W.Sagievnensekek.gom Sınıflandırma (Classification) Şadi Evren ŞEKER www.SadiEvrenSiwww.SadiEvrenSEKER.com Youtube : Bilgisayar Kavramlar Data Mining: Concepts and **Techniques** 02/11/16

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.corKaynaklar.renSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE





<u>Data Mining: Concepts and</u>
 <u>Techniques, Third Edition,</u> Jiawei Han,
 Micheline Kamber, Jian Pei

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE

Chapter 8. Classification (Siniflandirma)

Sınıflandırmaya Giriş **K-NN Algoritması** H Karar Ağaçları (Decision Trees) Bayes Sınıflandırma Yöntemi Kural Tabanlı Sınıflandırmalar (Rule-based Classification) Doğru Modelin seçilmesi Başarıyı arttıran bazı yöntemler (Ensemble Techniques)

Supervised vs. Unsupervised Learning (Gözetimli ve Gözetimsiz Öğrenme)

- Supervised learning (classification) [EvrenSEKER.com www.Sadil
 - Supervision: The training data (observations,
- measurements, etc.) are accompanied by **labels** indicating
- the class of the observations
 - New data is classified based on the training set
- Unsupervised learning (clustering)

The class labels of training data is unknown

 Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

Prediction Problems: Classification vs. Numeric Prediction

Classification

predicts categorical class labels (discrete or nominal)

 classifies data (constructs a model) based on the training set and the values (class labels) in a classifying attribute and uses it in classifying new data

Numeric Prediction

 models continuous-valued functions, i.e., predicts unknown or missing values

Typical applications

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

Classification—A Two-Step Process

www.saulevrensekek.com www.saulevrensekek.com www.Sadi

- Model construction: describing a set of predetermined classes
 - Each tuple/sample is assumed to belong to a predefined class, as determined by the class label attribute
 - The set of tuples used for model construction is training set
 - The model is represented as classification rules, decision trees, or mathematical formulae
- Model usage: for classifying future or unknown objects
 - Estimate accuracy of the model
 - The known label of test sample is compared with the classified result from the model
 - Accuracy rate is the percentage of test set samples that are correctly classified by the model
 - Test set is independent of training set (otherwise overfitting)
 - If the accuracy is acceptable, use the model to classify new data
 - Note: If the test set is used to select models, it is called validation (test) set

Process (1): Model Construction





Chapter 8. Classification: Basic Concepts

Classification: Basic Concepts **Decision Tree Induction Bayes Classification Methods Rule-Based Classification** Model Evaluation and Selection Techniques to Improve Classification Accuracy: **Ensemble Methods** Summary

Decision Tree Induction: An Example



Algorithm for Decision Tree Induction

- Basic algorithm (a greedy algorithm)
 - Tree is constructed in a top-down recursive divide-and-
- conquer manner
 - At start, all the training examples are at the root
 - Attributes are categorical (if continuous-valued, they are
 - discretized in advance)
 - Examples are partitioned recursively based on selected attributes
 - Test attributes are selected on the basis of 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 majority voting is employed for classifying the leaf
 There are no samples left

Brief Review of Entropy

- Entropy (Information Theory)
 - A measure of uncertainty associated with a random variable
 - Calculation: For a discrete random variable Y taking m distinct values {y₁, ..., y_m},

1.0

 $\stackrel{(X)}{\overset{H}{\overset{}}_{H}}_{H}$

 $0.5 \\ Pr(X = 1)$

- $H(Y) = -\sum_{i=1}^{m} p_i \log(p_i)$, where $p_i = P(Y = y_i)$
- Interpretation:
 - Higher entropy => higher uncertainty
 - Lower entropy => lower uncertainty
- Conditional Entropy

• $H(Y|X) = \sum_{x} p(x)H(Y|X = x)$ m = 2

Attribute Selection Measure: Information Gain (ID3/C4.5)

Select the attribute with the highest information gain

- 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) = -*∑_{i=1}^m p_i log₂(p_i)

 Information needed (after using A to split D into v partitions) to classify D: *Info_A(D) =* ∑_{i=1}^v | D_j | / D_j | × *Info(D_j)*

• Information gained by branching on attribute A $Gain(A) = Info(D) - Info_A(D)$ 13

Attribute Selection: Information Gain

Class P: buys_computer = "yes" Class N: buys computer = "no" $Info(D) = I(9,5) = -\frac{9}{14}\log_2(\frac{9}{14}) - \frac{5}{14}\log_2(\frac{5}{14}) = 0.940$ I(p_i, n_i) age n_i **p**_i 0.971 2 3 <=30 31...40 4 $\mathbf{0}$ 3 0.971 >40 2

age	income	student	credit_rating	buys_computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
3140	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
3140	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
3140	medium	no	excellent	yes
3140	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,

Gain(income) = 0.029Gain(student) = 0.151 $Gain(credit_rating) = 0.048$

Computing Information-Gain for Continuous-Valued Attributes

Let attribute A be a continuous-valued attribute Must determine the *best split point* for A Sort the value A in increasing order Typically, the midpoint between each pair of adjacent values is considered as a possible *split point* • $(a_i+a_{i+1})/2$ is the midpoint between the values of a_i and a_{i+1} The point with the *minimum expected information requirement* for A is selected as the split-point for A Split:

D1 is the set of tuples in D satisfying A ≤ split-point, and D2 is the set of tuples in D satisfying A > split-point

Gain Ratio for Attribute Selection (C4.5)

- Information gain measure is biased towards attributes with a large number of values
- C4.5 (a successor of ID3) uses gain ratio to overcome the problem (normalization to information gain)

$$SplitInfo_A(D) = -\sum_{j=1}^{\nu} \frac{|D_j|}{|D|} \times \log_2(\frac{|D_j|}{|D|})$$

GainRatio(A) = Gain(A)/SplitInfo(A)

Ex.

$$SplitInfo_{income}(D) = -\frac{4}{14} \times \log_2\left(\frac{4}{14}\right) - \frac{6}{14} \times \log_2\left(\frac{6}{14}\right) - \frac{4}{14} \times \log_2\left(\frac{4}{14}\right) = 1.557$$

gain_ratio(income) = 0.029/1.557 = 0.019

 The attribute with the maximum gain ratio is selected as the splitting attribute

Gini Index (CART, IBM IntelligentMiner)

• If a data set *D* contains examples from *n* classes, gini index, gini(D) is defined as $gini(D) = 1 - \sum_{j=1}^{n} p_j^2$

where p_i 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)
17

Computation of Gini Index

• Ex. 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 D₁: {low, medium} and 4 in D₂ gini_{income \in \{low,medium\}}(D) = $\left(\frac{10}{14}\right)Gini(D_1) + \left(\frac{4}{14}\right)Gini(D_2)$ $= \frac{10}{14}\left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) + \frac{4}{14}\left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2\right)$ = 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
- ¹⁸ Can be modified for categorical attributes

Comparing Attribute Selection Measures

The three measures, in general, return good results but

- Information gain:
 - biased towards multivalued attributes
- Gain ratio:
 - - tends to prefer unbalanced splits in which one partition is much smaller than the others
- Gini index:
 - biased to multivalued attributes
 - has difficulty when # of classes is large

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

Other Attribute Selection Measures

- <u>CHAID</u>: a popular decision tree algorithm, measure based on χ² test for independence
- <u>C-SEP</u>: performs better than info. gain and gini index in certain cases
- <u>G-statistic</u>: has a close approximation to χ^2 distribution
- MDL (Minimal Description Length) principle (i.e., 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
- Multivariate splits (partition based on multiple variable combinations)
 - <u>CART</u>: finds multivariate splits based on a linear comb. of attrs.
- Which attribute selection measure is the best?
 - Most give good results, none is significantly superior than others

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
 - <u>Prepruning</u>: 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
 - <u>Postpruning</u>: *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"

Enhancements to Basic Decision Tree Induction

- Allow for continuous-valued attributes
 - Dynamically define new discrete-valued attributes that partition the continuous attribute value into a discrete set of intervals
- Handle missing attribute values
 - Assign the most common value of the attribute
 - Assign probability to each of the possible values
- Attribute construction
 - Create new attributes based on existing ones that are sparsely represented

This reduces fragmentation, repetition, and replication

Classification in Large Databases

- Classification—a classical problem extensively studied by statisticians and machine learning researchers
 - Scalability: Classifying data sets with millions of examples and hundreds of attributes with reasonable speed
 - Why is decision tree induction popular?
 - relatively faster learning speed (than other classification methods)
 - convertible to simple and easy to understand classification rules
 - can use SQL queries for accessing databases
 - comparable classification accuracy with other methods
 - RainForest (VLDB' 98 Gehrke, Ramakrishnan & Ganti)
 - Builds an AVC-list (attribute, value, class label)

Scalability ⊢rameworк тог RainForest

- Separates the scalability aspects from the criteria that
- determine the quality of the tree
- Builds an AVC-list: AVC (Attribute, Value, Class_label)
- AVC-set (of an attribute X)
- Projection of training dataset onto the attribute X and class label where counts of individual class label are aggregated
 - AVC-group (of a node n)

Set of AVC-sets of all predictor attributes at the node n

Rainforest: Training Set and Its AVC Sets

www.saulevrensekek.com www.saulevrensekek.com www.Sadil

Training Examples

age	income	student	redit_rating	com
<=30	high	no	fair	no
<=30	high	no	excellent	no
3140	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
3140	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
3140	medium	no	excellent	yes
3140	high	yes	fair	yes
>40	medium	no	excellent	no

AVC-set on Age

Age	Buy_Computer		
Sad	yes	no	
<=30	2	3	
3140	4	0	
>40	3	2	

AVC-set	on	income
---------	----	--------

income	Buy_Computer		
	yes	no	
high	2	2	
medium	4	2	
low	3	1	

AVC-set on *Student*

AVC-set on credit_rating

student	Buy_Computer		Que d'it	Buy_Computer	
10001	yes	no	rating	yes	no
yes	6	nsfkei	fair	6	2
no	3	4	excellent	3	3

BOAT (Bootstrapped Optimistic Algorithm for Tree Construction)

 Use a statistical technique called *bootstrapping* to create several smaller samples (subsets), each fits in memory

 Each subset is used to create a tree, resulting in several trees

These trees are examined and used to construct a new tree T'

It turns out that T' is very close to the tree that would be generated using the whole data set together

Adv: requires only two scans of DB, an incremental alg.

Presentation of Classification Results



www.SadiEvrenSEKER.coMineSet 3.0 MiseKER.com www.SadiE



Interactive Visual Mining by Perception-Based Classification (PBC)



Chapter 8. Classification: Basic Concepts

Classification: Basic Concepts **Decision Tree Induction Bayes Classification Methods Rule-Based Classification** Model Evaluation and Selection Techniques to Improve Classification Accuracy: **Ensemble Methods** Summary 30

Bayesian Classification: Why?

- <u>A statistical classifier</u>: performs *probabilistic prediction, i.e.,* predicts class membership probabilities
- Foundation: Based on Bayes' Theorem.
- <u>Performance</u>: 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
- <u>Standard</u>: Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured

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 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|X), (i.e., *posteriori probability):* the probability that the hypothesis holds given the observed data sample X
- P(H) (*prior probability*): the initial probability
 - E.g., X will buy computer, regardless of age, income, ...
- P(X): probability that sample data is observed
- P(X|H) (likelihood): the probability of observing the sample X, given that the hypothesis holds
 - E.g., Given that X will buy computer, the prob. that X is 31..40, medium income

Prediction Based on Bayes' Theorem

- Given training data X, posteriori probability of a hypothesis H, P(H|X), follows the Bayes' theorem
- $P(H|\mathbf{X}) = \frac{P(\mathbf{X}|H)P(H)}{P(\mathbf{X})} = P(\mathbf{X}|H) \times P(H)/P(\mathbf{X})$ Informally, this can be viewed as
 - posteriori = likelihood x prior/evidence
- Predicts X belongs to C_i iff the probability P(C_i|X) is the highest among all the P(C_k|X) for all the k classes
- Practical difficulty: It requires initial knowledge of many probabilities, involving significant computational cost

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 X = (x₁, x₂, ..., x_n)
- Suppose there are *m* classes C₁, C₂, ..., C_m.
- Classification is to derive the maximum posteriori, i.e., the maximal P(C_i | 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(X) is constant for all classes, only $P(C_i | \mathbf{X}) = P(\mathbf{X} | C_i) P(C_i)$

needs to be maximized

Naïve Bayes Classifier

A simplified assumption: attributes are conditionally independent (i.e., no dependence relation between attributes):

 $P(\mathbf{X}|C_i) = \prod_{k=1}^{n} P(x_k|C_i) = P(x_1|C_i) \times P(x_2|C_i) \times \dots \times P(x_n|C_i)$ This greatly reduces the computation cost: Only counts the "

■ This greatly reduces the €0mputation cost: Only counts the class distribution

If A_k is categorical, P(x_k|C_i) is the # of tuples in C_i having value x_k for A_k divided by |C_{i, D}| (# of tuples of C_i in D)

 If A_k is continous-valued, P(x_k|C_i) is usually computed based on Gaussian distribution with a mean μ and standard deviation σ

and P(x_k|C_i) is $g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

www.SadiEvrenSEKER.com $P(\mathbf{X} | C_i) = g(x_k, \mu_{C_i}, \sigma_{C_i})$ on www.Sadil

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	_com
<=30	high	no	fair	no
<=30	high	no	excellent	no
3140	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
3140	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
3140	medium	no	excellent	yes
3140	high	yes	fair	yes
>40	medium	no	excellent	no
Naïve Bayes Classifier: An Example

no

no

no

excellent

1...40 hiah no fair ves • $P(C_i)$: P(buys computer = "yes") = 9/14 = 0.643no fair ves ves ves low no P(buys computer = "no") = 5/14 = 0.357excellen ves no ves fair ves Compute $P(X|C_i)$ for each class 40 medium ves fair ves 40 medium no excellen P(age = "<=30" | buys computer = "yes") = 2/9 = 0.222 $P(age = "<= 30" | buys_computer = "no") = 3/5 = 0.6$ P(income = "medium" | buys_computer = "yes") = 4/9 = 0.444 P(income = ``medium'' | buys computer = ``no'') = 2/5 = 0.4P(student = "yes" | buys_computer = "yes) = 6/9 = 0.667 P(student = "yes" | buys computer = "no") = 1/5 = 0.2P(credit_rating = "fair" | buys_computer = "yes") = 6/9 = 0.667 P(credit_rating = "fair" | buys_computer = "no") = 2/5 = 0.4 X = (age <= 30, income = medium, student = yes, credit_rating = fair) $P(X|C_i) : P(X|buys_computer = "yes") = 0.222 \times 0.444 \times 0.667 \times 0.667 = 0.044$ $P(X|buys_computer = "no") = 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$ $P(X|C_i)*P(C_i): P(X|buys_computer = "yes") * P(buys_computer = "yes") = 0.028$ P(X|buys_computer = "no") * P(buys_computer = "no") = 0.007 Therefore, X belongs to class ("buys_computer = yes") 37

Avoiding the Zero-Probability Problem

 Naïve Bayesian prediction requires each conditional prob. be non-zero. Otherwise, the predicted prob. will be zero

 $P(X \mid C_i) = \prod_{k=1}^{n} P(x_k \mid C_i)$

- Ex. 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

Prob(income = low) = 1/1003

Prob(income = medium) = 991/1003

Prob(income = high) = 11/1003

 The "corrected" prob. estimates are close to their "uncorrected" counterparts

Naïve Bayes Classifier: Comments

- Advantages
 - Easy to implement
 - Good results obtained in most of the cases
- Disadvantages
 - Assumption: class conditional independence, therefore loss of accuracy
 - Practically, dependencies exist among variables
 - E.g., hospitals: 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? Bayesian Belief Networks (Chapter 9)

Chapter 8. Classification: Basic Concepts

Classification: Basic Concepts **Decision Tree Induction Bayes Classification Methods Rule-Based Classification** Model Evaluation and Selection **Techniques to Improve Classification Accuracy: Ensemble Methods** Summary

Using IF-THEN Rules for Classification

- Represent the knowledge in the form of IF-THEN rules
 - R: IF *age* = youth AND *student* = yes THEN *buys_computer* = yes
 - Rule antecedent/precondition vs. rule consequent
- Assessment of a rule: coverage and accuracy
 - n_{covers} = # of tuples covered by R
 - n_{correct} = # of tuples correctly classified by R
 coverage(R) = n_{covers} / |D| /* D: training data set */
 accuracy(R) = n_{correct} / n_{covers}
 - If more than one rule are triggered, need **conflict resolution**
 - Size ordering: assign the highest priority to the triggering rules that has the "toughest" requirement (i.e., with the most attribute tests)
 - Class-based ordering: decreasing order of prevalence or misclassification cost per class
 - Rule-based ordering (decision list): rules are organized into one long priority list, according to some measure of rule quality or by experts

Rule Extraction from a Decision Tree

- Rules are easier to understand than large trees
- One rule is created *for each path* from the root to a leaf
- Each attribute-value pair along a path forms a conjunction: the leaf holds the class
 prediction
- Rules are mutually exclusive and exhaustive
 - Example: Rule extraction from our *buys_computer* decision-tree
 - IF age = young AND student = noTHEN buys_computer = noIF age = young AND student = yesTHEN buys_computer = yesIF age = mid-ageTHEN buys_computer = yesIF age = old AND credit_rating = excellentTHEN buys_computer = noIF age = old AND credit_rating = fairTHEN buys_computer = yes

age?

31..40

yes

yes

>40

excellent

no

credit rating?

fair

yes

<=30

Kule Induction: Sequential Covering Method

- Sequential covering algorithm: Extracts rules directly from training data
- Typical sequential covering algorithms: FOIL, AQ, CN2, RIPPER
- Rules are learned sequentially, each for a given class C_i will cover many tuples of C_i but none (or few) of the tuples of other classes
 - Steps:
 - Rules are learned one at a time
- Each time a rule is learned, the tuples covered by the rules are removed
 - Repeat the process on the remaining tuples until *termination* condition, e.g., when no more training examples or when the quality of a rule returned is below a user-specified threshold

 Comp. w. decision-tree induction: learning a set of rules simultaneously

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadai



Rule Generation

To generate a rule
 while(true)
 find the best predicate p
 if foil-gain(p) > threshold then add p to current rule
 else break

www.SadiEvi www.SadiEvi www.SadiEvi www.SadiEvi www.SadiEvi www.SadiEvi



www.Sadif www.Sadif www.Sadif www.Sadif www.Sadif www.Sadif www.Sadif

www.SadiEvreneexampl

How to Learn-One-Rule?

- Start with the most general rule possible: condition = empty
- Adding new attributes by adopting a greedy depth-first strategy
 - Picks the one that most improves the rule quality
- Rule-Quality measures: consider both coverage and accuracy
- Foil-gain (in FOIL & RIPPER): assesses info_gain by extending condition
 FOIL_Gain = pos'×(log₂ <u>pos'</u> - log₂ <u>pos</u>)
 pos + neg)

pos + neg
 favors rules that have high accuracy and cover many positive tuples
 Rule pruning based on an independent set of test tuples

 $FOIL_Prune(R) = \frac{pos - neg}{pos + neg}$ Pos/neg are # of positive/negative tuples covered by R.
If *FOIL_Prune* is higher for the pruned version of R, prune R

Chapter 8. Classification: Basic Concepts

Classification: Basic Concepts **Decision Tree Induction Bayes Classification Methods Rule-Based Classification** Model Evaluation and Selection Techniques to Improve Classification Accuracy: **Ensemble Methods** Summary

Model Evaluation and Selection

- Evaluation metrics: How can we measure accuracy? Other metrics to consider?
- Use validation test set of class-labeled tuples instead of training set when assessing accuracy
- Methods for estimating a classifier's accuracy:
 - Holdout method, random subsampling
 - Cross-validation
- Bootstrap KER.com WWW.SadiEvrenSEKER.com WWW.SadiE
- Comparing classifiers:
- Confidence intervals
- Cost-benefit analysis and ROC Curves

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sada

Classifier Evaluation Metrics: Confusion

www.SadiEvrenSEKER.com MatrixEvrenSEKER.com www.SadiE

Confusion Matrix:

Actual class\Predicted class	C ₁	¬ C ₁
Sadiev C ₁ SEKER.cc	True Positives (TP)	False Negatives (FN)
SadiEv¬C ₁ SEKER.co	False Positives (FP)	S True Negatives (TN)

Example of Confusion Matrix:

Actual class\Predicted class	buy_computer = yes	buy_computer = no	Total	/w.Sadil
buy_computer = yes	6954	46	7000	/w.Sadil
buy_computer = no	412	2588	3000	/w.Sadil
Total	7366	2634	10000	vw.Sadi

Given *m* classes, an entry, *CM*_{i,j} in a confusion matrix indicates
 # of tuples in class *i* that were labeled by the classifier as class *j*

May have extra rows/columns to provide totals

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadiev

Accuracy, Error Rate, Sensitivity and Specificity



 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
- Sensitivity: True Positive recognition rate

Sensitivity = TP/P

Specificity: True Negative recognition rate
 Specificity = TN/N

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

Precision and Recall, and Fmeasures

- Precision: exactness what % of tuples that the classifier labeled as positive are actually positive
 TP
- Recall: completeness what % of positive tuples did the classifier label as positive?
 TP + FP
 TP = TP
 - Perfect score is 1.0
 - Inverse relationship between precision & recall
- F measure (F_1 or F-score): harmonic mean of precision and recall, $F = \frac{2 \times precision \times recall}{2 \times precision \times recall}$

precision + recall

TP + FN

precision

• F_{β} : weighted measure of precision and recall • assigns ß times as much weight to recall as to precision $F_{\beta} = \frac{(1+\beta^2) \times precision \times recall}{\beta^2 \times precision + recall}$

Classifier Evaluation Metrics: Example

www.sadievrensekek.com_www.sadievrensekek.com_www.Sadii

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

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)	
www.SadiETotalnSEKER.co	230	9770	10000	96.40 (accuracy)	

Precision = 90/230 = 39.13%
Recall = 90/300 = 30.00%

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE

Holdout & Cross-Validation Methods

Holdout method

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

Evaluating Classifier Accuracy: Bootstrap

- Bootstrap
 - Works well with small data sets
 - Samples the given training tuples uniformly with replacement
 - i.e., each time a tuple is selected, it is equally likely to be selected again and re-added to the training set
- Several bootstrap methods, and a common one is .632 boostrap
 - A data set with *d* tuples is sampled *d* times, with replacement, resulting in a training set of *d* samples. The data tuples that did not make it into the training set end up forming the test set. About 63.2% of the original data end up in the bootstrap, and the remaining 36.8% form the test set (since $(1 - 1/d)^d \approx e^{-1} = 0.368$)
- Repeat the sampling procedure k times, overall accuracy of the model: $Acc(M) = \frac{1}{k} \sum_{i=1}^{k} (0.632 \times Acc(M_i)_{test_set} + 0.368 \times Acc(M_i)_{train_set})$ 54

Estimating Confidence Intervals: Classifier Models M₁ vs. M₂

- Suppose we have 2 classifiers, M₁ and M₂, which one is better?
- Use 10-fold cross-validation to obtain err(M₁) and err(M₂)
 These mean error rates are just estimates of error on the true population of future data cases
- What if the difference between the 2 error rates is just attributed to chance?
 - Use a test of statistical significance

Obtain confidence limits for our error estimates

Estimating Confidence Intervals: Null Hypothesis

- Perform 10-fold cross-validation
- Assume samples follow a t distribution with k-1 degrees of freedom (here, k=10)
- Use t-test (or Student's t-test) SadiEvrenSEKER.com
- Null Hypothesis: M₁ & M₂ are the same
- If we can reject null hypothesis, then
 - we conclude that the difference between M₁ & M₂ is statistically significant

Chose model with lower error rate

Estimating Confidence Intervals: t-test

If only 1 test set available: pairwise comparison

- For ith round of 10-fold cross-validation, the same cross partitioning is used to obtain err(M₁)_i and err(M₂)_i
- Average over 10 rounds to get $\overline{err}(M_1)$ and $\overline{err}(M_2)$
- t-test computes t-statistic with k-1 degrees of

freedom:

$$t = \frac{\overline{err}(M_1) - \overline{err}(M_2)}{\sqrt{var(M_1 - M_2)/k}} \quad \text{where}$$

$$var(M_1 - M_2) = \frac{1}{k} \sum_{i=1}^k \left[err(M_1)_i - err(M_2)_i - (\overline{err}(M_1) - \overline{err}(M_2)) \right]^2$$

If two test sets available: use **non-paired t-test** whe $var(M_1 - M_2) = \sqrt{\frac{var(M_1)}{k_1} + \frac{var(M_2)}{k_2}},$ re where $k_1 \& k_2$ are # of cross-validation samples used for M_1

Estimating Confidence Intervals: Table for t-distribution



TABLE B: 1-DISTRIBUTION CRITICAL VALUES

		1532-5	12532.55	1	Tai	l probabi	lity p					
'df	.25	.20	.15	.10	.05	.025	.02	.01	.005	.0025	.001	.0005
1	1.000	1,376	1.963	3.078	6.314	12.71	15.89	31.82	63.66	127.3	318.3	636.6
2	.816	1.061	1.386	1.886	2.920	4.303	4.849	6.965	9.925	14.09	22.33	31.60
3	.765	.978	1.250	1.638	2.353	3.182	3.482	4.541	5.841	7.453	10.21	12.92
4	.741	.941	1.190	1.533	2.132	2.776	2.999	3.747	4.604	5.598	7.173	8.610
5	.727	.920	1.156	1.476	2.015	2.571	2.757	3.365	4.032	4.773	5.893	6.869
6	718	.906	1.134	1.440	1.943	2.447	2.612	3.143	3.707	4.317	5.208	5.959
7	.711	.896	1.119	1.415	1.895	2.365	2.517	2.998	3.499	4.029	4.785	5.408
8	.706	.889	1.108	1.397	1.860	2.306	2.449	2.896	3.355	3.833	4.501	5:041
9	.703	.883	1.100	1.383	1.833	2.262	2.398	2.821	3.250	3.690	4.297	4.781
10	.700	.879	1.093	1.372	1.812	2.228	2,359	2.764	3.169	3.581	4.144	4.587
11	.697	.876	1.088	1.363	1.796	2,201	2.328	2.718	3.106	3.497	4.025	4.437
12	.695	.873	1.083	1.356	1.782	2.179	2.303	2.681	3.055	3.428	3.930	4.318
13	.694	.870	1.079	1.350	1.771	2.160	2.282	2.650	3.012	3.372	3.852	4.221
14	.692	.868	1.076	1.345	1.761	2.145	2.264	2.624	2.977	3.326	3.787	4.140
15	.691	.866	1.074	1.341	1.753	2.131	2.249	2.602	2.947	3.286	3.733	4.073
16	.690	.865	1.071	1.337	1.746	2.120	2.235	2.583	2.921	3.252	3.686	4.015
17	.689	.863	1.069	1.333	1.740	2.110	2.224	2.567	2.898	3.222	3.646	3.965
18	.688	.862	1.067	1.330	1.734	2.101	2.214	2.552	2.878	3.197	3.611	3.922
19	.688	.861	1.066	1.328	1.729	2.093	2.205	2.539	2.861	3.174	3.579	3.883
20	.687	.860	1.064	1.325	1.725	2.086	2.197	2.528	2.845	3.153	3.552	3.850
21	.686	.859	1.063	1.323	1.721	2.080	2.189	2.518	2.831	3.135	3.527	3.819
22	.686	.858	1.061	1.321	1.717	2.074	2.183	2.508	2.819	3.119	3.505	3.792
23	.685	.858	1.060	1.319	1.714	2.069	2.177	2.500	2.807	3.104	3.485	3.768
24	685	.857	1.059	1.318	1.711	2.064	2.172	2.492	2.797	3.091	3.467.	3.745
25	.684	.856	1.058	1.316	1.708	2.060	2.167	2.485	2.787	3.078	3.450	3.725
26	.684	.856	1.058	1.315	1.706	2.056	2.162	2.479	2.779	3.067	3.435	3.707
27	.684	.855	1.057	1.314	1.703	2.052	2.158	2.473	2.771	3.057	3.421	3.690
28	.683	.855	1.056	1.313	1.701	2.048	2.154	2.467	2.763	3.047	3.408	3.674
29	.683	.854	1.055	1.311	1.699	2.045	2.150	2.462	2.756	3.038	3.396	3.659
30	.683	.854	1.055	1.310	1.697	2.042	2.147	2:457	2.750	3.030	3.385	3.646
40	.681	.851	1.050	1.303	1.684	2.021	2.123	2.423	2.704	2.971	3.307	3.551
50	.679	.849	1.047	1.299	1.676	2.009	2.109	2.403	2.678	2.937	3.261	3.496
60	.679	.848	1.045	1.296	1.671	2.000	2.099	2.390	2.660	2.915	3.232	3,460
80	.678	.846	1.043	1.292	1.664	1.990	2.088	2.374	2.639	2.887	3.195	3.416
100	.677	.845	1.042	1.290	1.660	1.984	2.081	2.364	2.626	2.871	3.174	3.390
000	.675	.842	1.037	1.282	1.646	1.962	2.056	2.330	2.581	2.813	3.098	3,300
••	.674	.841	1.036	1.282	1.645	1.960	2.054	2.326	2.576	2.807	3.091	3.291
	50%	60%	70%	80%	90%	95%	96%	98%	99%	99.5%	99.8%	99.9%

- are *significantly* different for 95% of population
- **Confidence limit**, *z* = sig/2

Confidence level C

Estimating Confidence Intervals: Statistical Significance

- Are M₁ & M₂ significantly different?
 - Compute t. Select significance level (e.g. sig = 5%)
 - Consult table for t-distribution: Find t value corresponding to k-1 degrees of freedom (here, 9)
- t-distribution is symmetric: typically upper % points of distribution shown → look up value for confidence limit z=sig/2 (here, 0.025)
 - If t > z or t < -z, then t value lies in rejection region:</p>
 - Reject null hypothesis that mean error rates of M₁ & M₂ are same
- Conclude: <u>statistically significant</u> difference between M₁
 & M₂
 - Otherwise, conclude that any difference is chance

www.SadiEvrenSEKER.com www.SadiEvr

Model Selection: ROC Curves⁴⁴

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



ue positive raix

- Vertical axis represents the true positive rate
- Horizontal axis rep. the false positive rate
- The plot also shows a diagonal line
- A model with perfect accuracy will have an area of 1.0

Issues Affecting Model Selection

- Accuracy SEKER.com www.SadiEvrenSEKER.com www.Sadil
 - classifier accuracy: predicting class label
- Speed Speed Seker.com www.SadiEvrenSeker.com www.Sadil
 - time to construct the model (training time)
 - time to use the model (classification/prediction time)
- **Robustness**: handling noise and missing values
- Scalability: efficiency in disk-resident databases
- Interpretability
 - understanding and insight provided by the model
- Other measures, e.g., goodness of rules, such as decision tree size or compactness of classification rules

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiI

Chapter 8. Classification: Basic Concepts Classification: Basic Concepts **Decision Tree Induction Bayes Classification Methods** Rule-Based Classification Model Evaluation and Selection Techniques to Improve Classification Accuracy: **Ensemble Methods** Summary

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

Ensemble wethods: increasing the

www.SadiEvrenSEKER.co

www.SadiEvrenSEKER.com w www.SadiEvrenSEKER.com w www.SadiEvrenSEKER.com w



- Ensemble methods
 - Use a combination of models to increase accuracy
- Combine a series of k learned models, M₁, M₂, ..., M_k, with the aim of creating an improved model M*
 - Popular ensemble methods
 - Bagging: averaging the prediction over a collection of classifiers
 - Boosting: weighted vote with a collection of classifiers
- Ensemble: combining a set of heterogeneous classifiers

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sada

Bagging: Boostrap Aggregation

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

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE

www.SadiEvrenSEKER.com_www.SadiEvrenSEKER.com_www.SadiE www.SadiEvrenSEKER.cor<mark>Boosting</mark>renSEKER.com_www.SadiE

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

Adaboost (Freund and Schapire, 1997)

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE

- Given a set of *d* class-labeled tuples, (X₁, y₁), ..., (X_d, y_d)
- Initially, all the weights of tuples are set the same (1/d)
- Generate k classifiers in k rounds. At round i,
 - Tuples from D are sampled (with replacement) to form a training set D_i of the same size
 - Each tuple's chance of being selected is based on its weight
 - A classification model M_i is derived from D_i
 - Its error rate is calculated using D_i as a test set
 - If a tuple is misclassified, its weight is increased, o.w. it is decreased
 - Error rate: err(X_j) is the misclassification error of tuple X_j. Classifier M_i error rate is the sum of the weights of the misclassified tuples:

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

• The weight of classifier M_i 's vote is $\log \frac{1 - error(M_i)}{error(M_i)}$

Random Forest (Breiman 2001)

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

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

Classification of Class-Imbalanced Data Sets

- Class-imbalance problem: Rare positive example but numerous negative ones, e.g., medical diagnosis, fraud, oil-spill, fault, etc.
- Traditional methods assume a balanced distribution of classes and equal error costs: not suitable for class-imbalanced data
 - Typical methods for imbalance data in 2-class classification:
 - **Oversampling**: re-sampling of data from positive class
 - Under-sampling: randomly eliminate tuples from negative class
 - Threshold-moving: moves the decision threshold, t, so that the rare class tuples are easier to classify, and hence, less chance of costly false negative errors
 - Ensemble techniques: Ensemble multiple classifiers introduced above

Still difficult for class imbalance problem on multiclass tasks

Chapter 8. Classification: Basic Concepts Classification: Basic Concepts **Decision Tree Induction Bayes Classification Methods Rule-Based Classification** Model Evaluation and Selection Techniques to Improve Classification Accuracy:

69

Ensemble Methods



www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER **Summary (I)** SEKER.com www.SadiE

- Classification is a form of data analysis that extracts models describing important data classes.
- Effective and scalable methods have been developed for decision tree induction, Naive Bayesian classification, rule-based classification, and many other classification methods.
- Evaluation metrics include: accuracy, sensitivity, specificity, precision, recall, *F* measure, and *F_β* measure.
- Stratified k-fold cross-validation is recommended for accuracy estimation. Bagging and boosting can be used to increase overall accuracy by learning and combining a series of individual models.

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER Summary (II) nSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadievrenSEKER.com www.SadiE

- Significance tests and ROC curves are useful for model selection.
- There have been numerous comparisons of the different classification methods; the matter remains a research topic
- No single method has been found to be superior over all others for all data sets
- Issues such as accuracy, training time, robustness, scalability, and interpretability must be considered and can involve tradeoffs, further complicating the quest for an overall superior method

www.SadiEvrenSEKER.com www.SadiEvrenSEK

www.saulevrensekek.com www.saulevrensekek.com www.sauli

- C. Apte and S. Weiss. Data mining with decision trees and decision rules. Future Generation Computer Systems, 13, 1997
- C. M. Bishop, Neural Networks for Pattern Recognition. Oxford University Press, 1995
- L. Breiman, J. Friedman, R. Olshen, and C. Stone. Classification and Regression Trees.
 Wadsworth International Group, 1984
 - C. J. C. Burges. A Tutorial on Support Vector Machines for Pattern Recognition. Data Mining and Knowledge Discovery, 2(2): 121-168, 1998
- P. K. Chan and S. J. Stolfo. Learning arbiter and combiner trees from partitioned data for scaling machine learning. KDD'95
- H. Cheng, X. Yan, J. Han, and C.-W. Hsu,
 <u>Discriminative Frequent Pattern Analysis for Effective Classification</u>, ICDE'07
- H. Cheng, X. Yan, J. Han, and P. S. Yu,
 <u>Direct Discriminative Pattern Mining for Effective Classification</u>, ICDE'08
 - W. Cohen. Fast effective rule induction. ICML'95
- G. Cong, K.-L. Tan, A. K. H. Tung, and X. Xu. Mining top-k covering rule groups for gene expression data. SIGMOD'05

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi
References (2)

- A. J. Dobson. An Introduction to Generalized Linear Models. Chapman & Hall, 1990.
- G. Dong and J. Li. Efficient mining of emerging patterns: Discovering trends and differences. KDD'99.
- R. O. Duda, P. E. Hart, and D. G. Stork. Pattern Classification, 2ed. John Wiley, 2001
- U. M. Fayyad. Branching on attribute values in decision tree generation. AAAI' 94.
- Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. J. Computer and System Sciences, 1997.
- J. Gehrke, R. Ramakrishnan, and V. Ganti. Rainforest: A framework for fast decision tree construction of large datasets. VLDB' 98.
- J. Gehrke, V. Gant, R. Ramakrishnan, and W.-Y. Loh, **BOAT -- Optimistic Decision Tree Construction**. SIGMOD'99.
- T. Hastie, R. Tibshirani, and J. Friedman. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer-Verlag, 2001.
- D. Heckerman, D. Geiger, and D. M. Chickering. Learning Bayesian networks: The combination of knowledge and statistical data. Machine Learning, 1995.
- W. Li, J. Han, and J. Pei, CMAR: Accurate and Efficient Classification Based on Multiple Class-Association Rules, ICDM'01.

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

References (3)

www.sadievrensekek.com www.sadievrensekek.com www.sadii

- 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 and R. M. Cameron-Jones. **FOIL: A midterm report**. ECML' 93.
- J. R. Quinlan. **C4.5: Programs for Machine Learning**. Morgan Kaufmann, 1993.
- J. R. Quinlan. Bagging, boosting, and c4.5. AAAI'96.

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

References (4)

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

- J. Shafer, R. Agrawal, and M. Mehta. SPRINT : A scalable parallel classifier for data mining. VLDB' 96.
- J. W. Shavlik and T. G. Dietterich. Readings in Machine Learning. Morgan Kaufmann, 1990.
- P. Tan, M. Steinbach, and V. Kumar. Introduction to Data Mining. Addison Wesley, 2005.
- S. M. Weiss and C. A. Kulikowski. Computer Systems that Learn: Classification and Prediction Methods from Statistics, Neural Nets, Machine Learning, and Expert Systems. Morgan Kaufman, 1991.
- S. M. Weiss and N. Indurkhya. **Predictive Data Mining**. Morgan Kaufmann, 1997.
- I. H. Witten and E. Frank. Data Mining: Practical Machine Learning Tools and Techniques, 2ed. Morgan Kaufmann, 2005.
- X. Yin and J. Han. CPAR: Classification based on predictive association rules. SDM'03
 - H. Yu, J. Yang, and J. Han. Classifying large data sets using SVM with hierarchical clusters. KDD'03.

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadj



CS412 Midterm Exam Statistics

Opinion Question Answering: Like the style: 70.83%, dislike: 29.16% Exam is hard: 55.75%, easy: 0.6%, just right: 43.63% Time: plenty:3.03%, enough: 36.96%, not: 60% Score distribution: # of students (Total: 180) <40:2 >=90: 24 60-69:37 **80-89:54** 50-59:15 **70-79:46** 40-49:2 Final grading are based on overall score accumulation and relative class distributions

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE

Issues: Evaluating Classification Methods

- Accuracy
 - classifier accuracy: predicting class label
 - predictor accuracy: guessing value of predicted attributes
- Speed
- time to construct the model (training time) ER.com www.Sadil
 - time to use the model (classification/prediction time)
 - Robustness: handling noise and missing values
 - Scalability: efficiency in disk-resident databases
 - Interpretability
 - understanding and insight provided by the model
- Other measures, e.g., goodness of rules, such as decision tree size or compactness of classification rules

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sada

Predictor Error Measures

- www.sadievrensekek.com/www.sadievrensekek.com/www.sadił
 - Measure predictor accuracy: measure how far off the predicted value is from the actual known value
 - Loss function: measures the error betw. y_i and the predicted value y_i'
 - Absolute error: $|y_i y_i'|$
 - Squared error: $(y_i y_i')^2$
 - Test error (generalization error): the average loss over the test set • Mean absolute error: $\sum_{i=1}^{d} |y_i - y_i|^2$ Mean squared error: $\sum_{i=1}^{d} (y_i - y_i')^2$
 - Mean absolute error: $\sum_{i=1}^{d} |y_i y_i|^d$ ean squared error: $\frac{d}{d}$

• Relative absolute error: $\frac{\sum_{i=1}^{n} |y_i| \mathbf{R} \vec{e} \mathbf{I} \mathbf{a}$ tive squared error: $\frac{\sum_{i=1}^{d} |y_i - \overline{y}|}{\sum_{i=1}^{d} |y_i - \overline{y}|}$

The mean squared-error exaggerates the presence of outliers Popularly use (square) root mean-square error, similarly, root relative squared error

 $\sum_{i=1}^{n} (y_i - y_i')^2$

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sara

Methods

SLIQ (EDBT' 96 — Mehta et al.)

- Builds an index for each attribute and only class list and the current attribute list reside in memory
- SPRINT (VLDB' 96 J. Shafer et al.)
 - Constructs an attribute list data structure
- PUBLIC (VLDB' 98 Rastogi & Shim)
- Integrates tree splitting and tree pruning: stop growing the tree earlier

RainForest (VLDB' 98 — Gehrke, Ramakrishnan & Ganti)

- Builds an AVC-list (attribute, value, class label)
- BOAT (PODS' 99 Gehrke, Ganti, Ramakrishnan & Loh)

Uses bootstrapping to create several small samples

Data Cube-Based Decision-Tree Induction

- Integration of generalization with decision-tree induction (Kamber et al.' 97)
 - Classification at primitive concept levels
- E.g., precise temperature, humidity, outlook, etc.
- Low-level concepts, scattered classes, bushy classificationtrees
 - Semantic interpretation problems
 - Cube-based multi-level classification
 - Relevance analysis at multi-levels
- Information-gain analysis with dimension + level

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE

Veri Ön İşleme (Data Preprocessing) Data Preprocessing: Giriş Veri Kalitesi Veri Ön işlemedeki ana işlemler Veri Temizleme (Data Cleaning) Veri Uyumu (Data Integration) Veri Küçültme (Data Reduction) Veri Dönüştürme ve Verinin Ayrıklaştırılması (Data Transformation and Data Discretization) 82 www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE

Data Warehouse: A Multi-Tiered Architecture



www.SadiEvrenSEKER.com Way SadiEvrenSEKER.com Www.SadiEvrenSEKER.com SadiEvrenSEKER.com SadiEvrenSEKER.com SadiE

Çok boyutlu olarak veri kalitesi kriterleri : Neden Ön işlem yapılır?

- Kesinlik (Accuracy) doğru ve yanlış veriler
- Tamamlık (Completeness) : kaydedilmemiş veya ulaşılamayan veriler
- Tutarlılık (Consistency) verilerin bir kısmının güncel olmaması, sallantıda veriler (dangling)
- Güncellik (Timeliness)
- İnandırıcılık (Believability)
- Yorumlanabilirlik (Interpretability): Verinin ne kadar kolay
- anlaşılacağı

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadal

www.SadiEvrenSEKER.com Www.SadiEvrenSEKER.com Www.SadiE www.SadiEvren**Veri Ön İşleme İşlemleri**

Veri Temizleme (Data cleaning)

- Eksik verilerin doldurulması, gürültülü verilerin düzeltilmesi, aykırı verilerin (outlier) temizlenmesi, uyuşmazlıkların (inconsistencies) çözümlenmesi
- Veri Entegrasyonu (Data integration)
 - Farklı veri kaynaklarının, Veri Küplerinin veya Dosyaların entegre olması

Verinin Küçültülmesi (Data reduction)

- Boyut Küçültme (Dimensionality reduction)
- Sayısal Küçültme (Numerosity reduction)
- Verinin Sıkıştırılması (Data compression)
- Verinin Dönüştürülmesi ve Ayrıklaştırılması (Data transformation and data discretization)
 - Normalleştirme (Normalization)
 - Kavram Hiyerarşisi (Concept hierarchy generation)

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

Veri Ön İşleme (Data Preprocessing) Data Preprocessing: Giriş Veri Kalitesi Veri Ön işlemedeki ana işlemler Veri Temizleme (Data Cleaning) Veri Uyumu (Data Integration) Veri Küçültme (Data Reduction) Veri Dönüştürme ve Verinin Ayrıklaştırılması (Data Transformation and Data Discretization) 86

Veri Temizleme (Data Cleaning)

- Gerçek hayattaki veriler kirlidir: Çok sayıda makine, insan veya bilgisayar hataları, iletim bozulmaları yaşanabilir.
 - <u>Eksik Veri (incomplete</u>) bazı özelliklerin eksik olması (missing
 - data), sadece birleşik verinin (aggregate) bulunması
 - örn., Meslek="" (girilmemiş)
 - Gülrültülü Veri (noisy): Gürültü, hata veya aykırı veriler bulunması
 - örn., *Maaş*="-10" (hata)
 - <u>Tutarsız Veri (inconsistent)</u>: farklı kaynaklardan farklı veriler gelmesi
 - Yaş="42", Doğum Tarihi="03/07/2010"
 - Eski notlama "1, 2, 3", yeni notlama "A, B, C"
 - Tekrarlı kayıtlarda uyuşmazlık
 - Kasıtlı Problemler (Intentional)
 - Doğum tarihi bilinmeyen herkese 1 Ocak yazılması

Eksik Veriler (Incomplete (Missing) Data)

Veriye her zaman erişilmesi mümkün değildir

- Örn., bazı kayıtların alın(a)mamış olması. Satış sırasında müşterilerin gelir düzeyinin yazılmamış olması.
- Eksik veriler genelde aşağıdaki durumlarda olur:
 - Donanımsal bozukluklardan
 - Uyuşmazlık yüzünden silinen veriler
 - Anlaşılamayan verilerin girilmemiş olması
 - Veri girişi sırasında veriye önem verilmemiş olması
 - Verideki değişikliklerin kaydedilmemiş olması

Eksik verilerin çözülmesi gerekir
 88

Eksik veriler nasıl çözülür?

- İhmal etme: Eksik veriler işleme alınmaz, yokmuş gibi davranılır. Kullanılan VM yöntemine göre sonuca etkileri bilinmelidir.
- Eksik verilerin elle doldurulması: her zaman mümkün değildir ve bazan çok uzun ve maliyetli olabilir
 - Otomatik olarak doldurulması
 - Bütün eksik veriler için yeni bir sınıf oluşturulması ("bilinmiyor" gibi)
 - Ortalamanın yazılması
 - Sınıf bazında ortalamaların yazılması
 - Bayesian formül ve karar ağacı uygulaması

Gürültülü Veri (Noisy Data)

- Gürültü (Noise): ölçümdeki rasgele oluşan değerler
 - Yanlış özellik değerleri aşağıdaki durumlarda oluşabilir:
 - Veri toplama araçlarındaki hatalar
 - Veri giriş problemleri
 - Veri iletim problemleri
 - Teknoloji sınırları
 - İsimlendirmedeki tutarsızlıklar
- Veri temizlemesini gerektiren diğer durumlar
 - Tekrarlı kayıtlar
 - Eksik veriler
 - Tutarsız veriler

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sage

Gürültülü Veri Nasıl Çözülür?

- Paketleme (Binning)
 - Veri sıralanır ve eşit frekanslarda paketlere bölünür.
 - Eksik veriler farklı yöntemlerle doldurulur:
 - Mean
- Median
 Boundary
- Regrezisyon (Regression)
 - Regrezisyon fonksiyonlarına tabi tutularak eksik verilerin girilmesi
 - Bölütleme (Kümeleme , Clustering)
 - Aykırı verilerin bulunması ve temizlenmesi
- Bilgisayar ve insan bilgisinin ortaklaşa kullanılması
- detect suspicious values and check by human (e.g., 91)

Veri Temizleme Süreci

- Verideki farklılıkların yakalanması
 - Üst verinin (metadata) kullanılması (örn., veri alanı (domain, range), bağlılık (dependency), dağılım (distribution)
 - Aşırı yüklü alanlar (Field Overloading)
 - Veri üzerinde kural kontrolleri (unique, consecutive, null)
 - Ticari yazılımların kullanılması
 - Bilgi Ovalaması (Data scrubbing): Basit alan bilgileri kurallarla kontrol otmok (o.g. postal codo, spoll chock)
 - kontrol etmek (e.g., postal code, spell-check)
 - Veri Denetimi (Data auditing): veriler üzerinden kural çıkarımı ve kurallara uymayanların bulunması (örn., correlation veya clustering ile aykırıların (outliers) bulunması)
- Veri Göçü ve Entegrasyonu (Data migration and integration)
 - Data migration Araçları: Verinin dönüştürülmesine izin verir
 - ETL (Extraction/Transformation/Loading) Araçları: Genelde grafik arayüzü ile dönüşümü yönetme imkanı verir

92

- İki farklı işin entegre yürütülmesi
 - Iterative / interactive (Örn.., Potter's Wheels)



www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.con Örnekler.renSEKER.com www.SadiE



Vijayshankar Raman and Joseph M. Hellerstein , **Potter's Wheel: An Interactive Data Cleaning System** berkeley

Chapter 3: Data Preprocessing Data Preprocessing: An Overview Data Quality Major Tasks in Data Preprocessing Data Cleaning **Data Integration** 44 **Data Reduction** Data Transformation and Data Discretization Summary 95

Data Integration

- Data integration:
 - Combines data from multiple sources into a coherent store
 - Schema integration: e.g., A.cust-id = B.cust-#
 - Integrate metadata from different sources
- Entity identification problem:
 - Identify real world entities from multiple data sources, e.g., Bill
- Clinton = William Clinton
 - Detecting and resolving data value conflicts
 - For the same real world entity, attribute values from different sources are different

 Possible reasons: different representations, different scales, e.g., metric vs. British units Handling Redundancy in Data Integration

 Redundant data occur often when integration of multiple databases

 Object identification: The same attribute or object may have different names in different databases

 Derivable data: One attribute may be a "derived" attribute in another table, e.g., annual revenue

- Redundant attributes may be able to be detected by correlation analysis and covariance analysis
- Careful integration of the data from multiple sources may help reduce/avoid redundancies and inconsistencies and improve mining speed and quality

97

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE

www.SCorrelation Analysis (Nominal Data) w Sadi

X² (chi-square) test

- $\chi^{2} = \sum \frac{(Observed Expected)^{2}}{Expected}$ The larger the X² value, the more likely the variables are related
- The cells that contribute the most to the X² value are those whose actual count is very different from the expected count
- Correlation does not imply causality
- # of hospitals and # of car-theft in a city are correlated

Both are causally linked to the third variable: population
 98

Chi-Square Calculation: An Example

Image: Play chessNot play chessSum (row)Like science fiction250(90)200(360)450Not like science fiction50(210)1000(840)1050Sum(col.)30012001500

www.Sadi www.Sadi

 X² (chi-square) calculation (numbers in parenthesis are expected counts calculated based on the data distribution in the two categories)

 $\chi^{2} = \frac{(250-90)^{2}}{90} + \frac{(50-210)^{2}}{210} + \frac{(200-360)^{2}}{360} + \frac{(1000-840)^{2}}{840} = 507.93$ It shows that like_science_fiction and play_chess are correlated in the group

Correlation Analysis (Numeric Data)

 Correlation coefficient (also called Pearson's product moment coefficient)

 $r_{A,B} = \frac{\sum_{i=1}^{n} (a_i - \overline{A})(b_i - \overline{B})}{(n-1)\sigma_A \sigma_B} = \frac{\sum_{i=1}^{n} (a_i b_i) - n\overline{AB}}{(n-1)\sigma_A \sigma_B}$

where n is the number of tuples, \overline{A} and \overline{B} are the respective means of A and B, σ_A and σ_B are the respective standard deviation of A and B, and $\Sigma(a_ib_i)$ is the sum of the AB cross-product.

If r_{A,B} > 0, A and B are positively correlated (A's values increase as B's). The higher, the stronger correlation.

r_{A,B} = 0: independent; r_{AB} < 0: negatively correlated

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa100

Visually Evaluating Correlation



KER.com www.Sadif KER.com www.Sadif KER.com www.Sadif KER.com www.Sadif

Scatter plots showing the similarity from -1 to 1.

KER.com www.Sadił KER.com www.Sadił KER.com www.Sadił KER.com www.Sadił

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa₁₀₁

Correlation (viewed as linear relationship)

- Correlation measures the linear relationship between objects
- To compute correlation, we standardize data objects, A and B, and then take their dot product

 $a'_{k} = (a_{k} - mean(A)) / std(A)$ $b'_{k} = (b_{k} - mean(B)) / std(B)$ $correlation(A, B) = A' \bullet B'$

Covariance (Numeric Data)

Covariance is similar to correlation

$$Cov(A, B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^{n} (a_i - A)(b_i - B)}{n}$$

Correlation coefficient: $r_{A,B} = \frac{COV(A,B)}{\sigma_A \sigma_B}$

- where n is the number of tuples, \overline{A} and \overline{B} are the respective mean or **expected values** of A and B, σ_A and σ_B are the respective standard deviation of A and B.
 - Positive covariance: If Cov_{A,B} > 0, then A and B both tend to be larger than their expected values.
 - Negative covariance: If Cov_{A,B} < 0 then if A is larger than its expected value, B is likely to be smaller than its expected value.
 - Independence: Cov_{A,B} = 0 but the converse is not true:
 - Some pairs of random variables may have a covariance of 0 but are not independent. Only under some additional assumptions (e.g., the data follow multivariate normal distributions) does a covariance of 0 imply independence 103

Co-Variance: An Example

$$Cov(A,B) = E((A - \bar{A})(B - \bar{B})) = \frac{\sum_{i=1}^{n} (a_i - \bar{A})(b_i - \bar{B})}{n}$$

It can be simplified in computation as

$$Cov(A,B) = E(A \cdot B) - \overline{A}\overline{B}$$

- Suppose two stocks A and B have the following values in one week:
 (2, 5), (3, 8), (5, 10), (4, 11), (6, 14).
- Question: If the stocks are affected by the same industry trends, will their prices rise or fall together?
 - E(A) = (2 + 3 + 5 + 4 + 6)/5 = 20/5 = 4

• E(B) = (5 + 8 + 10 + 11 + 14) / 5 = 48 / 5 = 9.6

• $Cov(A,B) = (2 \times 5 + 3 \times 8 + 5 \times 10 + 4 \times 11 + 6 \times 14)/5 - 4 \times 9.6 = 4$

Thus, A and B rise together since Cov(A, B) > 0.

Chapter 3: Data Preprocessing Data Preprocessing: An Overview Data Quality Major Tasks in Data Preprocessing Data Cleaning **Data Integration** 44 Data Reduction Data Transformation and Data Discretization Summary 105

Data Reduction Strategies

- Data reduction: Obtain a reduced representation of the data set that is much smaller in volume but yet produces the same (or almost the same) analytical results
- Why data reduction? A database/data warehouse may store terabytes of data. Complex data analysis may take a very long time to run on the complete data set.
- Data reduction strategies
 - Dimensionality reduction, e.g., remove unimportant attributes
 - Wavelet transforms
 - Principal Components Analysis (PCA)
 - Feature subset selection, feature creation
 - Numerosity reduction (some simply call it: Data Reduction)
 - Regression and Log-Linear Models
 - Histograms, clustering, sampling
 - Data cube aggregation
 - Data compression
- www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa₁₀₆I

Data Reduction 1: Dimensionality Reduction

Curse of dimensionality

- When dimensionality increases, data becomes increasingly sparse
- Density and distance between points, which is critical to clustering, outlier analysis, becomes less meaningful
 - The possible combinations of subspaces will grow exponentially

Dimensionality reduction

- Avoid the curse of dimensionality
- Help eliminate irrelevant features and reduce noise
- Reduce time and space required in data mining
- Allow easier visualization
- **Dimensionality reduction techniques**
 - Wavelet transforms
 - Principal Component Analysis
 - Supervised and nonlinear techniques (e.g., feature selection)

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

Mapping Data to a New Space

Fourier transformWavelet transform



www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com
What Is Wavelet Transform?

- Decomposes a signal into different frequency subbands
 Applicable to ndimensional signals
 Data are transformed to
 - Data are transformed to preserve relative distance between objects at different levels of resolution
 - Allow natural clusters to become more distinguishable
 - Used for image compression



www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa₁₀₉

www.SadiEvrenSEKER.com www.SadiEvrenSEK

Wavelet Transformation



- Discrete wavelet transform (DWT) for linear signal processing, multi-resolution analysis
- Compressed approximation: store only a small fraction of the strongest of the wavelet coefficients
- Similar to discrete Fourier transform (DFT), but better lossy compression, localized in space
- Method:
 - Length, L, must be an integer power of 2 (padding with 0's, when necessary)
 - Each transform has 2 functions: smoothing, difference
 - Applies to pairs of data, resulting in two set of data of length L/2
 - Applies two functions recursively, until reaches the desired length

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Satio

Wavelet Decomposition

- Wavelets: A math tool for space-efficient hierarchical decomposition of functions
 - S = [2, 2, 0, 2, 3, 5, 4, 4] can be transformed to S_{$^$} = [2³/₄, -1¹/₄, ¹/₂, 0, 0, -1, -1, 0]
- Compression: many small detail coefficients can be replaced by 0's, and only the significant coefficients are retained

WWV	Resolution	Averages	Detail Coefficients	w.Sadil
WWV	8	[2, 2, 0, 2, 3, 5, 4, 4]		w.Sadil
	4	[2,1,4,4]	$[0,\ -1,\ -1,\ 0]$	w SadiF
	2	$[1\frac{1}{2}, 4]$	$[\frac{1}{2}, 0]$	TY I DUUT
WWV	1	$[\overline{2}\frac{3}{4}]$	$[-1\frac{1}{4}]$	w.Sadil

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi



Why Wavelet Transform?

- Use hat-shape filters
 - Emphasize region where points cluster
 - Suppress weaker information in their boundaries
 - Effective removal of outliers
 - Insensitive to noise, insensitive to input order
- Multi-resolution
- Efficient enseker.com www.SadiEvrenSEKER.com www.SadiE
 - Complexity O(N)
 - Only applicable to low dimensional data

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi3 www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE

Principal Component Analysis (PCA)

- Find a projection that captures the largest amount of variation in data
- The original data are projected onto a much smaller space, resulting in dimensionality reduction. We find the eigenvectors of the covariance matrix, and these eigenvectors define the new space

 X_2 e \mathbf{X}_{1} 114

Principal Component Analysis (Steps)

- Given N data vectors from *n*-dimensions, find $k \le n$ orthogonal vectors (*principal components*) that can be best used to represent data
 - Normalize input data: Each attribute falls within the same range
 - Compute *k* orthonormal (unit) vectors, i.e., *principal components*
- Each input data (vector) is a linear combination of the k principal component vectors
 - The principal components are sorted in order of decreasing "significance" or strength
- Since the components are sorted, the size of the data can be reduced by eliminating the *weak components*, i.e., those with low variance (i.e., using the strongest principal components, it is possible to reconstruct a good approximation of the original data)

Works for numeric data only

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Safis

Attribute Subset Selection

- Another way to reduce dimensionality of data
- Redundant attributes
 - Duplicate much or all of the information contained in one or more other attributes
 - E.g., purchase price of a product and the amount of sales tax paid
- Irrelevant attributes
 - Contain no information that is useful for the data mining task at hand
 - E.g., students' ID is often irrelevant to the task of predicting students' GPA

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa116

Heuristic Search in Attribute Selection

- There are 2^d possible attribute combinations of d attributes
- Typical heuristic attribute selection methods:
 - Best single attribute under the attribute independence assumption: choose by significance tests
- Best step-wise feature selection:
 - The best single-attribute is picked first
 Then next best attribute condition to the first, ...
 - Step-wise attribute elimination:
 - Repeatedly eliminate the worst attribute
 - Best combined attribute selection and elimination
 - Optimal branch and bound:
 - Use attribute elimination and backtracking

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Saiij

Attribute Creation (Feature Generation)

- Create new attributes (features) that can capture the important information in a data set more effectively than the original ones
- Three general methodologies
 - Attribute extraction
 - Domain-specific
 - Mapping data to new space (see: data reduction)
 - E.g., Fourier transformation, wavelet transformation, manifold approaches (not covered)
 - Attribute construction
- Combining features (see: discriminative frequent patterns in Chapter 7)
 Data discretization

Data Reduction 2: Numerosity Reduction

- Reduce data volume by choosing alternative, smaller forms of data representation
 - Parametric methods (e.g., regression)
- Assume the data fits some model, estimate model parameters, store only the parameters, and discard the data (except possible outliers)
- Ex.: Log-linear models—obtain value at a point in *m*-D space as the product on appropriate marginal subspaces
- Non-parametric methods
 - Do not assume models

Major families: histograms, clustering, sampling, ...

Parametric Data Reduction: Regression and Log-Linear Models

Linear regression

- Data modeled to fit a straight line
- Often uses the least-square method to fit the line

Multiple regression

 Allows a response variable Y to be modeled as a linear function of multidimensional feature vector

Log-linear model

Approximates discrete multidimensional probability distributions

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa₁₂₀E

Regression Analysis

- Regression analysis: A collective name for techniques for the modeling and analysis of numerical data consisting of values of a *dependent variable* (also called *response variable* or *measurement*) and of one or more *independent variables* (aka.
 explanatory variables or *predictors*)
- The parameters are estimated so as to give a "best fit" of the data
- Most commonly the best fit is evaluated by using the *least squares method*, but other criteria have also been used



SEKER.com www.Sadil

Used for prediction (including forecasting of time-series data), inference, hypothesis testing, and modeling of causal relationships

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Saiil

Regress Analysis and Log-Linear Models

- Linear regression: Y = w X + b
 - Two regression coefficients, w and b, specify the line and are to be estimated by using the data at hand
 - Using the least squares criterion to the known values of Y₁, Y₂, ...,
 X₁, X₂,
- <u>Multiple regression</u>: $Y = b_0 + b_1 X_1 + b_2 X_2$
 - Many nonlinear functions can be transformed into the above
- Log-linear models:
 - Approximate discrete multidimensional probability distributions
 - Estimate the probability of each point (tuple) in a multi-dimensional space for a set of discretized attributes, based on a smaller subset of dimensional combinations

172

Useful for dimensionality reduction and data smoothing



www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.coClusteringenSEKER.com www.SadiE

- Partition data set into clusters based on similarity, and store cluster representation (e.g., centroid and diameter) only
 - Can be very effective if data is clustered but not if data is "smeared"
- Can have hierarchical clustering and be stored in multidimensional index tree structures
- There are many choices of clustering definitions and clustering algorithms

Cluster analysis will be studied in depth in Chapter 10

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE

- www.sadievrensekek.com www.sadievrensekek.com www.sadi
- Sampling: obtaining a small sample s to represent the whole data set N
- Allow a mining algorithm to run in complexity that is potentially sub-linear to the size of the data
 - Key principle: Choose a representative subset of the data
- Simple random sampling may have very poor performance in the presence of skew
- Develop adaptive sampling methods, e.g., stratified sampling:
- Note: Sampling may not reduce database I/Os (page at a time)

Types of Sampling

Simple random sampling

- There is an equal probability of selecting any particular item
- Sampling without replacement
- Once an object is selected, it is removed from the population
 - Sampling with replacement
 - A selected object is not removed from the population
- Stratified sampling:
 - Partition the data set, and draw samples from each partition (proportionally, i.e., approximately the same percentage of the data)
- Used in conjunction with skewed data

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

Sampling: With or without Replacement www.SadiEvrenSEKER.com_www.SadiE SRSWOR (simple random) sample without replacement) SRSWR Raw Data 127



Data Cube Aggregation

- The lowest level of a data cube (base cuboid)
- The aggregated data for an individual entity of interest
- E.g., a customer in a phone calling data warehouse
- Multiple levels of aggregation in data cubes
- Further reduce the size of data to deal with
- Reference appropriate levels
- Use the smallest representation which is enough to solve the task
- Queries regarding aggregated information should be answered using data cube, when possible

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa129

Data Reduction 3: Data Compression

- String compression
 - There are extensive theories and well-tuned algorithms
 - Typically lossless, but only limited manipulation is possible without expansion
- Audio/video compression
 - Typically lossy compression, with progressive refinement
 - Sometimes small fragments of signal can be reconstructed without reconstructing the whole
 - Time sequence is not audio
 - Typically short and vary slowly with time
- Dimensionality and numerosity reduction may also be considered as forms of data compression

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa₁₃₀



Chapter 3: Data Preprocessing Data Preprocessing: An Overview Data Quality Major Tasks in Data Preprocessing Data Cleaning **Data Integration** YW **Data Reduction** Data Transformation and Data Discretization Summary 132

Data Transformation

- A function that maps the entire set of values of a given attribute to a new set of replacement values s.t. each old value can be identified with one of the new values
 - Methods
 - Smoothing: Remove noise from data
 - Attribute/feature construction
 - New attributes constructed from the given ones
 - Aggregation: Summarization, data cube construction
 - Normalization: Scaled to fall within a smaller, specified range
 - min-max normalization
 - z-score normalization
 - Sa Inormalization by decimal scaling
- Discretization: Concept hierarchy climbing

www.SadiEvrenSEKERNormalizationSEKER.com www.SadiE **Min-max normalization**: to [new_min_A, new_max_A] $v' = \frac{v - min_A}{.} (new_max_A - new_min_A) + new_min_A$ $max_{A} - min_{A}$ • Ex. Let income range \$12,000 to \$98,000 normalized to [0.0, 1.0]. Then \$73,000 is mapped to $\frac{73,600-12,000}{98,000-12,000}(1.0-0)+0=0.716$ **Z-score normalization** (μ : mean, σ : standard deviation): $v' = \frac{v - \mu_A}{v}$ \mathcal{O}_A $\frac{73,600 - 54,000}{16,000} = 1.225$ • Ex. Let $\mu = 54,000$, $\sigma = 16,000$. Then Normalization by decimal scaling Where *j* is the smallest integer such that Max(|v'|) < 1

Discretization

- Three types of attributes
 - Nominal—values from an unordered set, e.g., color, profession
 - Ordinal—values from an ordered set, e.g., military or academic
- www.SirankiirenSEKER.com/v

- Numeric—real numbers, e.g., integer or real numbers
- Discretization: Divide the range of a continuous attribute into intervals
 - Interval labels can then be used to replace actual data values
 - Reduce data size by discretization
 - Supervised vs. unsupervised
 - Split (top-down) vs. merge (bottom-up)
- Discretization can be performed recursively on an attribute
- Prepare for further analysis, e.g., classification

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

Data Discretization Methods

- Typical methods: All the methods can be applied recursively Binning Top-down split, unsupervised Histogram analysis Top-down split, unsupervised Clustering analysis (unsupervised, top-down split or bottom-up merge) Decision-tree analysis (supervised, top-down split)
- Correlation (e.g., χ²) analysis (unsupervised, bottom-up merge)

Simple Discretization: Binning

Equal-width (distance) partitioning

- Divides the range into N intervals of equal size: uniform grid
- if A and B are the lowest and highest values of the attribute, the width of intervals will be: W = (B A)/N.
- The most straightforward, but outliers may dominate presentation
- Skewed data is not handled well
- Equal-depth (frequency) partitioning
 - Divides the range into *N* intervals, each containing approximately
 - same number of samples
 - Good data scaling
 - Managing categorical attributes can be tricky

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sadi

Binning Methods for Data Smoothing

- Sorted data for price (in dollars): 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34
- * Partition into equal-frequency (equi-depth) bins:
 - Bin 1: 4, 8, 9, 15 - Bin 2: 21, 21, 24, 25
 - Bin 3: 26, 28, 29, 34
- * Smoothing by **bin means**:
 - Bin 1: 9, 9, 9, 9 - Bin 2: 23, 23, 23, 23
 - Bin 3: 29, 29, 29, 29
- * Smoothing by **bin boundaries**:
 - Bin 1: 4, 4, 4, 15
 - Bin 2: 21, 21, 25, 25
 - Bin 3: 26, 26, 26, 34
- www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa₁₃₈E

Labels (Binning vs. Clustering)







Equal frequency (binning)



K-means clustering leads to better results

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa

Discretization by Classification & Correlation Analysis

- Classification (e.g., decision tree analysis)
 - Supervised: Given class labels, e.g., cancerous vs. benign
 - Using *entropy* to determine split point (discretization point)
- Top-down, recursive split
 - Details to be covered in Chapter 7
- Correlation analysis (e.g., Chi-merge: χ²-based discretization)
 - Supervised: use class information
- Bottom-up merge: find the best neighboring intervals (those having similar distributions of classes, i.e., low χ² values) to merge
 Merge performed recursively, until a predefined stopping condition

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa₁₄₀

Concept Hierarchy Generation

- Concept hierarchy organizes concepts (i.e., attribute values) hierarchically and is usually associated with each dimension in a data warehouse
- Concept hierarchies facilitate <u>drilling and rolling</u> in data warehouses to view data in multiple granularity
- Concept hierarchy formation: Recursively reduce the data by collecting and replacing low level concepts (such as numeric values for *age*) by higher level concepts (such as *youth, adult,* or *senior*)
- Concept hierarchies can be explicitly specified by domain experts and/ or data warehouse designers
- Concept hierarchy can be automatically formed for both numeric and nominal data. For numeric data, use discretization methods shown.

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa141

Concept Hierarchy Generation for Nominal Data

- Specification of a partial/total ordering of attributes explicitly at the schema level by users or experts
 - street < city < state < country</pre>
 - Specification of a hierarchy for a set of values by explicit data grouping
 - {Urbana, Champaign, Chicago} < Illinois</p>
- Specification of only a partial set of attributes
 - E.g., only street < city, not others</p>
- Automatic generation of hierarchies (or attribute levels) by the analysis of the number of distinct values
- E.g., for a set of attributes: { street, city, state, country}
- www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa₁₄₂

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE

Automatic Concept Hierarchy Generation

- Some hierarchies can be automatically generated based on the analysis of the number of distinct values per attribute in the data set
 - The attribute with the most distinct values is placed at the lowest level of the hierarchy
 - Exceptions, e.g., weekday, month, quarter, year



Chapter 3: Data Preprocessing Data Preprocessing: An Overview Data Quality Major Tasks in Data Preprocessing Data Cleaning **Data Integration Data Reduction** Data Transformation and Data Discretization Summary 144
www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.SadiE www.SadiEvrenSEKER.co**Summary**enSEKER.com www.SadiE

- Data quality: accuracy, completeness, consistency, timeliness, believability, interpretability
- Data cleaning: e.g. missing/noisy values, outliers
 - **Data integration** from multiple sources:
 - Entity identification problem
- Remove redundancies
 - Detect inconsistencies
 - **Data reduction**
 - Dimensionality reduction
 - Numerosity reduction
 - Data compression
 - Data transformation and data discretization
 - Normalization
- Concept hierarchy generation

www.SadiEvrenSEKER.com www.SadiEvrenSEKER.com www.Sa145

Kaynaklar www.salolievicensienek.com **Data Mining: Concepts and Techniques, Third Edition** (The Morgan Third Edition Kaufmann Series in Data Management Systems) 3rd Edition by Jiawei Han (Author), Micheline Kamber (Author), Jian Pei (Author) DATA MINING Concepts and Techniques M< Jiawei Han | Micheline Kamber | Jian Pei

www.SadiEvrenSEKER.com www.SadiEvrenSEK

 D. P. Ballou and G. K. Tayi. Enhancing data quality in data warehouse environments. Comm. of ACM, 42:73-78, 1999

- A. Bruce, D. Donoho, and H.-Y. Gao. Wavelet analysis. *IEEE Spectrum*, Oct 1996
- T. Dasu and T. Johnson. Exploratory Data Mining and Data Cleaning. John Wiley, 2003
- J. Devore and R. Peck. *Statistics: The Exploration and Analysis of Data*. Duxbury Press, 1997.
- H. Galhardas, D. Florescu, D. Shasha, E. Simon, and C.-A. Saita. Declarative data cleaning: Language, model, and algorithms. *VLDB'01*
- M. Hua and J. Pei. Cleaning disguised missing data: A heuristic approach. KDD'07
- H. V. Jagadish, et al., Special Issue on Data Reduction Techniques. Bulletin of the Technical Committee on Data Engineering, 20(4), Dec. 1997
- H. Liu and H. Motoda (eds.). Feature Extraction, Construction, and Selection: A Data Mining Perspective. Kluwer Academic, 1998
 - J. E. Olson. Data Quality: The Accuracy Dimension. Morgan Kaufmann, 2003
- D. Pyle. Data Preparation for Data Mining. Morgan Kaufmann, 1999

- V. Raman and J. Hellerstein. Potters Wheel: An Interactive Framework for Data Cleaning and Transformation, VLDB' 2001
- T. Redman. Data Quality: The Field Guide. Digital Press (Elsevier), 2001
- R. Wang, V. Storey, and C. Firth. A framework for analysis of data quality research. IEEE Trans.
 Knowledge and Data Engineering, 7:623-640, 1995