


Introduction to Data Mining

Chapter 8. Classification: Basic Concepts
Jiawei Han, Computer Science, Univ. Illinois at
Urbana-Champaign, 2017

1

Chapter 8. Classification: Basic Concepts

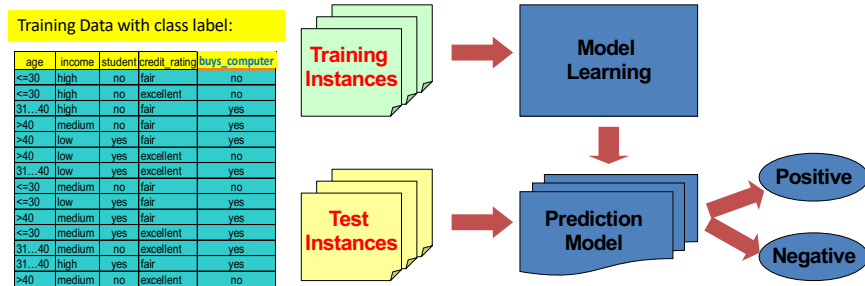
- Classification: Basic Concepts 
- Decision Tree Induction
- Bayes Classification Methods
- Linear Classifier
- Model Evaluation and Selection
- Techniques to Improve Classification Accuracy: Ensemble Methods
- Additional Concepts on Classification
- Summary

2

Supervised vs. Unsupervised Learning (1)

- Supervised learning (classification)

- Supervision: The training data such as observations or measurements are accompanied by **labels** indicating the classes which they belong to
- New data is classified based on the models built from the training set

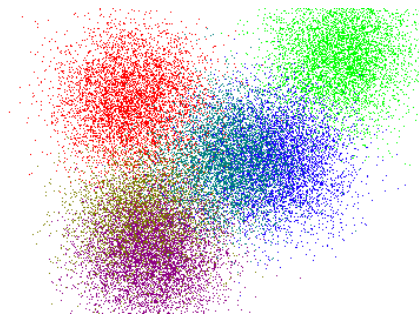


3

Supervised vs. Unsupervised Learning (2)

- Unsupervised learning (clustering)

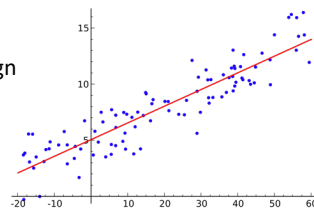
- The class labels of training data are unknown
- Given a set of observations or measurements, establish the possible existence of classes or clusters in the data



4

Prediction Problems: Classification vs. Numeric Prediction

- **Classification**
 - Predict categorical class labels (discrete or nominal)
 - Construct a model based on the training set and the **class labels** (the values in a classifying attribute) and use it in classifying new data
- **Numeric prediction**
 - Model continuous-valued functions (i.e., predict unknown or missing values)
- Typical applications of classification
 - Credit/loan approval
 - Medical diagnosis: if a tumor is cancerous or benign
 - Fraud detection: if a transaction is fraudulent
 - Web page categorization: which category it is



5

Classification—Model Construction, Validation and Testing

- **Model construction**
 - Each sample is assumed to belong to a predefined class (shown by the **class label**)
 - The set of samples used for model construction is **training set**
 - Model: Represented as decision trees, rules, mathematical formulas, or other forms
- **Model Validation and Testing:**
 - **Test:** Estimate accuracy of the model
 - The known label of test sample is compared with the classified result from the model
 - **Accuracy:** % of test set samples that are correctly classified by the model
 - Test set is independent of training set
 - **Validation:** If *the test set* is used to select or refine models, it is called **validation** (or development) (**test**) **set**
- **Model Deployment:** If the accuracy is acceptable, use the model to classify new data

6

Chapter 8. Classification: Basic Concepts

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



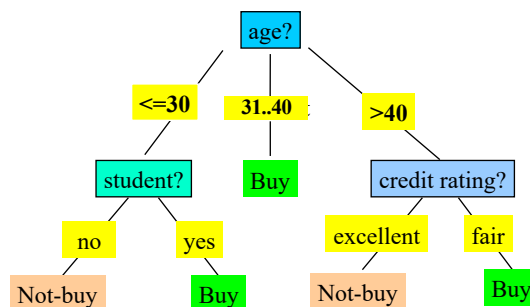
7

Decision Tree Induction: An Example

Decision tree construction:

- A top-down, recursive, divide-and-conquer process

Resulting tree:



Training data set: Who buys computer?

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Note: The data set is adapted from "Playing Tennis" example of R. Quinlan

8

From Entropy to Info Gain: A Review of Entropy

- Entropy (Information Theory)
 - A measure of uncertainty associated with a random number
 - Calculation: For a discrete random variable Y taking m distinct values $\{y_1, y_2, \dots, y_m\}$

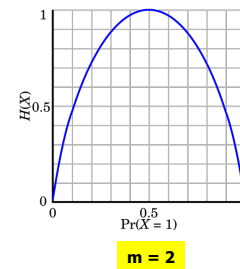
$$H(Y) = - \sum_{i=1}^m p_i \log(p_i) \text{ where } p_i = P(Y = y_i)$$

- Interpretation

- Higher entropy → higher uncertainty
- Lower entropy → lower uncertainty

- Conditional entropy

$$H(Y|X) = \sum_x p(x) H(Y|X = x)$$



9

Information Gain: An Attribute Selection Measure

- Select the attribute with the highest information gain (used in typical decision tree induction algorithm: ID3/C4.5)
- Let p_i be the probability that an arbitrary tuple in D belongs to class C_i , estimated by $|C_{i,D}|/|D|$
- Expected information (entropy) needed to classify a tuple in D :

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

- Information needed (after using A to split D into v partitions) to classify D :

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

- Information gained by branching on attribute A

$$Gain(A) = Info(D) - Info_A(D)$$

I : the expected information needed to classify a given sample

E (entropy) : expected information based on the partitioning into subsets by A

10

Example: Attribute Selection with Information Gain

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

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

age	p_i	n_i	$I(p_i, n_i)$
<=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	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

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

$\frac{5}{14} I(2,3)$ means "age <=30" has 5 out of 14 samples, with 2 yes'es and 3 no's.

Hence

$$Gain(age) = Info(D) - Info_{age}(D) = 0.246$$

Similarly, we can get

$$Gain(income) = 0.029$$

$$Gain(student) = 0.151$$

$$Gain(credit_rating) = 0.048$$

11

Decision Tree Induction: Algorithm

- Basic algorithm
 - Tree is constructed in a **top-down, recursive, divide-and-conquer manner**
 - At start, all the training examples are at the root
 - Examples are partitioned recursively based on selected attributes
 - On each node, attributes are selected based on the training examples on that node, and a heuristic or statistical measure (e.g., **information gain**)
- Conditions for stopping partitioning
 - All samples for a given node belong to the same class
 - There are no remaining attributes for further partitioning
 - There are no samples left
- Prediction
 - Majority voting** is employed for classifying the leaf

12

How to Handle Continuous-Valued Attributes?

- Method 1: Discretize continuous values and treat them as categorical values
 - E.g., age: < 20, 20..30, 30..40, 40..50, > 50
- Method 2: Determine the **best split point** for continuous-valued attribute A
 - Sort the value A in increasing order; e.g. 15, 18, 21, 22, 24, 25, 29, 31, ...
 - *Possible split point*: the midpoint between *each pair of adjacent values*
 - $(a_i + a_{i+1})/2$ is the midpoint between the values of a_i and a_{i+1}
 - e.g., $(15+18)/2 = 16.5$, 19.5, 21.5, 23, 24.5, 27, 30, ...
 - The point with the *maximum information gain* for A is selected as the **split-point** for A
- Split: Based on split point P
 - The set of tuples in D satisfying $A \leq P$ vs. those with $A > P$

13

Gain Ratio: A Refined Measure for Attribute Selection

- Information gain measure is biased towards attributes with a large number of values
- Gain ratio: Overcomes the problem (as a normalization to information gain)

$$SplitInfo_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left(\frac{|D_j|}{|D|} \right)$$
 - $GainRatio(A) = Gain(A)/SplitInfo(A)$
 - The attribute with the maximum gain ratio is selected as the splitting attribute
 - Gain ratio is used in a popular algorithm C4.5 (a successor of ID3) by R. Quinlan
- Example
 - $SplitInfo_{income}(D) = -\frac{4}{14} \log_2 \frac{4}{14} - \frac{6}{14} \log_2 \frac{6}{14} - \frac{4}{14} \log_2 \frac{4}{14} = 1.557$
 - $GainRatio(income) = 0.029/1.557 = 0.019$

14

Another Measure: Gini Index

- Gini index: Used in CART, and also in IBM IntelligentMiner
- If a data set D contains examples from n classes, gini index, $gini(D)$ is defined as
 - $gini(D) = 1 - \sum_{j=1}^n p_j^2$
 - p_j is the relative frequency of class j in D
- If a data set D is split on A into two subsets D_1 and D_2 , the gini index $gini(D)$ is defined as
 - $gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$
- Reduction in Impurity:
 - $\Delta gini(A) = gini(D) - gini_A(D)$
- The attribute provides the smallest $gini_{split}(D)$ (or the largest reduction in impurity) is chosen to split the node (*need to enumerate all the possible splitting points for each attribute*)

15

Computation of Gini Index

- Example: D has 9 tuples in buys_computer = “yes” and 5 in “no”

$$gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

- Suppose the attribute income partitions D into 10 in D_1 : {low, medium} and 4 in D_2
 - $gini_{income \in \{low, medium\}}(D) = \frac{10}{14} gini(D_1) + \frac{4}{14} gini(D_2)$
 - $= \frac{10}{14} \left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) + \frac{4}{14} \left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2\right) = 0.443 = gini_{income \in \{high\}}(D)$
 - $Gini_{\{low, high\}}$ is 0.458; $Gini_{\{medium, high\}}$ is 0.450
 - Thus, split on the {low, medium} (and {high}) since it has the lowest Gini index

- All attributes are assumed continuous-valued
- May need other tools, e.g., clustering, to get the possible split values
- Can be modified for categorical attributes

16

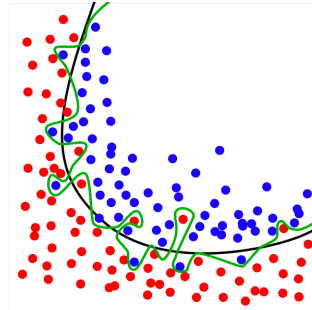
Other Attribute Selection Measures

- Minimal Description Length (MDL) principle
 - Philosophy: The simplest solution is preferred
 - The best tree as the one that requires the fewest # of bits to both (1) encode the tree, and (2) encode the exceptions to the tree
- CHAID: a popular decision tree algorithm, measure based on χ^2 test for independence
- Multivariate splits (partition based on multiple variable combinations)
 - CART: finds multivariate splits based on a linear combination of attributes
- There are many other measures proposed in research and applications
 - E.g., G-statistics, C-SEP
- Which attribute selection measure is the best?
 - Most give good results, none is significantly superior than others

17


Overfitting and Tree Pruning

- Overfitting: An induced tree may overfit the training data
 - Too many branches, some may reflect anomalies due to noise or outliers
 - Poor accuracy for unseen samples
- Two approaches to avoid 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
 - Postpruning: *Remove branches* from a “fully grown” tree—get a sequence of progressively pruned trees
 - Use a set of data different from the training data to decide which is the “best pruned tree”



18

Chapter 8. Classification: Basic Concepts

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

19

What Is Bayesian Classification?

- A statistical classifier
 - Perform *probabilistic prediction* (i.e., predict class membership probabilities)
- Foundation—Based on Bayes' Theorem
- Performance
 - A simple Bayesian classifier, *naïve Bayesian classifier*, has comparable performance with decision tree and selected neural network classifiers
- Incremental
 - Each training example can incrementally increase/decrease the probability that a hypothesis is correct—prior knowledge can be combined with observed data
- Theoretical Standard
 - Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured

20

Bayes' Theorem: Basics

- Total probability Theorem: $P(B) = \sum_{i=1}^M P(B|A_i)P(A_i)$
- Bayes' Theorem: $P(H|\mathbf{X}) = \frac{P(\mathbf{X}|H)P(H)}{P(\mathbf{X})} = P(\mathbf{X}|H) \times P(H) / P(\mathbf{X})$
 - Let \mathbf{X} be a data sample ("evidence"): class label is unknown
 - Let H be a *hypothesis* that X belongs to class C
 - Classification is to determine $P(H|\mathbf{X})$, (i.e., *posteriori probability*): the probability that the hypothesis holds given the observed data sample \mathbf{X}
 - $P(H)$ (*prior probability*): the initial probability
 - E.g., \mathbf{X} will buy computer, regardless of age, income, ...
 - $P(\mathbf{X})$: probability that sample data is observed
 - $P(\mathbf{X}|H)$ (*likelihood*): the probability of observing the sample \mathbf{X} , given that the hypothesis holds
 - E.g., Given that \mathbf{X} will buy computer, the prob. that X is 31..40, medium income

21

Bayes' Theorem: Basics

- Total probability Theorem:

$$p(B) = \sum_i p(B|A_i)p(A_i)$$

- Bayes' Theorem:

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

What we should choose posteriori probability likelihood prior probability

What we just see What we knew previously

- \mathbf{X} : a data sample ("evidence")
- H : X belongs to class C

Prediction can be done based on Bayes' Theorem:

Classification is to derive the maximum posteriori

22

Bayes' Theorem Example: Picnic Day

- ❑ The morning is cloudy ☹
- ❑ What is the chance of rain? $P(\text{Rain} \mid \text{Cloud}) = ?$
- ❑ 50% of all rainy days start off cloudy. $P(\text{Cloud} \mid \text{Rain}) = 50\%$
- ❑ Cloudy mornings are common (40% of days start cloudy) $P(\text{Cloud}) = 40\%$
- ❑ This is usually a dry month (only 3 of 30 days tend to be rainy) $P(\text{Rain}) = 10\%$

$$P(\text{Rain} \mid \text{Cloud}) = P(\text{Rain}) P(\text{Cloud} \mid \text{Rain}) / P(\text{Cloud}) = 10\% * 50\% / 40\% = 12.5\%$$
- ❑ The chance of rain is probably not as high as expected ☺
- ❑ Bayes' Theorem allows us to tell back and forth between posterior and likelihood (e.g., $P(\text{Rain} \mid \text{Cloud})$ and $P(\text{Cloud} \mid \text{Rain})$), tests the reality, which is the most important trick in Bayesian Inference

23

Naïve Bayes Classifier: Making a Naïve Assumption

- Practical difficulty of Naïve Bayes inference: It requires initial knowledge of many probabilities, which may not be available or involving significant computational cost
- A Naïve Special Case
 - Make an additional assumption to simplify the model, but achieve comparable performance.

attributes are conditionally independent
(i.e., no dependence relation between
attributes)

$$p(X|C_i) = \prod_k p(x_k|C_i) = p(x_1|C_i) \cdot p(x_2|C_i) \cdots p(x_n|C_i)$$

- Only need to count the class distribution w.r.t. features

24

Naïve Bayes Classifier: Training Dataset

Class:

C1:buys_computer = 'yes'

C2:buys_computer = 'no'

Data to be classified:

X = (age <=30, Income = medium,
Student = yes, Credit_rating = Fair)

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

25

Naïve Bayes Classifier: An Example

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

- $P(C_i)$: $P(\text{buys_computer} = \text{"yes"}) = 9/14 = 0.643$
 $P(\text{buys_computer} = \text{"no"}) = 5/14 = 0.357$
- Compute $P(X|C_i)$ for each class
 $P(\text{age} = \text{"<=30"} | \text{buys_computer} = \text{"yes"}) = 2/9 = 0.222$
 $P(\text{age} = \text{"<=30"} | \text{buys_computer} = \text{"no"}) = 3/5 = 0.6$
 $P(\text{income} = \text{"medium"} | \text{buys_computer} = \text{"yes"}) = 4/9 = 0.444$
 $P(\text{income} = \text{"medium"} | \text{buys_computer} = \text{"no"}) = 2/5 = 0.4$
 $P(\text{student} = \text{"yes"} | \text{buys_computer} = \text{"yes"}) = 6/9 = 0.667$
 $P(\text{student} = \text{"yes"} | \text{buys_computer} = \text{"no"}) = 1/5 = 0.2$
 $P(\text{credit_rating} = \text{"fair"} | \text{buys_computer} = \text{"yes"}) = 6/9 = 0.667$
 $P(\text{credit_rating} = \text{"fair"} | \text{buys_computer} = \text{"no"}) = 2/5 = 0.4$
- $X = (\text{age} \leq 30, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit_rating} = \text{fair})$
 $P(X|C_1) : P(X | \text{buys_computer} = \text{"yes"}) = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$
 $P(X|C_2) : P(X | \text{buys_computer} = \text{"no"}) = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$
 $P(C_1|X) = P(X|C_1) \cdot P(C_1) : P(X | \text{buys_computer} = \text{"yes"}) \cdot P(\text{buys_computer} = \text{"yes"}) = 0.028$
 $P(C_2|X) = P(X|C_2) \cdot P(C_2) : P(X | \text{buys_computer} = \text{"no"}) \cdot P(\text{buys_computer} = \text{"no"}) = 0.007$
Therefore, X belongs to class ("buys_computer = yes")

26

Avoiding the Zero-Probability Problem

- Naïve Bayesian prediction requires each conditional probability be **non-zero**

- Otherwise, the predicted probability will be zero

$$p(X|C_i) = \prod_k p(x_k|C_i) = p(x_1|C_i) \cdot p(x_2|C_i) \cdots p(x_n|C_i)$$

- Example. Suppose a dataset with 1000 tuples:

income = low (0), income = medium (990), and income = high (10)

- Use **Laplacian correction** (or Laplacian estimator)

- Adding 1 to each case

$$\text{Prob}(\text{income} = \text{low}) = 1/(1000 + 3)$$

$$\text{Prob}(\text{income} = \text{medium}) = (990 + 1)/(1000 + 3)$$

$$\text{Prob}(\text{income} = \text{high}) = (10 + 1)/(1000 + 3)$$

- The “corrected” probability estimates are close to their “uncorrected” counterparts


27

Naïve Bayes Classifier: Strength vs. Weakness

- Strength
 - Easy to implement
 - Good results obtained in most of the cases
- Weakness
 - Assumption: attributes conditional independence, therefore loss of accuracy
 - Practically, dependencies exist among variables
 - E.g., Patients: Profile: age, family history, etc.
Symptoms: fever, cough etc.
Disease: lung cancer, diabetes, etc.
 - Dependencies among these cannot be modeled by Naïve Bayes Classifier
- How to deal with these dependencies?
 - Use Bayesian Belief Networks (to be covered in the next chapter)

28

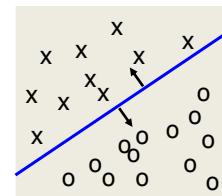
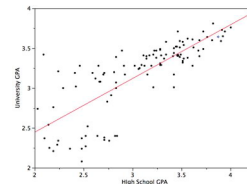
Chapter 8. Classification: Basic Concepts

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

29

Linear Regression vs. Linear Classifier

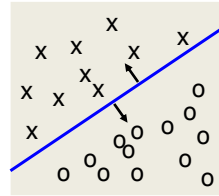
- **Linear regression**
 - Data modeled to fit a straight line
 - *Linear equation: $Y = wX + b$*
 - Often uses the least-square method to fit the line
 - Used to **predict continuous values**
- **Linear Classifier**
 - Built a classification model using a straight line
 - Used for **(categorical data) binary classification**



30

Linear Classifier: General Ideas

- Binary Classification
- $f(x)$ is a linear function based on the example's attribute values
 - The prediction is based on the value of $f(x)$
 - Data above the blue line belongs to class 'x' (i.e., $f(x) > 0$)
 - Data below blue line belongs to class 'o' (i.e., $f(x) < 0$)
- Classical Linear Classifiers
 - Linear Discriminant Analysis (LDA) (not covered)
 - Logistic Regression
 - Perceptron (later)
 - SVM (later)



31

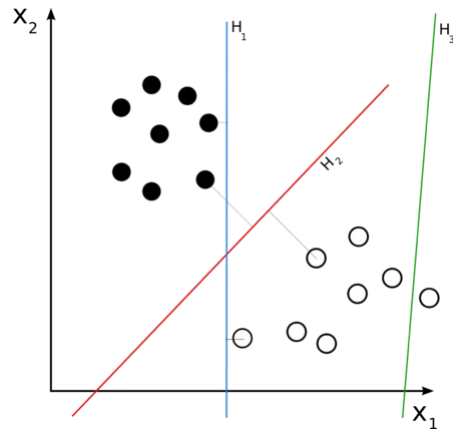
Linear Classifier: An Example

- A toy rule to determine whether a faculty member has tenure
 - $\text{Year} \geq 6 \text{ or Title} = \text{"Professor"} \Leftrightarrow \text{Tenure}$
- How to express the rule as a linear classifier?
- Features
 - x_1 ($x_1 \geq 0$) is an integer denoting the year
 - x_2 is a Boolean denoting whether the title is "Professor"
- A feasible linear classifier: $f(x) = (x_1 - 5) + 6 \cdot x_2$
 - When x_2 is True, because $x_1 \geq 0$, $f(x)$ is always greater than 0
 - When x_2 is False, because $f(x) > 0 \Leftrightarrow x_1 \geq 6$
- There are many more feasible classifiers
 - $f(x) = (x_1 - 5.5) + 6 \cdot x_2$
 - $f(x) = 2 \cdot (x_1 - 5) + 11 \cdot x_2$
 -

32

Key Question: Which Line Is Better?

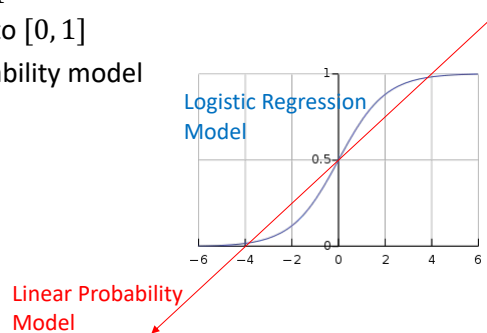
- There might be many feasible linear functions
 - Both H1 and H2 will work
- Which one is better?
 - H2 looks “better” in the sense that it is also furthest from both groups
 - We will introduce more in the SVM section



33

Logistic Regression: General Ideas

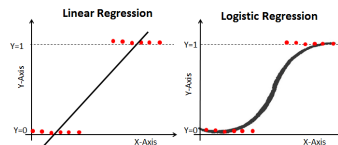
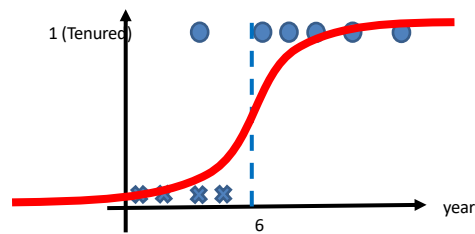
- Key Idea: Turns linear predictions into probabilities
- Sigmoid function:
 - $S(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{e^x+1}$
 - Projects $(-\infty, +\infty)$ to $[0, 1]$
- Compare to linear probability model
 - More smooth



34

Logistic Regression: An Example

- Suppose we only consider the year as feature



35

Logistic Regression: Maximum Likelihood

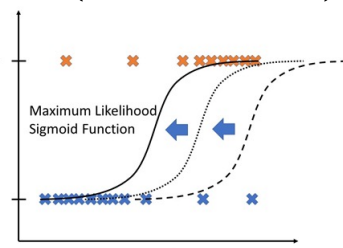
- How to select the correct sigmoid function that best fits the training data?
- The answer is maximum likelihood.
- The prediction function to learn
 - $p(Y = 1 | X = x; \mathbf{w}) = S(w_0 + \sum_{i=1}^n w_i \cdot x_i)$
 - $\mathbf{w} = (w_0, w_1, w_2, \dots, w_n)$ are the parameters of the model

- Maximum Likelihood

- Log likelihood:

$$l(\mathbf{w}) = \sum_{i=1}^N y_i \log p(Y = 1 | X = x_i; \mathbf{w}) + (1 - y_i) \log(1 - p(Y = 1 | X = x_i; \mathbf{w}))$$

- There's no close form solution
 - Gradient Descent
 - Update \mathbf{w} based on training data.



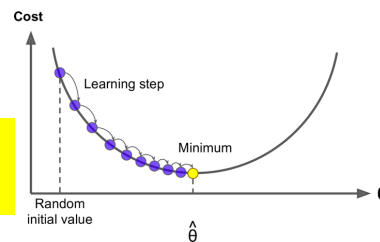
36

Gradient Descent

- Gradient Descent is an iterative optimization algorithm for finding the minimum of a function (e.g., the negative log likelihood)
- For a function $F(x)$ at a point \mathbf{a} , $F(x)$ *decreases fastest* if we go in the direction of the negative gradient of \mathbf{a}


$$\mathbf{a}_{n+1} = \mathbf{a}_n - \gamma \nabla F(\mathbf{a}_n)$$

When the gradient is zero, we arrive at the local minimum



37

Chapter 8. Classification: Basic Concepts

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

38

Model Evaluation and Selection

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

39

Classifier Evaluation Metrics: Confusion Matrix

- **Confusion Matrix:**

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

- In a confusion matrix w. m classes, $CM_{i,j}$ indicates # of tuples in class i that were labeled by the classifier as class j
 - May have extra rows/columns to provide totals
- **Example of Confusion Matrix:**

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

40

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

A\P	C	¬C	
C	TP	FN	P
¬C	FP	TN	N
	P'	N'	All

Predictions

Real-world truth

- **Classifier accuracy**, or recognition rate
 - Percentage of test set tuples that are correctly classified
 - **Accuracy** = $(TP + TN)/All$
- **Error rate**: $1 - accuracy$, or
- **Error rate** = $(FP + FN)/All$

Class imbalance problem

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

41

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

A\P	C	¬C	
C	TP	FN	P
¬C	FP	TN	N
	P'	N'	All

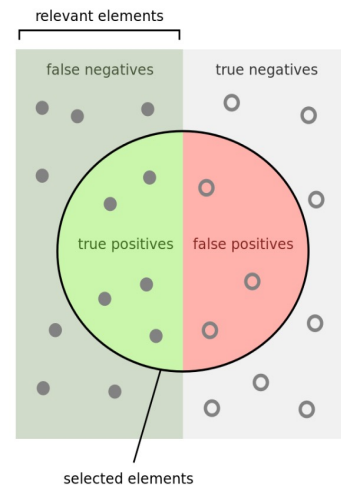
- **Precision**: Exactness: what % of tuples that the classifier labeled as positive are actually positive?

$$P = \text{Precision} = \frac{TP}{TP + FP} \quad \text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

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

$$R = \text{Recall} = \frac{TP}{TP + FN} \quad \text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

- Range: $[0, 1]$



https://en.wikipedia.org/wiki/Precision_and_recall

42

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

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

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

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

- **F1-measure (balanced F-measure)**
 - » That is, when $\beta = 1$,

$$F_1 = \frac{2PR}{P + R}$$

43

Classifier Evaluation Metrics: Example

- Use the same confusion matrix, calculate the measure just introduced

Actual Class\Predicted class	cancer = yes	cancer = no	Total	Recognition(%)
cancer = yes	90	210	300	30.00 (sensitivity)
cancer = no	140	9560	9700	98.56 (specificity)
Total	230	9770	10000	96.50 (accuracy)

- Sensitivity = $TP/P = 90/300 = 30\%$
- Specificity = $TN/N = 9560/9700 = 98.56\%$
- Accuracy = $(TP + TN)/All = (90+9560)/10000 = 96.50\%$
- Error rate = $(FP + FN)/All = (140 + 210)/10000 = 3.50\%$
- Precision = $TP/(TP + FP) = 90/(90 + 140) = 90/230 = 39.13\%$
- Recall = $TP/(TP + FN) = 90/(90 + 210) = 90/300 = 30.00\%$
- $F1 = 2 \times P \times R / (P + R) = 2 \times 39.13\% \times 30.00\% / (39.13\% + 30\%) = 33.96\%$

44

Classifier Evaluation: Holdout

- **Holdout method**
 - Given data is randomly partitioned into two independent sets
 - Training set (e.g., 2/3) for model construction
 - Test set (e.g., 1/3) for accuracy estimation
 - Repeated random sub-sampling validation: a variation of holdout
 - Repeat holdout k times, accuracy = avg. of the accuracies obtained

45

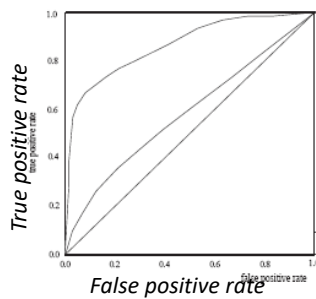
Classifier Evaluation: Cross-Validation

- **Cross-validation** (k -fold, where $k = 10$ is most popular)
 - Randomly partition the data into k *mutually exclusive* subsets, each approximately equal size
 - At i -th iteration, use D_i as test set and others as training set
 - Leave-one-out: k folds where $k = \#$ of tuples, for small sized data
 - ***Stratified cross-validation***: folds are stratified so that class distribution, in each fold is approximately the same as that in the initial data

46

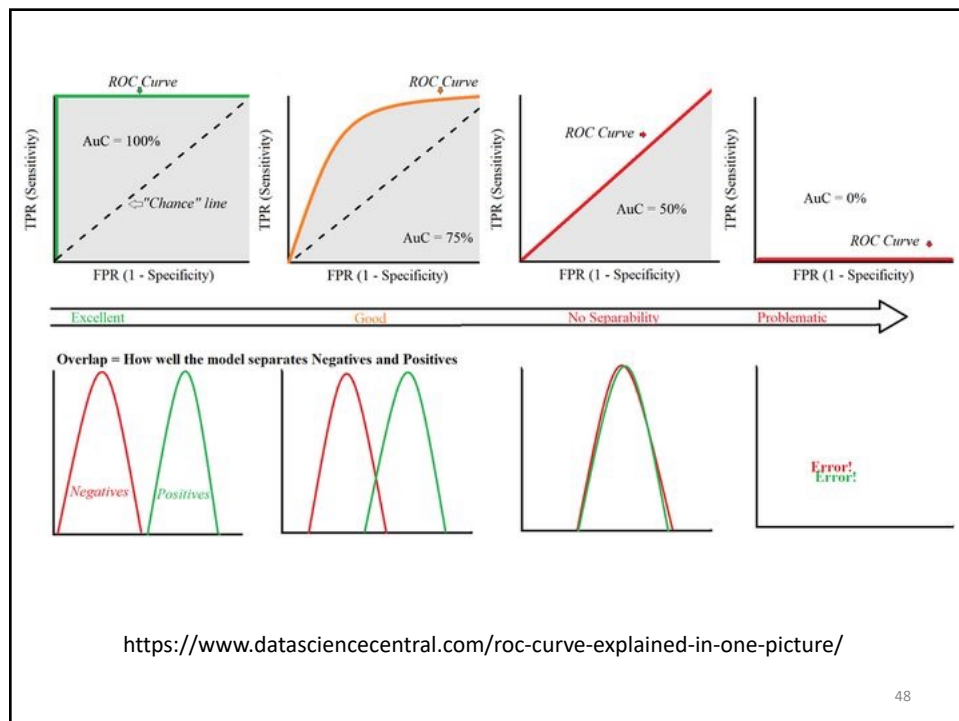
Model Selection: ROC Curves

- **ROC** (Receiver Operating Characteristics) curves: for visual comparison of classification models
- Originated from signal detection theory
- Shows the trade-off between the TP rate and the FP rate
- The area under the ROC curve (**AUC**: Area Under Curve) is a measure of the accuracy of the model
- Rank the test tuples in decreasing order: the one that is most likely to belong to the positive class appears at the top of the list
- The closer to the diagonal line (i.e., the closer the area is to 0.5), the less accurate is the model



- Vertical axis represents true positive rate (TP/P)
- Horizontal axis rep. false positive rate (FP/N)
- The plot also shows a diagonal line
- A model with perfect accuracy will have AUC = 1.0

47



48

48

Classifier Evaluation: Bootstrap

- **Bootstrap** is a statistical technique used to estimate the accuracy of a classifier (or any model) by resampling the data with replacement. It helps assess the performance of a model without needing a separate test set, and is particularly useful when you have limited data. Here's how it works for estimating a classifier's accuracy:
- **How Bootstrap Works**
- 1. **Create Resampled Datasets:**
 1. From the original dataset of size NNN, create BBB bootstrap samples by randomly sampling NNN instances **with replacement**.
 2. Each bootstrap sample will have some data points repeated and others omitted.
- 2. **Train and Test:**
 1. Train the classifier on each bootstrap sample.
 2. Test the classifier on the data points that were **not included** in the sample (called the "out-of-bag" data).
- 3. **Evaluate Metrics:**
 1. Calculate the performance metric (e.g., accuracy) on the out-of-bag data for each bootstrap iteration.
- 4. **Aggregate Results:**
 1. Compute the average and standard deviation of the performance metric across all bootstrap samples to estimate accuracy and variability.


49

Classifier Evaluation: Bootstrap

- **Example of Bootstrap Process**
- For a dataset with 100 instances:
 1. Randomly sample 100 instances with replacement to create a bootstrap sample.
Example: Original dataset = [1,2,3,...,100]; bootstrap sample = [2,5,3,5,97,...].
 2. Train the classifier on this sample.
 3. Test it on the out-of-bag data (instances not selected in the sample).
 4. Repeat the process B times (e.g., B=1000).
- **Limitations of Bootstrap**
 - **Overfitting:** Training on overlapping samples may lead to an optimistic bias if the model overfits.
 - **Computational Cost:** Requires training and testing the model many times, which can be expensive for complex models or large datasets.

50

Chapter 8. Classification: Basic Concepts

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

51

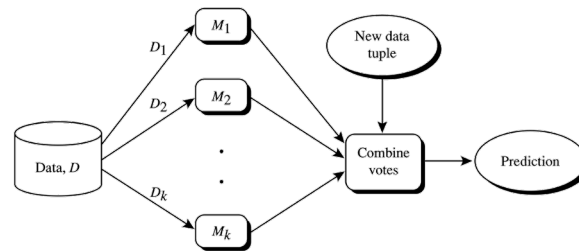
Techniques to Improve Classification Accuracy

- Introducing Ensemble Methods
- Bagging
- Boosting
- Random Forests
- Imbalanced Classification

52

Ensemble Methods: Increasing the Accuracy

- Ensemble methods
 - Use a combination of models to increase accuracy
 - Combine a series of k learned models, M_1, M_2, \dots, M_k , with the aim of creating an **improved** model M^*



53

Ensemble Methods: Increasing the Accuracy

- What are the requirements to generate an improved model?
 - Example: majority voting

		X1	X2	X3
Base model performance	M1	✓	✓	X
	M2	X	✓	✓
	M3	✓	X	✓
Ensemble performance	Voting Ensemble	✓	✓	✓
		Case 1: Ensemble has positive effect		

		X1	X2	X3
Base model performance	M1	✓	✓	X
	M2	✓	✓	X
	M3	✓	✓	X
Ensemble performance	Voting Ensemble	✓	✓	X
		Case 2: Ensemble has no effect		

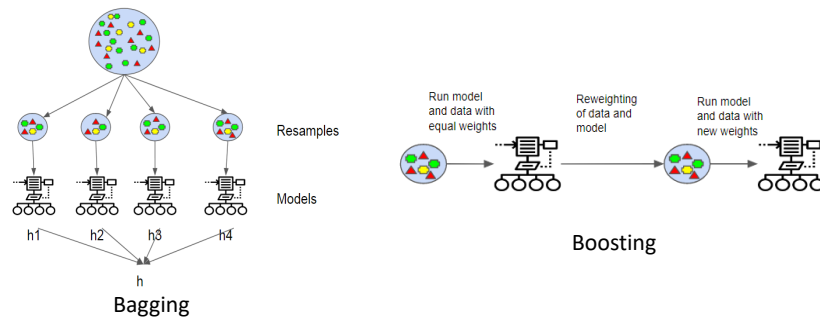
		X1	X2	X3
Base model performance	M1	✓	X	X
	M2	X	✓	X
	M3	X	X	✓
Ensemble performance	Voting Ensemble	X	X	X
		Case 3: Ensemble has negative effect		

- ❑ Base models should be
 - ❑ Accurate
 - ❑ Diverse

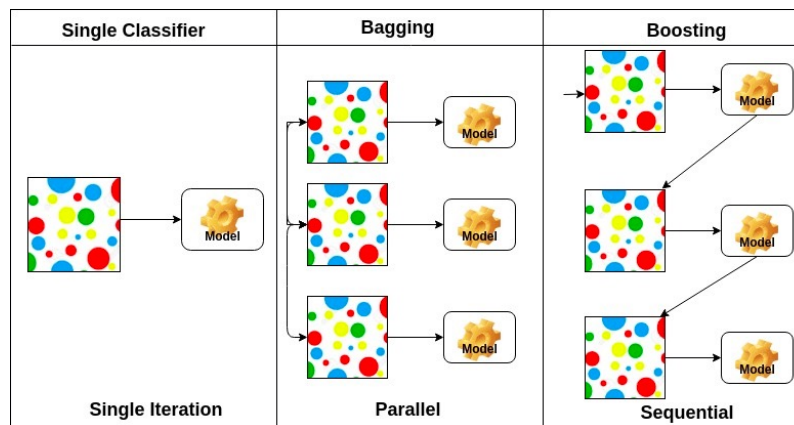
54

Ensemble Methods: Increasing the Accuracy

- Popular ensemble methods
 - Bagging:** Trains each model using a subset of the training set, and models learned in parallel
 - Boosting:** Trains each new model instance to emphasize the training instances that previous models misclassified, and models learned in order



55



<https://www.datacamp.com/community/tutorials/adaboost-classifier-python>

56

56

Bagging: Bootstrap Aggregation

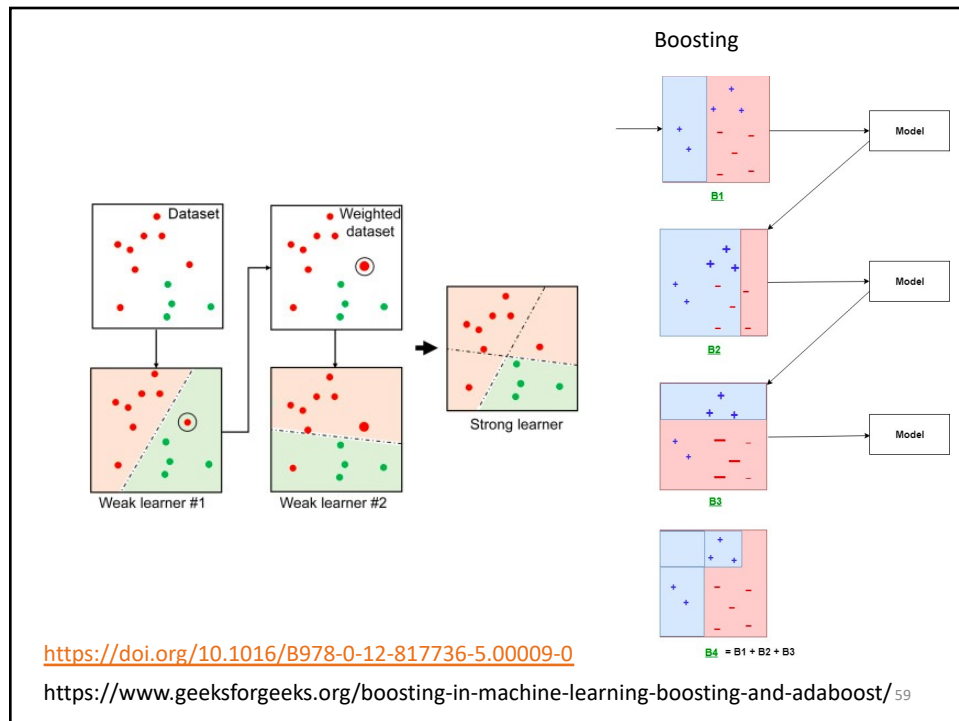
- Analogy: Diagnosis based on multiple doctors' majority vote
- **Training:** Given a set D of d tuples, at each iteration i
 - create bootstrap sample (training set), D_i , by sampling D with replacement
 - learn a classifier model M_i for each training set D_i
- **Classification:** classify an unknown sample X
 - Each classifier M_i returns its class prediction
 - The bagged classifier M^* counts the votes and assigns the class with the most votes to X
- **Prediction:** It can be applied to the prediction of continuous values by taking the average value of each prediction instead of voting
- Accuracy: Improved accuracy in prediction
 - Often significantly better than a single classifier derived from D
 - For noise data: Not considerably worse, more robust

57

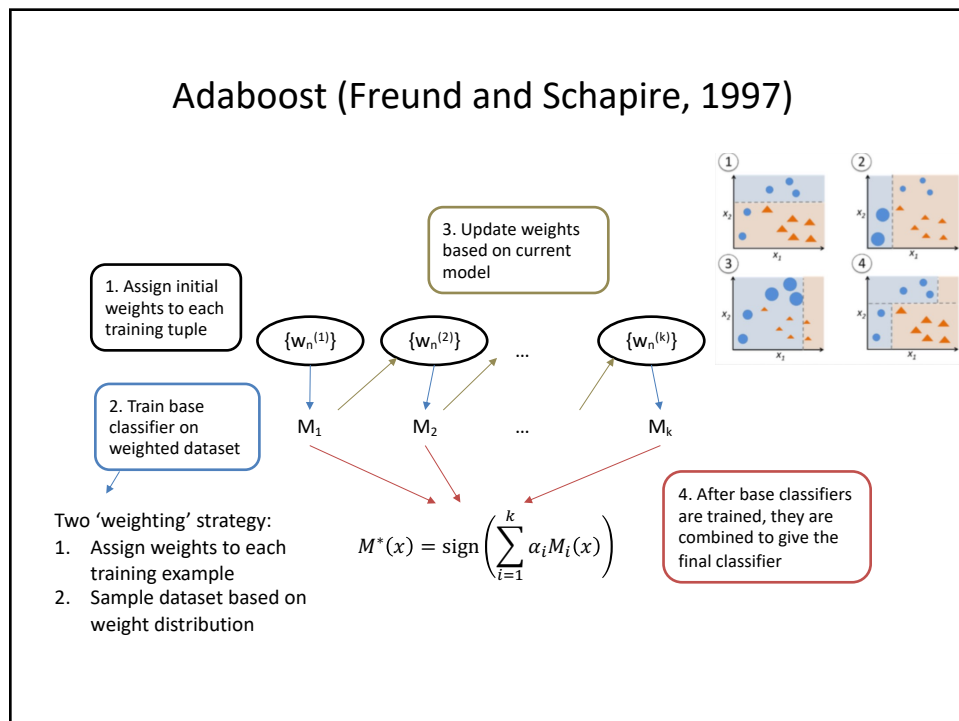
Boosting

- Analogy: Consult several doctors, based on a combination of weighted diagnoses—weight assigned based on the previous diagnosis accuracy
- How boosting works?
 - **Weights** are assigned to each training tuple
 - A series of k classifiers is iteratively learned
 - After a classifier M_i is learned, the weights are updated to allow the subsequent classifier, M_{i+1} , to **pay more attention to the training tuples that were misclassified** by M_i
 - The final **M^* combines the votes** of each individual classifier, where the weight of each classifier's vote is a function of its accuracy
- Boosting algorithm can be extended for numeric prediction
- Comparing with bagging: Boosting tends to have greater accuracy, but it also risks overfitting the model to misclassified data

58



59



60

Adaboost (Freund and Schapire, 1997)

- Given a set of d class-labeled tuples, $(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_d, y_d)$
- Initially, all the weights of tuples are set the same ($1/d$)
- Generate k classifiers in k rounds. At round i ,
 - Tuples from D are sampled (with replacement) to form a training set D_i of the same size
 - Each tuple's chance of being selected is based on its weight
 - A classification model M_i is derived from D_i
 - Its error rate is calculated using D_i as a test set
 - If a tuple is misclassified, its weight is increased; otherwise, it is decreased
- Error rate: $\text{err}(\mathbf{X}_i)$ is the misclassification error of tuple \mathbf{X}_i . Classifier M_i error rate is the sum of the weights of the misclassified tuples:

- The weight of classifier M_i 's vote is

$$\text{error}(M_i) = \sum_j w_j \times \text{err}(\mathbf{X}_j)$$

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

61

Gradient Boosting

- Operates on:
 - A differentiable loss function
 - A weak learner to make predictions (usually trees)
 - An additive model to add weak learners to minimize the loss function
- Each time adds an additional weak learner

$$\begin{aligned} \hat{y}_i^{(0)} &= 0 \\ \hat{y}_i^{(1)} &= f_1(x_i) = \hat{y}_i^{(0)} + f_1(x_i) \\ \hat{y}_i^{(2)} &= f_1(x_i) + f_2(x_i) = \hat{y}_i^{(1)} + f_2(x_i) \quad \text{Previous model} \\ &\dots \\ \hat{y}_i^{(t)} &= \sum_{k=1}^t f_k(x_i) = \hat{y}_i^{(t-1)} + f_t(x_i) \quad \text{New weak learner} \end{aligned}$$
- Scalable implementation: XGBoost

62

Random Forest: Basic Concepts

- Random Forest (first proposed by L. Breiman in 2001)
 - Bagging with **decision trees** as base models
 - *Data bagging*
 - Use a **subset of training data** by sampling with replacement for each tree
 - *Feature bagging* ← Advantage of decision trees – more diversity
 - At each node use a **random selection of attributes** as candidates and split by the best attribute among them
 - During classification, each tree votes and the most popular class is returned

63

Random Forest

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

64

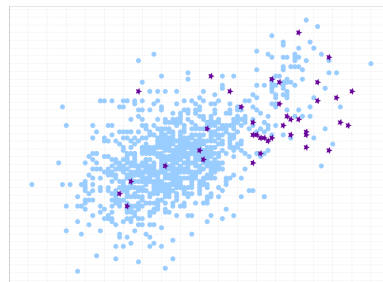
Ensemble Methods Recap

- Random forest and XGBoost are the most commonly used algorithms for tabular data
- Pros
 - Good performance for tabular data, requires no data scaling
 - Can scale to large datasets
 - Can handle missing data to some extent
- Cons
 - Can overfit to training data if not tuned properly
 - Lack of interpretability (compared to decision trees)

65

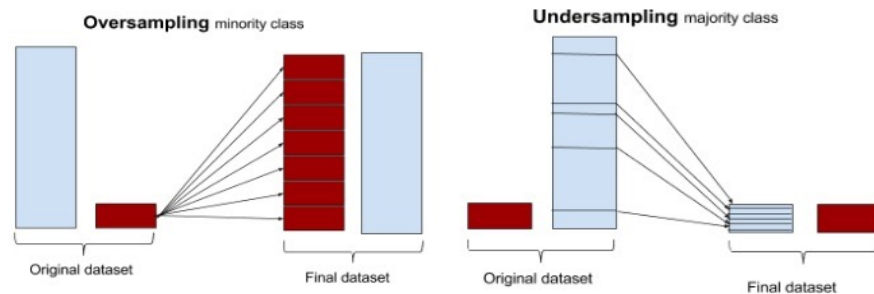
Classification of Class-Imbalanced Data Sets

- Traditional methods assume a balanced distribution of classes and equal error costs. But in real world situations, we may face imbalanced data sets, which has rare positive examples but numerous negative ones.
- Medical diagnosis: Medical screening for a condition is usually performed on a large population of people without the condition, to detect a small minority with it (e.g., HIV prevalence in the USA is $\sim 0.4\%$)
- Fraud detection: About 2% of credit card accounts are defrauded per year. (Most fraud detection domains are heavily imbalanced.)
- Product defect, accident (oil-spill), disk drive failures, etc.



66

- Typical methods on imbalanced data (Balance the training set)
 - **Oversampling:** Oversample the minority class.
 - **Under-sampling:** Randomly eliminate tuples from majority class
 - **Synthesizing:** Synthesize new minority classes

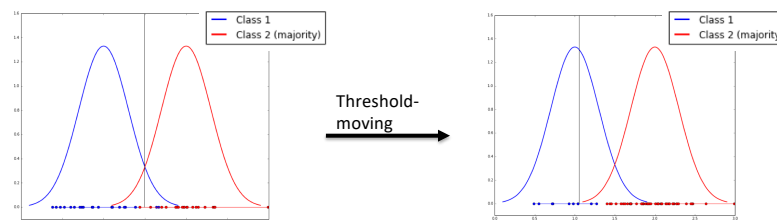


67

67

Classification of Class-Imbalanced Data Sets

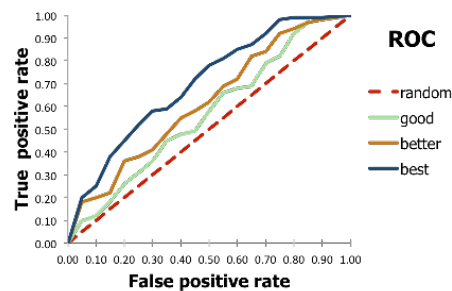
- Typical methods on imbalanced data (At the algorithm level)
 - **Threshold-moving:** Move the decision threshold, t , so that the rare class tuples are easier to classify, and hence, less chance of costly false negative errors
 - **Class weight adjusting:** Since false negative costs more than false positive, we can give larger weight to false negative
- Ensemble multiple classifiers introduced in the following chapter



68


Evaluate imbalanced data classifier

- Can we use Accuracy to evaluate imbalanced data classifier?
- Accuracy simply counts the number of errors. If a data set has 2% positive labels and 98% negative labels, a classifier that map all inputs to negative class will get an accuracy of 98%!
- ROC Curve



69

Chapter 8. Classification: Basic Concepts

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

70

Multiclass Classification

- Classification involving more than two classes (i.e., > 2 Classes)
- Method 1. **One-versus-all**
 - Given m classes, train m binary classifiers, one for each class
 - Classifier j : treat tuples in class j as *positive* & **all the rest** as *negative*
 - To classify a tuple \mathbf{X} , the set of classifiers vote as an ensemble
- Method 2. **All-versus-all**: Train a binary classifier for any two classes
 - Given m classes, construct $m(m - 1)/2$ binary classifiers
 - A classifier is trained using tuples of the two classes
 - To classify a tuple \mathbf{X} , each classifier votes
 - \mathbf{X} is assigned to the class with maximal vote
- Comparison: One-vs.-one tends to perform better than one-vs.-rest
- **A related problem**: multi-label classification (each data tuple could belong to one or more classes)

71

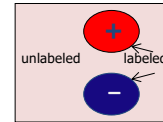
Semi-Supervised Learning

- **Definition**: Learning from a mix of labeled and unlabeled data, typically with far more unlabeled data.
- **Goal**: Improve learning performance by leveraging the unlabeled data alongside the labeled data.
- **Example**: Training a model for email spam classification using a small labeled dataset and a large unlabeled dataset.
- **Data**: Combines both labeled and unlabeled data.
- **Similarity**: Shares characteristics with supervised (labeled) and unsupervised (unlabeled) learning.
- **Difference**: Bridges the gap between supervised and unsupervised learning by combining both types of data.

72

Semi-Supervised Classification

- Semi-supervised: Uses labeled and unlabeled data to build a classifier
- Self-training
 - Build a classifier using the labeled data
 - Use it to label the unlabeled data, and those with the most confident label prediction are added to the set of labeled data
 - Repeat the above process
 - Adv.: easy to understand; Disadv.: may reinforce errors
- Co-training: Use two or more classifiers to teach each other
 - Each learner uses a mutually independent set of features of each tuple to train a good classifier, say f_1 and f_2
 - Then f_1 and f_2 are used to predict the class label for unlabeled data X
 - Teach each other: The tuple having the most confident prediction from f_1 is added to the set of labeled data for f_2 & vice versa
- Other methods include joint probability distribution of features and labels



73

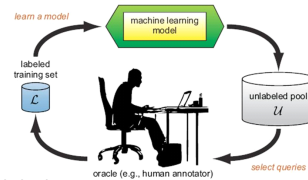
Self-Supervised Learning

- **Definition:** A form of unsupervised learning where the model generates its own supervisory signal by creating pseudo-labels from the data.
- **Goal:** Pre-train models by solving tasks derived from data itself, often for use in transfer learning.
- **Example:** Predicting missing words in a sentence (used in training language models like GPT).
- **Data:** Uses unlabeled data but creates internal labels or tasks.
- **Similarity:** Shares use of unlabeled data with unsupervised learning.
- **Difference:** Employs specific pretext tasks to create supervisory signals.

74

Active Learning

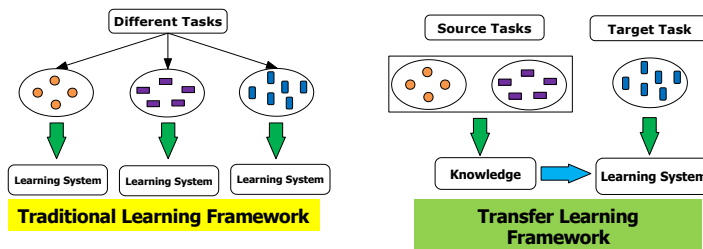
- A special case of semi-supervised learning
 - Unlabeled data: Abundant
 - Class labels are expensive to obtain
- Active learner: Interactively query teachers (oracle) for labels
- Pool-based approach: Uses a pool of unlabeled data
 - L : a small subset of D is labeled, U : a pool of unlabeled data in D
 - Use a query function to carefully select one or more tuples from U and request labels from an oracle (a human annotator)
 - The newly labeled samples are added to L , and learn a model
 - Goal: **Achieve high accuracy using as few labeled data as possible**
- Evaluated using *learning curves*: Accuracy as a function of the number of instances queried (# of tuples to be queried should be small)
- A lot of algorithms have been developed for active learning



75

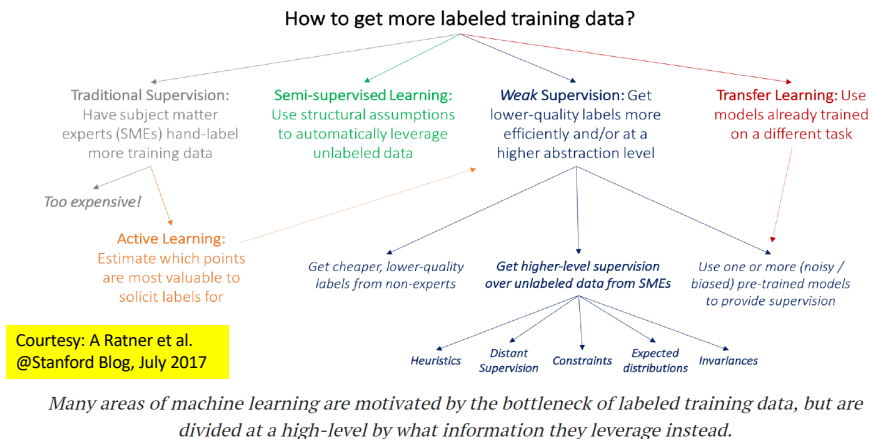
Transfer Learning: Conceptual Framework

- Transfer learning: Extract knowledge from one or more source tasks (e.g., recognizing cars) and apply the knowledge to a target task (e.g., recognizing trucks)
- Traditional learning: Build a new classifier for each new task
- Transfer learning: Build new classifier by applying existing knowledge learned from source tasks
- Many algorithms are developed, applied to text classification, spam filtering, etc.




76

Relationships Among Different Kinds of Supervisions



77

Chapter 8. Classification: Basic Concepts

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

78

Summary

- Classification: Model construction from a set of training data
- Effective and scalable methods
 - Decision tree induction, Bayes classification methods, linear classifier, ...
 - No single method has been found to be superior over all others for all data sets
- Evaluation metrics: Accuracy, sensitivity, specificity, precision, recall, F measure
- Model evaluation: Holdout, cross-validation, bootstrapping, ROC curves (AUC)
- Improve Classification Accuracy: Bagging, boosting
- Additional concepts on classification: Multiclass classification, semi-supervised classification, active learning, transfer learning, weak supervision

79

References (1)

- C. Apte and S. Weiss. **Data mining with decision trees and decision rules**. Future Generation Computer Systems, 13, 1997
- P. K. Chan and S. J. Stolfo. **Learning arbiter and combiner trees from partitioned data for scaling machine learning**. KDD'95
- A. J. Dobson. **An Introduction to Generalized Linear Models**. Chapman & Hall, 1990.
- R. O. Duda, P. E. Hart, and D. G. Stork. **Pattern Classification**, 2ed. John Wiley, 2001
- U. M. Fayyad. **Branching on attribute values in decision tree generation**. AAAI'94.
- Y. Freund and R. E. Schapire. **A decision-theoretic generalization of on-line learning and an application to boosting**. J. Computer and System Sciences, 1997.
- J. Gehrke, R. Ramakrishnan, and V. Ganti. **Rainforest: A framework for fast decision tree construction of large datasets**. VLDB'98.
- J. Gehrke, V. Ganti, R. Ramakrishnan, and W.-Y. Loh, **BOAT -- Optimistic Decision Tree Construction**. SIGMOD'99.
- T. Hastie, R. Tibshirani, and J. Friedman. **The Elements of Statistical Learning: Data Mining, Inference, and Prediction**. Springer-Verlag, 2001.

80

References (2)

- T.-S. Lim, W.-Y. Loh, and Y.-S. Shih. **A comparison of prediction accuracy, complexity, and training time of thirty-three old and new classification algorithms.** Machine Learning, 2000
- J. Magidson. **The Chaid approach to segmentation modeling: Chi-squared automatic interaction detection.** In R. P. Bagozzi, editor, Advanced Methods of Marketing Research, Blackwell Business, 1994
- M. Mehta, R. Agrawal, and J. Rissanen. **SLIQ : A fast scalable classifier for data mining.** EDBT'96
- T. M. Mitchell. **Machine Learning.** McGraw Hill, 1997
- S. K. Murthy, **Automatic Construction of Decision Trees from Data: A Multi-Disciplinary Survey,** Data Mining and Knowledge Discovery 2(4): 345-389, 1998
- J. R. Quinlan. **Induction of decision trees.** Machine Learning, 1:81-106, 1986.
- J. R. Quinlan. **C4.5: Programs for Machine Learning.** Morgan Kaufmann, 1993.
- J. R. Quinlan. **Bagging, boosting, and c4.5.** AAAI'96.

81

References (3)

- R. Rastogi and K. Shim. **Public: A decision tree classifier that integrates building and pruning.** VLDB'98
- J. Shafer, R. Agrawal, and M. Mehta. **SPRINT : A scalable parallel classifier for data mining.** VLDB'96
- J. W. Shavlik and T. G. Dietterich. **Readings in Machine Learning.** Morgan Kaufmann, 1990
- P. Tan, M. Steinbach, and V. Kumar. **Introduction to Data Mining.** Addison Wesley, 2005
- S. M. Weiss and C. A. Kulikowski. **Computer Systems that Learn: Classification and Prediction Methods from Statistics, Neural Nets, Machine Learning, and Expert Systems.** Morgan Kaufman, 1991
- S. M. Weiss and N. Indurkha. **Predictive Data Mining.** Morgan Kaufmann, 1997
- I. H. Witten and E. Frank. **Data Mining: Practical Machine Learning Tools and Techniques,** 2ed. Morgan Kaufmann, 2005

82

Classification Is to Derive the Maximum Posteriori

- Let D be a training set of tuples and their associated class labels, and each tuple is represented by an n -D attribute vector $\mathbf{X} = (x_1, x_2, \dots, x_n)$
- Suppose there are m classes C_1, C_2, \dots, C_m .
- Classification is to derive the maximum posteriori, i.e., the maximal $P(C_i|\mathbf{X})$
- This can be derived from Bayes' theorem

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

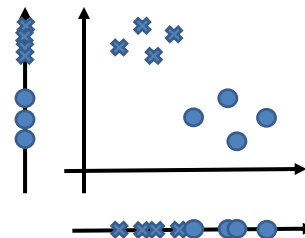
- Since $P(\mathbf{X})$ is constant for all classes, only

needs to be maximized $P(C_i|\mathbf{X}) \propto P(\mathbf{X}|C_i)P(C_i)$

83

Linear Discriminant Analysis (LDA)

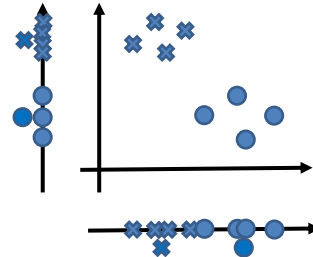
- Linear Discriminant Analysis (LDA) works when the attributes are all continuous
 - For the categorical attributes, discriminant correspondence analysis is the equivalent technique
- Basic Ideas: Project all samples on a line such that different classes are well separated
- Example: Suppose we have 2 classes and 2-dimensional samples x_1, \dots, x_n
 - n_1 samples come from class 1
 - n_2 samples come from class 2
- Let the line direction be given by unit vector \mathbf{v}
- There are two candidates of projections
 - Vertical: $\mathbf{v} = (0,1)$
 - Horizontal: $\mathbf{v} = (1,0)$
- Which one looks better?
- How to mathematically measure it?



84

Fisher's LDA (Linear Discriminant Analysis)

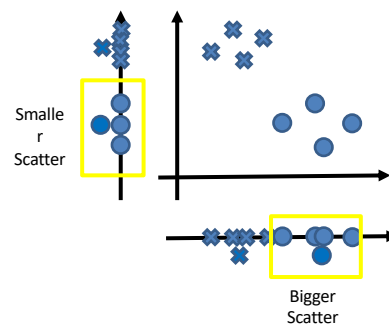
- $v^T x_i$ is the distance of projection of x_i from the origin
- Let μ_1 and μ_2 be the means of class 1 and class 2 in the original space
 - $\mu_1 = \frac{1}{n_1} \sum_{i \in \text{class 1}} x_i$
 - $\mu_2 = \frac{1}{n_2} \sum_{i \in \text{class 2}} x_i$
- The distance between the means of the projected points
 - $|v^T \mu_1 - v^T \mu_2|$
 - Good? No. Horizontal one may have larger distance



85

Fisher's LDA (con't)

- Normalization needed
- Scatter: Sample variance multiplied by n
 - $s_1 = \sum_{i \in \text{class 1}} (v^T x_i - v^T \mu_1)^2$
 - $s_2 = \sum_{i \in \text{class 2}} (v^T x_i - v^T \mu_2)^2$
- Fisher's LDA
 - Maximize $J(v) = \frac{(v^T \mu_1 - v^T \mu_2)^2}{s_1 + s_2}$
 - Closed-form optimal solution



86

Fisher's LDA: Summary

- Advantages
 - Useful for dimension reduction
 - Easy to extend to multi-classes
- Fisher's LDA will fail
 - When $\mu_1 = \mu_2$, $J(v)$ is always 0.
 - When classes have large overlap when projected to any line