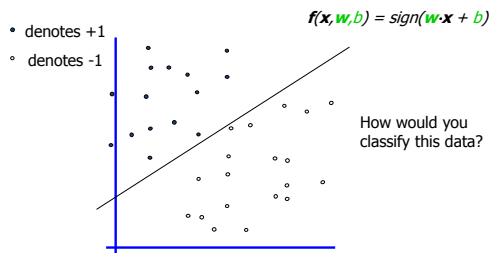
174

Linear Classifiers



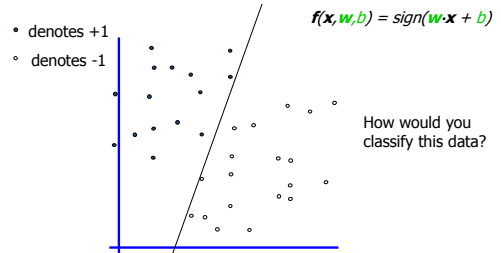
175

Linear Classifiers



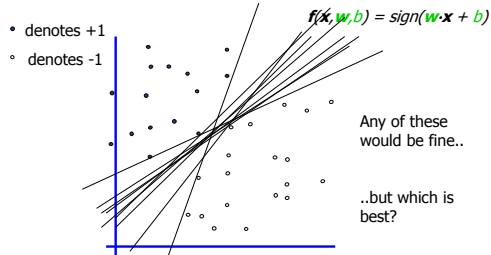
176

Linear Classifiers



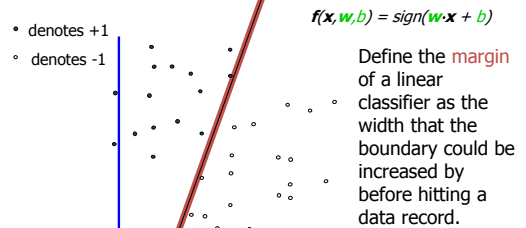
177

Linear Classifiers



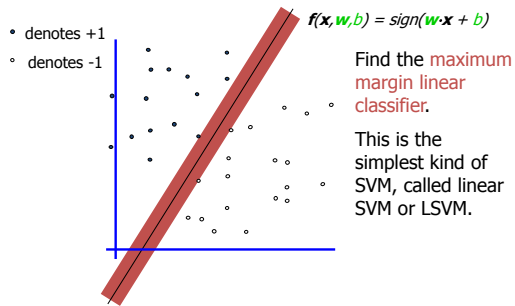
178

Classifier Margin



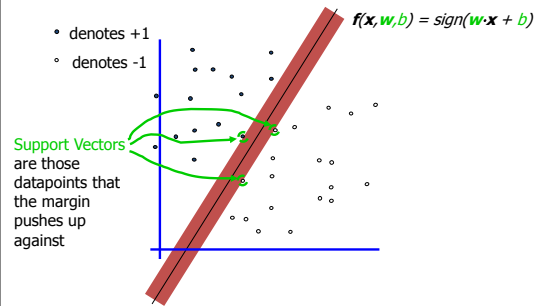
179

Maximum Margin



180

Maximum Margin



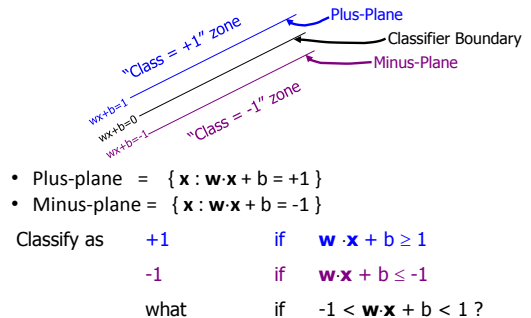
181

Why Maximum Margin?

- If we made a small error in the location of the boundary, this gives us the least chance of causing a misclassification.
- Model is immune to removal of any non-support-vector data records.
- There is some theory (using VC dimension) that is related to (but not the same as) the proposition that this is a good thing.
- Empirically it works very well.

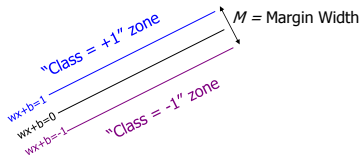
182

Specifying a Line and Margin



183

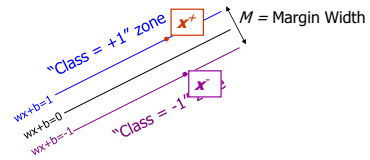
Computing Margin Width



- Plus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$
- Goal: compute M in terms of \mathbf{w} and b
 - Note: vector \mathbf{w} is perpendicular to plus-plane
 - Consider two vectors \mathbf{u} and \mathbf{v} on plus-plane and show that $\mathbf{w} \cdot (\mathbf{u} - \mathbf{v}) = 0$
 - Hence it is also perpendicular to the minus-plane

184

Computing Margin Width



- Choose arbitrary point \mathbf{x}^* on minus-plane
- Let \mathbf{x}^+ be the point in plus-plane closest to \mathbf{x}^*
- Since vector \mathbf{w} is perpendicular to these planes, it holds that $\mathbf{x}^+ = \mathbf{x}^* + \lambda \mathbf{w}$, for some value of λ

185

Putting It All Together

- We have so far:
 - $\mathbf{w} \cdot \mathbf{x}^+ + b = +1$ and $\mathbf{w} \cdot \mathbf{x}^* + b = -1$
 - $\mathbf{x}^+ = \mathbf{x}^* + \lambda \mathbf{w}$
 - $|\mathbf{x}^+ - \mathbf{x}^*| = M$
- Derivation:
 - $\mathbf{w} \cdot (\mathbf{x}^* + \lambda \mathbf{w}) + b = +1$, hence $\mathbf{w} \cdot \mathbf{x}^* + b + \lambda \mathbf{w} \cdot \mathbf{w} = 1$
 - This implies $\lambda \mathbf{w} \cdot \mathbf{w} = 2$, i.e., $\lambda = 2 / \mathbf{w} \cdot \mathbf{w}$
 - Since $M = |\mathbf{x}^+ - \mathbf{x}^*| = |\lambda \mathbf{w}| = \lambda |\mathbf{w}| = \lambda (\mathbf{w} \cdot \mathbf{w})^{0.5}$
 - We obtain $M = 2 (\mathbf{w} \cdot \mathbf{w})^{0.5} / \mathbf{w} \cdot \mathbf{w} = 2 / (\mathbf{w} \cdot \mathbf{w})^{0.5}$

186

Finding the Maximum Margin

- How do we find \mathbf{w} and b such that the margin is maximized and *all training records are in the correct zone for their class*?
- Solution: Quadratic Programming (QP)
- QP is a well-studied class of optimization algorithms to maximize a **quadratic function** of some real-valued variables subject to **linear constraints**.
 - There exist algorithms for finding such constrained quadratic optima efficiently and reliably.

187

Quadratic Programming

Find $\arg \max_{\mathbf{u}} c + \mathbf{d}^T \mathbf{u} + \frac{\mathbf{u}^T \mathbf{R} \mathbf{u}}{2}$ ← Quadratic criterion

Subject to

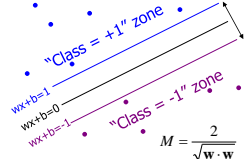
$$\left. \begin{aligned} a_{11}u_1 + a_{12}u_2 + \dots + a_{1m}u_m &\leq b_1 \\ a_{21}u_1 + a_{22}u_2 + \dots + a_{2m}u_m &\leq b_2 \\ &\vdots \\ a_{n1}u_1 + a_{n2}u_2 + \dots + a_{nm}u_m &\leq b_n \end{aligned} \right\} n \text{ additional linear inequality constraints}$$

And subject to

$$\left. \begin{aligned} a_{(n+1)1}u_1 + a_{(n+1)2}u_2 + \dots + a_{(n+1)m}u_m &= b_{(n+1)} \\ a_{(n+2)1}u_1 + a_{(n+2)2}u_2 + \dots + a_{(n+2)m}u_m &= b_{(n+2)} \\ &\vdots \\ a_{(n+e)1}u_1 + a_{(n+e)2}u_2 + \dots + a_{(n+e)m}u_m &= b_{(n+e)} \end{aligned} \right\} e \text{ additional linear equality constraints}$$

188

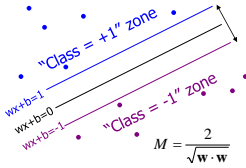
What Are the SVM Constraints?



- Consider n training records $(\mathbf{x}(k), y(k))$, where $y(k) = +/- 1$
- How many constraints will we have?
- What should they be?
- What is the quadratic optimization criterion?

189

What Are the SVM Constraints?



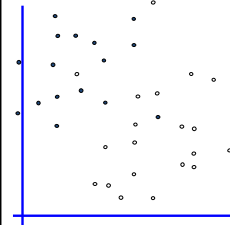
- What is the quadratic optimization criterion?
 - Minimize $\mathbf{w} \cdot \mathbf{w}$

- Consider n training records $(\mathbf{x}(k), y(k))$, where $y(k) = +/- 1$
 - How many constraints will we have? n .
 - What should they be?
- For each $1 \leq k \leq n$:
- $\mathbf{w} \cdot \mathbf{x}(k) + b \geq 1$, if $y(k)=1$
 - $\mathbf{w} \cdot \mathbf{x}(k) + b \leq -1$, if $y(k)=-1$

190

Problem: Classes Not Linearly Separable

- denotes +1
- denotes -1

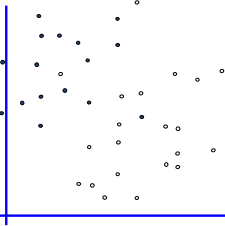


- Inequalities for training records are not satisfiable by any \mathbf{w} and b

191

Solution 1?

- denotes +1
- denotes -1

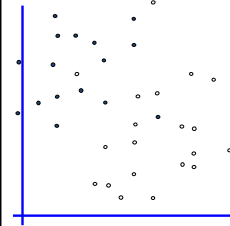


- Find minimum $\mathbf{w} \cdot \mathbf{w}$, while also minimizing number of training set errors
 - Not a well-defined optimization problem (cannot optimize two things at the same time)

192

Solution 2?

- denotes +1
- denotes -1

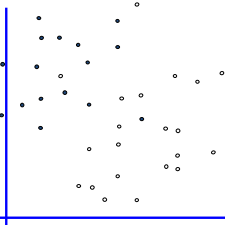


- Minimize $\mathbf{w} \cdot \mathbf{w} + C \cdot (\text{\#trainSetErrors})$
 - C is a tradeoff parameter
- Problems:
 - Cannot be expressed as QP, hence finding solution might be slow
 - Does not distinguish between disastrous errors and near misses

193

Solution 3

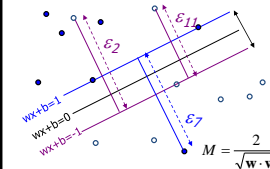
- denotes +1
- denotes -1



- Minimize $\mathbf{w} \cdot \mathbf{w} + C \cdot (\text{distance of error records to their correct place})$
- This works!
- But still need to do something about the unsatisfiable set of inequalities

194

What Are the SVM Constraints?



- What is the quadratic optimization criterion?
 - Minimize

$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^n \varepsilon_k$$

- Consider n training records $(\mathbf{x}(k), y(k))$, where $y(k) = +/- 1$
 - How many constraints will we have? n .
 - What should they be?
- For each $1 \leq k \leq n$:
- $\mathbf{w} \cdot \mathbf{x}(k) + b \geq 1 - \varepsilon_k$, if $y(k)=1$
 - $\mathbf{w} \cdot \mathbf{x}(k) + b \leq -1 + \varepsilon_k$, if $y(k)=-1$
 - $\varepsilon_k \geq 0$

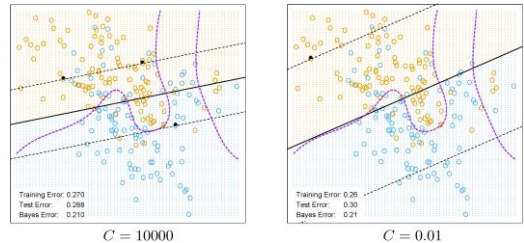
195

Facts About the New Problem Formulation

- Original QP formulation had $d+1$ variables
 - w_1, w_2, \dots, w_d and b
- New QP formulation has $d+1+n$ variables
 - w_1, w_2, \dots, w_d and b
 - $\epsilon_1, \epsilon_2, \dots, \epsilon_n$
- C is a new parameter that needs to be set for the SVM
 - Controls tradeoff between paying attention to margin size versus misclassifications

196

Effect of Parameter C



Source: Hastie, Tibshirani, and Friedman. The Elements of Statistical Learning

197

An Equivalent QP (The “Dual”)

$$\text{Maximize } \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k=1}^n \sum_{l=1}^n \alpha_k \alpha_l \cdot y(k) \cdot y(l) \cdot \mathbf{x}(k) \cdot \mathbf{x}(l)$$

$$\text{Subject to these constraints: } \forall k : 0 \leq \alpha_k \leq C \quad \sum_{k=1}^n \alpha_k y(k) = 0$$

Then define:

$$\mathbf{w} = \sum_{k=1}^n \alpha_k \cdot y(k) \cdot \mathbf{x}(k)$$

Then classify with:

$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

$$b = \text{AVG}_{k:0 < \alpha_k < C} \left\{ \frac{1}{y(k)} - \mathbf{x}(k) \cdot \mathbf{w} \right\}$$

198

Important Facts

- Dual formulation of QP can be optimized more quickly, but result is equivalent
- Data records with $\alpha_k > 0$ are the **support vectors**
 - Those with $0 < \alpha_k < C$ lie on the plus- or minus-plane
 - Those with $\alpha_k = C$ are on the wrong side of the classifier boundary (have $\epsilon_k > 0$)
- Computation for \mathbf{w} and b only depends on those records with $\alpha_k > 0$, i.e., the support vectors
- Alternative QP has another major advantage, as we will see now...

199

Easy To Separate

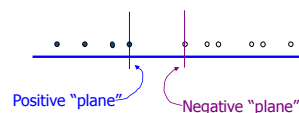
What would SVMs do with this data?



200

Easy To Separate

Not a big surprise



201

Harder To Separate

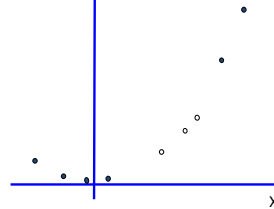
What can be done about this?



202

Harder To Separate

$X' (= X^2)$



Non-linear basis functions:

Original data: (X, Y)

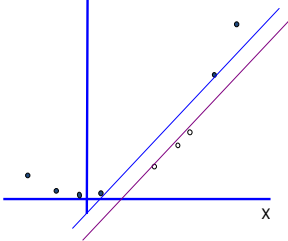
Transformed: (X, X^2, Y)

Think of X^2 as a new attribute, e.g., X'

203

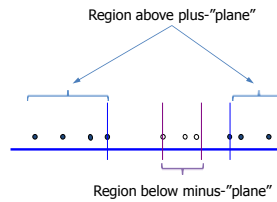
Now Separation Is Easy Again

$X' (= X^2)$



204

Corresponding "Planes" in Original Space



205

Common SVM Basis Functions

- Polynomial of attributes X_1, \dots, X_d of certain max degree, e.g., X_4^2
- Radial basis function
 - Symmetric around center, i.e., $\text{KernelFunction}(|\mathbf{X} - \mathbf{c}| / \text{kernelWidth})$
- Sigmoid function of \mathbf{X} , e.g., hyperbolic tangent
- Let $\Phi(\mathbf{x})$ be the transformed input record
 - Previous example: $\Phi(x) = (x, x^2)$

206

Quadratic Basis Functions

$$\Phi(\mathbf{x}) = \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \vdots \\ \sqrt{2}x_d \\ x_1^2 \\ x_2^2 \\ \vdots \\ x_d^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2}x_1x_3 \\ \vdots \\ \sqrt{2}x_1x_d \\ \sqrt{2}x_2x_3 \\ \vdots \\ \sqrt{2}x_1x_d \\ \vdots \\ \sqrt{2}x_{d-1}x_d \end{pmatrix}$$

Constant Term

Linear Terms

Pure Quadratic Terms

Quadratic Cross-Terms

Number of terms (assuming d input attributes):
 $(d+2)\text{-choose-2}$
 $= (d+2)(d+1)/2$
 $\approx d^2/2$

Why did we choose this specific transformation?

207

Dual QP With Basis Functions

$$\text{Maximize } \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k=1}^n \sum_{l=1}^n \alpha_k \alpha_l \cdot y(k) \cdot y(l) \cdot \Phi(\mathbf{x}(k)) \cdot \Phi(\mathbf{x}(l))$$

$$\text{Subject to these constraints: } \forall k : 0 \leq \alpha_k \leq C \quad \sum_{k=1}^n \alpha_k y(k) = 0$$

Then define:

$$\mathbf{w} = \sum_{k=1}^n \alpha_k \cdot y(k) \cdot \Phi(\mathbf{x}(k))$$

Then classify with:

$$\mathbf{f}(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)$$

$$b = \text{AVG}_{k:0 < \alpha_k < C} \left\{ \frac{1}{y(k)} - \Phi(\mathbf{x}(k)) \cdot \mathbf{w} \right\}$$

208

Computation Challenge

- Input vector \mathbf{x} has d components (its d attribute values)
- The transformed input vector $\Phi(\mathbf{x})$ has $d^2/2$ components
- Hence computing $\Phi(\mathbf{x}(k)) \cdot \Phi(\mathbf{x}(l))$ now costs order $d^2/2$ instead of order d operations (additions, multiplications)
- ...or is there a better way to do this?
 - Take advantage of properties of certain transformations

209

Quadratic Dot Products

$$\Phi(\mathbf{a}) \cdot \Phi(\mathbf{b}) = \begin{pmatrix} 1 \\ \sqrt{2}a_1 \\ \sqrt{2}a_2 \\ \vdots \\ \sqrt{2}a_d \\ a_1^2 \\ a_2^2 \\ \vdots \\ a_d^2 \\ \sqrt{2}a_1a_2 \\ \sqrt{2}a_1a_3 \\ \vdots \\ \sqrt{2}a_1a_d \\ \sqrt{2}a_2a_3 \\ \vdots \\ \sqrt{2}a_2a_d \\ \vdots \\ \sqrt{2}a_{d-1}a_d \end{pmatrix} \cdot \begin{pmatrix} 1 \\ \sqrt{2}b_1 \\ \sqrt{2}b_2 \\ \vdots \\ \sqrt{2}b_d \\ b_1^2 \\ b_2^2 \\ \vdots \\ b_d^2 \\ \sqrt{2}b_1b_2 \\ \sqrt{2}b_1b_3 \\ \vdots \\ \sqrt{2}b_1b_d \\ \sqrt{2}b_2b_3 \\ \vdots \\ \sqrt{2}b_2b_d \\ \vdots \\ \sqrt{2}b_{d-1}b_d \end{pmatrix} = \underbrace{1}_{1} + \underbrace{\sum_{i=1}^d 2a_i b_i}_{\sum_{i=1}^d 2a_i b_i} + \underbrace{\sum_{i=1}^d a_i^2 b_i^2}_{\sum_{i=1}^d a_i^2 b_i^2} + \underbrace{\sum_{i=1}^d \sum_{j=i+1}^d 2a_i a_j b_i b_j}_{\sum_{i=1}^d \sum_{j=i+1}^d 2a_i a_j b_i b_j}$$

210

Quadratic Dot Products

Now consider another function of \mathbf{a} and \mathbf{b} :

$$\begin{aligned} \Phi(\mathbf{a}) \cdot \Phi(\mathbf{b}) &= \\ 1 + 2 \sum_{i=1}^d a_i b_i + \sum_{i=1}^d a_i^2 b_i^2 + \sum_{i=1}^d \sum_{j=i+1}^d 2a_i a_j b_i b_j &= (\mathbf{a} \cdot \mathbf{b} + 1)^2 \\ &= (\mathbf{a} \cdot \mathbf{b})^2 + 2\mathbf{a} \cdot \mathbf{b} + 1 \\ &= \left(\sum_{i=1}^d a_i b_i \right)^2 + 2 \sum_{i=1}^d a_i b_i + 1 \\ &= \sum_{i=1}^d \sum_{j=1}^d a_i a_j b_i b_j + 2 \sum_{i=1}^d a_i b_i + 1 \\ &= \sum_{i=1}^d (a_i b_i)^2 + 2 \sum_{i=1}^d \sum_{j=i+1}^d a_i a_j b_i b_j + 2 \sum_{i=1}^d a_i b_i + 1 \end{aligned}$$

211

Quadratic Dot Products

- The results of $\Phi(\mathbf{a}) \cdot \Phi(\mathbf{b})$ and of $(\mathbf{a} \cdot \mathbf{b} + 1)^2$ are identical
- Computing $\Phi(\mathbf{a}) \cdot \Phi(\mathbf{b})$ costs about $d^2/2$, while computing $(\mathbf{a} \cdot \mathbf{b} + 1)^2$ costs only about $d+2$ operations
- This means that we can work in the high-dimensional space ($d^2/2$ dimensions) where the training records are more easily separable, but pay about the same cost as working in the original space (d dimensions)
- Savings are even greater when dealing with higher-degree polynomials, i.e., degree $q > 2$, that can be computed as $(\mathbf{a} \cdot \mathbf{b} + 1)^q$

212

Any Other Computation Problems?

$$\mathbf{w} = \sum_{k=1}^n \alpha_k \cdot y(k) \cdot \Phi(\mathbf{x}(k)) \quad b = \text{AVG}_{k:0 < \alpha_k < C} \left\{ \frac{1}{y(k)} - \Phi(\mathbf{x}(k)) \cdot \mathbf{w} \right\}$$

- What about computing \mathbf{w} ?
 - Finally need $\mathbf{f}(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)$:

$$\mathbf{w} \cdot \Phi(\mathbf{x}) = \sum_{k=1}^n \alpha_k \cdot y(k) \cdot \Phi(\mathbf{x}(k)) \cdot \Phi(\mathbf{x})$$
 - Can be computed using the same trick as before
- Can apply the same trick again to b , because

$$\Phi(\mathbf{x}(k)) \cdot \mathbf{w} = \sum_{j=1}^n \alpha_j \cdot y(j) \cdot \Phi(\mathbf{x}(k)) \cdot \Phi(\mathbf{x}(j))$$

213

SVM Kernel Functions

- For which transformations, called kernels, does the same trick work?
- Polynomial: $K(a,b)=(a \cdot b + 1)^q$
- Radial-Basis-style (RBF):

$$K(\mathbf{a}, \mathbf{b}) = \exp\left(-\frac{(\mathbf{a} - \mathbf{b})^2}{2\sigma^2}\right)$$

$q, \sigma, \kappa,$ and δ are magic parameters that must be chosen by a model selection method.

- Neural-net-style sigmoidal:

$$K(\mathbf{a}, \mathbf{b}) = \tanh(\kappa \cdot \mathbf{a} \cdot \mathbf{b} - \delta)$$

214

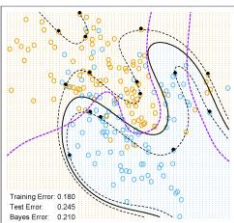
Overfitting

- With the right kernel function, computation in high dimensional transformed space is no problem
- But what about overfitting? There seem to be so many parameters...
- Usually not a problem, due to maximum margin approach
 - Only the support vectors determine the model, hence SVM complexity depends on number of support vectors, not dimensions (still, in higher dimensions there might be more support vectors)
 - Minimizing $\mathbf{w} \cdot \mathbf{w}$ discourages extremely large weights, which smoothes the function (recall weight decay for neural networks!)

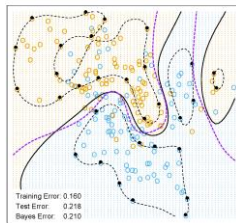
215

Different Kernels

SVM - Degree-4 Polynomial in Feature Space



SVM - Radial Kernel in Feature Space



Source: Hastie, Tibshirani, and Friedman. The Elements of Statistical Learning

216

Multi-Class Classification

- SVMs can only handle two-class outputs (i.e. a categorical output variable with arity 2).
- With output arity N , learn N SVM's
 - SVM 1 learns "Output==1" vs "Output != 1"
 - SVM 2 learns "Output==2" vs "Output != 2"
 - :
 - SVM N learns "Output== N " vs "Output != N "
- Predict with each SVM and find out which one puts the prediction the furthest into the positive region.

217

Why Is SVM Effective on High Dimensional Data?

- Complexity of trained classifier is characterized by the number of support vectors, not dimensionality of the data
- If all other training records are removed and training is repeated, the same separating hyperplane would be found
- The number of support vectors can be used to compute an upper bound on the expected error rate of the SVM, which is independent of data dimensionality
- Thus, an SVM with a small number of support vectors can have good generalization, even when the dimensionality of the data is high

218

SVM vs. Neural Network

- SVM
 - Relatively new concept
 - Deterministic algorithm
 - Nice Generalization properties
 - Hard to train – learned in batch mode using quadratic programming techniques
 - Using kernels can learn very complex functions
- Neural Network
 - Relatively old
 - Nondeterministic algorithm
 - Generalizes well but doesn't have strong mathematical foundation
 - Can easily be learned in incremental fashion
 - To learn complex functions—use multilayer perceptron (not that trivial)

219

Classification and Prediction Overview

- Introduction
- Decision Trees
- Statistical Decision Theory
- Nearest Neighbor
- Bayesian Classification
- Artificial Neural Networks
- Support Vector Machines (SVMs)
- Prediction
- Accuracy and Error Measures
- Ensemble Methods

221

What Is Prediction?

- Essentially the same as classification, but output is continuous, not discrete
 - Construct a model, then use model to predict continuous output value for a given input
- Major method for prediction: **regression**
 - Many variants of regression analysis in statistics literature; not covered in this class
- Neural network and k-NN can do regression “out-of-the-box”
- SVMs for regression exist
- What about trees?

222

Regression Trees and Model Trees

- Regression tree: proposed in CART system (Breiman et al. 1984)
 - CART: Classification And Regression Trees
 - Each leaf stores a continuous-valued prediction
 - Average output value for the training records in the leaf
- Model tree: proposed by Quinlan (1992)
 - Each leaf holds a regression model—a multivariate linear equation
- Training: like for classification trees, but uses variance instead of purity measure for selecting split predicates

223

Classification and Prediction Overview

- Introduction
- Decision Trees
- Statistical Decision Theory
- Nearest Neighbor
- Bayesian Classification
- Artificial Neural Networks
- Support Vector Machines (SVMs)
- Prediction
- Accuracy and Error Measures
- Ensemble Methods

224

Classifier Accuracy Measures

		Predicted class		total
		buy_computer = yes	buy_computer = no	
True class	buy_computer = yes	6954	46	7000
	buy_computer = no	412	2588	3000
total		7366	2634	10000

- Accuracy of a classifier M, $acc(M)$: percentage of test records that are correctly classified by M
 - Error rate (misclassification rate) of M = $1 - acc(M)$
 - Given m classes, $CM[i,j]$, an entry in a **confusion matrix**, indicates # of records in class i that are labeled by the classifier as class j

	C_1	C_2
C_1	True positive	False negative
C_2	False positive	True negative

225

Precision and Recall

- Precision: measure of exactness
 - $t\text{-pos} / (t\text{-pos} + f\text{-pos})$
- Recall: measure of completeness
 - $t\text{-pos} / (t\text{-pos} + f\text{-neg})$
- F-measure: combination of precision and recall
 - $2 * \text{precision} * \text{recall} / (\text{precision} + \text{recall})$
- Note: Accuracy = $(t\text{-pos} + t\text{-neg}) / (t\text{-pos} + t\text{-neg} + f\text{-pos} + f\text{-neg})$

226

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
- Always predicting the majority class defines the **baseline**
 - A good classifier should do better than baseline

227

Cost-Sensitive Measures: Cost Matrix

		PREDICTED CLASS	
		C(i j)	Class=Yes
ACTUAL CLASS	Class=Yes	C(Yes Yes)	C(No Yes)
	Class=No	C(Yes No)	C(No No)

$C(i|j)$: Cost of misclassifying class j example as class i

228

Computing Cost of Classification

Cost Matrix	PREDICTED CLASS		
	C(i j)	+	-
ACTUAL CLASS	+	-1	100
	-	1	0

Model M_1	PREDICTED CLASS		
ACTUAL CLASS	+	-	
	+	150	40
-	60	250	

Model M_2	PREDICTED CLASS		
ACTUAL CLASS	+	-	
	+	250	45
-	5	200	

Accuracy = 80%
Cost = 3910

Accuracy = 90%
Cost = 4255

229

Prediction Error Measures

- Continuous output: it matters how far off the prediction is from the true value
- Loss function**: distance between y and predicted value y'
 - Absolute error: $|y - y'|$
 - Squared error: $(y - y')^2$
- Test error (generalization error): average loss over the test set
- Mean absolute error: $\frac{1}{n} \sum_{i=1}^n |y(i) - y'(i)|$ Mean squared error: $\frac{1}{n} \sum_{i=1}^n (y(i) - y'(i))^2$
- Relative absolute error: $\frac{\sum_{i=1}^n |y(i) - y'(i)|}{\sum_{i=1}^n |y(i) - \bar{y}|}$ Relative squared error: $\frac{\sum_{i=1}^n (y(i) - y'(i))^2}{\sum_{i=1}^n (y(i) - \bar{y})^2}$
- Squared-error exaggerates the presence of outliers

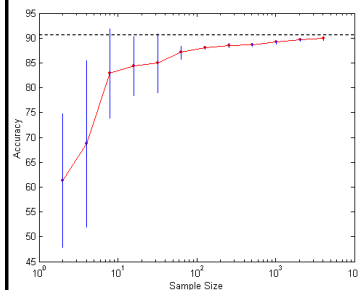
230

Evaluating a Classifier or Predictor

- Holdout method**
 - The given data set is randomly partitioned into two sets
 - Training set (e.g., 2/3) for model construction
 - Test set (e.g., 1/3) for accuracy estimation
 - Can repeat holdout multiple times
 - Accuracy = avg. of the accuracies obtained
- Cross-validation** (k-fold, where $k = 10$ is most popular)
 - Randomly partition data into k mutually exclusive subsets, each approximately equal size
 - In i -th iteration, use D_i as test set and others as training set
 - Leave-one-out: k folds where $k = \#$ of records
 - Expensive, often results in high variance of performance metric

231

Learning Curve



- Accuracy versus sample size
- Effect of small sample size:
 - Bias in estimate
 - Variance of estimate
- Helps determine how much training data is needed
 - Still need to have enough test and validation data to be representative of distribution

232

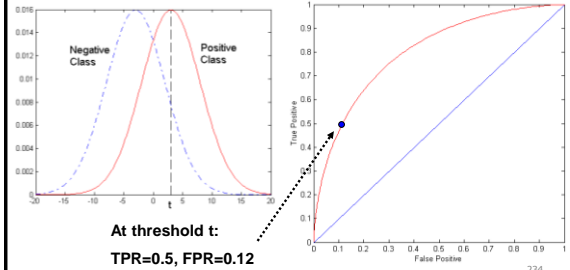
ROC (Receiver Operating Characteristic)

- Developed in 1950s for signal detection theory to analyze noisy signals
 - Characterizes trade-off between positive hits and false alarms
- ROC curve plots T-Pos rate (y-axis) against F-Pos rate (x-axis)
- Performance of each classifier is represented as a point on the ROC curve
 - Changing the threshold of the algorithm, sample distribution or cost matrix changes the location of the point

233

ROC Curve

- 1-dimensional data set containing 2 classes (positive and negative)
 - Any point located at $x > t$ is classified as positive

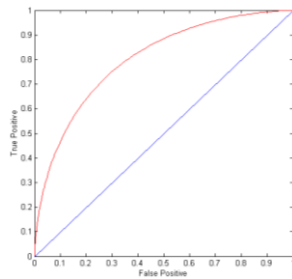


234

ROC Curve

(TPR, FPR):

- (0,0): declare everything to be negative class
- (1,1): declare everything to be positive class
- (1,0): ideal
- Diagonal line:
 - Random guessing



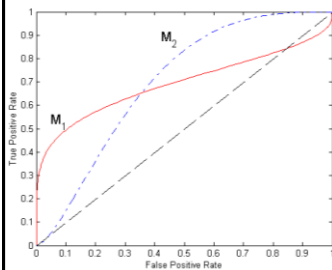
235

Diagonal Line for Random Guessing

- Classify a record as positive with fixed probability p , irrespective of attribute values
- Consider test set with a positive and b negative records
- True positives: $p*a$, hence true positive rate = $(p*a)/a = p$
- False positives: $p*b$, hence false positive rate = $(p*b)/b = p$
- For every value $0 \leq p \leq 1$, we get point (p,p) on ROC curve

236

Using ROC for Model Comparison



- Neither model consistently outperforms the other
 - M1 better for small FPR
 - M2 better for large FPR
- Area under the ROC curve
 - Ideal: area = 1
 - Random guess: area = 0.5

237

How to Construct an ROC curve

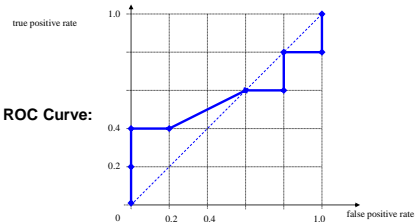
record	$P(+ \mathbf{x})$	True Class
1	0.95	+
2	0.93	+
3	0.87	-
4	0.85	-
5	0.85	-
6	0.85	+
7	0.76	-
8	0.53	+
9	0.43	-
10	0.25	+

- Use classifier that produces posterior probability $P(+|\mathbf{x})$ for each test record \mathbf{x}
- Sort records according to $P(+|\mathbf{x})$ in decreasing order
- Apply threshold at each unique value of $P(+|\mathbf{x})$
 - Count number of TP, FP, TN, FN at each threshold
 - TP rate, $TPR = TP/(TP+FN)$
 - FP rate, $FPR = FP/(FP+TN)$

238

How To Construct An ROC Curve

Class	+	-	+	-	+	-	+	-	+	-
TP	3	4	4	3	3	2	2	1	0	0
FP	5	5	4	4	3	1	0	0	0	0
TN	0	0	1	1	2	4	5	5	5	5
FN	0	5	1	2	2	3	2	4	5	5
TFR	1	0.8	0.8	0.6	0.6	0.4	0.4	0.2	0	0
FFR	1	1	0.8	0.8	0.6	0.2	0	0	0	0



239

Test of Significance

- Given two models:
 - Model M1: accuracy = 85%, tested on 30 instances
 - Model M2: accuracy = 75%, tested on 5000 instances
- Can we say M1 is better than M2?
 - How much confidence can we place on accuracy of M1 and M2?
 - Can the difference in accuracy be explained as a result of random fluctuations in the test set?

240

Confidence Interval for Accuracy

- Classification can be regarded as a Bernoulli trial
 - A Bernoulli trial has 2 possible outcomes, "correct" or "wrong" for classification
 - Collection of Bernoulli trials has a Binomial distribution
 - Probability of getting c correct predictions if model accuracy is p (=probability to get a single prediction right):

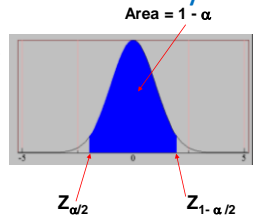
$$\binom{n}{c} p^c (1-p)^{n-c}$$
- Given c , or equivalently, $ACC = c/n$ and n (#test records), can we predict p , the **true accuracy** of the model?

241

Confidence Interval for Accuracy

- Binomial distribution for X ="number of correctly classified test records out of n "
 - $E(X)=pn$, $Var(X)=p(1-p)n$
 - Accuracy = X/n
 - $E(ACC) = p$, $Var(ACC) = p(1-p)/n$
 - For large test sets ($n>30$), Binomial distribution is closely approximated by normal distribution with same mean and variance
 - ACC has a normal distribution with mean= p , variance= $p(1-p)/n$
- $$P\left(Z_{\alpha/2} < \frac{ACC-p}{\sqrt{p(1-p)/n}} < Z_{1-\alpha/2}\right) = 1-\alpha$$
- Confidence Interval for p :

$$p = \frac{2n \cdot ACC + Z_{\alpha/2}^2 \pm \sqrt{Z_{\alpha/2}^2 + 4n \cdot ACC - 4n \cdot ACC^2}}{2(n + Z_{\alpha/2}^2)}$$



242

Confidence Interval for Accuracy

- Consider a model that produces an accuracy of 80% when evaluated on 100 test instances
 - $n = 100$, $ACC = 0.8$
 - Let $1-\alpha = 0.95$ (95% confidence)
 - From probability table, $Z_{\alpha/2} = 1.96$

N	50	100	500	1000	5000
p(lower)	0.670	0.711	0.763	0.774	0.789
p(upper)	0.888	0.866	0.833	0.824	0.811

1- α	Z
0.99	2.58
0.98	2.33
0.95	1.96
0.90	1.65

$$p = \frac{2n \cdot ACC + Z_{\alpha/2}^2 \pm \sqrt{Z_{\alpha/2}^2 + 4n \cdot ACC - 4n \cdot ACC^2}}{2(n + Z_{\alpha/2}^2)}$$

243

Comparing Performance of Two Models

- Given two models M1 and M2, which is better?
 - M1 is tested on D_1 (size= n_1), found error rate = e_1
 - M2 is tested on D_2 (size= n_2), found error rate = e_2
 - Assume D_1 and D_2 are independent
 - If n_1 and n_2 are sufficiently large, then

$$eR_1 \sim N(\mu_1, \sigma_1)$$

$$eR_2 \sim N(\mu_2, \sigma_2)$$
 - Estimate: $\hat{\mu}_i = e_i$ and $\hat{\sigma}_i^2 = \frac{e_i(1-e_i)}{n_i}$

244

Testing Significance of Accuracy Difference

- Consider random variable $d = \text{err}_1 - \text{err}_2$
 - Since $\text{err}_1, \text{err}_2$ are normally distributed, so is their difference
 - Hence $d \sim N(d_t, \sigma_t)$ where d_t is the true difference
- Estimator for d_t :
 - $E[d] = E[\text{err}_1 - \text{err}_2] = E[\text{err}_1] - E[\text{err}_2] \approx e_1 - e_2$
 - Since D_1 and D_2 are independent, variance adds up:

$$\hat{\sigma}_d^2 = \hat{\sigma}_1^2 + \hat{\sigma}_2^2 = \frac{e_1(1-e_1)}{n_1} + \frac{e_2(1-e_2)}{n_2}$$
 - At $(1-\alpha)$ confidence level, $d_t = E[d] \pm Z_{\alpha/2} \hat{\sigma}_d$

245

An Illustrative Example

- Given: M1: $n_1 = 30, e_1 = 0.15$
M2: $n_2 = 5000, e_2 = 0.25$
- $E[d] = |e_1 - e_2| = 0.1$
- 2-sided test: $d_t = 0$ versus $d_t \neq 0$

$$\hat{\sigma}_d^2 = \frac{0.15(1-0.15)}{30} + \frac{0.25(1-0.25)}{5000} = 0.0043$$
- At 95% confidence level, $Z_{\alpha/2} = 1.96$

$$d_t = 0.100 \pm 1.96 \sqrt{0.0043} = 0.100 \pm 0.128$$
- Interval contains zero, hence difference may not be statistically significant
- But: may reject null hypothesis ($d_t \neq 0$) at lower confidence level

246

Significance Test for K-Fold Cross-Validation

- Each learning algorithm produces k models:
 - L1 produces $M11, M12, \dots, M1k$
 - L2 produces $M21, M22, \dots, M2k$
 - Both models are tested on the same test sets D_1, D_2, \dots, D_k
 - For each test set, compute $d_j = e_{1,j} - e_{2,j}$
 - For large enough k , d_j is normally distributed with mean d_t and variance σ_t
 - Estimate:

$$\hat{\sigma}_d^2 = \frac{\sum_{j=1}^k (d_j - \bar{d})^2}{k(k-1)}$$
- t-distribution: get t coefficient $t_{1-\alpha, k-1}$ from table by looking up confidence level $(1-\alpha)$ and degrees of freedom $(k-1)$
- $$d_t = \bar{d} \pm t_{1-\alpha, k-1} \hat{\sigma}_d$$

247

Classification and Prediction Overview

- Introduction
- Decision Trees
- Statistical Decision Theory
- Nearest Neighbor
- Bayesian Classification
- Artificial Neural Networks
- Support Vector Machines (SVMs)
- Prediction
- Accuracy and Error Measures
- Ensemble Methods

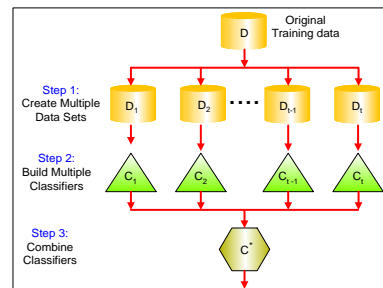
248

Ensemble Methods

- Construct a set of classifiers from the training data
- Predict class label of previously unseen records by aggregating predictions made by multiple classifiers

249

General Idea



250

Why Does It Work?

- Consider 2-class problem
- Suppose there are 25 base classifiers
 - Each classifier has error rate $\epsilon = 0.35$
 - Assume the classifiers are independent
- Return majority vote of the 25 classifiers
 - Probability that the ensemble classifier makes a wrong prediction:

$$\sum_{i=13}^{25} \binom{25}{i} \epsilon^i (1-\epsilon)^{25-i} = 0.06$$

251

Base Classifier vs. Ensemble Error

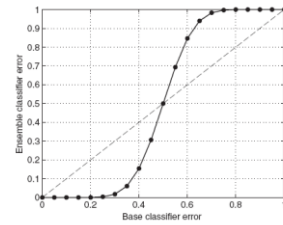


Figure 5.30. Comparison between errors of base classifiers and errors of the ensemble classifier.

252

Model Averaging and Bias-Variance Tradeoff

- Single model: lowering bias will usually increase variance
 - “Smoother” model has lower variance but might not model function well enough
- Ensembles can overcome this problem
 1. Let models overfit
 - Low bias, high variance
 2. Take care of the variance problem by averaging many of these models
- This is the basic idea behind **bagging**

253

Bagging: Bootstrap Aggregation

- Given training set with n records, sample n records randomly with replacement

Original Data	1	2	3	4	5	6	7	8	9	10
Bagging (Round 1)	7	8	10	8	2	5	10	10	5	9
Bagging (Round 2)	1	4	9	1	2	3	2	7	3	2
Bagging (Round 3)	1	8	5	10	5	5	9	6	3	7

- Train classifier for each bootstrap sample
- Note: each training record has probability $1 - (1 - 1/n)^n$ of being selected at least once in a sample of size n

254

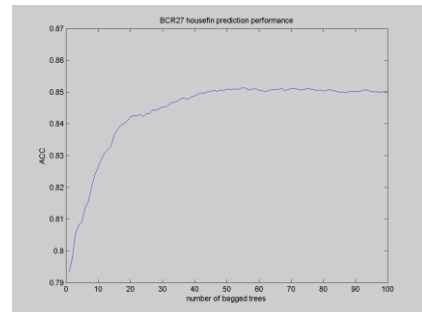
Bagged Trees

- Create k trees from training data
 - Bootstrap sample, grow large trees
- Design goal: independent models, high variability between models
- Ensemble prediction = average of individual tree predictions (or majority vote)
- Works the same way for other classifiers

$$(1/k) \cdot \text{Tree}_1 + (1/k) \cdot \text{Tree}_2 + \dots + (1/k) \cdot \text{Tree}_k$$

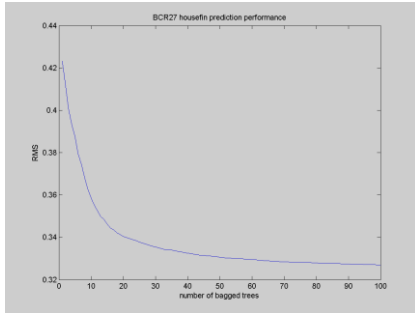
255

Typical Result



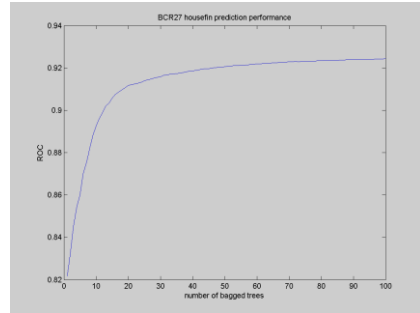
256

Typical Result



257

Typical Result



258

Bagging Challenges

- Ideal case: all models independent of each other
- Train on independent data samples
 - Problem: limited amount of training data
 - Training set needs to be representative of data distribution
 - Bootstrap sampling allows creation of many “almost” independent training sets
- Diversify models, because similar sample might result in similar tree
 - Random Forest: limit choice of split attributes to small random subset of attributes (new selection of subset for each node) when training tree
 - Use different model types in same ensemble: tree, ANN, SVM, regression models

259

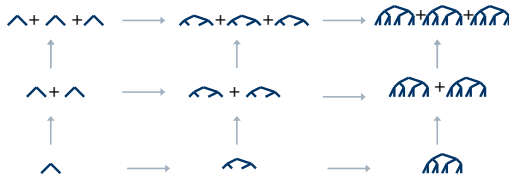
Additive Grove

- Ensemble technique for predicting continuous output
- Instead of individual trees, train additive models
 - Prediction of single Grove model = sum of tree predictions
- Prediction of ensemble = average of individual Grove predictions
- Combines large trees and additive models
 - Challenge: how to train the additive models without having the first trees fit the training data too well
 - Next tree is trained on residuals of previously trained trees in same Grove model
 - If previously trained trees capture training data too well, next tree is mostly trained on noise

$$(1/k) \cdot \left[\text{Tree}_1 + \dots + \text{Tree}_k \right] + (1/k) \cdot \left[\text{Tree}_1 + \dots + \text{Tree}_k \right] + \dots + (1/k) \cdot \left[\text{Tree}_1 + \dots + \text{Tree}_k \right]$$

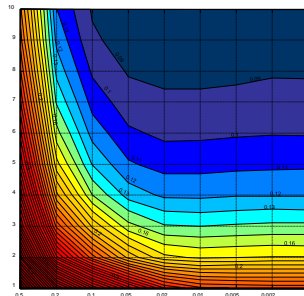
260

Training Groves



261

Typical Grove Performance



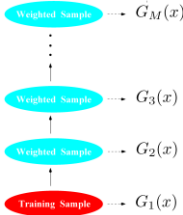
- Root mean squared error
 - Lower is better
- Horizontal axis: tree size
 - Fraction of training data when to stop splitting
- Vertical axis: number of trees in each single Grove model
- 100 bagging iterations

262

Boosting

FINAL CLASSIFIER

$$G(x) = \text{sign} \left[\sum_{m=1}^M \alpha_m G_m(x) \right]$$



- Iterative procedure to adaptively change distribution of training data by focusing more on previously misclassified records
- Initially, all n records are assigned equal weights
- Record weights may change at the end of each boosting round

263

Boosting

- Records that are wrongly classified will have their weights increased
- Records that are classified correctly will have their weights decreased

Original Data	1	2	3	4	5	6	7	8	9	10
Boosting (Round 1)	7	3	2	8	7	9	4	10	6	3
Boosting (Round 2)	5	4	9	4	2	5	1	7	4	2
Boosting (Round 3)	4	4	8	10	4	5	4	6	3	4

- Assume record 4 is hard to classify
- Its weight is increased, therefore it is more likely to be chosen again in subsequent rounds

264

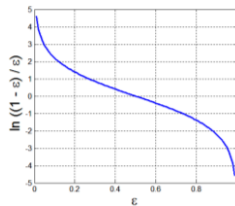
Example: AdaBoost

- Base classifiers: C_1, C_2, \dots, C_T
- Error rate (n training records, w_j are weights that sum to 1):

$$\epsilon_i = \sum_{j=1}^n w_j \delta(C_i(x_j) \neq y_j)$$

- Importance of a classifier:

$$\alpha_i = \ln \left(\frac{1 - \epsilon_i}{\epsilon_i} \right)$$



265

AdaBoost Details

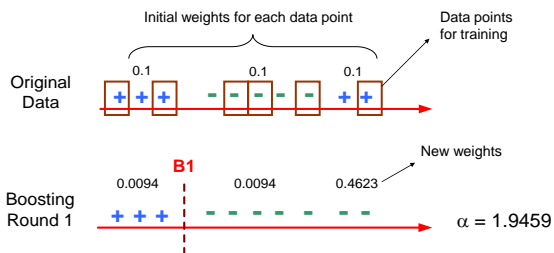
- Weight update:
$$w_j^{(i+1)} = \frac{w_j^{(i)}}{Z_i} \begin{cases} \epsilon_i & \text{if } C_i(x_j) = y_j \\ 1 - \epsilon_i & \text{if } C_i(x_j) \neq y_j \end{cases}$$
 where Z_i is the normalization factor

- Weights initialized to $1/n$
- Z_i ensures that weights add to 1
- If any intermediate rounds produce error rate higher than 50%, the weights are reverted back to $1/n$ and the resampling procedure is repeated
- Final classification:

$$C^*(x) = \arg \max_y \sum_{i=1}^T \alpha_i \delta(C_i(x) = y)$$

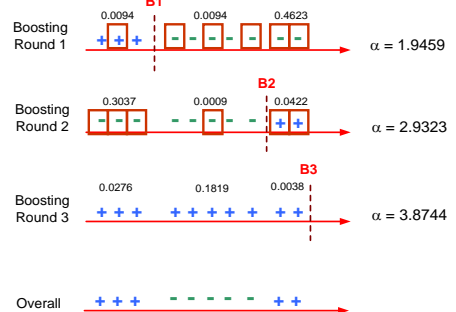
266

Illustrating AdaBoost



Note: The numbers appear to be wrong, but they convey the right idea... 267

Illustrating AdaBoost



Note: The numbers appear to be wrong, but they convey the right idea... 268

Bagging vs. Boosting

- Analogy
 - Bagging: diagnosis based on multiple doctors' majority vote
 - Boosting: weighted vote, based on doctors' previous diagnosis accuracy
- Sampling procedure
 - Bagging: records have same weight; easy to train in parallel
 - Boosting: weights record higher if model predicts it wrong; inherently sequential process
- Overfitting
 - Bagging robust against overfitting
 - Boosting susceptible to overfitting: make sure individual models do not overfit
- Accuracy usually significantly better than a single classifier
 - Best boosted model often better than best bagged model
- Additive Grove
 - Combines strengths of bagging and boosting (additive models)
 - Shown empirically to make better predictions on many data sets
 - Training more tricky, especially when data is very noisy

269

Classification/Prediction Summary

- Forms of data analysis that can be used to train models from data and then make predictions for new records
- Effective and scalable methods have been developed for decision tree induction, Naive Bayesian classification, Bayesian networks, rule-based classifiers, Backpropagation, Support Vector Machines (SVM), nearest neighbor classifiers, and many other classification methods
- Regression models are popular for prediction. Regression trees, model trees, and ANNs are also used for prediction.

270

Classification/Prediction Summary

- K-fold cross-validation is a popular method for accuracy estimation, but determining accuracy on large test set is equally accepted
 - If test sets are large enough, a significance test for finding the best model is not necessary
- Area under ROC curve and many other common performance measures exist
- Ensemble methods like bagging and boosting can be used to increase overall accuracy by learning and combining a series of individual models
 - Often state-of-the-art in prediction quality, but expensive to train, store, use
- No single method is superior over all others for all data sets
 - Issues such as accuracy, training and prediction time, robustness, interpretability, and scalability must be considered and can involve trade-offs

271