

CSCI577 Final

1 CLASSIFICATION

Definition. Constructing a method of classifying new instances using information in a training set

- Naive Bayes (Conditional Probabilities)
- Decision Trees
- Logistic Regression
- Neural Networks

1.1 TDIDT

Definition. Top-Down Induction of Decision Trees

1.2 Adequacy

Definition. No two instances with the same values of all the attributes may belong to different classes. Naive bayes can still be used when this doesn't obtain, as it will still be able to obtain the probabilities of each class. kNN can still be used, as long as then multiple datapoints in the same location in euclidean space would still function as expected.

Application. TDIDT algorithm is negatively affected

Algorithm. Until no more splitting is possible:

- IF all the instances in the training set belong to the same class THEN return the value of the class
- ELSE
 - (1) (a) Select an attribute A to split on
 - (2) (b) Sort the instances in the training set into subsets, one for each value of attribute A
 - (3) (c) Return a tree with one branch for each non-empty subset
 - Each branch having a descendant subtree or a class value produced by applying the algorithm recursively

1.3 Overfitting

Understand the concept of overfitting and be able to tell how you would know that a classification system overfit

Definition. If classifier generates a decision tree (or other mechanism) too well adapted to the training set Performs well on training set, not well on other data. Some overfitting inevitable.

Remedy:

- Adjust a decision tree while it is being generated: Pre-pruning
- Modify the tree after creation: Post-pruning

1.3.1 Clashes. Two (or more) instances of a training set have identical attribute values, but different classification. Especially a problem for TDIDT's 'Adequacy condition'.

Stems from.

- Classification incorrectly recorded
- Recorded attributes insufficient - Would need more attributes, normally impossible

Solutions.

- Discard the branch to the clashing node from the node above
- Of clashing instances, assume the majority label

1.3.2 Prepruning. When pre-pruning, may reduce accuracy on training set but may be better on test set (and subsequent) data than unpruned classifier.

- (1) test whether a termination condition applies.
 - If so, current subset is treated as a 'clash set'
 - Resolve by 'delete branch,' 'majority voting,' etc.
- (2) two methods methods
 - Size Cutoff – prune of subset has fewer than X instances

- Maximum depth: prune if length of branch exceed Y

1.3.3 PostPruning.

- (1) look for a non-leaf nodes that have descendants of length 1.
- (2) In this tree, only node G and D are candidates for pruning (consolidation).

1.4 Discretizing

1.4.1 Equal Width Intervals.

1.4.2 Pseudo Attributes.

1.4.3 Processing Sorted Instance Table.

1.4.4 ChiMerge.

Rationalization. Initially, each distinct value of a numerical attribute A is considered to be one interval. χ^2 tests are performed for every pair of adjacent intervals. Adjacent intervals with the least χ^2 values are merged together, because χ^2 low values for a pair indicates similar class distributions. This merging process proceeds recursively until a predefined stopping criterion is met. For two adjacent intervals, if χ^2 test concludes that the class is independent intervals should be merged. If χ^2 test concludes that they are not independent, i.e. the difference in relative class frequency is statistically significant, the two intervals should remain separate.

Calculation. To calculate expected value for any combination of row and class:

- (1) Take the product of the corresponding row sum and column sum
- (2) Divided by the grand total of the observed values for the two rows.

Then:

- (1) Using observed and expected values, calculate, for each of the cells: $\frac{(O-E)^2}{E}$
- (2) Sum each cell's χ^2

When exceeds χ^2 threshold, hypothesis is rejected. Small value for supports hypothesis. Important adjustment, when $E < 0.5$ replace it with 0.5.

- (1) Select the smallest value
- (2) Compare it to the threshold
- (3) If it falls below the threshold, merge it with the row immediately below it
- (4) recalculate χ^2 , Only need to do this for rows adjacent to the recently merged one.

Large numbers of intervals does little to solve the problem of discretization. Just one interval cannot contribute to a decision making process. Modify significance level hypothesis of independence must pass, triggering interval merge. Set a minimum and a maximum number of intervals

1.5 Entropy

Definition. Entropy is the measure of the presence of there being more than one possible classification. Used for splitting attributes in decision trees Entropy minimizes the complexity, number of branches, in the decision tree. No guarantee that using entropy will always lead to a small decision Tree Used for feature reduction: Calculate the value of information gain for each attribute in the original dataset. Discard all attributes that do not meet a specified criterion. Pass the revised dataset to the preferred classification algorithm

Entropy has bias towards selecting attributes with a large number of values

Calculation. To decide if you split on an attribute:

- (1) find the entropy of the data in each of the branches after the split
- (2) then take the average of those and use it to find information gain.

(3) The attribute split with the highest information gain (lowest entropy) is selected.

- Entropy is always positive or zero
- Entropy is zero when $p_i = 1$, aka when all instances have the same class
- Entropy is at its max value for the of classes when all classes are evenly distributed

If there are classes, we can denote the proportion of instances with classification i by p_i for $i = 1$ to K . $p_i = \frac{\text{instances of class } i}{\text{total number of instances}}$

$$\text{Entropy} = E = - \sum_{i=1}^K p_i \log_2 p_i$$

where $K =$ non-empty classes and $p_i = \frac{|i|}{N}$, instances in class i over total number of instances N .

1.6 GINI

Calculation.

- (1) For each non-empty column, form the sum of the squares of the values in the body of the table and divide by the column sum.
- (2) Add the values obtained for all the columns and divide by N (the number of instances).
- (3) Subtract the total from 1.

1.7 Information Gain

Definition. The difference between the entropy before and after splitting on a given attribute in a decision tree. Maximizing information gain is the same as minimizing E_{new} .

Calculation.

$$\text{Information Gain} = E_{start} - E_{new}$$

Starting node:

$$E_{start} = -\frac{4}{24} \log_2 \frac{4}{24} - \frac{5}{24} \log_2 \frac{5}{24} - \frac{15}{24} \log_2 \frac{15}{24} \quad (1)$$

After splitting on attribute:

$$E_{new} = \frac{8}{24} E_1 + \frac{8}{24} E_2 + \frac{8}{24} E_3 \quad (2)$$

Uses.

2 CLUSTERING

Definition. Grouping data into separate groups. Use distance metric between two datapoints. Groups should be distinct from another and composed of items similar to one another, and different from items in other groups.

2.1 Naïve Bayes

$$P(c_i|v) = P(c_i) \prod_{j=1}^n P(a_j = v_j | \text{class} = c_i)$$

2.2 Nearest Neighbors

Mainly used when all attribute values are continuous

General strategy:

- (1) Find the k training instances that are closest to the unseen instance.
- (2) Take the most commonly occurring classification for these instances.
 - KMeans
 - DBSCAN

3 SEQUENCE MINING

TODO

Definition. Finding meaningful, recurring sequences of events

4 ASSOCIATION RULE ANALYSIS

Definition. Given a collection of collections (database of transactions of food items), find items with high co-occurrence.

Let m be the number possible items that can be bought. Let I denote the set of all possible items. Possible itemsets: $s^{|I|}$ An itemset S matches a transaction T (itself an itemset) if $S \subset T$.

4.1 Support

Definition. $support(S)$: proportion of itemsets matched by S . Proportion of transactions that contain all the items in S . Frequency with which the items in S occur together in the database.

Calculation.

$$support(S) = \frac{count(S)}{n}$$

where n is the number of transactions in the database.

4.1.1 Uses.

4.2 Confidence

Calculation. Confidence of a rule can be calculated either by

$$Confidence(L \rightarrow R) = \frac{count(L \cup R)}{count(L)}$$

or

$$Confidence(L \rightarrow R) = \frac{support(L \cup R)}{support(L)}$$

Reject rules where

$$support < minsup \approx 0.01 = 1\%$$

Also called a frequent|large|supported itemset.

Reject rules where

$$confidence < minconf \approx 0.8 = 80\%$$

Uses.

4.3 Lift

Definition. Lift measures how many more times the items in and occur together than would be expected if they were statistically independent.

Calculation.

$$\begin{aligned} \text{Lift}(L \rightarrow R) &= \frac{count(L \cup R)}{count(L) \times support(R)} \\ &= \frac{support(L \cup R)}{support(L) \times support(R)} \\ &= \frac{confidence(L \rightarrow R)}{support(R)} \\ &= \frac{N \times confidence(L \rightarrow R)}{count(R)} \\ &= \frac{N \times confidence(R \rightarrow L)}{count(R)} \\ &= \text{Lift}(R \rightarrow L) \end{aligned} \quad (3)$$

Uses.

4.4 Frequent Itemsets

- (1) Find itemsets of size k made from 2 supported itemsets of size $k - 1$
- (2) For each new itemset:
 - (a) check if every sub-itemset in it also exists in the supported itemsets of size $k - 1$.
 - (b) If not every sub-itemset does, then prune it
 - (i) Now with the final candidates, determine if they have minimum support

- (ii) To determine association rules, find which itemsets have at least minimum confidence

4.5 Rules Possible

$${}_k C_i$$

or

$$2^k - 2$$