

UNIVERSITY OF STAVANGER, DAT550

Summary of lecture notes



Håvard Godal
Spring 2021

Contents

1	Introduction to Data Mining	5
1.1	Data Mining	5
1.2	Machine Learning	5
1.3	Data Mining or Machine Learning	6
1.4	Properties of Data Mining	6
1.4.1	Commercial viewpoint	6
1.4.2	Scientific Viewpoint	6
1.4.3	KDD Process	7
1.5	Data Mining tasks	7
1.6	Regression	7
1.7	Unsupervised Learning	8
1.8	Association Rule Mining	8
1.9	Challenges in Data Mining	8
2	Data	9
2.1	Attributes	9
2.1.1	Values	9
2.1.2	Types	9
2.1.3	Discrete and Continuous Attributes	10
2.2	Datasets	10
2.2.1	Record Data	10
2.2.2	Graph Data	11
2.2.3	Ordered Data	11
2.3	Data Quality	11
2.3.1	Noise	11
2.3.2	Outliers	11
2.3.3	Missing Values	11
2.3.4	Duplicate Data	12
2.4	Distance/Similarity Functions	12
2.4.1	Similarity/Dissimilarity for Simple Attributes	12
2.4.2	Euclidean distance	12
2.4.3	Minkowski Distance	13
2.4.4	Properties of Distance Functions	13
2.4.5	Properties of Similarity Functions	14
2.5	Similarity/Coefficient	14
2.5.1	Simple Matching Coefficient	14
2.5.2	Jaccard Similarity/Coefficient	14

2.5.3	Cosine Similarity	14
2.6	Dot Product	15
2.7	Data Preprocessing	15
2.7.1	Aggregation	15
2.7.2	Data Sampling	15
2.8	Dimensionality Reduction	16
3	Exploring Data	17
3.1	Summary Statistics	17
3.2	Visualisation	18
3.2.1	Representation	18
3.2.2	Arrangement	18
3.2.3	Selection	18
3.2.4	Visualisation Techniques	18
4	Decision Trees	20
4.1	Classification Definition	20
4.1.1	Decision Tree	21
4.2	Hunt's Algorithm	22
4.3	Building Trees	22
4.3.1	Tree Induction	22
4.4	Measures of Node Impurity	23
4.4.1	Gini	24
4.4.2	Entropy	25
4.4.3	Misclassification error	26
4.4.4	Comparison among Splitting Criterias	26
4.5	Stopping Criteria for Tree Induction	26
4.6	Advantages of Decision Tree based Classification	27
4.7	Classification Issues	27
4.7.1	Underfitting and Overfitting	27
4.7.2	Address Overfitting	28
4.7.3	Generalization Errors	28
4.7.4	Handling Missing Attribute Values	29
4.7.5	Other Issues	30
5	Model Evaluation	31
5.1	Metrics for Performance Evaluation	31
5.1.1	Confusion Matrix	31
5.1.2	Cost Matrix	32
5.1.3	Cost-Sensitive Measures	32
5.2	Methods for Performance Evaluation	33
5.2.1	Learning Curve	33
5.2.2	Other Methods of Estimation	34

5.3	Methods for Model Comparison	35
5.3.1	Receiver Operating Characteristic (ROC)	35
6	Regression Trees and Ensemble methods	36
6.1	Regression Trees	36
6.1.1	Constructing a Regression Tree	36
6.2	Ensemble Methods	37
6.3	Manipulating the Training set	37
6.4	Manipulating the Features	39

1 Introduction to Data Mining

1.1 Data Mining

Data mining can be described as:

*"Non-trivial extraction of implicit,
previously unknown and potential useful information from data"*

*"Exploration and analysis, by automatic or semi-automatic means,
of large quantities of data in order to discover meaningful patterns"*

1.2 Machine Learning

Machine learning can be described as:

*"The field of study that gives computers the ability to learn
without being explicitly programmed"*

*"A computer program is said to learn from experience E
with respect to some class of tasks T and performance measure P,
if its performance at tasks in T, as measured by P, improves with experience E"*

1.3 Data Mining or Machine Learning

Data mining is the process of discovering patterns in large data sets involving methods at the intersection of machine learning, statistics, and database systems.

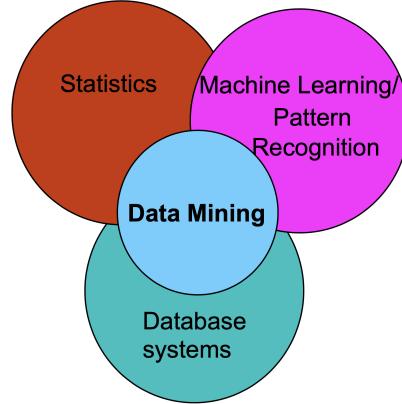


Figure 1.1: Diagram of Data Mining

1.4 Properties of Data Mining

1.4.1 Commercial viewpoint

Lots of data is being collected and warehoused, such as web data, purchases, and transactions. Computation power has become cheap and powerful, and the competitive pressure is strong. It is therefore necessary to provide better, customized services for an edge.

1.4.2 Scientific Viewpoint

Data is collected and stored at enormous speeds, up to multiple terabytes per hour. Such data can be generated from the large hadron collider and from scanning the universe. Such data can not be interpreted by traditional techniques, and requires data mining to classify and segment the data, and form hypothesis formations.

1.4.3 KDD Process

Knowledge Discovery in Databases (KDD) refers to the overall process of discovering useful knowledge from data.

KDD is an integration of multiple technologies for data management such as database management and data warehousing, statistic machine learning, decision support, and others such as visualisation and parallel computing.

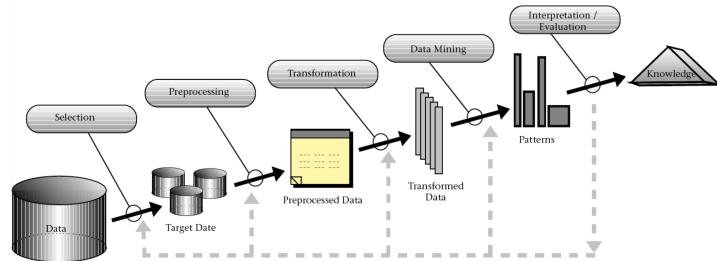


Figure 1.2: The KDD Process

1.5 Data Mining tasks

- Prediction methods
 - Use some variables to predict unknown or future values of other variables.
 - * Supervised (classification or regression)
 - * Unsupervised (clustering)
 - * Semi-supervised
- Descriptive methods
 - Find human-interpretable patterns that describe the data.
 - * Rule mining
 - * Frequent patterns
 - * Anomaly detection

1.6 Regression

1. Given a set of data points/instances along with “correct” answers or labels (training data)
2. Feed it to an algorithm which can learn a model and predict the correct value for an unseen data point (test data)
3. The learned model is an approximate representation of the training data

4. Predict continuous valued output (price in the previous example)
5. The algorithm should produce more right answers (goal)

1.7 Unsupervised Learning

- Unsupervised learning includes unlabeled data
- One way of doing this would be to cluster data into groups
 - Group data points which are similar together, while separating dissimilar items as much as possible
 - This is a clustering algorithm

1.8 Association Rule Mining

- Given a set of records each of which contain some number of items from a given collection
- Produce dependency rules which will predict occurrence of an item based on occurrences of other items

Such rules can be used for marketing and sales promotion. Rule mining can also be used for inventory management, and much more.

1.9 Challenges in Data Mining

- Scalability
- Dimensionality
- Heterogenous or complex data
- Data ownership and distribution
- Privacy concern
- Data quality
- Evolving/streaming data

2 Data

2.1 Attributes

- An attribute is a property or characteristic of an object
 - Attribute are also known as variables, fields, characteristics, or features
- A collection of attributes describe an object
 - Objects are also known as records, points, cases, samples, entries, or instances

2.1.1 Values

Attribute values are numbers or symbols assigned to an attribute.

The distinction between attributes and attribute values are:

- Same attribute can be mapped to different attribute values
- Different attributes can be mapped to the same set of values

2.1.2 Types

Attribute Type	Operations	Transformations
Nominal	Mode, Entropy, Contingency, Correlation, χ^2 -test	Any permutation of values
Ordinal	Median, Rank correlation, Percentiles, Run tests, Sign tests	Order-preserving change of values $new_value = f(old_value)$ where f is a monotonic function.
Interval	Mean, Standard deviation, t -tests Pearson's correlation, F -tests	$new_value = a * old_value + b$ where a and b are constants
Ratio	Geometric mean, Harmonic mean, Percent variation	$new_value = a * old_value$

2.1.3 Discrete and Continuous Attributes

Discrete attributes:

- Finite or countably infinite set of values
- Often integer or binary variables

Continuous attributes:

- Real numbers as attribute values
- Typically floating-point variables

2.2 Datasets

Some important characteristics of structured data are:

- Dimensionality
 - Curse of dimensionality
 - * As the number of features or dimensions grow, the amount of data needed to generalize accurately grows exponentially.
 - * When data moves from one dimensions to i.e. three dimensions, the given data fills less and less of the data space. In order to maintain an accurate representation of the space, the data for analysis grows exponentially.
 - * When sorting or classifying data, low dimensional spaces tend to show the data as very similar, but in higher dimensions, the data might be further away from each other.
- Sparsity
 - Only presence counts
- Resolution
 - Patterns depend on the scale

2.2.1 Record Data

Data that consists of a collection of records, each of which consists of a fixed set of attributes.

Data Matrix

If data objects have the same fixed set of numeric attributes, then the data objects can be thought of as points in a multi-dimensional space, where each dimension represents a distinct attribute.

Such data set can be represented by an m by n matrix, where there are m rows, one for each object, and n columns, one for each attribute.

Document Data

Each document becomes a *term* vector. Each term is a component (attribute) of the vector.

The value of each component is the number of times the corresponding term occurs in the document.

Transaction Data

A special type of record data, where each record (transaction) involves a set of items.

For example, consider a grocery store. The set of products purchased by a customer during one shopping trip constitute a transaction, while the individual products that were purchased are the items.

2.2.2 Graph Data

I.e. chemical structures.

2.2.3 Ordered Data

I.e. sequential data or gene sequences.

2.3 Data Quality

2.3.1 Noise

Noise refers to modification of original values.

2.3.2 Outliers

Outliers are data objects with characteristics that are considerably different than most of the other data objects in the data set.

2.3.3 Missing Values

Reasons for missing values may be that some information is not collected or applicable to all cases.

Handling missing values can be done by:

- Eliminating data objects
- Estimating the missing values
- Ignoring the missing values during analysis
- Replacing the missing values with all possible values weighted by their probabilities

2.3.4 Duplicate Data

Data set may include data objects that are duplicates, or almost duplicates of one another. This is a major issue when merging data from heterogenous sources.

2.4 Distance/Similarity Functions

- Similarity
 - Numerical measure of how alike two data objects are
 - Higher value means more alike
 - Often falls in the range [0, 1]
- Dissimilarity
 - Numerical measure of how different two data objects are
 - Lower when objects are more alike
 - Minimum dissimilarity is often 0
 - Upper limit varies
- Proximity refers to a similarity or dissimilarity

2.4.1 Similarity/Dissimilarity for Simple Attributes

Given that p and q are the attribute values for two data objects:

Attribute Type	Dissimilarity	Similarity
Nominal	$d = \begin{cases} 0 & \text{if } p = q \\ 1 & \text{if } p \neq q \end{cases}$	$s = \begin{cases} 1 & \text{if } p = q \\ 0 & \text{if } p \neq q \end{cases}$
Ordinal	$d = \frac{ p-q }{n-1}$	$s = 1 - \frac{ p-q }{n-1}$
Interval/Ratio	$d = p - q $	$s = -d, s = \frac{1}{1+d}, s = 1 - \frac{d - \min(d)}{\max(d) - \min(d)}$

2.4.2 Euclidean distance

Euclidean distance is defined in the following theorem, where n is the number of dimensions (attributes) and p_k and q_k are, respectively, the k th attributes (components) or data objects p and q .

Theorem 2.1: Euclidean Distance

$$\text{distance} = \sqrt{\sum_{k=1}^n (p_k - q_k)^2}$$

Standardization is necessary if scales differ.

2.4.3 Minkowski Distance

Minkowski Distance is a generalization of Euclidean Distance where r is a parameter, n is the number of dimensions (attributes) and p_k and q_k are, respectively, the k^{th} attributes (components) or data objects p and q .

Theorem 2.2: Minkowski Distance

$$\text{distance} = \sqrt[r]{\sum_{k=1}^n |p_k - q_k|^r}$$

- $r = 1$: Manhattan distance
- $r = 2$: Euclidean distance
- $r \rightarrow \infty$: "Supremum" distance
 - This is the maximum difference between any component of the vectors

2.4.4 Properties of Distance Functions

Distances, such as the Euclidean distance, have some well known properties.

- Positive definiteness
 - $d(p, q) \geq 0$ for all p and q
 - $d(p, q) = 0$ only if $p = q$
- Symmetry
 - $d(p, q) = d(q, p)$ for all p and q
- Triangle Inequality
 - $d(p, r) \leq d(p, q) + d(q, r)$ for all p , q and r

Important 2.3: Metric

A distance that satisfies the properties mentioned above is a metric.

2.4.5 Properties of Similarity Functions

Well known properties of similarities.

- Maximum Similarity
 - $s(p, q) = 1$ (or maximum similarity) only if $p = q$
- Symmetry
 - $s(p, q) = s(q, p)$ for all p and q

2.5 Similarity/Coefficient

A common situation is that objects, p and q , have only binary attributes.

Let's define:

- M_{01} - Number of attributes where p is 0 and q is 1
- M_{10} - Number of attributes where p is 1 and q is 0
- M_{00} - Number of attributes where p is 0 and q is 0
- M_{11} - Number of attributes where p is 1 and q is 1

2.5.1 Simple Matching Coefficient

$$SMC = \frac{\text{Number of matches}}{\text{Number of attributes}} \quad (2.1)$$

$$SMC = \frac{M_{00} + M_{11}}{M_{01} + M_{10} + M_{00} + M_{11}} \quad (2.2)$$

2.5.2 Jaccard Similarity/Coefficient

The Jaccard Similarity/Coefficient is used for categorical attributes and sets.

$$J = \frac{\text{Number of } M_{11}}{\text{Number of non-both-zero attributes}} \quad (2.3)$$

$$J = \frac{M_{11}}{M_{01} + M_{10} + M_{11}} \quad (2.4)$$

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cup B|} \quad (2.5)$$

2.5.3 Cosine Similarity

If \mathbf{A} and \mathbf{B} are two documented vectors, then

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}} \quad (2.6)$$

2.6 Dot Product

Given two vectors a and b :

$$\vec{a} \cdot \vec{b} = \sum_{i=1}^n a_i b_i = a_1 b_1 + a_2 b_2 + \cdots + a_n b_n \quad (2.7)$$

2.7 Data Preprocessing

2.7.1 Aggregation

Aggregation is to combine two or more attributes (or objects) into a single attribute (or object).

Purpose:

- Data reduction
- Reduce the number of attributes or objects
- Change of scale
- More "stable" data
- Often reduce variability

2.7.2 Data Sampling

Data sampling is the main technique employed for data selection. It is often used for both the preliminary investigation of the data and the final data analysis. Statisticians sample because obtaining the entire set of data of interest is too expensive or time consuming. Sampling is used in data mining because processing the entire set of data of interest is too expensive or time consuming.

A good strategy for choosing the sample size is to aim for a sample size that is 10% of the population, as long as the sample size is smaller than 1000. It is also important to have a large enough sample size so that all attributes/groups are represented.

Sampling is effective when:

- The sample is representative
- The sample has approximately the same properties (of interest) as the original set of data

Types of sampling:

- Simple random sampling
 - Equal probability of selecting any item.

- Sampling without replacement
 - As each item is selected, it is removed from the population.
- Sampling with replacement
 - Objects are not removed from the population as they are selected for the sample.
- Stratified sampling
 - Split the data into several partitions, then draw random samples from each partition.

Reservoir Sampling

Theorem 2.4: Reservoir Sampling

Keep a reservoir of r samples

1. Keep the first r items in memory
2. When the i^{th} item arrives ($i > r$)
 - a) Keep the new item with probability $\frac{r}{i}$,
or discard the new item with probability $1 - \frac{r}{i}$
 - b) Discard one of the items in the reservoir at random
if the new item was kept.

This means that as i increases, the probability of a new item being kept in the reservoir reduces.

2.8 Dimensionality Reduction

- Purpose
 - Avoid curse of dimensionality
 - Reduce amount of time and memory required by data mining
 - Allow data to be more easily visualized
 - May help eliminate irrelevant features or reduce noise
- Techniques
 - Principle Component analysis
 - Singular Value Decomposition

3 Exploring Data

3.1 Summary Statistics

- Frequency
- Mode
- Percentiles
- Mean
- Median
- Range
- Variance

Theorem 3.1: Mean

$$\text{mean}(x) = \bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$$

Theorem 3.2: Median

$$\text{median}(x) = \begin{cases} x_{(r+1)} & \text{if } m = 2r + 1 \\ \frac{1}{2} (x_{(r)} + x_{(r+1)}) & \text{if } m = 2r \end{cases}$$

Theorem 3.3: Variance

$$\text{variance}(x) = s_x^2 = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2$$

Theorem 3.4: Average Absolute Deviation

$$AAD(x) = \frac{1}{m} \sum_{i=1}^m |x_i - \bar{x}|$$

Theorem 3.5: Mean Absolute Deviation

$$MAD(x) = \text{median} (\{|x_1 - \bar{x}| \cdots |x_m - \bar{x}|\})$$

3.2 Visualisation

Visualisation is the conversion of data into a visual or tabular format so that the characteristics of the data and the relationships among data items or attributes can be analyzed or reported. Visualisation allows humans to detect general patterns and trends, as well as detect outliers and unusual patterns.

3.2.1 Representation

The representation is the mapping of information to a visual format. Data objects, their attributes, and the relationships among data objects are translated into graphical elements such as points, lines, shapes, and colors.

3.2.2 Arrangement

Arrangement is the placement of visual elements within a display. This can make a large difference in how easy it is to understand the data.

3.2.3 Selection

Selection is the elimination or the de-emphasis of certain objects and attributes. Such selections may involve choosing a subset of attributes, or choosing a subset of objects. Only selecting some attributes can be done through dimensionality reduction, while selecting some objects can be done through stratified sampling. Some things to keep in mind when using data selection is to contain the diversity of the objects.

3.2.4 Visualisation Techniques

- Histogram plot
- Box plot
- Scatter plot

- Matrix plot
- Parallel coordinates plot
- Star plot

4 Decision Trees

4.1 Classification Definition

Given a collection of records (a training set) where each record contains a set of attributes, the classification has to find a model for one particular attribute (the class) as a function of the values of other attributes.

The goal is for previously unseen records to have the class attribute assigned as accurately as possible.

A test set is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with the training set, used to build the model, and the test set used to validate it.

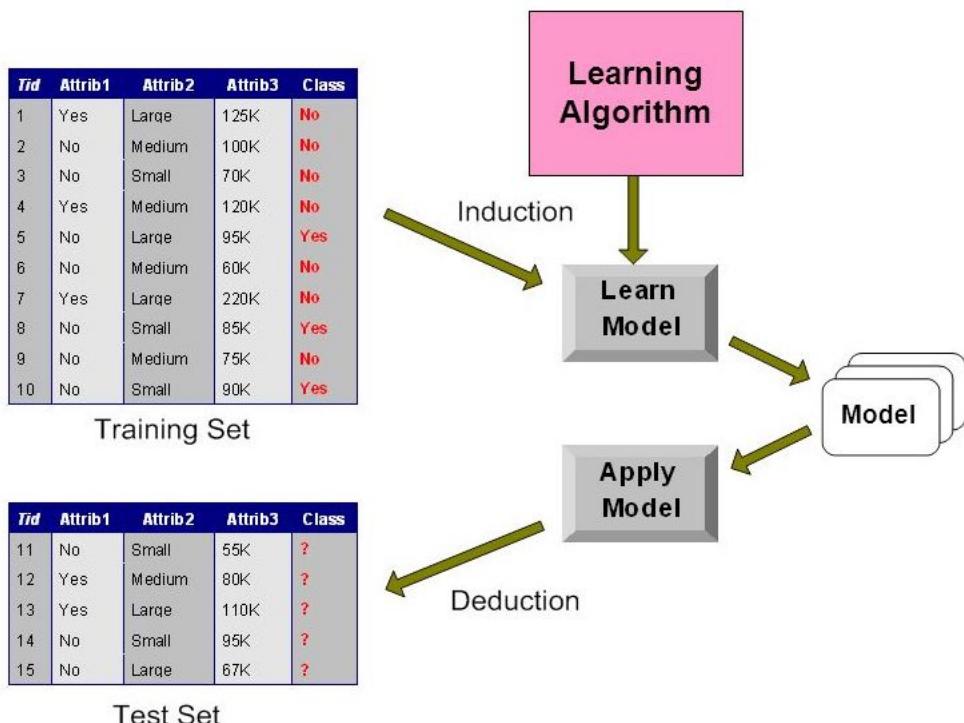


Figure 4.1: Typical Classification Task

4.1.1 Decision Tree

A tree is built by splitting the source set, constituting the root node of the tree, into subsets, which constitute the successor children. The splitting is based on a set of splitting rules based on classification features. This process is repeated on each derived subset in a recursive manner called recursive partitioning. The recursion is completed when the subset at a node has all the same values of the target variable, or when splitting no longer adds value to the predictions. This process of top-down induction of decision trees is an example of a greedy algorithm, and it is by far the most common strategy for learning decision trees from data.

Decision trees used in data mining are of two main types:

- Classification tree analysis - When the predicted outcome is the class (discrete) to which the data belongs.
- Regression tree analysis - When the predicted outcome can be considered a real number.

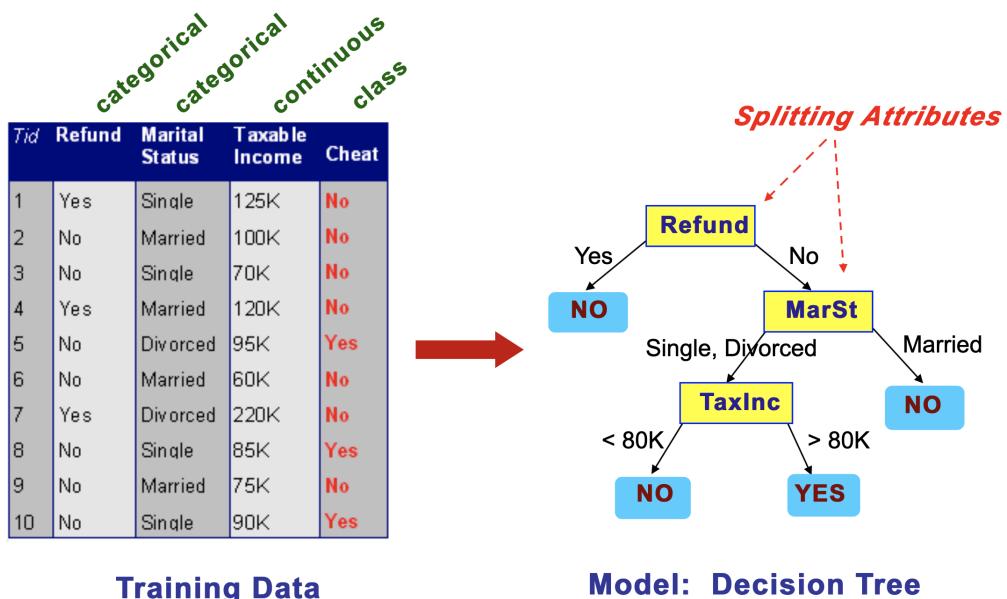


Figure 4.2: Typical Decision Tree

4.2 Hunt's Algorithm

- Let D_t be the set of training records that reach a node t .
- General procedure:
 - If D_t contains records that belong to the same class as y_t , then t is a leaf node labeled as y_t .
 - * This is a pure node.
 - If D_t is an empty set, then t is a leaf node labeled by the default class y_d .
 - If D_t contains records that belong to more than one class, use an attribute test to split the data into smaller subsets. Recursively apply the procedure to each subset.
 - * This is an impure node.

4.3 Building Trees

4.3.1 Tree Induction

The greedy strategy splits the records based on an attribute test that optimizes a certain criterion.

Issues with this approach include:

- How to split the records. How should the attribute test condition be specified, and how is the best split determined?
- Determining when to stop.

Specifying the Attribute Test Condition

A proper test condition depends on the attribute type, which can be nominal, ordinal, or continuous.

The condition also depends on the number of ways to split. This can be either a two-way split (binary) or a multi-way split. The binary split is often preferred as of the bushiness of the multi-way split. This means that very few records will match the leaf-nodes, and will therefore be overfitted.

Ordinal attributes can be split the same way as nominal attributes, with some limitations if the order is to be preserved.

For continuous attributes, there are two ways of handling splitting:

- Discretization to form an ordinal categorical attribute.
 - Static - discretize once at the beginning.
 - Dynamic - ranges can be found by equal interval bucketing, equal frequency bucketing, or clustering.

- Binary decision.
 - Consider all possible splits and find the best cut.
 - Can be more compute intensive.

Determine Best Split

When it comes to greedy approaches, nodes with a homogeneous class distribution are preferred.



Figure 4.3: Class distributions

4.4 Measures of Node Impurity

Finding the best split follows closely to this process:

1. Define the set of records as M_0 , where some records belong to class C_0 , while the rest belongs to C_1 .
2. Define some splitting criterias A and B .
3. Use each criteria to generate new sets of records. The new sets of records now has a different (or equal) distribution of C_0 - and C_1 -records.
 - M_1 and M_2 from A , collected into M_{12}
 - M_3 and M_4 from B , collected into M_{34}
4. The gain in purity can be measured as such: $Gain = M_0 - M_{12}$ vs. $M_0 - M_{34}$
5. The greedy algorithm can then choose the attribute that gives the greatest gain in purity.

4.4.1 Gini

The Gini Index measures the degree of probability of a particular variable being wrongly classified when it is randomly chosen.

Theorem 4.1: Gini Index

$$GINI(t) = 1 - \sum_j^{n_c} [p(j|t)]^2$$

$$\text{Maximum: } (1 - \frac{1}{n_c}) \quad \text{Minimum: } 0.0 \text{ (pure)}$$

Note: $p(j|t)$ is the relative frequency of class j at node t .

Gini can be used when splitting nodes. When a node p is split into k partitions (children), the quality of split follows the theorem:

Theorem 4.2: Gini Split

$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

n_i = Number of records at child i .

n = Number of records at node p .

Compute the Gini Index for continuous attributes can often be quite expensive as all values have to be evaluated for splitting. One process (to be done for each attribute) for finding a good split is as such:

1. Sort the attribute on values.
2. Linearly scans the values, each time choosing the mean of the current value and the next value as the splitting point position and compute the Gini Index.
3. Choose the split position that has the least Gini Index.

Cheat	No	No	No	Yes	Yes	Yes	No	No	No	No											
Annual Income																					
Sorted Values →	60	70	75	85	90	95	100	120	125	220											
Split Positions →	55	65	72	80	87	92	97	110	122	172	230										
<= > <= > <= > <= > <= > <= > <= >																					
Yes	0	3	0	3	0	3	1	2	2	1	3	0	3	0	3	0	3	0			
No	0	7	1	6	2	5	3	4	3	4	3	4	4	4	3	5	2	6	1	7	0
Gini	0.420	0.400	0.375	0.343	0.417	0.400	0.300	0.343	0.375	0.400	0.420										

Figure 4.4: Gini Index of Continuous Attributes

4.4.2 Entropy

Entropy is a measure of the uncertainty associated with a random variable. Entropy is a splitting criterion based on information gain and measures the homogeneity of a node.

Theorem 4.3: Entropy

$$\text{Entropy}(t) = - \sum_j^{n_c} p(j|t) \log_2 p(j|t)$$

Maximum: $(\log_2 n_c)$ Minimum: 0.0

Note: $p(j|t)$ is the relative frequency of class j at node t .

In the same way Gini can be used to split nodes, so can Entropy:

Theorem 4.4: Entropy Split

$$GAIN_{split} = \text{Entropy}(p) - \left(\sum_{i=1}^k \frac{n_i}{n} \text{Entropy}(i) \right)$$

Parent node p is split into k partitions.

n_i is the number of records in partition i

The gain measures reduction in entropy achieved from the split. Choose the split that achieves the most reduction, that is, maximizes gain.

The disadvantage is that Entropy tends to prefer splits that result in a large number of partitions, each being small but pure (overfitting).

This is fixed by including the number of partitions in the equation.

Theorem 4.5: Entropy Split Ratio

$$GainRatio_{split} = \frac{Gain_{split}}{SplitInfo}$$

$$SplitInfo = - \sum_{i=1}^k \frac{n_i}{n} \log \frac{n_i}{n}$$

Parent node p is split into k partitions.

n_i is the number of records in partition i

4.4.3 Misclassification error

The Classification error measures the misclassification error made by a node.

Classification error at node t :

Theorem 4.6: Classification Error

$$\text{Error}(t) = 1 - \max P(i|t)$$

$$\text{Maximum: } (1 - \frac{1}{n_c}) \quad \text{Minimum: } 0.0$$

4.4.4 Comparison among Splitting Criterias

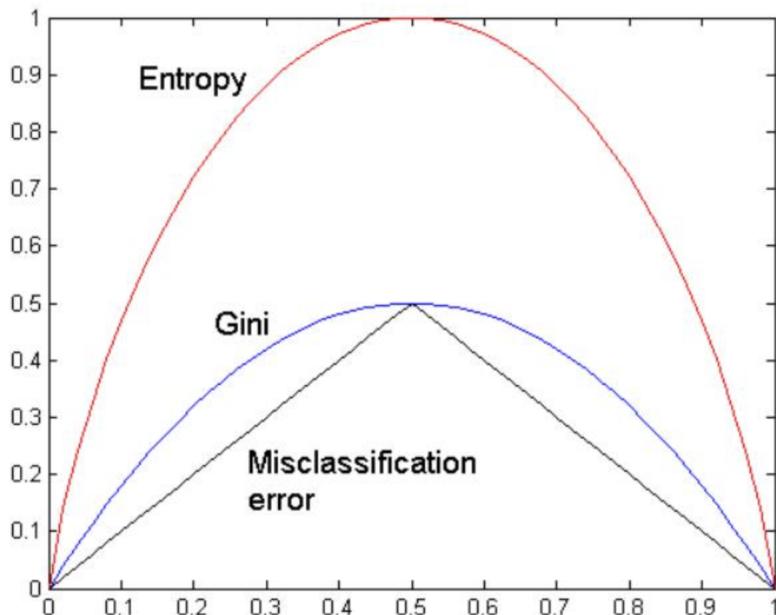


Figure 4.5: Splitting Criterias for a two-class problem

4.5 Stopping Criteria for Tree Induction

- Stop expanding a node when all the records belong to the same class.
- Stop expanding a node when all the records have similar attribute values.
- Early termination (probably for limiting tree depth).

4.6 Advantages of Decision Tree based Classification

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

4.7 Classification Issues

4.7.1 Underfitting and Overfitting

- Underfitting: When a model cannot capture the underlying trend of the data. Intuitively, underfitting occurs when the model or the algorithm does not fit the data well enough. Specifically, underfitting occurs if the model or algorithm shows low variance but high bias. Underfitting is often a result of an excessively simple model. Overfitting can also arise from insufficient data points.
- Overfitting: When a model captures the noise of the data. Intuitively, overfitting occurs when the model or the algorithm fits the data too well. Specifically, overfitting occurs if the model or algorithm shows low bias but high variance. Overfitting is often a result of an excessively complicated model, and it can be prevented by fitting multiple models and using validation or cross-validation to compare their predictive accuracies on test data.

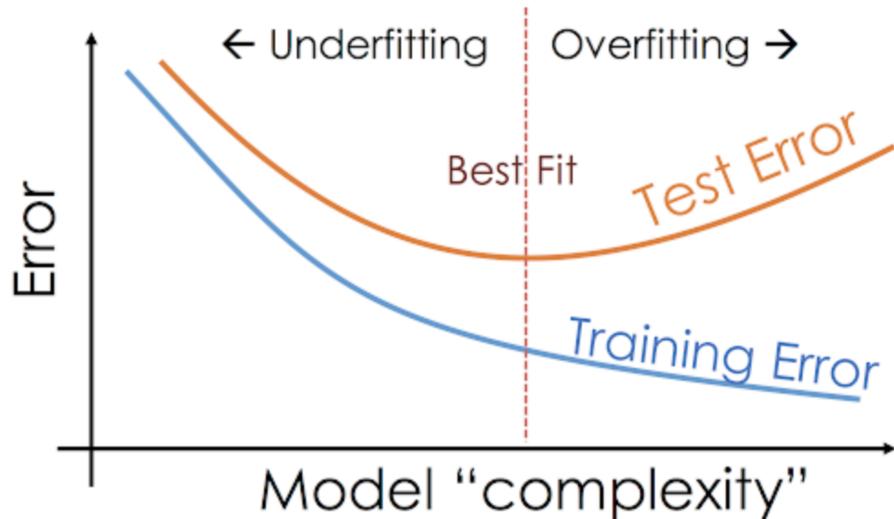


Figure 4.6: Overfitting and underfitting

4.7.2 Address Overfitting

Decision Tree Pruning

Pre-Pruning (Early Stopping Rule):

- Stop the algorithm before it becomes a fully-grown tree.
 - Stop if all instances belong to the same class.
 - Stop if all the attribute values are the same.
 - Stop if the number of records in a node are lower than the user-specified threshold.
 - Stop if splitting a node does not improve the impurity measure.

Post-Pruning:

1. Grow decision tree to its entirety.
2. Trim the nodes of the decision tree in a bottom-up fashion.
3. If generalization error improves after trimming, replace sub-tree by a leaf node.
4. The class label of a leaf node is determined from the majority class of instances in the sub-tree.

Occam's Razor

- Given two models of similar generalization errors, one should prefer the simpler model over the more complex model.
- For complex models, there is a greater chance that it was fitted accidentally by the error in data.
- Therefore, one should include model complexity when evaluating a model.

4.7.3 Generalization Errors

Generalization error is a measure of how accurately a model can predict outcome values for previously unseen data. The measure is calculated from the ratio of misclassifications that can occur when choosing the majority class for all records in the node.

- Training errors: Error on training ($\sum_{t=1}^n e(t)$)
- Generalization errors: Error on validation ($\sum_{t=1}^n e'(t)$)

Methods for estimating generalization errors:

- Optimistic approach: $e'(t) = e(t)$

- Pessimistic approach:
 - For leaf each leaf node, $e'(t_i) = e(t_i) + \frac{1}{2}$
 - Total error: $e'(T) = e(T) + N\frac{1}{2}$ (N = number of leaf nodes)
- Reduced error pruning (REP)
 - Uses validation data set to estimate generalization error.

4.7.4 Handling Missing Attribute Values

Missing attribute values affect the decision tree in three major ways:

- How impurity measures are computed.
- How to distribute instance with missing value to child nodes.
- How a test instance with missing value is classified.

Computing Impurity Measure (Impurity Gain)

One way of computing the information gain is as normal, without including the instances with the missing attribute. This means that the entropy-computation for the parent is as normal, but for the children, instances with the missing attributes should be omitted, but included in the denominator. Finally, the gain is computed as normal but is multiplied by the fraction of "valid" records.

Distributing Instances

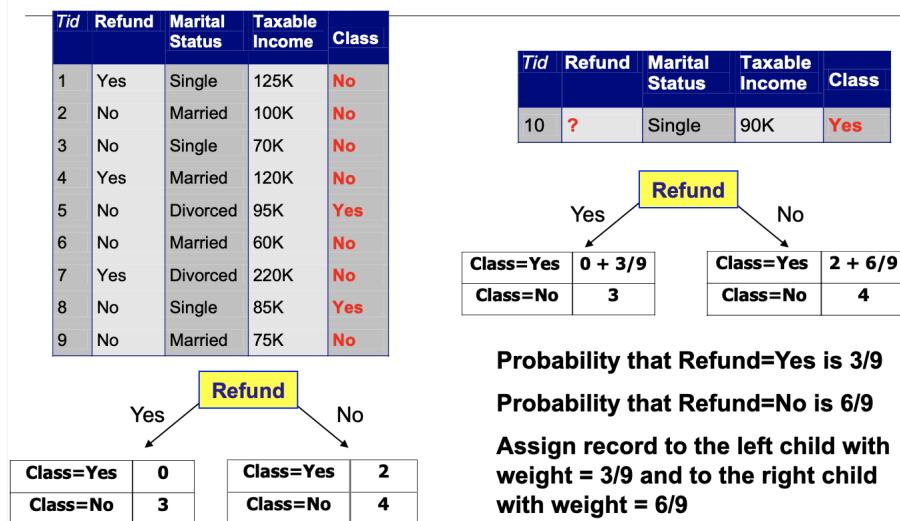


Figure 4.7: Distributing Instances

The figure explains that since the instance with the unknown attribute corresponds to the class yes, it is distributed on both nodes, affecting the yes-class in each node. The distribution ratio is based on the number of instances of the same class, divided by the total number of instances.

Classify Instances

To classify an instance with an unknown class, as well as having the current attribute unknown, can be done by understanding what option is more likely. If most of the instances are of attribute 1 the instance with the unknown attribute is assumed to also have attribute 1 and follows the tree accordingly.

4.7.5 Other Issues

Data Fragmentation

- The number of instances gets smaller as you traverse down the tree.
- Number of instances at the leaf node could be too small to make any statistically significant decision.

Search Strategy

- Finding an optimal decision tree is NP-hard (exponentially).
- The algorithm presented so far uses a greedy, top-down, recursive partitioning strategy to induce a reasonable solution.
- Could use bottom-up or bi-directional algorithms.

Expressiveness

Decision trees provide expressive representation for learning discrete-valued functions, but they do not generalize well to certain types of boolean functions like parity functions. Also, decision trees are not expressive enough for modeling continuous variables.

Decision Boundary

The borderline between two neighboring regions of different classes is known as a decision boundary.

A decision boundary is parallel to the axes because the test conditions only involve a single attribute at a time.

This could be fixed with oblique decision trees, but includes more expressive representation and can be computationally expensive.

Tree Replication

The same subtree appears in multiple branches.

5 Model Evaluation

5.1 Metrics for Performance Evaluation

Metrics for performance evaluation focuses on the predictive capability of a model, that is, how well a model can predict the class of the test data, and evaluate the performance of the training-data. This focus is instead of focusing on how fast it takes to classify or build models, scalability, etc.

5.1.1 Confusion Matrix

This can be done through a confusion matrix:

		Predicted		Precision
		FALSE	TRUE	
Actual	FALSE	True Negative (TN)	False Positive (FP)	
	TRUE	False Negative (FN)	True Positive (TP)	
				Recall

Figure 5.1: Confusion Matrix

One of the most common and widely-used metric for evaluating the performance is accuracy:

Theorem 5.1

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

The accuracy metric has some limitations, among those are the ability to evaluate two-class problems where almost all instances are of class 0, while only a few are of class 1. The accuracy metric will result in 99.9% accuracy, which is misleading if the model predicts all instances to belong to class 0.

5.1.2 Cost Matrix

One way of solving the issue is to use a cost matrix instead:

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

Figure 5.2: Cost Matrix

Here, $C(i|j)$ is the cost of misclassifying class j as class i . The cost equals the weighted sum of all predictions.

If true positives and true negatives are equally weighted, as well as the false positives and the false negatives, the accuracy is proportional to the cost. This results in the following equation:

$$Cost = N[q - (q - p) * Accuracy] \quad (5.1)$$

Where $N = TP + TN + FP + FN$, $q = FP$, and $p = TP$.

5.1.3 Cost-Sensitive Measures

Some cost-sensitive measures that are preferred over accuracy:

$$Precision(p) = \frac{TP}{TP + FP} \quad (5.2)$$

$$Recall(r) = \frac{TP}{TP + FN} \quad (5.3)$$

$$F-measure(F) = \frac{2TP}{2TP + FP + FN} \quad (5.4)$$

- Precision is biased towards TP and FP
- Recall is biased towards TP and FN
- F-measure is biased towards all except TN

$$WeightedAccuracy = \frac{\omega_1 TP + \omega_4 TN}{\omega_1 TP + \omega_2 FN + \omega_3 FP + \omega_2 TN} \quad (5.5)$$

5.2 Methods for Performance Evaluation

Methods on how to obtain a reliable estimate of performance.

The performance of a model may depend on other factors besides the learning algorithm:

- Class distribution.
- Cost of misclassification.
- Size of training and test sets.

5.2.1 Learning Curve

Learning curves show how accuracy changes with varying sample sizes.

Learning curves require a sampling schedule:

- Arithmetic sampling. (Linear)
- Geometric sampling. (Exponential)

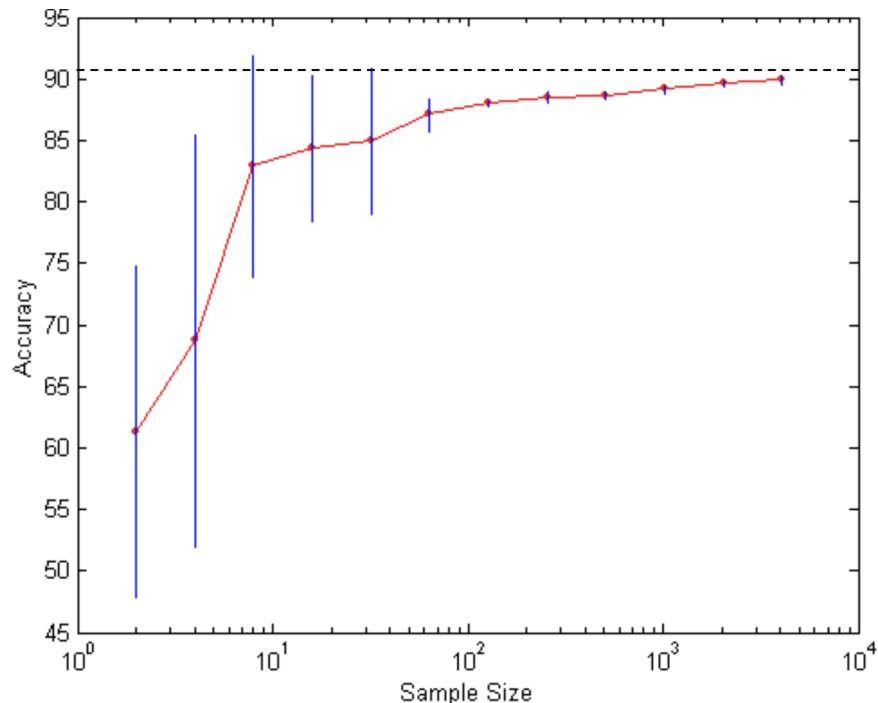


Figure 5.3: A Learning Curve

The effects of small sample sizes includes:

- Bias in the estimate (Underfitting)
- Variance of estimate (Overfitting)

5.2.2 Other Methods of Estimation

- Holdout
 - Reserve 80% for training and 20% for testing.
- Random subsampling
 - Repeated holdouts
- Cross validation (popular)
 - Partition data into k disjoint subsets.
 - k -fold: Train on $k - 1$ partitions, test on the remaining one.
 - Leave-one-out: $k = n$
- Stratified sampling
 - Oversampling
 - Undersampling
- Bootstrap
 - Sampling with replacement

5.3 Methods for Model Comparison

Methods for model comparison explains how to compare the relative performance among competing models.

5.3.1 Receiver Operating Characteristic (ROC)

ROC was developed in the 1950s to characterize the trade-off between true positives and false positives.

A ROC curve plots TP on the y -axis against FP on the x -axis, where the performance of each classifier is represented as a point on the ROC curve. The point changes based on the threshold of the algorithm, the sample distribution, and/or the cost matrix.

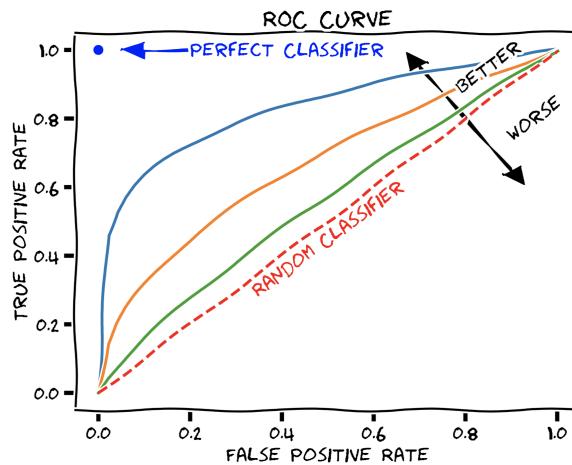


Figure 5.4: ROC Curve

In some cases when comparing the ROC curve for different models, no models consistently outperform the others. Evaluating such models can be based on the area under the ROC curve (AUC), or which model fits better to the desired outcome.

Constructing a ROC curve

1. Use a classifier that produces posterior probability for each test instance $P(+|A)$.
 2. Sort the instances according to $P(+|A)$ in descending order.
 3. Apply threshold at each unique value of $P(+|A)$.
 4. Count the number of TP , FP , TN and FN at each threshold.
- $TPR = \frac{TP}{TP+FN}$
 - $FPR = \frac{FP}{FP+TN}$

6 Regression Trees and Ensemble methods

6.1 Regression Trees

In a decision tree, there is typically a data set with N features (X_1, X_2, \dots, X_N). Each of the features has a domain, with the domain type being dependent on the type of data the feature contains (categorical or numerical). The output variable Y will have a domain D_y .

- Categorical - Classification
- Numerical - Regression

6.1.1 Constructing a Regression Tree

1. Find a split (X_i, v) that creates two child nodes D_L and D_R from the parent node D . The split should maximize the weighted sum of the variances.
 - For ordered domains, sort X_i and consider a split between each pair of adjacent values.
 - For categorical X_i find the best split based on subsets (Breiman's algorithm).

$$\text{Weighted_Variance} = |D| * \text{Var}(D) - (|D_L| * \text{Var}(D_L) + |D_R| * \text{Var}(D_R)) \quad (6.1)$$

$$\text{Var}(D) = \frac{1}{n} \sum_{i=1}^{|D|} (y_i - \bar{y})^2 \quad (6.2)$$

Stopping a Tree

- When the leaf is pure, i.e. $\text{Var}(D) < \epsilon$
- When the number of instances in the leaf ($|D|$) is too small.

Predictor

- Regression - Average value of y_i of the instances in the leaf.
- Classification - Choose the most common y_i in the leaf.

6.2 Ensemble Methods

Ensemble methods construct a set of classifiers from a given training data. The classification is done by predicting class values from previously unseen records by aggregating predictions made by multiple classifiers.

The general idea of ensemble methods is to have the original training data and then:

1. Create multiple data sets from the training data.
2. Build multiple classifiers corresponding to the data sets.
3. Combine the predictions of the classifiers.

Ensemble methods work the sum of the error rate of multiple independent classifiers are lower than the error rate itself. The following equation calculates the probability that the ensemble classifier makes a wrong prediction, having ϵ as the error rate.

$$\sum_{i=N/2(+1)}^N \binom{N}{i} \epsilon^i (1-\epsilon)^{N-i} \quad (6.3)$$

6.3 Manipulating the Training set

Bagging

Bagging is sampling with replacement

1. Let k be the number of bootstrap samples.
 2. `for i in range(k):`
 - a) Create a bootstrap sample of size N , D_i .
 - b) Train a base classifier C_i on the bootstrap sample D_i .
- $$C^*(x) = \operatorname{argmax} (\sum_i \delta(C_i(x) = y))$$

$\delta = 1$ if the argument is true and 0 otherwise (the majority class is chosen).

Bagging can also increase the complexity of simple classifiers such as decision stumps.

Boosting

Boosting is an iterative procedure to adaptively change the distribution of training data by focusing more on previously misclassified records. Initially, all N records are assigned equal weights.

Unlike bagging, the weights of the records may change at the end of each boosting round.

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

One way of boosting is to use the AdaBoost algorithm:

Algorithm 5.7 AdaBoost algorithm.

```

1:  $\mathbf{w} = \{w_j = 1/N \mid j = 1, 2, \dots, N\}$ . {Initialize the weights for all  $N$  examples.}
2: Let  $k$  be the number of boosting rounds.
3: for  $i = 1$  to  $k$  do
4:   Create training set  $D_i$  by sampling (with replacement) from  $D$  according to  $\mathbf{w}$ .
5:   Train a base classifier  $C_i$  on  $D_i$ .
6:   Apply  $C_i$  to all examples in the original training set,  $D$ .
7:    $\epsilon_i = \frac{1}{N} [\sum_j w_j \delta(C_i(x_j) \neq y_j)]$  {Calculate the weighted error.}
8:   if  $\epsilon_i > 0.5$  then
9:      $\mathbf{w} = \{w_j = 1/N \mid j = 1, 2, \dots, N\}$ . {Reset the weights for all  $N$  examples.}
10:    Go back to Step 4.
11:   end if
12:    $\alpha_i = \frac{1}{2} \ln \frac{1-\epsilon_i}{\epsilon_i}$ .
13:   Update the weight of each example according to Equation 5.69.
14: end for
15:  $C^*(\mathbf{x}) = \underset{y}{\operatorname{argmax}} \sum_{j=1}^T \alpha_j \delta(C_j(\mathbf{x}) = y))$ .
```

Figure 6.1: AdaBoost Algorithm

6.4 Manipulating the Features

Random Forest

A random forest is a class of ensemble methods specifically designed for decision tree classifiers. It combines the predictions made by multiple decision trees where each tree is generated based on the values of an independent sample set. The samples are generated from a fixed probability distribution.

As in bagging, we build decision trees on bootstrapped training samples. Each time a split in a tree is considered, a random sample of F features is chosen as split candidates from the full set of m features.

(If $F = m$, we have bagging)

A random forest algorithm can look like this:

- `for i in range(B):`
1. Draw a bootstrap sample D^* of size N from the training data.
 - a) Grow a random forest tree to the bootstrapped data.
 - i. Select F features at random from the m variables.
 - ii. Pick the best variable/split point among the F s.
 - iii. Split the node into two child nodes.
 - Output the ensemble of trees.

Prediction at a new point x can be done:

- For regression - Average of the results.
- For classification - Majority vote.

Random Forest Recommendations

- For classification, the default value for F is \sqrt{m} and the minimum node size is one.
- For regression, the default value for F is $\frac{m}{3}$ and the minimum node size is five.