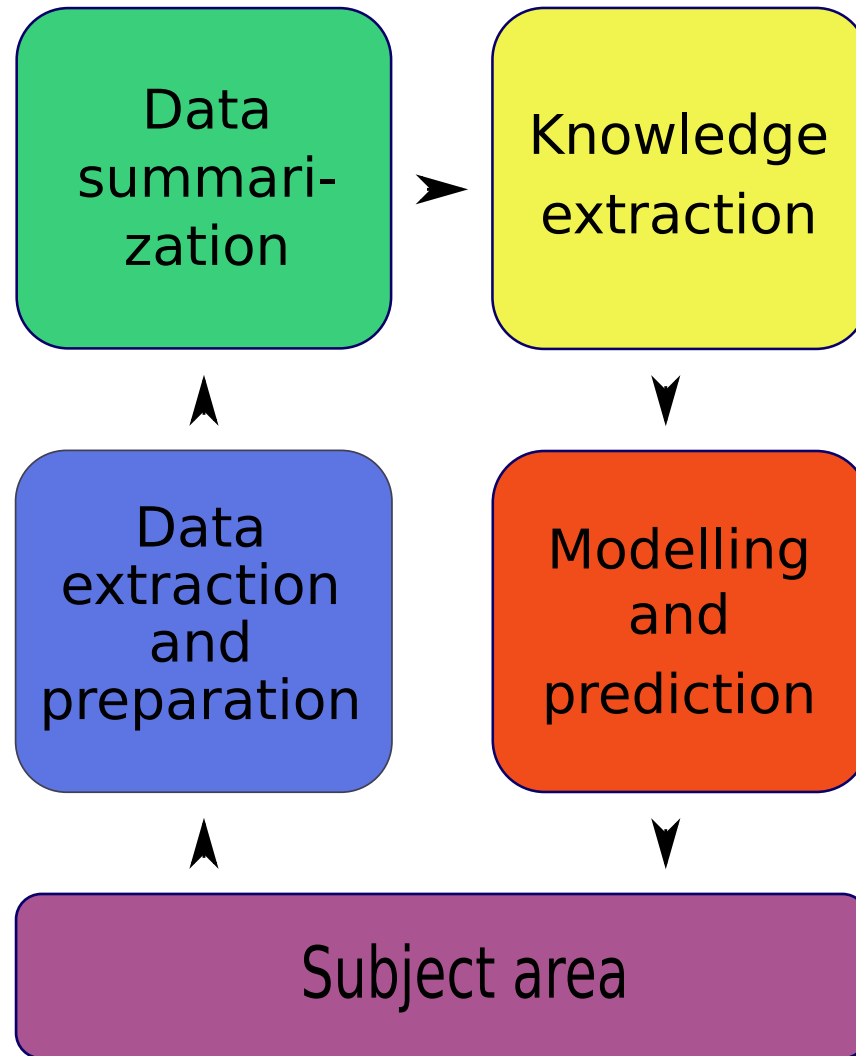


The data analysis cycle



Data Analysis: Building Models and Predicting Outcomes

Institute of Technology Tallaght

Department of Computing

Models

[S] statistics [M] machine learning		Type of dependent [S] / target [M] variable		
		categorical	numeric	probability of categorical
Independent variable [S]/ feature [M] type	categorical	<ul style="list-style-type: none"> ◆ supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	<ul style="list-style-type: none"> ◆ supervised segmentation using variance measures [M] 	<ul style="list-style-type: none"> ◆ probability estimation tree [M]
	numeric	<ul style="list-style-type: none"> ◆ linear classifiers <ul style="list-style-type: none"> ◇ logistic regression [S] ◇ support vector machines [M] ◆ non-linear classifiers [M] 	<ul style="list-style-type: none"> ◆ statistical regression (linear and non-linear, simple and multivariate) [S] 	

- Models are used for predicting a variable, called
 - *dependent variable* in the context of statistics
 - *target variable* in the context of machine learning
- They are built using properties of data called
 - *independent variables* in the context of statistics
 - *features* in the context of machine learning

		Type of dependent [S] / target [M] variable		
		categorical	numeric	probability of categorical
Independent variable [S] / feature [M] type	categorical	<ul style="list-style-type: none"> supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	<ul style="list-style-type: none"> supervised segmentation using variance measures [M] 	<ul style="list-style-type: none"> probability estimation tree [M]
	numeric	<ul style="list-style-type: none"> linear classifiers logistic regression [S] support vector machines [M] non-linear classifiers [M] 	<ul style="list-style-type: none"> statistical regression (linear and non-linear, simple and multivariate) [S] 	

Naïve Bayes

- A method for estimating the probabilities of target variable values, based on the assumption of true independence among the attributes of the data
- Complete independence among all attributes is rarely possible in real data sets but the predictions made by the Naïve Bayes family of models are nevertheless very good for practical purposes
- The Naïve Bayes formula for the probability of a data instance belonging to class C_j , given a certain combination of attribute values, X_k , is:

$$P(C_j|X_k) = \frac{P(C_j) \prod_{i=1}^d P(x_{ik}|C_j)}{P(X_k)}$$

where the class C_j is a possible value of the categorical target variable, X_k is the vector of given attribute values (i.e. a data instance, identified as k), d is the number of attributes used in the model and x_{ik} is the value of attribute x_i in the vector of given attribute values (i.e. in data instance k). $P(C_j)$ is the estimated probability of any data instance in the set belonging to class C_j (i.e. the prior probability of C_j). $P(x_{ik}|C_j)$ is the estimated

probability of attribute x_i having value x_{ik} , given that the class is C_j . $P(X_k)$ is the estimated probability of data instance X_k occurring. The formula assumes independence among the attributes, $x_1, x_2 \dots x_d$.

- The model can be used both with numeric and categorical attributes. Probabilities for categorical attributes can be estimated using counts. For numeric attributes probability densities for a suitable distribution (e.g. normal or Poisson) are used.

Example

Outlook	Temperature	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Source: [DM]

First, let's look at the table and see what we can say about whether the game will be played if it is **Sunny**. There are five instances that have the value of **Sunny** for `Outlook` and of those 3 are associated with a **No** for `Play` and 2 with a **Yes**. In terms of probabilities, we can write this as:

$$P(\mathbf{Yes}|\mathbf{Sunny}) = \frac{2}{5} \quad P(\mathbf{No}|\mathbf{Sunny}) = \frac{3}{5}$$

Based on these values our prediction might be that if it is sunny, the game is more likely not to go ahead than to go ahead. However, we have used only one of the four attributes of the weather available to us. It would make more sense to use all data available i.e. all four attributes to make the prediction. What if one day the attribute values were: `Outlook=Sunny`, `Temperature=Cool`, `Humidity=High`, `Windy=True` and we wanted to see how likely it was that the game would go ahead? We don't even have a case like that in the table, let alone several cases to use for counting up incidences of **Yes** and **No**. This is where Naïve Bayes comes in:

$$P(\mathbf{Yes}|\mathbf{E}) = \frac{P(\mathbf{E}|\mathbf{Yes})P(\mathbf{Yes})}{P(\mathbf{E})} \quad P(\mathbf{No}|\mathbf{E}) = \frac{P(\mathbf{E}|\mathbf{No})P(\mathbf{No})}{P(\mathbf{E})} \quad (1)$$

where **E** is used as shorthand for the 'evidence', consisting of the set of values {**Sunny**, **Cool**, **High**, **True**}, hence $P(\mathbf{E}) = P(\mathbf{Sunny}, \mathbf{Cool}, \mathbf{High}, \mathbf{True})$.

The assumption of independence of attributes allows us to simply multiply the individual probabilities, whether conditional or not:

$$P(\mathbf{E}|\mathbf{Yes}) = P(\mathbf{Sunny}|\mathbf{Yes})P(\mathbf{Cool}|\mathbf{Yes})P(\mathbf{High}|\mathbf{Yes})P(\mathbf{True}|\mathbf{Yes}) \quad P(\mathbf{E}) = P(\mathbf{Sunny})P(\mathbf{Cool})P(\mathbf{High})P(\mathbf{True})$$

and these individual probabilities *can* be estimated using the data in the table.

Finally, because **Yes** and **No** are the only two possible outcomes, $P(\mathbf{Yes}|\mathbf{E}) + P(\mathbf{No}|\mathbf{E}) = P(\mathbf{E})$. From this and equation (1) it follows that

$$\frac{P(\mathbf{E}|\mathbf{Yes})P(\mathbf{Yes})}{P(\mathbf{E})} + \frac{P(\mathbf{E}|\mathbf{No})P(\mathbf{No})}{P(\mathbf{E})} = 1$$

and further that

$$P(\mathbf{E}) = P(\mathbf{E}|\mathbf{Yes})P(\mathbf{Yes}) + P(\mathbf{E}|\mathbf{No})P(\mathbf{No})$$

Substituting this into equations (1) we can reduce the probabilities that need to be estimated to the conditional probabilities:

$$P(\mathbf{Yes}|\mathbf{E}) = \frac{P(\mathbf{E}|\mathbf{Yes})P(\mathbf{Yes})}{P(\mathbf{E}|\mathbf{Yes})P(\mathbf{Yes}) + P(\mathbf{E}|\mathbf{No})P(\mathbf{No})} \quad P(\mathbf{No}|\mathbf{E}) = \frac{P(\mathbf{E}|\mathbf{No})P(\mathbf{No})}{P(\mathbf{E}|\mathbf{Yes})P(\mathbf{Yes}) + P(\mathbf{E}|\mathbf{No})P(\mathbf{No})}$$

Using the counts from the table to estimate the conditional probabilities we get:

$$P(\mathbf{E}|\mathbf{Yes})P(\mathbf{Yes}) = P(\mathbf{Sunny}|\mathbf{Yes})P(\mathbf{Cool}|\mathbf{Yes})P(\mathbf{High}|\mathbf{Yes})P(\mathbf{True}|\mathbf{Yes})P(\mathbf{Yes}) = \frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14} = 0.0053$$

$$P(\mathbf{E}|\mathbf{No})P(\mathbf{No}) = P(\mathbf{Sunny}|\mathbf{No})P(\mathbf{Cool}|\mathbf{No})P(\mathbf{High}|\mathbf{No})P(\mathbf{True}|\mathbf{No})P(\mathbf{No}) = \frac{3}{9} \times \frac{1}{9} \times \frac{4}{9} \times \frac{3}{9} \times \frac{5}{14} = 0.0206$$

and the *posterior** probabilities for **Yes** and **No**:

$$P(\mathbf{Yes}|\mathbf{E}) = \frac{0.0053}{0.0053 + 0.0206} = 0.205 \quad P(\mathbf{No}|\mathbf{E}) = \frac{0.0206}{0.0053 + 0.0206} = 0.795$$

So, with the evidence at hand (**Sunny**, **Cool**, **High** and **Yes**), we can say that it is nearly four times more probable that the game will not be played than that it will.

*The *a priori* probability of an outcome is the general probability of that outcome, while the posterior, or *a posteriori* probability is the probability of the outcome based on some information, i.e. 'evidence', in addition to the general probability.

		Type of dependent [S] / target [M] variable		
		categorical	numeric	probability of categorical
Independent variable [S] / feature [M] type	categorical	<ul style="list-style-type: none"> supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	<ul style="list-style-type: none"> supervised segmentation using variance measures [M] 	<ul style="list-style-type: none"> probability estimation tree [M]
	numeric	<ul style="list-style-type: none"> linear classifiers <ul style="list-style-type: none"> logistic regression [S] support vector machines [M] non-linear classifiers [M] 	<ul style="list-style-type: none"> statistical regression (linear and non-linear, simple and multivariate) [S] 	

Gaussian Naïve Bayes

- Used in Naïve Bayes probability estimations to include conditional probability contributions by normally distributed numeric variables
- If in the formula

$$P(C_j|X_k) = \frac{P(C_j) \prod_{i=1}^d P(x_{ik}|C_j)}{P(X_k)}$$

the value x_i is numeric the counting of appearances of x_{ik} would not make sense. Instead of the probability $P(x_{ik})$ a *probability density* (the quantity on the y-axis of a distribution graph for a numeric variable) value is used. Then the formula for the probability of membership of a class conditional on a particular set of values (the evidence) becomes:

$$P(C_j|X_k) = \frac{P(C_j) \prod_{i=1}^{d_{categorical}} P(x_{ik}|C_j) \prod_{i=1}^{d_{numeric}} f(x_{ik}|C_j)}{P(X_k)}$$

where $d_{categorical}$ and $d_{numeric}$ are the number of categorical and of numeric attributes in the data set, respectively, and $f(x_{ik}|C_j)$ is the conditional probability density of x at x_{ik} , given class C_j .

- In the particular case of a normally distributed variable x , the probability density is calculated using the following formula:

$$f_n(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x - \bar{x}}{\sigma} \right)^2}$$

where σ is the population standard deviation for x and \bar{x} is the mean value of x . For conditional probability density σ and \bar{x} are calculated for the values of x in the data subset of instances that meet the condition, for example of instances that are of a particular class.

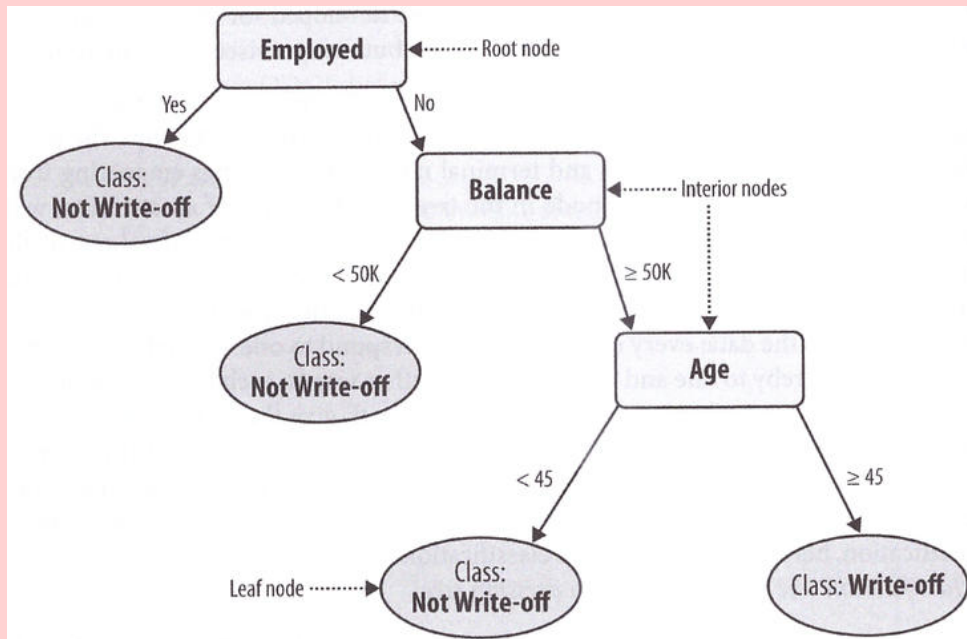
		Type of dependent [S] / target [M] variable		
		categorical	numeric	probability of categorical
Independent variable [S] / feature [M] type	categorical	<ul style="list-style-type: none"> supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	<ul style="list-style-type: none"> supervised segmentation using variance measures [M] 	<ul style="list-style-type: none"> probability estimation tree [M]
	numeric	<ul style="list-style-type: none"> linear classifiers logistic regression [S] support vector machines [M] non-linear classifiers [M] 	<ul style="list-style-type: none"> statistical regression (linear and non-linear, simple and multivariate) [S] 	

Classification trees

- Classification trees are models built using supervised segmentation
- A classification tree is built through a *tree learning* process using one of many existing *induction algorithms* (e.g. ID3 or CART) which
 - are all based on the *supervised segmentation* principles discussed in the previous lecture i.e. finding informative attributes for the target variable
 - progressively segment the data set, each instance of segmentation represented by a node in the tree and the branches stemming from it to 'child' nodes
 - result in leaf nodes that are associated with a particular value of the target variable
- A popular type of model because they are:
 - easy to understand
 - easy to use: at each step the problem is just a smaller version of the problem at the previous step (i.e. recursive)
 - relatively efficient
- Tree induction can continue until all attributes are used or leaf subsets are pure, but this is not done in practice (more about this in section on *overfitting*)

Examples

The picture below shows a simple classification tree for a target variable that has two values: **Write-off** and **Not Write-off**. The attribute **Employed** is used to split the entire dataset and is hence shown in the root node. As the value of **Yes** for this attribute is found to indicate the value **No Write-off** for the target variable, the node representing the subset grouped by value **Yes** for attribute **Employed** is a leaf node. Attributes **Balance** and **Age** are used to split the dataset further, until the leaf nodes are deemed to be informative enough regarding the value of the target variable (see section on overfitting also).



Source: [DSB]

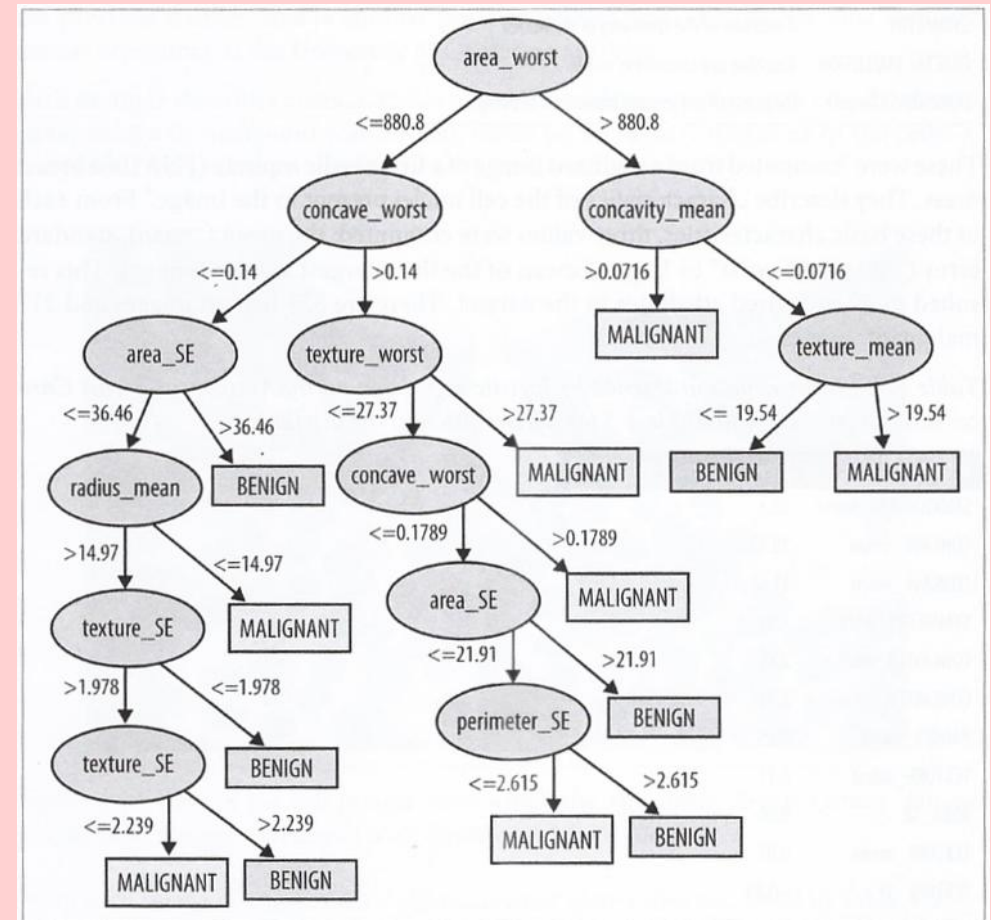


Figure 4-12. Decision tree learned from the Wisconsin Breast Cancer dataset.

Source: [DSB]

The picture above shows a classification tree for a real dataset relating to breast tissue scans examined for cancer. The target variable has values **MALIGNANT** and **BENIGN**. The attributes are numeric, categorized into ranges that produce a good predictive model.

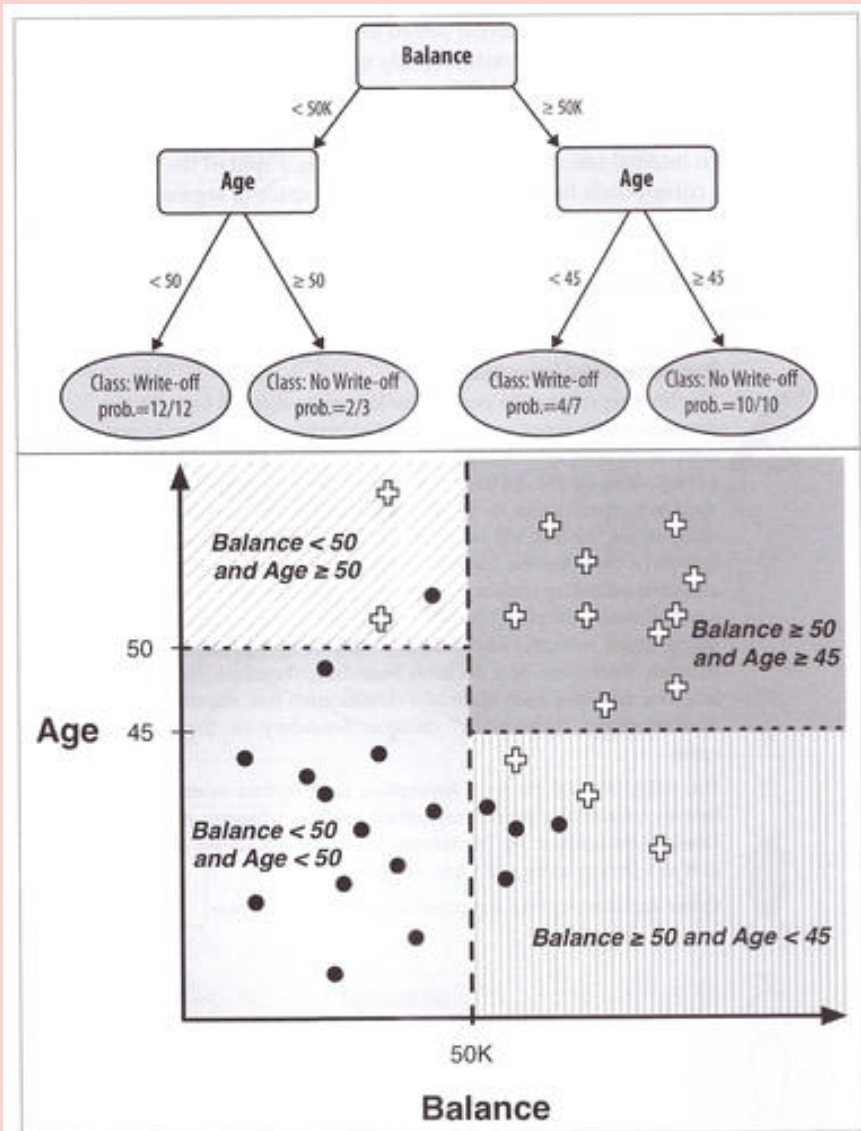


Figure 3-15. A classification tree and the partitions it imposes in instance space. The black dots correspond to instances of the class Write-off, the plus signs correspond to instances of class non-Write-off. The shading shows how the tree leaves correspond to segments of the population in instance space.

Source: [DSB]

Segment visualization

The picture shows a classification tree using two attributes: **Balance** and **Age**. With two attributes, it is possible to show the data in a two dimensional coordinate system, together with the lines along which the data set is split. The splitting for this particular set could be improved, for example if the age categorization boundary were changed to 55 for $Balance < 50K$, the two data subsets for $Balance < 50K$ would be pure.

Regression trees

[S] statistics [M] machine learning		Type of dependent [S] / target [M] variable		
Independent variable feature [M] type	[S]/ type	categorical	numeric	probability of categorical
		<ul style="list-style-type: none"> supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	<ul style="list-style-type: none"> supervised segmentation using variance measures [M] 	<ul style="list-style-type: none"> probability estimation tree [M]
	numeric	<ul style="list-style-type: none"> linear classifiers logistic regression [S] support vector machines [M] non-linear classifiers [M] 	<ul style="list-style-type: none"> statistical regression (linear and non-linear, simple and multivariate) [S] 	

- Regression trees are built for *continuous numeric target variables*
- Predictions made by a regression tree model are in the form of a numeric value
- Instead of measures of impurity, regression tree building algorithms use *variance* to define splitting criteria
- Variance is used in the context of a test such as ANOVA, which evaluates splits based on variance within and between subsets: the most informative split is that with the greatest difference between inter-subset and intra-subset variance

		Type of dependent [S] / target [M] variable		
		categorical	numeric	probability of categorical
Independent variable [S] / feature [M] type	categorical	<ul style="list-style-type: none"> supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	<ul style="list-style-type: none"> supervised segmentation using variance measures [M] 	<ul style="list-style-type: none"> probability estimation tree [M]
	numeric	<ul style="list-style-type: none"> linear classifiers logistic regression [S] support vector machines [M] non-linear classifiers [M] 	<ul style="list-style-type: none"> statistical regression (linear and non-linear, simple and multivariate) [S] 	

Probability estimation trees

- Apart from classification and regression, tree models can be used to determine the *probability* of a data instance belonging to a certain class
- There are many applications for which probability estimation is the more appropriate approach than classification:
 - Identification of segments with highest probability of class membership in situations where the overall probability is very low. For example, if a bank is trying to identify the customers that are likely to default on their mortgage and the probability of default is low in all segments of the customer population, a classification model would classify all the segments as **not likely to default** and would not be very useful. A probability estimation model would distinguish between probabilities in the different segments, e.g. 0.5, 0.3 and 0.1, all low but different!
 - Ranking of data instances by probability. For example, if a company has a budget to spend on customer retention, they may want to rank all the customers by probability of leaving and offer incentives to the number of highest ranking potential leavers that the budget allows.

- *Probability estimation trees* are induced using the same principles as classification trees
- Probability of class membership for a segment (corresponding to a leaf data subset in a tree), is estimated using the following formula:

$$p(c) = \frac{n + 1}{n + m + 2}$$

where c is the class, n is the number of instances belonging to class c and m is the number of instances not belonging to class c . The formula $p(c) = n/(m + n)$ intuitively seems more appropriate, however, the *Laplace correction* is included to account for effects of small data subsets.

Examples

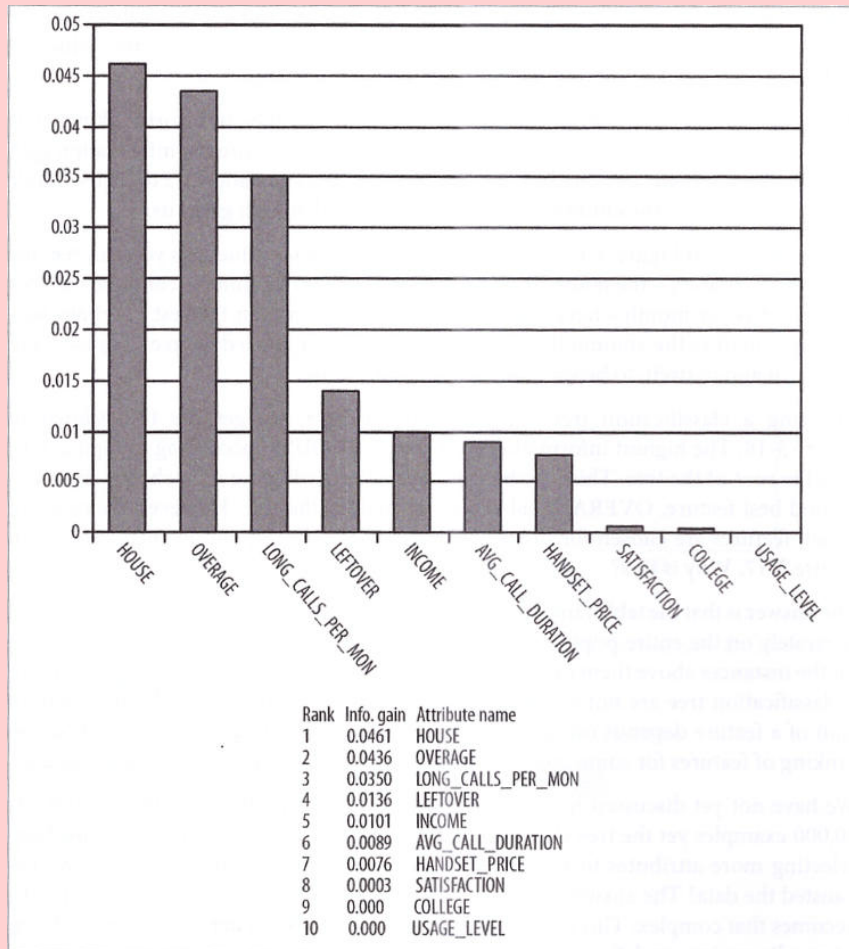


Figure 3-17. Churn attributes from Table 3-2 ranked by information gain.

Source: [DSB]

A probability estimation tree for target variable with values **CHURN** and **STAY**, regarding the likelihood of the customers of a mobile phone company to leave. The picture on the left lists the attributes that are used in the tree shown in the picture on the right.

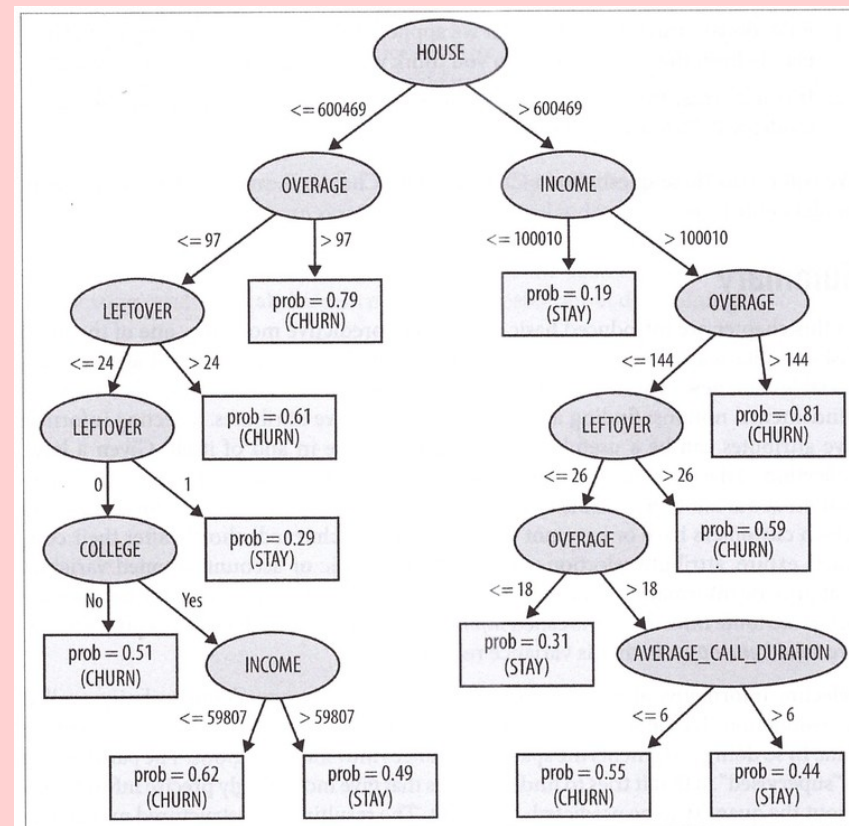


Figure 3-18. Classification tree learned from the cellular phone churn data. Rectangular leaves correspond to segments of the population, defined by the path from the root at the top. Probabilities at the leaves are the estimated probabilities of churning for the corresponding segment; in parentheses are shown the classifications resulting from applying a decision threshold of 0.5 to the probabilities (i.e., are the individuals in the segment more likely to CHURN or to STAY?).

Source: [DSB]

[S] statistics [M] machine learning		Type of dependent [S] / target [M] variable		
		categorical	numeric	probability of categorical
Independent variable [S]/ feature [M] type	categorical	<ul style="list-style-type: none"> supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	<ul style="list-style-type: none"> supervised segmentation using variance measures [M] 	<ul style="list-style-type: none"> probability estimation tree [M]
	numeric	<ul style="list-style-type: none"> linear classifiers logistic regression [S] support vector machines [M] non-linear classifiers [M] 	<ul style="list-style-type: none"> statistical regression (linear and non-linear, simple and multivariate) [S] 	

Linear classifiers

- A linear classifier as a predictive model has numeric inputs and a categorical value as output
- A linear classifier is most easily illustrated in a two-dimensional cartesian coordinate system, where data instances labelled by the value of the categorical target variable are shown at the points defined by their attribute values. The linear classifier is a straight line that splits the entire space into two parts, each containing as homogeneous as possible a data subset with regard to the target variable (see examples below).
- This idea can be extended to any number of dimensions, corresponding to the same number of attributes, with the dividing artefact being an $(n - 1)$ -dimensional plane, where n is the number of attributes.
- The form of the general linear model is:

$$f(x) = w_0 + w_1x_1 + w_2x_2... + w_{n-1}x_{n-1}$$

where n is the number of variables, x_i are the attributes and w_i are the linear function parameters. The function $f(x)$ is called the *linear discriminant function*.

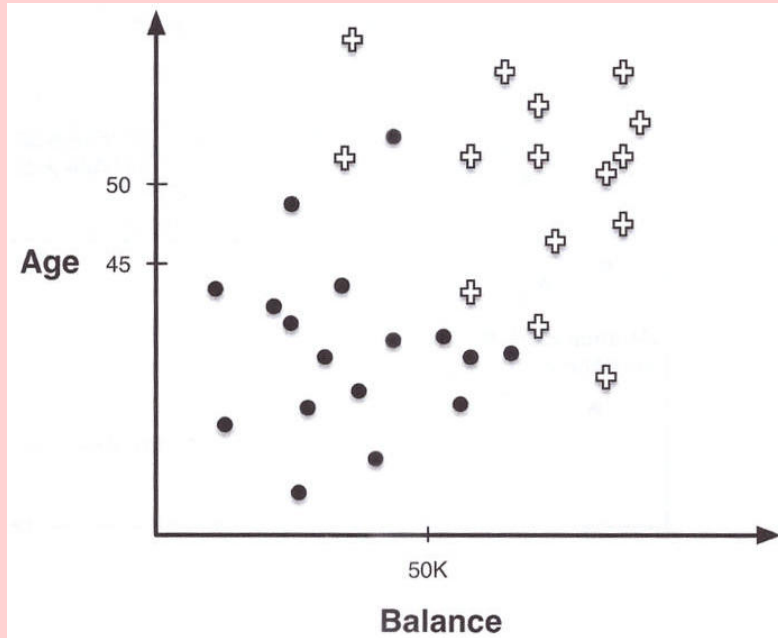
- The aim of linear classification is to find the geometrical artefact that best divides the attribute space with respect to the target classes (target variable values). This line, plane or hyperplane (the term used for a plane with more than 2 dimensions) is expressed as

$$f(x) = 0$$

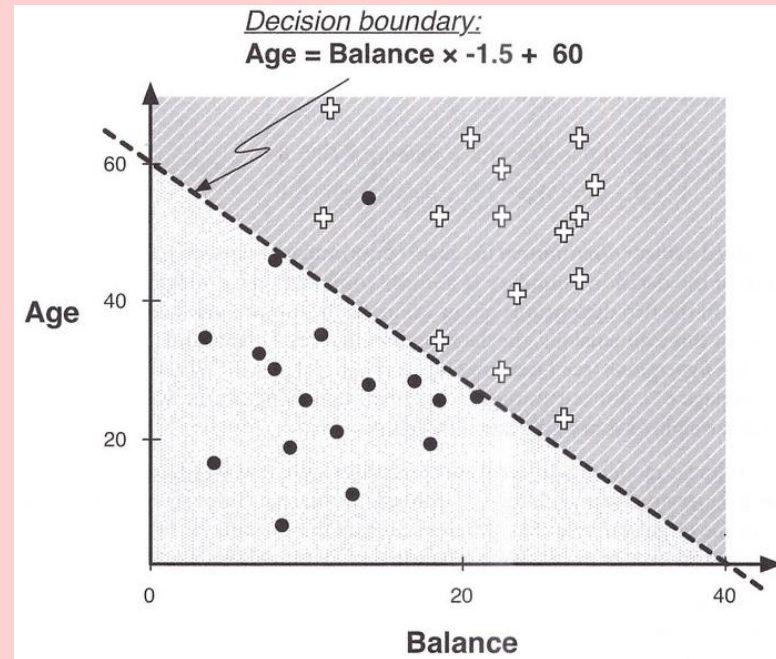
and its parameters parameters (weights) w_i are chosen by a process of optimization, or learning, with an exact objective that depends on the method used.

- Once the weights are determined i.e. learned, the linear discriminant function can be used to assign new data instances to classes, with $f(x) > 0$ placing data instance x in one class and $f(x) \leq 0$ in the other.
- The weights, w_i , can be interpreted as a measure of importance of the attributes to which they apply, provided the attribute values were normalised preceding the weight learning (optimization) process.
- The most common linear classifier methods are:
 - Logistic regression
 - Support Vector Machines

Examples



Source: [DSB]



Source: [DSB]

The picture on the left shows the same data set as that shown for segment visualisation, but without the splitting lines. A boundary line is shown alongside its linear formula using attributes **Balance** and **Age**.

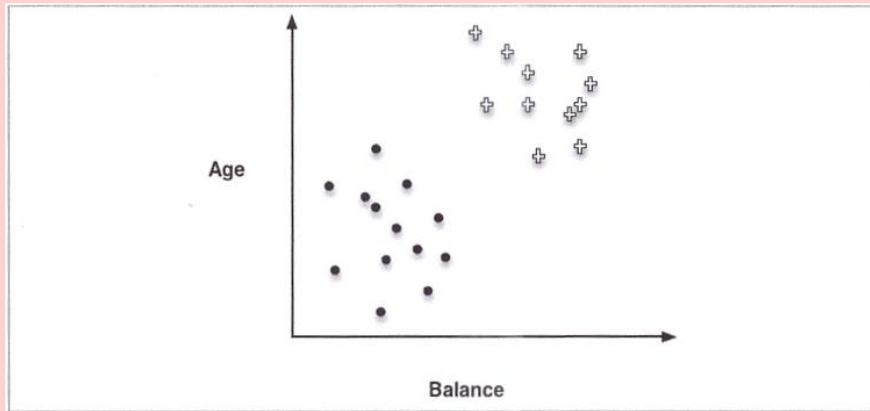


Figure 4-4. A basic instance space in two dimensions containing points of two classes.

Source: [DSB]

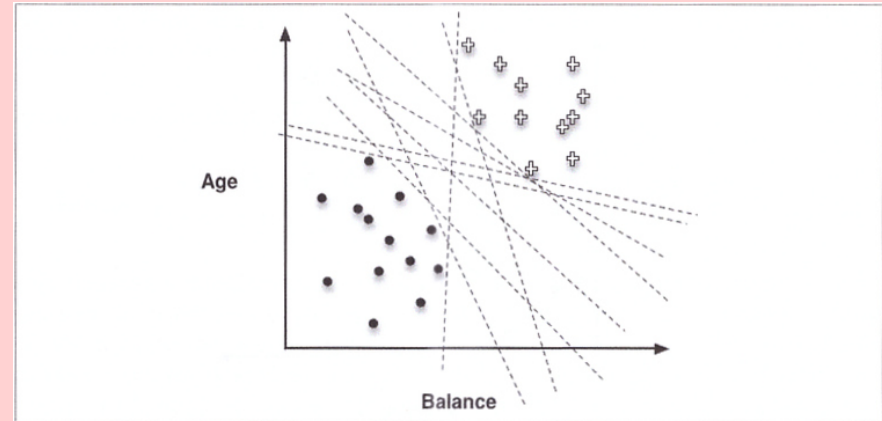


Figure 4-5. Many different possible linear boundaries can separate the two groups of points of Figure 4-4.

Source: [DSB]

The picture on the left shows another dataset in the two-attribute space, while the one on the right illustrates a few of the infinite number of ways in which the linear discriminant could be defined.

[S] statistics [M] machine learning		Type of dependent [S] / target [M] variable		
		categorical	numeric	probability of categorical
Independent variable feature [M] type	categorical	<ul style="list-style-type: none"> supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	<ul style="list-style-type: none"> supervised segmentation using variance measures [M] 	<ul style="list-style-type: none"> probability estimation tree [M]
	numeric	<ul style="list-style-type: none"> linear classifiers ◊ logistic regression [S] ◊ support vector machines [M] non-linear classifiers [M] 	<ul style="list-style-type: none"> statistical regression (linear and non-linear, simple and multivariate) [S] 	

Logistic regression

- Logistic regression is, in spite of its name, in fact a classification model. Its input attributes are numerical and the output is either a category or a probability of belonging to a category.

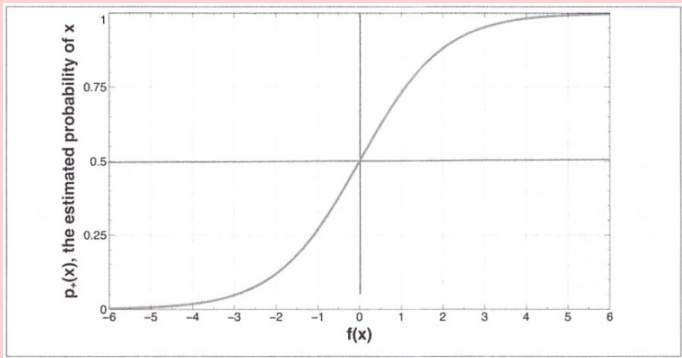


Figure 4-10. Logistic regression’s estimate of class probability as a function of $f(x)$, (i.e., the distance from the separating boundary). This curve is called a “sigmoid” curve because of its “S” shape, which squeezes the probabilities into their correct range (between zero and one).

Source: [DSB]

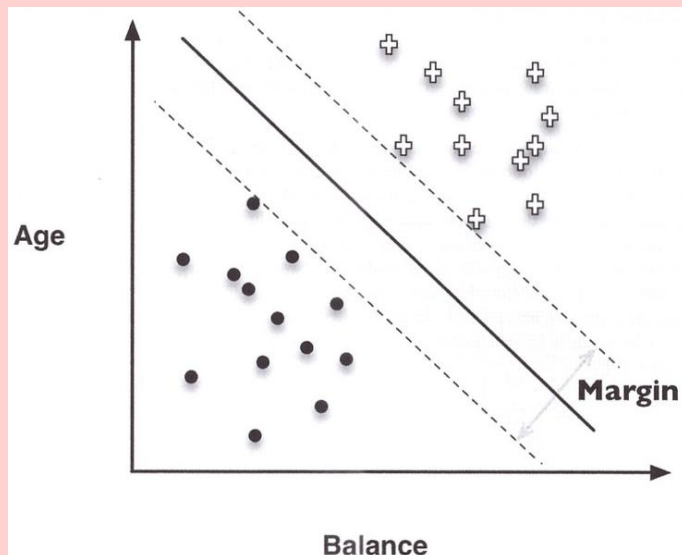
- Instead of using the linear discriminant function, it uses a transformation of it:

$$p_{+}(x) = \frac{1}{1 + e^{-f(x)}}$$
 where $p_{+}(x)$ is the probability of data instance x belonging to a particular class (having one of two target variable values), $f(x)$ and is a function identical to the linear discriminant function.
- The value of $p_{+}(x)$ is limited to the range between -1.0 and 1.0 and once the model is in place (i.e. the parameter values have been learned) and given a new data instance, the logistic model outputs a probability of the instance belonging to the class of interest
- Defining a probability threshold (e.g. 0.5), allows the use of the logistic model for classification.

[S] statistics [M] machine learning		Type of dependent [S] / target [M] variable		
		categorical	numeric	probability of categorical
Independent variable [S]/ feature [M] type	categorical	<ul style="list-style-type: none"> supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	<ul style="list-style-type: none"> supervised segmentation using variance measures [M] 	<ul style="list-style-type: none"> probability estimation tree [M]
	numeric	<ul style="list-style-type: none"> linear classifiers logistic regression [S] support vector machines [M] non-linear classifiers [M] 	<ul style="list-style-type: none"> statistical regression (linear and non-linear, simple and multivariate) [S] 	

Support Vector Machines

- Support Vector Machines are another type of linear classifier.



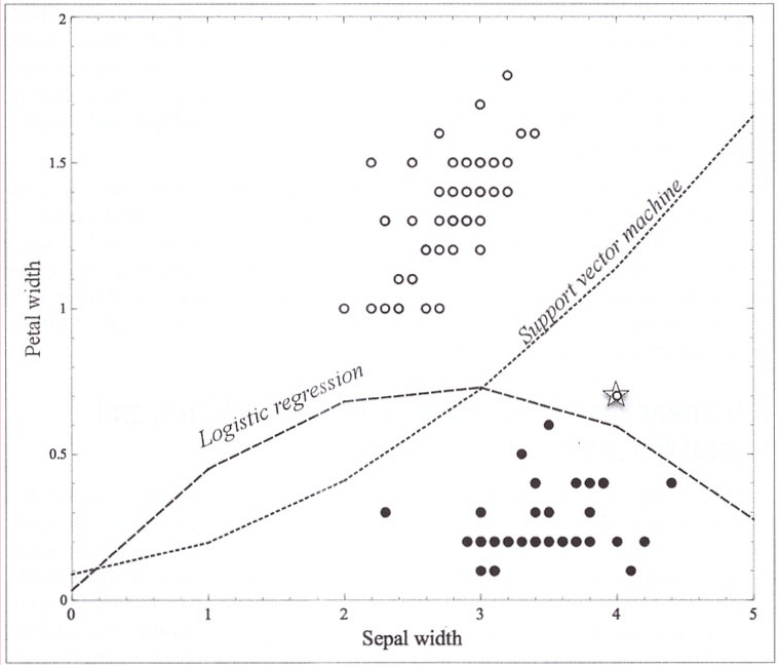
Source: [DSB]

- The boundary line (or plane or hyperplane) is determined by finding the maximal margin that can be achieved between data instances, with respect to attribute values, between instances belonging to one target class and those belonging to the other. The boundary is drawn along the middle of the margin-wide strip.
- If the data instances cannot be cleanly separated in this manner, a margin-wide strip is still defined, with the optimization including 'penalties' to account for data instances 'on the wrong side' of the boundary

[S] statistics [M] machine learning		Type of dependent [S] / target [M] variable		
Independent variable feature [M] type	[S] / type	categorical	numeric	probability of categorical
		<ul style="list-style-type: none"> supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	<ul style="list-style-type: none"> supervised segmentation using variance measures [M] 	<ul style="list-style-type: none"> probability estimation tree [M]
Independent variable feature [M] type	[S] / type	categorical	numeric	probability of categorical
		<ul style="list-style-type: none"> linear classifiers logistic regression [S] support vector machines [M] non-linear classifiers [M] 	<ul style="list-style-type: none"> statistical regression (linear and non-linear, simple and multivariate) [S] 	

Non-linear classifiers

Classification methods may use non-linear functions (where the attribute values contribute with a power other than 1). In this case the boundary line for classification is not straight.



Source: [MSD]

The picture shows data being classified using non-linear logistic regression and a non-linear support vectory machine.

		Type of dependent [S] / target [M] variable		
		categorical	numeric	probability of categorical
Independent variable [S] / feature [M] type	categorical	<ul style="list-style-type: none"> supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	<ul style="list-style-type: none"> supervised segmentation using variance measures [M] 	<ul style="list-style-type: none"> probability estimation tree [M]
	numeric	<ul style="list-style-type: none"> linear classifiers logistic regression [S] support vector machines [M] non-linear classifiers [M] 	<ul style="list-style-type: none"> statistical regression (linear and non-linear, simple and multivariate) [S] 	

Statistical regression

- Statistical regression models can be linear and non-linear, single variable and multi-variate.
- As a prediction model, regression takes numeric inputs and produces numeric outputs.
- A regression model is *fitted* using known data, then used for the prediction of the dependent variable in data where it is not known.
- We will be looking at linear regression only, but the principles can be extended to the other types of regression.
- A linear regression model can be generated if there is a linear relationship between the independent and dependent variable.
- The mathematical model for linear regression is the same linear function as used for classification:

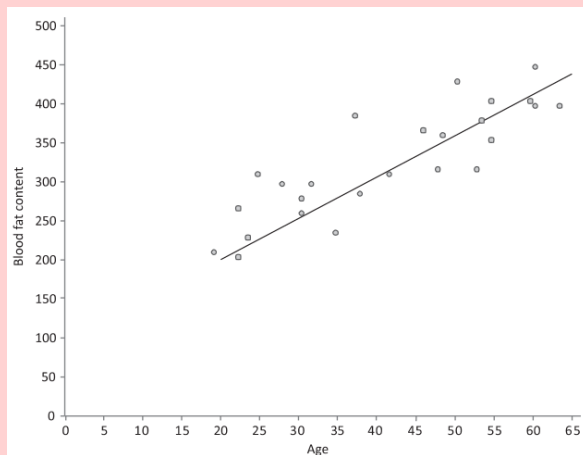
$$f(x) = w_0 + w_1x_1 + w_2x_2... + w_{n-1}x_{n-1}$$

where $x_1, x_2...x_n$ are the independent variables and $w_0, w_1, w_2... w_n$ are the coefficients of the linear function. With regression, $f(x)$ corresponds to the dependent variable and the coefficients are chosen on the principle of minimizing the average distance of the existing

data observations from the model line (or plane or hyperplane). The model *represents the relationship* between the independent variables on one side and the dependent variable on the other: given a new observation, the dependent variable can be deduced from the independent variables by using the model. In the case of classification the purpose of the line was to *divide the space* and hence define subsets of the data.

Fitting a linear regression model

- A linear regression model with one independent variable can be fitted 'manually'.



The picture shows an example of some data with one independent variable together with the fitted regression line.

Source: [MSD]

- The regression line formula is:

$$y = a + bx$$

where y is the dependent variable, x is the independent variable and a and b are the coefficients of the linear function.

- The coefficient b , which is in fact the slope of the line, is calculated as follows:

$$b = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

where x_i and y_i are the independent and dependent variable values for the i^{th} observation, n is the number of observations and \bar{x} and \bar{y} are the means of the independent and dependent variable, respectively.

- The coefficient a can be calculated by substituting the variables in the equation by the means:

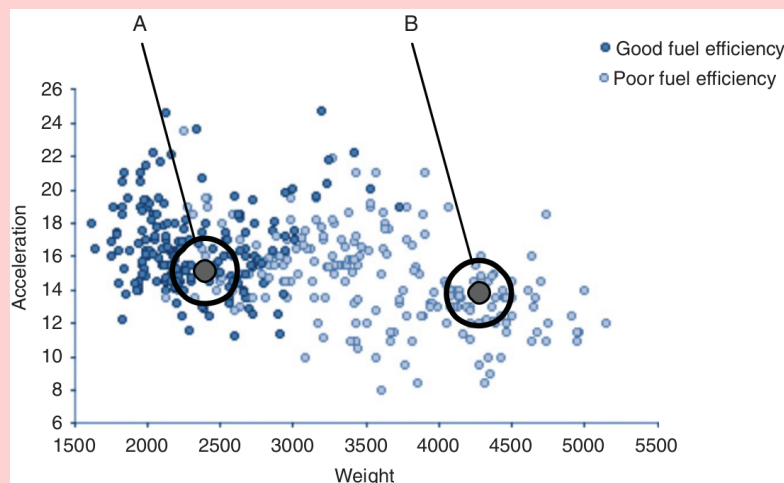
$$a = \bar{y} - b \times \bar{x}$$

- Fitting a model with more than one independent variable is more complicated and is generally done using a computer.

		Type of dependent [S] / target [M] variable		
		categorical	numeric	probability of categorical
Independent variable [S]/ feature [M] type	categorical	♦ supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] ■	♦ supervised segmentation using variance measures [M] ■	♦ probability estimation tree [M] ■
	numeric	♦ linear classifiers ♦ logistic regression [S] ♦ support vector machines [M] ♦ non-linear classifiers [M] ■	♦ statistical regression (linear and non-linear, simple and multivariate) [S] ■	■

k-Nearest Neighbours (k-NN)







- This method can be used with any types of attributes and target variables
- It is an **algorithm** for determining a target variable's value for a data instance, given a labelled data set (one in which the target variable's value is known for all instances)



Source: [MSD]

In the picture **fuel efficiency** (possible values: **Good** and **Bad**) is determined for two cars, based on existing fuel efficiency data as related to weight and acceleration, using k-NN.

- The distance is calculated between the new data instance and the instances in the dataset, based on the values of the attributes, which can be numeric (distance is Euclidean) or categorical (some measure of distance e.g. Jaccard)
- The target variable of the new instance is calculated based on the target variable values of the closest k instances (e.g. as average if numeric or most frequent value if categorical)
- Determining the best value for k is also part of the algorithm

[S] statistics [M] machine learning		Type of dependent [S] / target [M] variable		
		categorical	numeric	probability of categorical
Independent variable feature [M] type	categorical	♦ supervised segmentation using impurity measures (can be used for the numeric-categorical combination if numeric attributes converted to categories) [M] 	♦ supervised segmentation using variance measures [M] 	♦ probability estimation tree [M] 
	numeric	♦ linear classifiers ◊ logistic regression [S] ◊ support vector machines [M] ♦ non-linear classifiers [M] 	♦ statistical regression (linear and non-linear, simple and multivariate) [S] 	

Artificial neural networks (ANN)

- Computational models that loosely emulate the workings of a brain
- ANNs can be designed for any kind of input and output variable types
- They are non-linear models, with parameters (weights) that need to be fitted to data
- The structure of an ANN consists of nodes and connections between those nodes:
 - a node applies a weight to the value that enters it, despatching the new value further into the network
 - the connections define the flow of information between nodes
 - the input attributes are 'fed' into the first layer of nodes and the output values received from the last layer
- The learning, i.e. parameter fitting, in an ANN takes place through multiple cycles of data flowing through the network and weight adjustments using a suitable algorithm

Overfitting

- *Overfitting* is the application of model-building procedures to an extent that causes the created model to include properties specific to the training data, in addition to those actually pertinent to a general model
- For example:
 - a classification tree can be refined until all the leaf nodes are pure, but in most cases such a tree is over-fitted
 - with a mathematical function overfitting may occur through the addition of too many attributes: more attributes means more dimensions and a better fit, but to the detriment of general applicability
- Overfitting can be avoided by using *holdout* data to test that the accuracy of a model when used on new data matches, or is close to, the accuracy of the model when used on training data

Examples

The optimal complexity of a particular type of model built with some data set can be determined with the help of a **fitting graph**. A fitting graph plots the accuracy of models of increasing complexity (e.g. represented by the number of nodes if a model is a decision tree), both for the training set and for a test set. As the complexity increases, so does the accuracy for both the training and test (holdout) set, until a critical complexity value at which the accuracy for the two sets diverges. It is for complexity higher than this critical value that the model is overfitted.

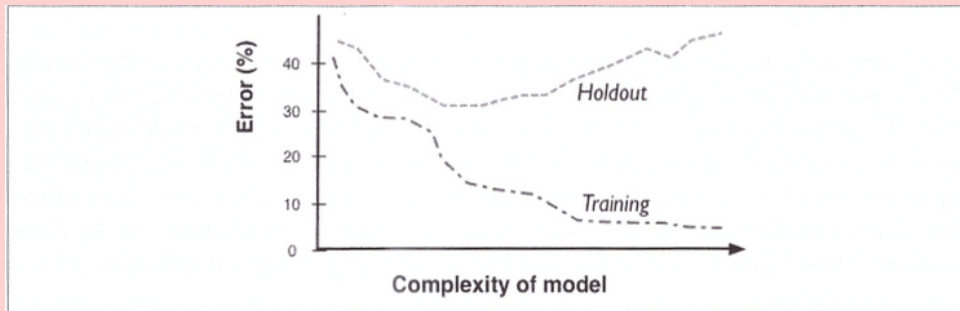


Figure 5-1. A typical fitting graph. Each point on a curve represents an accuracy estimation of a model with a specified complexity (as indicated on the horizontal axis). Accuracy estimates on training data and testing data vary differently based on how complex we allow a model to be. When the model is not allowed to be complex enough, it is not very accurate. As the models get too complex, they look very accurate on the training data, but in fact are overfitting—the training accuracy diverges from the holdout (generalization) accuracy.

Source: [DSB]

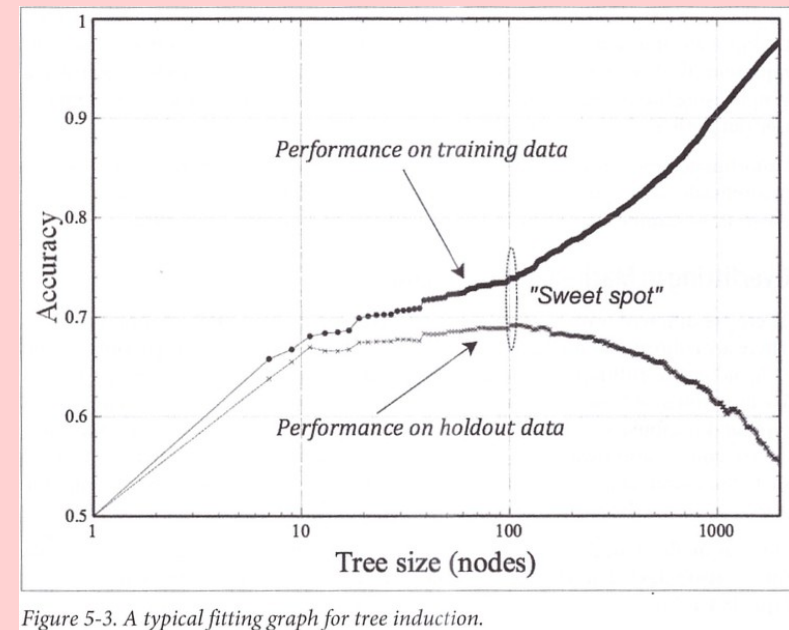


Figure 5-3. A typical fitting graph for tree induction.

Source: [DSB]

References

The pictures in this presentation were taken from the following books. The source for each picture is cited beside it.

[DSB] *Data Science for Business: What you need to know about data mining and data-analytic thinking*, by Foster Provost and Tom Fawcett, O'Reilly Media, 2013.

[MSD] *Making Sense of Data I: A Practical Guide to Exploratory Data Analysis and Data Mining*, by Glenn J. Myatt and Wayne P. Johnson, John Wiley & Sons, 2014.