

Module 4: Linear Classification

TMA4268 Statistical Learning V2020

Stefanie Muff, Department of Mathematical Sciences, NTNU

February xx, 2020

Introduction

Learning material for this module

- James et al (2013): An Introduction to Statistical Learning. Chapter 4.

We need more statistical theory than is presented in the textbook, which you find in this module page.

Some of the figures and slides in this presentation are taken (or are inspired) from “An Introduction to Statistical Learning, with applications in R” (Springer, 2013) with permission from the authors: G. James, D. Witten, T. Hastie and R. Tibshirani.

Topics in this module (todo)

What is classification?

- By now our responses Y was assumed *continuous*, while covariates were allowed to be *categorical*.
- Now we allow the response to be *categorical*.
- This is even more common than continuous responses. Examples:
 - Spam filters `email` $\in \{\text{spam}, \text{ham}\}$,
 - Eye color $\in \{\text{blue}, \text{brown}, \text{green}\}$.
 - Medical condition $\in \{\text{disease1}, \text{disease2}, \text{disease3}\}$.

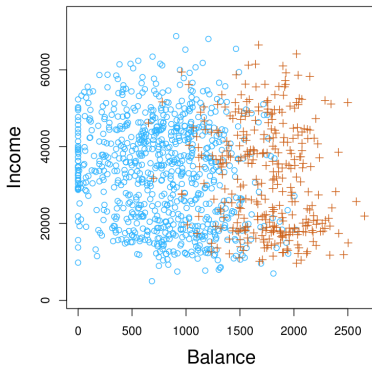
- Suppose we have a qualitative response value that can be a member in one of K classes $\mathcal{C} = \{c_1, c_2, \dots, c_K\}$.
- In classification we build a function $f(X)$ that takes a vector of input variables X and predicts its class membership, such that $Y \in \mathcal{C}$.
- We would also assess the *uncertainty* in this classification. Sometimes the role of the different predictors may be of main interest.
- We often build models that **predict probabilities of categories**, *given* certain covariates X .

Example: The credit card data

The `Default` dataset is available from the ISLR package.

Aim: to predict whether an individual will default on his or her credit card payment, given the annual income and credit card balance.

Orange: `default=yes`, blue: `default=no`.



General classification setup

Set-up: Training observations $\{(x_1, y_1), \dots, (x_n, y_n)\}$ where the response variable Y is categorical, e.g $Y \in \mathcal{C} = \{0, 1, \dots, 9\}$ or $Y \in \mathcal{C} = \{dog, cat, \dots, horse\}$.

Aim: To *build* a classifier $f(X)$ that assigns a class label from \mathcal{C} to a future unlabelled observation x and to assess the *uncertainty* in this classification.

Performance measure: Most popular is the misclassification error rate (training and test version).

What are the methods?

Three methods for classification are discussed here:

- Logistic regression
- K -nearest neighbours
- Linear and quadratic discriminant analysis

Linear regression for binary classification?

Suppose we have a binary outcome, for example whether a credit card user defaults $Y = \text{yes}$ or no , given covariates X to predict Y . We could use *dummy encoding* for Y like

$$Y = \begin{cases} 0 & \text{if no ,} \\ 1 & \text{if yes .} \end{cases}$$

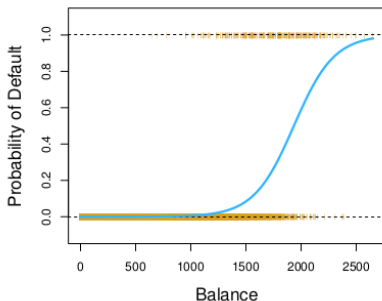
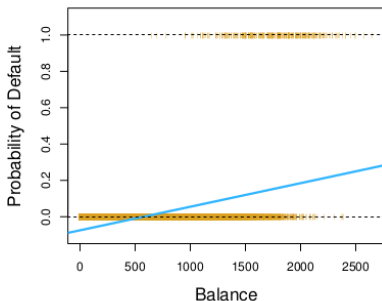
Can we simply perform a linear regression of Y on X and classify as **yes** if $\hat{Y} > 0.5$?

- In this case of a binary outcome, linear regression does a good job as a classifier, and is equivalent to linear discriminant analysis which we discuss later.
- Since in the population $E(Y \mid X = x) = \Pr(Y = 1 \mid X = x)$, we might think that regression is perfect for this task.
- However, linear regression might produce probabilities less than zero or bigger than one.

→ We need to use **logistic regression**.

Linear vs. logistic regression

Let's anticipate a bit, to see why linear regression does not so well. We estimate the probability that someone defaults, given the credit card balance as predictor:



Linear regression for categorical classification?

What when there are more than two possible outcomes? For example, a medical diagnosis Y , given predictors X can be categorized as

$$Y = \begin{cases} 1 & \text{if stroke ,} \\ 2 & \text{if drug overdose ,} \\ 3 & \text{if epileptic seizure .} \end{cases}$$

This suggests an ordering, but this is artificial.

- Linear and logistic regression are not appropriate here.
- We need *Multiclass logistic regression* and *Discriminant Analysis*.

However:

- It is still possible to use linear regression for classification problems with two classes. It is actually not even a bad idea, and works well under some conditions. If the conditional class densities are (multivariate) normal with equal covariance matrices then this linear regression (with 0 and 1 response) will in fact give the same classification as linear discriminant analysis (LDA).
- For categorical outcomes with more than two level, it requires some extra work (multivariate Y due to the dummy variable coding).
- We leave linear regression for now.
- For two classes *binary regression*, in particular *logistic regression*, is very popular - and is up next.

Logistic regression

- In logistic regression we consider a classification problem with two classes.
- Assume that Y is coded ($\mathcal{C} = \{1, 0\}$ or $\{\text{success}, \text{failure}\}$), and we focus on success ($Y = 1$).
- We may assume that Y_i follows a **Bernoulli distribution** with probability of success p_i .

$$Y_i = \begin{cases} 1 & \text{with probability } p_i, \\ 0 & \text{with probability } 1 - p_i. \end{cases}$$

- **Aim:** For covariates (X_1, \dots, X_p) , we want to estimate $p_i = \Pr(Y_i = 1 \mid X_1, \dots, X_p)$.

- We need a clever way to *link* our covariates X_1, \dots, X_p with this probability p_i . Aim: want to relate the *linear predictor*

$$\eta_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

to p_i . How?

- The idea is to use a so-called *link-function* to link p_i to the linear predictor.
- In logistic regression, we use the *logistic link function*

$$\log \left(\frac{p_i}{1 - p_i} \right) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} . \quad (1)$$

- **Q:** What is the rationale behind this? Other transformations that bring the linear predictor into a 0 to 1 range?

Logistic regression with one covariate

- Equation (1) can be rearranged and solved for p_i . Let's look at this for only one covariate:

$$p_i = \frac{e^{\beta_0 + \beta_1 x_i}}{1 + e^{\beta_0 + \beta_1 x_i}}.$$

- **Important:** These values p_i will always lie in the interval between 0 and 1, with an S-shaped curve.
- The parameter β_1 determines the rate of increase or decrease of the S-shaped curve, and the sign indicates whether the curve ascends or descends.

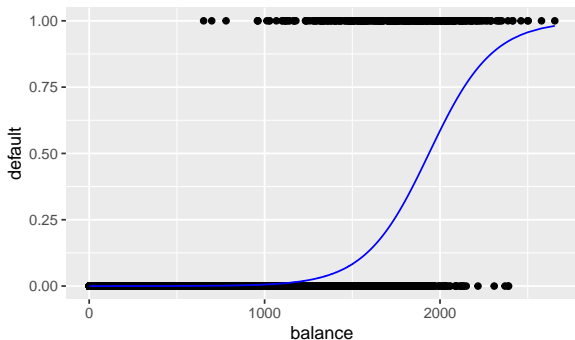
Example: Default credit card data

- The parameters are estimated using the method of maximum likelihood - we will look at that soon.
- Let's first do it! In R, this works with the `glm()` function, where we specify `family="binomial"`.

```
library(ISLR)
Default$default <- as.numeric(Default$default) - 1
glm_default = glm(default ~ balance, data = Default, family = "binomial")
summary(glm_default)$coef
```

```
##              Estimate Std. Error  z value    Pr(>|z|)
## (Intercept) -10.651330614 0.3611573721 -29.49221 3.623124e-191
## balance      0.005498917 0.0002203702  24.95309 1.976602e-137
```

Plotting the fitted line (in blue):



Default data: here $\hat{\beta}_0 = -10.65$ and $\hat{\beta}_1 = 0.005$.

Estimating the regression coefficients with ML

- The coefficients β_0, β_1, \dots are estimated with *maximum likelihood* (ML).
- We assume that pairs of covariates and responses $\{\mathbf{x}_i, y_i\}$ are measured independently of each other. Given n such observation pairs, the likelihood function of a logistic regression model can be written as:

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n L_i(\boldsymbol{\beta}) = \prod_{i=1}^n f(y_i; \boldsymbol{\beta}) = \prod_{i=1}^n (p_i)^{y_i} (1 - p_i)^{1-y_i},$$

where $\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2, \dots, \beta_p)^T$ enters into p_i

$$p_i = \frac{\exp(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip})}{1 + \exp(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip})}.$$

- The maximum likelihood estimates are found by maximizing the likelihood.
- To make the math easier, we usually work with the log-likelihood (the log is a monotone transform, thus it will give the same result as maximizing the likelihood).

$$\begin{aligned}\log(L(\beta)) &= l(\beta) = \sum_{i=1}^n \left(y_i \log p_i + (1 - y_i) \log(1 - p_i) \right) \\ &= \sum_{i=1}^n \left(y_i \log \left(\frac{p_i}{1 - p_i} \right) + \log(1 - p_i) \right) \\ &= \sum_{i=1}^n \left(y_i (\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip}) - \log(1 + e^{\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip}}) \right).\end{aligned}$$

- To maximize the log-likelihood function we find the $r + 1$ partial derivatives, and set equal to 0.
- This gives us a set of $p + 1$ non-linear equations in the β s.
- This set of equations does not have a closed form solution.
- The equation system is therefore solved numerically using the *Newton-Raphson algorithm* (or Fisher Scoring).

Qualitative interpretation of the coefficients

Let's again look at the regression output:

```
summary(glm_default)$coef
```

##	Estimate	Std. Error	z value	Pr(> z)
## (Intercept)	-10.651330614	0.3611573721	-29.49221	3.623124e-191
## balance	0.005498917	0.0002203702	24.95309	1.976602e-137

- The z -statistic is equal to $\frac{\hat{\beta}}{SE(\hat{\beta})}$, and is approximately $N(0, 1)$ distributed¹
- The p -value is $\Pr(|Z| > |z|)$ for a $Z \sim N(0, 1)$ random variable
- **Balance** seems an important predictor for **Default** (yes/no), with a higher balance leading to a higher probability of defaulting ($\beta_1 > 0$).

¹With this knowledge we can construct confidence intervals and test hypotheses about the β s, with the aim to understand which covariate(s) contribute to our posterior probabilities and classification.

Quantitative interpretation of the coefficients

Remember from equation (1) that

$$\log \left(\frac{p_i}{1 - p_i} \right) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} ,$$

thus

$$\frac{p_i}{1 - p_i} = e^{\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}} = e^{\eta_i} .$$

The quantity $p_i/(1 - p_i)$ is called the *odds*. Odds represent *chances* (e.g. in betting).

1. You think your football team will win tonight with a probability $p = 80\%$. What is the odds that it will win?
2. The odds for the best horse in a race is 9 : 1. What is the probability that this horse will win?

Why is the odds relevant?

Let's again rearrange the *odds* in the logistic regression model:

$$\begin{aligned}\frac{p_i}{1 - p_i} = \text{odds}(Y_i = 1 \mid X = x) &= \frac{P(Y_i = 1 \mid X = x)}{P(Y_i = 0 \mid X = x)} \\ &= \exp(\beta_0) \cdot \exp(\beta_1 x_{i1}) \cdot \dots \cdot \exp(\beta_p x_{ip}) .\end{aligned}$$

→ We have a *multiplicative model* for the odds - which can help us to interpret our β s.

The odds ratio

To understand the effect of a regression coefficient β_j , let's see what happens if we increase x_{ij} to $x_{ij} + 1$, while all other covariates are kept fixed.

Using simple algebra and the formula on the previous slide, you will see that

$$\frac{\text{odds}(Y_i = 1 \mid X_j = x_{ij} + 1)}{\text{odds}(Y_i = 1 \mid X_j = x_{ij})} = \exp(\beta_j) . \quad (2)$$

Interpretation: By increasing covariate x_{ij} by one unit, we change the odds ratio for $Y_i = 1$ by a factor $\exp(\beta_j)$.

Moreover: Taking the log on equation (2), it follows that β_j can be interpreted as a **log odds-ratio**.

Let's now fit the logistic regression model for **default**, given **balance**, **income** and the binary variable **student** as predictors:

```
glm_default2 = glm(default ~ balance + income + student, data = Default,  
  family = "binomial")  
  
summary(glm_default2)$coef
```

	Estimate	Std. Error	z value	Pr(> z)
## (Intercept)	-1.086905e+01	4.922555e-01	-22.080088	4.911280e-108
## balance	5.736505e-03	2.318945e-04	24.737563	4.219578e-135
## income	3.033450e-06	8.202615e-06	0.369815	7.115203e-01
## studentYes	-6.467758e-01	2.362525e-01	-2.737646	6.188063e-03

Questions:

- What happens with the odds-ratio when **income** increases by 10'000 dollars?
- What happens with the odds-ratio when **balance** increases by 100 dollars?

Predictions

Let's catch up. Why were we doing all this in the first place?

Answer: We wanted to build a model that **predict probabilities of categories** Y , *given* certain covariates X_1, \dots, X_p .

- For given parameter estimates $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ and a new observation \mathbf{x}_0 , we can estimate the probability $\hat{p}(\mathbf{x}_0)$ that the new observation belongs to the class defined by $Y = 1$

$$\hat{p}(\mathbf{x}_0) = \frac{e^{\hat{\eta}_0}}{1 + e^{\hat{\eta}_0}} ,$$

with linear predictor

$$\hat{\eta}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_{01} + \dots \hat{\beta}_p x_{0p} .$$

In the case of qualitative covariates, a dummy variable needs to be introduced. This is done as for linear regression.

So in the `Default` example, we can predict the probability that someone defaults.

For example: “What is the estimated probability that a student defaults with a balance of 2000, and an income of 40000?”

$$\hat{p}(X) = \frac{e^{\beta_0 + 2000 \cdot \beta_1 + 40000 \cdot \beta_2 + 1 \cdot \beta_3}}{1 + e^{\beta_0 + 2000 \cdot \beta_1 + 40000 \cdot \beta_2 + 1 \cdot \beta_3}} = 0.5196$$

Using R:

```
eta <- summary(glm_default2)$coef[, 1] %*% c(1, 2000, 40000, 1)
exp(eta)/(1 + exp(eta))
```

```
##           [,1]
## [1,] 0.5196218
```

Example: South African heart disease data set

- **SAhert** data set from the **ElemStatLearn** package, a retrospective sample of males in a heart-disease high-risk region in South Africa.
- 462 observations on 10 variables.
- All subjects are male in the age range 15-64.
- 160 cases (individuals who have suffered from a conorary heart disease) and 302 controls (individuals who have not suffered from a conorary heart disease).

The response value (**chd**) and covariates

- **chd** : conorary heart disease {yes, no} coded by the numbers {1, 0}
- **sbp** : systolic blood pressure
- **tobacco** : cumulative tobacco (kg)
- **ldl** : low density lipoprotein cholesterol
- **famhist** : family history of heart disease. Categorical variable with two levels: {Absent, Present}.
- **obesity** : a numerical value
- **alcohol** : current alcohol consumption
- **age** : age at onset

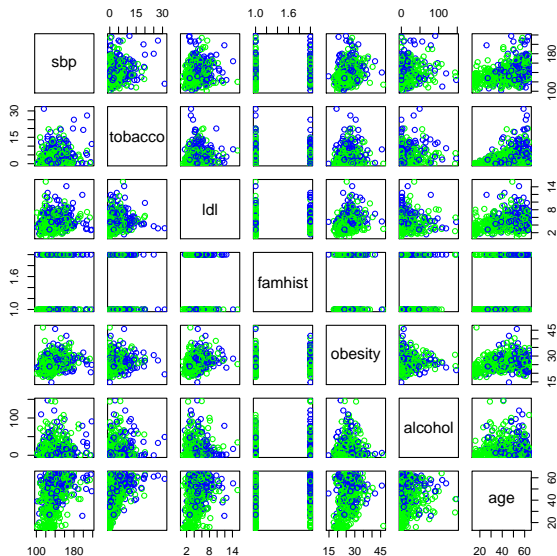
The goal is to identify important risk factors. We start by loading and looking at the data:

```
library(ElemStatLearn)

d.heart <- SAheart
d.heart$chd <- as.factor(d.heart$chd)
d.heart <- d.heart[, c("sbp", "tobacco", "ldl", "famhist", "obesity",
  "alcohol", "age", "chd")]
head(d.heart)
```

##	sbp	tobacco	ldl	famhist	obesity	alcohol	age	chd
## 1	160	12.00	5.73	Present	25.30	97.20	52	1
## 2	144	0.01	4.41	Absent	28.87	2.06	63	1
## 3	118	0.08	3.48	Present	29.14	3.81	46	0
## 4	170	7.50	6.41	Present	31.99	24.26	58	1
## 5	134	13.60	3.50	Present	25.99	57.34	49	1
## 6	132	6.20	6.47	Present	30.77	14.14	45	0

`pairs()` plot with $Y = 1$ (case) in green and $Y = 0$ (control) in blue:



Fitting the model using all predictors in R:

```
glm_heart = glm(chd ~ ., data = d.heart, family = "binomial")
summary(glm_heart)
```

```
##
## Call:
## glm(formula = chd ~ ., family = "binomial", data = d.heart)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.7517  -0.8378  -0.4552   0.9292   2.4434
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -4.1295997  0.9641558  -4.283 1.84e-05 ***
## sbp           0.0057607  0.0056326   1.023  0.30643
## tobacco      0.0795256  0.0262150   3.034  0.00242 **
## ldl           0.1847793  0.0574115   3.219  0.00129 **
## famhistPresent 0.9391855  0.2248691   4.177 2.96e-05 ***
## obesity      -0.0345434  0.0291053  -1.187  0.23529
## alcohol       0.0006065  0.0044550   0.136  0.89171
## age          0.0425412  0.0101749   4.181 2.90e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 596.11  on 461  degrees of freedom
## Residual deviance: 483.17  on 454  degrees of freedom
## AIC: 499.17
##
## Number of Fisher Scoring iterations: 4
```

- What are important predictors for CHD?
- How would you now calculate $\hat{p}(X)$ that someone will develop CHD, given covariates X ?

Multinomial logistic regression

- Generalization of the logistic regression model to a response variable with more than two (K) possible classes.
- The probability that Y belongs to class k , given an observation vector $\mathbf{x} = (x_1, x_2, \dots, x_p)^T$ is (usually) modelled by:

$$\log \frac{\Pr(Y = k|\mathbf{x})}{\Pr(Y = K|\mathbf{x})} = \beta_{0k} + \beta_{1k}x_1 + \dots + \beta_{pk}x_p,$$

where K is chosen as the *reference category* (arbitrary). We would thus (jointly) estimate $K - 1$ sets of regression parameters.

- The multinomial logistic regression model is implemented in the `glmnet` or `VGAM` package in R.
- We will not discuss this further since LDA is more popular (than logistic regression) in the multi-class setting. And, as we shall see soon - they are not that different.

Bayes classifier

- We were going for classification, but now estimated $\hat{p}(X) = \Pr(Y | X)$ using logistic regression. Idea: probability can be used for classification.
- For K classes $\mathcal{C} = \{c_1, c_2, \dots, c_K\}$, with elements numbered as $1, 2, \dots, K$, the probability that a new observation x_0 belongs to class k is²

$$p_k(x_0) = \Pr(Y = k | X = x_0), \quad k = 1, 2, \dots, K.$$

- The *Bayes classifier* assigns an observation to the most likely class, given its predictor values.
- **Example** for two groups $\{A, B\}$. A new observation x_0 will be classified to A if $\Pr(Y = A | X = x_0) > 0.5$ and to class B otherwise.

²This is the conditional probability: the probability that $Y = k$ given the observation x_0 .

Properties of the Bayes classifier

- It has the *smallest test error rate*.
- The class boundaries using the Bayes classifier is called the *Bayes decision boundary*.
- The overall Bayes error rate is given as

$$1 - E(\max_j \Pr(Y = j \mid X))$$

where the expectation is over X .

- The Bayes error rate is comparable to the *irreducible error* in the regression setting.

What is test error?

To understand what test error means, we need to understand the *training error* first.

Training error

- **Training error rate:** The proportion of mistakes that are made if we apply our estimator \hat{f} to the training observations, i.e.

$$\hat{y}_i = \hat{f}(x_i)$$

$$\frac{1}{n} \sum_{i=1}^n \mathbf{I}(y_i \neq \hat{y}_i) ,$$

with indicator function \mathbf{I} , which is defined as:

$$\mathbf{I}(a \neq \hat{a}) = \begin{cases} 1 & \text{if } a \neq \hat{a} , \\ 0 & \text{else.} \end{cases}$$

- The indicator function counts the number of times our model has made a wrong classification. The training error rate is the fraction of misclassifications made on our training set.
- A very low training error rate may imply overfitting.

Test error

- **Test error rate:** The fraction of misclassifications when our model is applied on a test set

$$\text{Ave}(I(y_0 \neq \hat{y}_0)) ,$$

where the average is over all the test observations (x_0, y_0) .

- Again, this gives a better indication of the true performance of the classifier (than the training error).
- We assume that a *good* classifier is a classifier that has a *low* test error.

Training error rate

- the proportion of mistakes that are made if we apply classifier \hat{f} to the training observations, i.e. $\hat{y}_i = \hat{f}(x_i)$.

$$\frac{1}{n} \sum_{i=1}^n \mathbf{I}(y_i \neq \hat{y}_i),$$

where \mathbf{I} is the indicator function (to give our 0/1 loss) which is defined as:

$$\mathbf{I}(a \neq \hat{a}) = \begin{cases} 1 & \text{if } a \neq \hat{a} \\ 0 & \text{else} \end{cases}$$

- The indicator function counts the number of times our model has made a wrong classification. The training error rate is the fraction of misclassifications made on our training set.
- A very low training error rate may imply overfitting.

Test error rate

- Here the fraction of misclassifications is calculated when our model is applied on a test set. From what we have learned about regression we can deduct that this gives a better indication of the true performance of the classifier (than the training error)

$$\text{Ave}(I(y_0 \neq \hat{y}_0)) ,$$

where the average is over all the test observations (x_0, y_0) .

- We assume that a *good* classifier is a classifier that has a *low* test error.

The K -nearest neighbour (KNN) classifier

- Caveat: We usually don't know the conditional distribution of Y given X for real data.
- We have discussed how to do it with logistic regression (for two categories).
- The K -nearest neighbor (KNN) classifier estimates this conditional distribution and chooses the most likely category (Bayes classifier).³

³Attention!! K refers to the number of neighbours used for the classifier, and *not* to the number of classes!! The latter is assumed known.

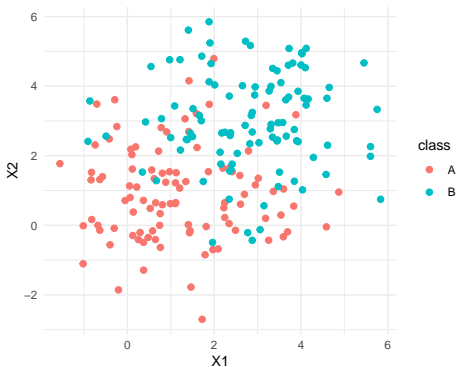
The K -nearest neighbour classifier (KNN) works in the following way:

- Given a new observation x_0 it searches for the K points in our training data that are closest to it (Euclidean distance).
- These points make up the neighborhood of x_0 , \mathcal{N}_0 .
- The point x_0 is classified by taking a majority vote of the neighbors.
- That means that x_0 is classified to the most occurring class among its neighbors (“majority vote”)

$$\Pr(Y = j|X = x_0) = \frac{1}{K} \sum_{i \in \mathcal{N}_0} I(y_i = j) .$$

A synthetic example

- Simulate 2×100 observations from a bivariate normal distribution with mean vectors $\mu_A = (1, 1)^T$, $\mu_B = (3, 3)^T$, and covariance matrix $\Sigma_A = \Sigma_B = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$.
- Aim: Find a rule to classify a new observation to A or B .

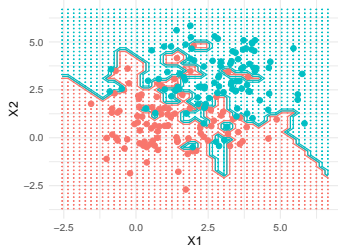


- Assume we have a new observation $X_0 = (x_{01}, x_{02})^T$ which we want to classify as belonging to the class A or B .
- To illustrate this problem we fit the K -nearest neighbor classifier to our simulated data set with $K = 1, 3, 10$ and 150 (next slide).

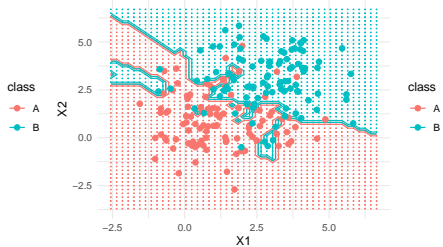
Interpretation:

- The small colored dots show the predicted classes for an evenly-spaced grid.
- The lines show the decision boundaries.
- If our new observation falls into the region within the red decision boundary, it will be classified as A . If it falls into the region within the green decision boundary, it will be classified as B .

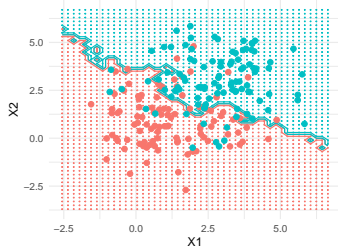
$k = 1$



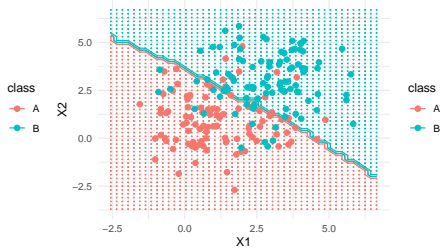
$k = 3$



$k = 10$

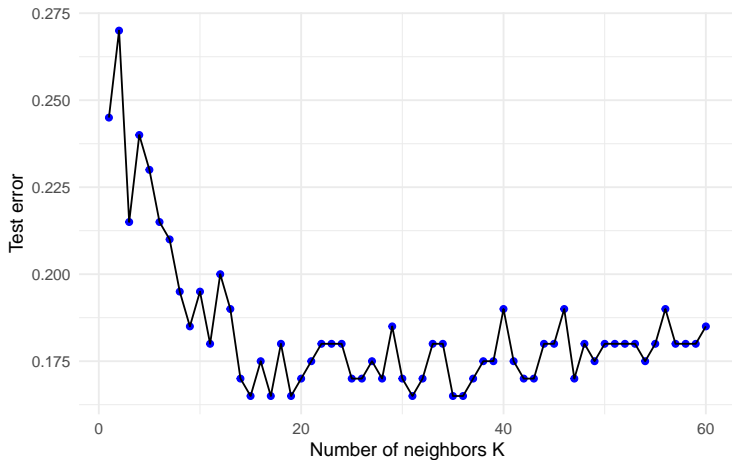


$k = 150$



- The choice of K has a big influence on the result of our classification. For $K = 1$ the classification is made to the same class as the one nearest neighbor. As K gets very large, the decision boundary tends towards a straight line (which is the Bayes boundary in this set-up).
- To find the optimal value of K the typical procedure is to try different values of K and then test the predictive power of the different classifiers, for example by cross-validation, which will be discussed in M5.
- After trying all choices for K between 1 and 50, a few choices of K gave the smallest misclassification error rate, estimating by leave-one out cross-validation (see M5). The smallest error rate is equal to 0.165 (16.5% misclassification).

Cross-validated test error rate for KNN



Bias-variance trade-off in a classification setting

Remember the bias-variance trade-off that we discussed in Module 2.

Considerations:

- A too low value of K will give a very flexible classifier (with high variance and low bias) which will fit the training set too well (it will overfit) and make poor predictions for new observations.
- Choosing a high value for K makes the classifier lose its flexibility and the classifier will have low variance but high bias.

→ Critical to the success of the classifier to choose the correct level of flexibility (K).

The curse of dimensionality

- The nearest neighbor classifier can be quite good if the number of predictor p is small and the number of observations n is large. We need enough close neighbors to make a good classification.
- The effectiveness of the KNN classifier falls quickly when the dimension of the predictor space is high.
- Why? Because the nearest neighbors tend to be far away in high dimensions and the method no longer is local. This is referred to as the *curse of dimensionality*.

Bayes decision rule - two paradigms

Two approaches to estimate $\Pr(Y = k \mid X = x)$:

Diagnostic paradigm

This is what we did so far: The focus was *directly* estimating the posterior distribution for the classes

$$\Pr(Y = k \mid X = x) .$$

Examples: Logistic regression. KNN classification.

Sampling paradigm

- The focus is on estimating the prior probabilities $\pi_k = \Pr(Y = k)$ for the classes and the class conditional distributions $f_k(x) = \Pr(X = x \mid Y = k)$.
- We classify to the class with the maximal product $\pi_k f_k(x)$. In this paradigm, we need to model the pdf for each class.
- Popular: the multivariate normal distribution!

Given a continuous X and categorical Y , and

- the probability *density* function $f_k(x) = \Pr(X = x \mid Y = k)$ for X in class k .
- the *prior* probability for class k $\pi_k = \Pr(Y = k)$ is the prior probability.

How do we get $\Pr(Y = k \mid X = x_0)$?

Bayes theorem

$$\begin{aligned} p_k(X) = \Pr(Y = k \mid X = x) &= \frac{\Pr(\mathbf{X} = \mathbf{x} \cap Y = k)}{f(\mathbf{x})} \\ &= \frac{f_k(x)\pi_k}{\sum_{l=1}^K f_l(x)\pi_l} \end{aligned}$$

Discriminant Analysis

Discriminant analysis is relying on the *sampling paradigm*.

- The approach is to model the distribution of X in each of the classes separately, and then use Bayes theorem to flip things around and obtain $\Pr(Y \mid X)$.
- When we use normal (Gaussian) distributions for each class, this leads to linear or quadratic discriminant analysis.
- However, this approach is quite general, and other distributions can be used as well. We will focus on normal distributions.

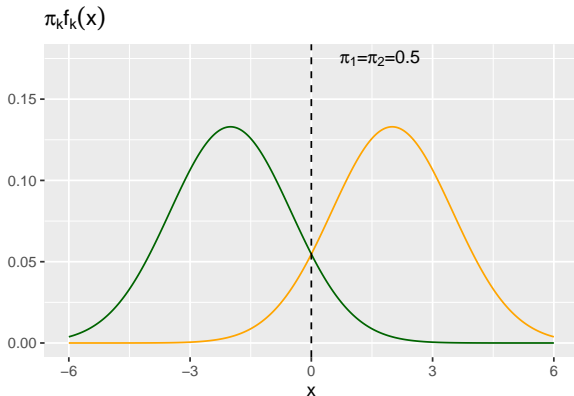
Classification to the highest density

Suppose we have observations coming from two classes:
 $\{\text{green}, \text{orange}\}$

$$X_{\text{green}} \sim \mathcal{N}(-2, 1.5^2) \text{ and } X_{\text{orange}} \sim \mathcal{N}(2, 1.5^2)$$

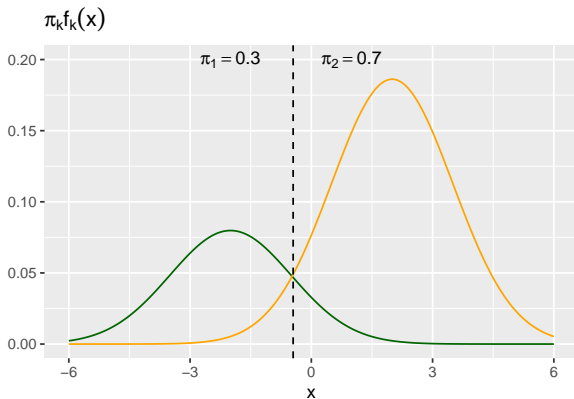
- Assume probabilities to be equal, $\pi_1 = \pi_2 = 0.5$.

We plot $\pi_k f_k(x)$ for the two classes:



The decision boundary is where the point of intersection of the two lines is, because here $\pi_1 f_1(x) = \pi_2 f_2(x)$. Thus all points to left of the decision boundary will be classified as green and similarly, all points to the right of the decision boundary will be classified as orange.

We now specify different priors, such that $\pi_1 = 0.3$ and $\pi_2 = 0.7$. We see that this has affected the decision boundary, which now has shifted to the left:



Why discriminant analysis?

- When the classes are well-separated, the parameter estimates for the logistic regression model are surprisingly unstable. Linear discriminant analysis does not suffer from this problem.
- If n is small and the distribution of the predictors X is approximately normal in each of the classes, the linear discriminant model is again more stable than the logistic regression model.
- Linear discriminant analysis is popular when we have more than two response classes, because it also provides low-dimensional views of the data.

Linear discriminant analysis (LDA) when $p = 1$

- Assume there is only one predictor X ($p = 1$).
- Class conditional distributions $f_k(X)$ are assumed normal (Gaussian) for $k = 1, \dots, K$, that is

$$f_k(x) = \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{1}{2}\left(\frac{x-\mu_k}{\sigma_k}\right)^2}$$

has parameters μ_k (mean) and σ_k (standard deviation).

- With LDA we assume that all of the classes have the *same standard deviation* $\sigma_k = \sigma$.
- In addition we have prior class probabilities $\pi_k = \Pr(Y = k)$, so that $\sum_{k=1}^K \pi_k = 1$.

We can insert the expression for each class distribution into Bayes formula to obtain the posterior probability $p_k(x) = \Pr(Y = k|X = x)$

$$p_k(x) = \frac{f_k(\mathbf{x})\pi_k}{f(\mathbf{x})} = \frac{\pi_k \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu_k}{\sigma}\right)^2}}{\sum_{l=1}^K \pi_l \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu_l}{\sigma}\right)^2}} .$$

Our rule is to classify to the class for which $p_k(x)$ is largest.

Taking logs, and discarding terms that do not depend on k , we see that this is equivalent to assigning x to the class with the largest *discriminant score* $\delta_k(x)$:

$$\delta_k(x) = x \cdot \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k).$$

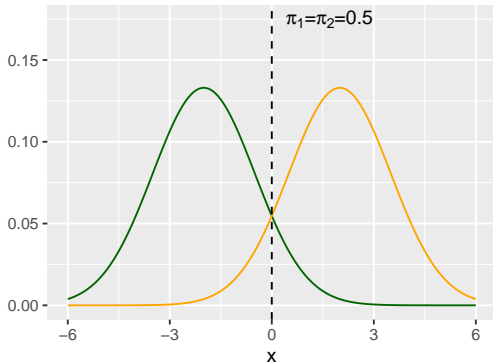
- This decision boundaries between the classes are *linear* in x .
- For $K = 2$ classes and $\pi_1 = \pi_2$, the decision boundary is at

$$x = \frac{\mu_1 + \mu_2}{2} .$$

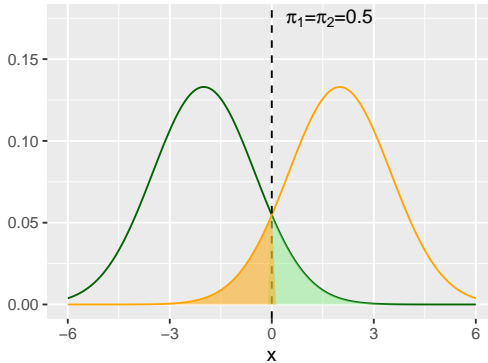
(Show this by setting $\delta_1(x) = \delta_2(x)$ and resolving for x .)

Back to our example

$$X_{\text{green}} \sim \mathcal{N}(-2, 1.5^2) \text{ and } X_{\text{orange}} \sim \mathcal{N}(2, 1.5^2)$$



- Where is the decision boundary?
- How can we calculate the error rate?
- Is it possible to construct a classifier and it has a lower error rate than the Bayes error rate?



- The Bayes decision boundary is at $x = 0$, and we would make a (Bayes) error of `round(2*0.5*pnorm(0,2,1.5))=0.09`.
- The Bayes classifier has the lowest test error rate.

In the above example we knew the true distributions $p_k(X)$ and the priors π_k . But typically we don't know these parameters, we only have the training data.

Idea: we simply estimate the parameters and plug them into the rule.

Parameter estimators

- Prior probability for class k is (often) estimated by taking the fraction of observations n_k (out of n) coming from class k :

$$\hat{\pi}_k = \frac{n_k}{n}.$$

- The mean value for class k is simply the sample mean of all observations from class k :

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i.$$

- The standard deviation: sample standard deviation across all classes:

$$\hat{\sigma}^2 = \frac{1}{n - K} \sum_{k=1}^K \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2 = \sum_{k=1}^K \frac{n_k - 1}{n - K} \cdot \hat{\sigma}_k^2.$$

$\hat{\sigma}_k$: estimated standard deviation of all observations from class k .

How can we then estimate the goodness of our estimator?

1. Use the training set to estimate parameters and class boundary:
2. Use the test set to estimate misclassification rate.

Simulate data and use training data to estimate μ_k and $\hat{\sigma}^2$:

```
n = 1000
pi1 = pi2 = 0.5
mu1 = -2
mu2 = 2
sigma = 1.5
set.seed(1)
n1train = rbinom(1, n, pi1)
n2train = n - n1train
n1test = rbinom(1, n, pi1)
n2test = n - n1test
train1 = rnorm(n1train, mu1, sigma)
train2 = rnorm(n2train, mu2, sigma)
test1 = rnorm(n1test, mu1, sigma)
test2 = rnorm(n2test, mu2, sigma)
var2.1 = var(train1)
var2.2 = var(train2)
var.pool = ((n1train - 1) * var2.1 + (n2train - 1) * var2.2)/(n - 2)
```

Then set

$$\hat{\delta}_1(x) = \hat{\delta}_2(x)$$

and resolve for x to obtain a decision rule (boundary).

Todo: Verify that the following code will give you the training and test error rates:

```
rule = 0.5 * (mean(train1) + mean(train2)) + var.pool * (log(n2train/n) -  
               log(n1train/n))/(mean(train1) - mean(train2))  
  
c((sum(train1 > rule) + sum(train2 < rule))/n, (sum(test1 > rule) + sum(test2 <  
               rule))/n)
```

```
## [1] 0.105 0.115
```

This is a rather good performance, compared to the minimal Bayes error rate. But keep in mind that we LDA classifier relies on the Normal assumption, and that $\sigma_k = \sigma$ for all classes is assumed.

The confusion matrix

- The confusion matrix is a table that can show the performance of classifier, given that the true values are known.
- We can make a confusion matrix from the training or test set, but will in most cases do so for the test set only.
- The rows represent the true classes, while the columns represent the predicted classes. Inside the table we have counts (just labelled “correct” and “wrong” below - but should be numbers).
- The sum of the diagonal is the total number of correct classifications. The sum of all elements off the diagonal is the total number of misclassifications.

	Predicted 1	Predicted 2	...	Predicted K
True 1	correct	wrong	...	wrong
True 2	wrong	correct	...	wrong
...
True K	wrong	wrong	...	correct

The confusion matrix can be obtained in R by using the `table` function with the predicted classes and the true classes need to be given as function arguments, or directly using the `caret` package.

We will soon come back to the special case of two classes - where we may think of this as “-” (non disease) and “+” (disease).

Multivariate LDA ($p > 1$)

Linear discriminant analysis can be generalized to situations when $p > 1$ covariates are used. The decision boundaries are still linear. The multivariate normal distribution function:

$$f(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

This gives the following expression for the discriminant function:

$$\delta_k(x) = \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \log \pi_k.$$

Some details on the derivation of the discriminant function - when $K = 2$ classes, and the classes are denoted 0 and 1

$$\Pr(Y = 0|X = \mathbf{x}) = \Pr(Y = 1|X = \mathbf{x})$$

$$\frac{\pi_0 f_0(x)}{\pi_0 f_0(x) + \pi_1 f_1(x)} = \frac{\pi_1 f_0(1)}{\pi_0 f_0(x) + \pi_1 f_1(x)}$$

$$\pi_0 f_0(\mathbf{x}) = \pi_1 f_1(\mathbf{x})$$

$$\pi_0 \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_0)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_0)} = \pi_1 \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_1)}$$

$$\log(\pi_0) - \frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_0)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_0) = \log(\pi_1) - \frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_1)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_1)$$

$$\log(\pi_0) - \frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x} + \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_0 - \frac{1}{2}\boldsymbol{\mu}_0^T \Sigma^{-1} \boldsymbol{\mu}_0 = \log(\pi_1) - \frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x} + \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_1 - \frac{1}{2}\boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1$$

$$\log(\pi_0) + \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_0 - \frac{1}{2}\boldsymbol{\mu}_0^T \Sigma^{-1} \boldsymbol{\mu}_0 = \log(\pi_1) + \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_1 - \frac{1}{2}\boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1$$

$$\delta_0(\mathbf{x}) = \delta_1(\mathbf{x})$$

Estimators for $p > 1$:

- Prior probability for class k (unchanged from $p=1$): $\hat{\pi}_k = \frac{n_k}{n}$.
- The mean value for class k is simply the sample mean of all observations from class k (but now these are vectors):

$$\hat{\boldsymbol{\mu}}_k = \frac{1}{n_k} \sum_{i:y_i=k} \mathbf{X}_i.$$

- The covariance matrices for each class:

$$\hat{\boldsymbol{\Sigma}}_k = \frac{1}{n_k - 1} \sum_{i:y_i=k} (\mathbf{X}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{X}_i - \hat{\boldsymbol{\mu}}_k)^T$$

- Pooled version:

$$\hat{\boldsymbol{\Sigma}} = \sum_{k=1}^K \frac{n_k - 1}{n - K} \cdot \hat{\boldsymbol{\Sigma}}_k.$$

Optional: **Proof that the estimator $\hat{\boldsymbol{\Sigma}}_k$ is unbiased for each class (from TMA4267)**. Proof for pooled version not provided.

Posterior probabilities

Sometimes the probability that an observation comes from a class k is more interesting than the actual classification itself. These class probabilities can be estimated from the priors and class conditional distributions, or from the discriminant functions:

$$\begin{aligned}\hat{P}(Y = k|X = x) &= \frac{\hat{\pi}_k \cdot \frac{1}{(2\pi)^{p/2}|\hat{\Sigma}|^{1/2}} \exp(-\frac{1}{2}(\mathbf{x} - \hat{\mu}_k)^T \hat{\Sigma}^{-1}(\mathbf{x} - \hat{\mu}_k))}{\sum_{l=1}^K \hat{\pi}_l \frac{1}{(2\pi)^{p/2}|\hat{\Sigma}|^{1/2}} \exp(-\frac{1}{2}(\mathbf{x} - \hat{\mu}_l)^T \hat{\Sigma}^{-1}(\mathbf{x} - \hat{\mu}_l))} \\ &= \frac{e^{\hat{\delta}_k(x)}}{\sum_{l=1}^K e^{\hat{\delta}_l(x)}}.\end{aligned}$$

Quadratic Discriminant Analysis (QDA)

In quadratic discriminant analysis we assume that the distributions of the classes is multivariate normal (Gaussian), but with covariance matrix Σ_k for each class.

The discriminant functions are now given by:

$$\begin{aligned}\delta_k(x) &= -\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) - \frac{1}{2} \log |\Sigma_k| + \log \pi_k \\ &= -\frac{1}{2} x^T \Sigma_k^{-1} x + x^T \Sigma_k^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma_k^{-1} \mu_k - \frac{1}{2} \log |\Sigma_k| + \log \pi_k.\end{aligned}$$

These decision boundaries are *quadratic* functions of x .

LDA vs QDA

In QDA we are allowed for the covariance matrices to be different for the classes, but for LDA they are assumed to be the same, so QDA is more flexible than LDA.

Q:

- But, if the covariance matrices in theory are equal - will they not be estimated equal?
- Should we not always prefer QDA to LDA?

A:

Explanation similar to a “Bias-variance trade-off”:

If the assumption of equal covariance matrices is wrong

- then LDA may suffer from high bias for the parameter estimators.
- and QDA is better off. But, for small sample sizes the covariance matrices might be poorly estimated (high variance of estimators).

If the number of covariates is high:

- then QDA requires estimating $K \cdot p \cdot (p + 1)/2$ parameters,
- while LDA only requires $p \cdot (p + 1)/2$.

Therefore, LDA is less flexible than QDA and might therefore have much less variance.

Naive Bayes (optional)

This is method that is popular when p is large.

In Naive Bayes (Idiot's Bayes) we assume that each class density is the product of marginal densities - i.e. inputs are conditionally independent in each class.

$$f_k(x) = \prod_{j=1}^p f_{kj}(x_j)$$

This is generally not true, but it simplifies the estimation dramatically.

The original naive Bayes used univariate normal marginal distributions, but generalizations can be made.

Instead of estimating a full covariance matrix, only the diagonal elements are estimated.

This method often produces good results, even though the joint pdf is not the product of the marginal pdf. This might be because we are not focussing on estimation of class pdfs, but class boundaries.

Example: Which type of iris species?

The **iris** flower data set was introduced by the British statistician and biologist Ronald Fisher in 1936.

- Three plant species {setosa, virginica, versicolor} (50 observation of each), and
- four features: Sepal.Length, Sepal.Width, Petal.Length and Petal.Width.

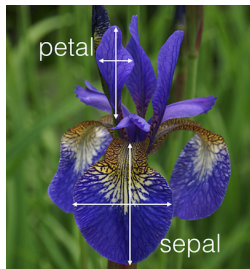


Figure 1: Iris plant with sepal and petal leaves

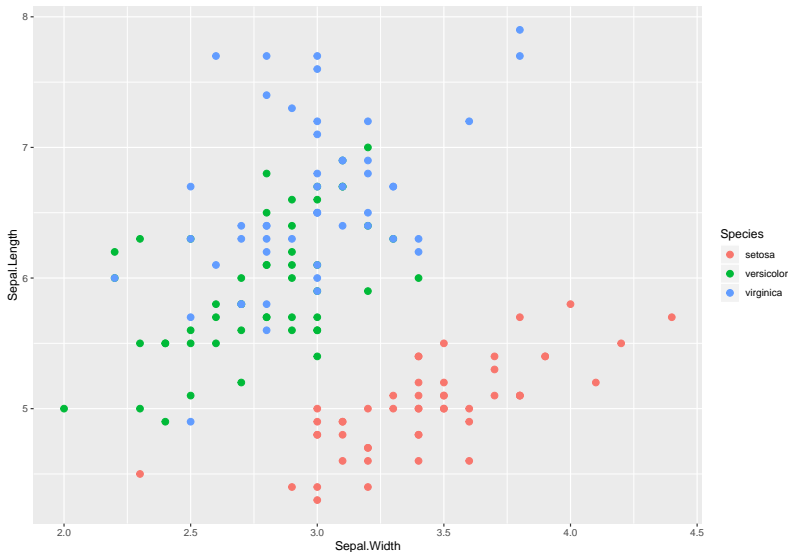
Example: Classification of iris plants

We will use **sepal width** and **sepal length** to build a classifier, here 50 observations from each class available.

```
attach(iris)
library(class)
library(MASS)
library(ggplot2)
library(dplyr)
kable(head(iris), format = whichformat)
```

Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
4.7	3.2	1.3	0.2	setosa
4.6	3.1	1.5	0.2	setosa
5.0	3.6	1.4	0.2	setosa
5.4	3.9	1.7	0.4	setosa

```
iris0_plot = ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length,  
  geom_point(size = 2.5)  
iris0_plot
```



Iris: LDA

reporting on parameter estimates

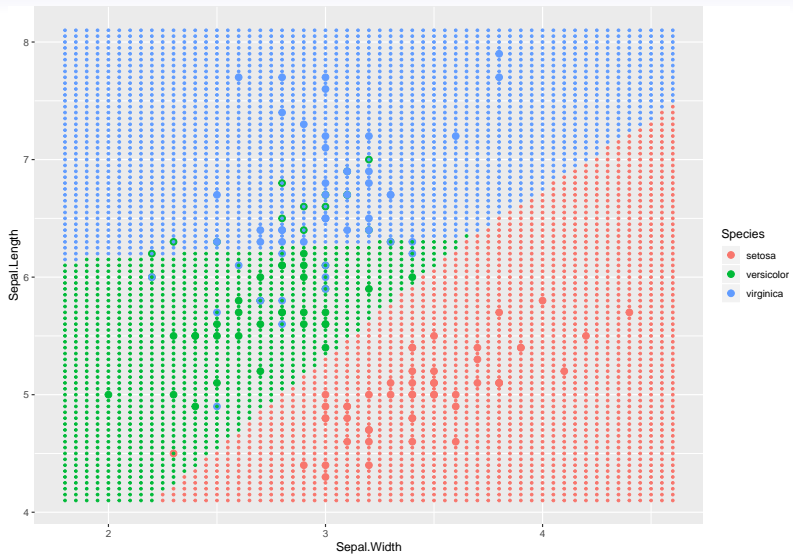
```
grS = cbind(iris[iris$Species == "setosa", 1:2])
grVe = cbind(iris[iris$Species == "versicolor", 1:2])
grVi = cbind(iris[iris$Species == "virginica", 1:2])
nk = dim(grS)[1]
apply(grS, 2, mean)
```

```
## Sepal.Length  Sepal.Width
##           5.006           3.428
```

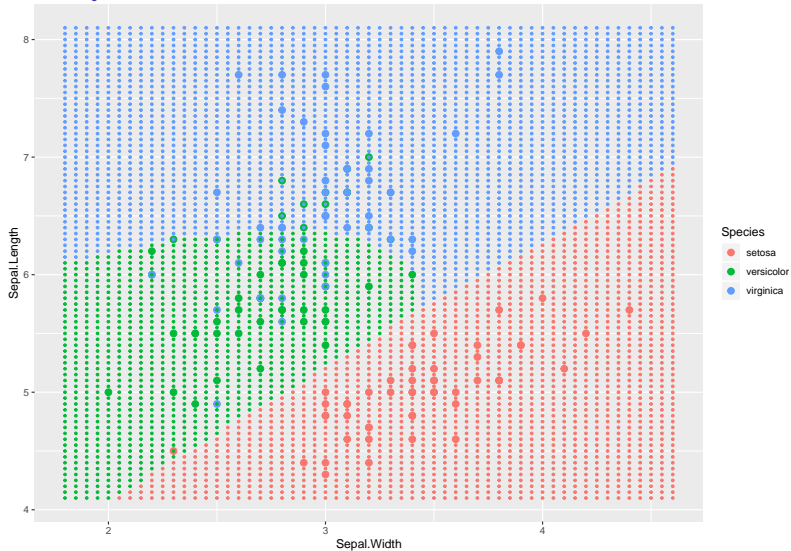
```
cS = cov(grS)
apply(grVe, 2, mean)
```

```
## Sepal.Length  Sepal.Width
##           5.936           2.770
```

```
cVe = cov(grVe)
apply(grVi, 2, mean)
```



Iris: QDA



Iris: compare LDA and QDA

We now want to compare the predictive performance of our two classifiers. We do this by dividing the original `iris` data set, randomly, into train and test samples of equal size:

```
set.seed(1)
train = sample(1:150, 75)

iris_train = iris[train, ]
iris_test = iris[-train, ]

iris_lda2 = lda(Species ~ Sepal.Length + Sepal.Width, data = iris,
  prior = c(1, 1, 1)/3)
```

Training error

```
table(predict(iris_lda2, newdata = iris_train)$class, iris_train
```

```
##
```

```
##           setosa versicolor virginica
```

```
##  setosa         27           0         0
```

```
##  versicolor     1          15         8
```

```
##  virginica      0           5        19
```

Test error

```
iris_lda2_predict = predict(iris_lda2, newdata = iris_test)
```

```
table(iris_lda2_predict$class, iris$Species[-train])
```

```
##
```

```
##           setosa versicolor virginica
```

```
##  setosa         22           0         0
```

```
##  versicolor     0          22        11
```

```
##  virginica      0           8        12
```

The LDA classifier has a training error rate of $15/75$, that is 20 %. When tested on a new set it correctly classified $24+13+23$ times and misclassified 15 times. This gives a misclassification rate of

$$\text{Test error}_{\text{LDA}} = \frac{15}{75} = 0.2.$$

Using a different division into training and test set will give (small) changes to these numbers.

Iris: training and test error for QDA

```
iris_qda2 = qda(Species ~ Sepal.Length + Sepal.Width, data = iris,
  prior = c(1, 1, 1)/3)
```

Important: use the same division into training and test set for methods we want to compare.

Training error

```
table(predict(iris_qda2, newdata = iris_train)$class, iris_train
```

```
##
```

```
##           setosa versicolor virginica
```

```
##  setosa         28           0         0
```

```
##  versicolor     0          16         9
```

```
##  virginica      0           4        18
```

Test error

```
iris_qda2_predict = predict(iris_qda2, newdata = iris_test)
```

```
table(iris_qda2_predict$class, iris$Species[-train])
```

```
##
```

```
##           setosa versicolor virginica
```

```
##  setosa         22           0         0
```

```
##  versicolor     0          18        12
```

```
##  virginica      0          12        11
```


The QDA classifier has a training error rate of 16%. When tested on a new set, the misclassification error rate was

$$\text{Test error}_{\text{QDA}} = \frac{18}{75} = .24$$

The LDA classifier has given the smallest test error for classifying iris plants based on sepal width and sepal length for our test set and should be preferred in this case.

1. Would another division of the data into training and test set give the same conclusion (that LDA is better than QDA for this data set)? (A: Not necessarily, but probably.)

We will look into other choice than dividing into one training and one test set in Module 5 (crossvalidation).

2. What about the other two covariates? Would adding them to the model (4 covariates) give a better classification rule? (A: Probably. Try if you want.)

Fishers idea (Optional)

In 1936 R. A. Fisher developed LDA.

- His aim was to find a linear combination of the explanatory variables which *maximized the ratio of its between class to within class variance*.
- In this way the observations are transformed so that they are separated as much as possible.
- His approach has the advantage that it is suited for visual inspection and graphical description , it “separates” the populations.

Let the between-class variance be denoted

$\mathbf{B} = \sum_{m=1}^M (\boldsymbol{\mu}_m - \bar{\boldsymbol{\mu}})(\boldsymbol{\mu}_m - \bar{\boldsymbol{\mu}})^T$, where $\boldsymbol{\mu}_m$ denotes the mean of class ω_m and $\bar{\boldsymbol{\mu}}$ the overall mean.

The within-class variance is assumed equal for all classes and is denoted $\boldsymbol{\Sigma}$ ($\boldsymbol{\Sigma}$ is assumed to have full rank).

The linear combination $\mathbf{l}^T \mathbf{X}$ that maximize $\mathbf{l}^T \mathbf{B} \mathbf{l} / \mathbf{l}^T \mathbf{\Sigma} \mathbf{l}$ under the constraint that $\mathbf{l}^T \mathbf{\Sigma} \mathbf{l} = 1$ is found to be the scaled eigenvectors of $\mathbf{\Sigma}^{-1} \mathbf{B}$ corresponding to the nonzero eigenvalues of $\mathbf{\Sigma}^{-1} \mathbf{B}$. The eigenvector corresponding to the largest eigenvalue defines the first discriminant $\mathbf{l}_1^T \mathbf{X}$. The second linear discriminant $\mathbf{l}_2^T \mathbf{X}$ is constructed from the eigenvector corresponding to the second largest eigenvalue and so on. (We also have $Cov(\mathbf{l}_j^T \mathbf{X}, \mathbf{l}_i^T \mathbf{X}) = 0$ for $i \neq j$ and $Var(\mathbf{l}_j^T \mathbf{X}) = 1$.)

The number of linear discriminants equals the number of nonzero eigenvalues. Observations are assigned to the class of the nearest (Euclidean distance) class mean in the discriminant space.

This equals classification to the nearest Mahalanobis distance population mean in the input space. Again, the parameters μ_i and Σ and the between-class variance \mathbf{B} are usually unavailable. Replacing the parameters by estimates from the training set leads to Fisher's sample linear discriminants.

Diagnostic paradigm

Remember:

- The **diagnostic paradigm**: We focus on *directly* estimating the posterior distribution for the classes $\Pr(Y = k \mid X = x)$.
- The **sampling paradigm**: There focus is on estimating the prior probabilities for the classes and the class conditional distributions. We classify to the class with the maximal product $\pi_k f_k(x)$.

Now we move to the *diagnostic paradigm* and the K -nearest neighbor classifier.

The K -nearest neighbour classifier estimates $\Pr(Y = k \mid X = x)$ and classifies a new observation based on this estimated probability

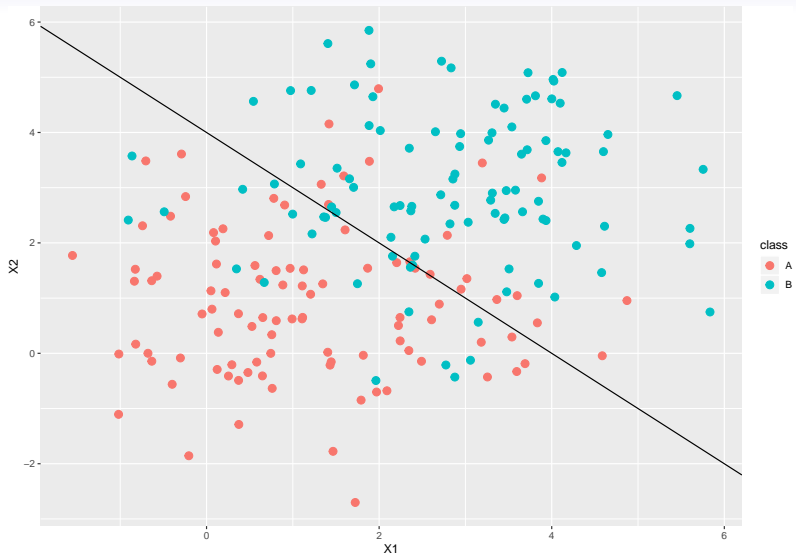
Synthetic example for KNN classification

Consider a two-class example, and equal class prior probabilities.

A new observation x_0 will be classified to A if

$\Pr(Y = A|X = x_0) > 0.5$ and to class B otherwise.

- The figure below shows a plot of 100 observations from two classes A (red dots) and B (turquoise dots),
- simulated from a bivariate normal distribution with mean vectors $\mu_A = (1, 1)^T$ and $\mu_B = (3, 3)^T$ and a covariance matrix $\Sigma_A = \Sigma_B = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$.
- We want to find a rule to classify a new observation to class A or B .



Remark: since the truth is known here we can calculate the Bayes boundary and the Bayes error.

Since we have bivariate normal class distributions with common covariance matrix, the optimal boundary is given by LDA. The boundary will be at $\delta_A(\mathbf{x}) = \delta_B(\mathbf{x})$, where $\delta_A(\mathbf{x}) = \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_A - \frac{1}{2} \boldsymbol{\mu}_A^T \Sigma^{-1} \boldsymbol{\mu}_A + \log \pi_A$, and for $\delta_B(\mathbf{x})$ with $\boldsymbol{\mu}_B$.

$$\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_A - \frac{1}{2} \boldsymbol{\mu}_A^T \Sigma^{-1} \boldsymbol{\mu}_A + \log \pi_A = \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_B - \frac{1}{2} \boldsymbol{\mu}_B^T \Sigma^{-1} \boldsymbol{\mu}_B + \log \pi_B$$

$$\mathbf{x}^T \Sigma^{-1} (\boldsymbol{\mu}_A - \boldsymbol{\mu}_B) - \frac{1}{2} \boldsymbol{\mu}_A^T \Sigma^{-1} \boldsymbol{\mu}_A + \frac{1}{2} \boldsymbol{\mu}_B^T \Sigma^{-1} \boldsymbol{\mu}_B + \log \pi_A - \log \pi_B = 0$$

Inserting numerical values gives: $-x_1 - x_2 + 4 = 0$, and boundary $x_2 = 4 - x_1$.

```
muA = matrix(c(1, 1), ncol = 1)
muB = matrix(c(3, 3), ncol = 1)
sigmainv = diag(2)/2
sigmainv %*% (muA - muB)
```


K-nearest neighbour classifier

(warning: K is not the number of classes, but neighbours...)

The K -nearest neighbour classifier (KNN) works in the following way:

- Given a new observation x_0 it searches for the K points in our training data that are closest to it.
- These points make up the neighborhood of x_0 , \mathcal{N}_0 .
- The point x_0 is classified by taking a majority vote of the neighbors.
- This means that KNN estimate the posterior class probability as:

$$\hat{P}(Y = j|X = x_0) = \frac{1}{K} \sum_{i \in \mathcal{N}_0} I(y_i = j).$$

What did we cover today - and what is left?

- What is classification and discrimination?
- What is the Bayes classifier and the Bayes risk?
- What is the sampling paradigm and what is modelled then?
- Linear discriminant analysis: model, method, results.
- Quadratic discriminant analysis: model, method, results.
- Naive Bayes - when and why?
- That is the diagnostic paradigm and what is modelled then?
- KNN - majority vote or estimate posterior class probability?

Part B: more on methods that model and estimate the posterior class probabilities.

So far - and moving on

Part B: Modelling posterior probabilities, ROC/AUC and comparisons

What to remember from Part A?

Aim: Discrimination and classification

Today - data from:

- Default: will a new customer **default** or not based on his/her status (**student** or not), **balance** and **income**?
- South African heart disease data set: classify to coronary heart disease or not, based on 9 covariates.

Notation

Training set: observations (independent pairs) $\{(x_1, y_1), \dots, (x_n, y_n)\}$ where the response variable Y is *qualitative* and labelled $1, 2, \dots, K$.

The training set is used to construct the classification rule (by estimating parameters in class densities or posterior probabilities).

Test set: observations (independent pairs), same format as the training set.

The test set is used to evaluate the classification rule.

Loss function: The misclassifications are given the loss 1 and the correct classifications loss 0 - this is called *0/1-loss*.

Bayes classifier

- Assume that we know or can estimate the probability that a new observation x_0 belongs to class k :

$$p_k(x_0) = \Pr(Y = k|X = x_0), \quad k = 1, 2, \dots, K.$$

This is the probability that $Y = k$ given the observation x_0 . The *Bayes classifier assigns an observation to the most likely class*, given its predictor values.

Two paradigms

- The **sampling paradigm**: There focus is on estimating the prior probabilities for the classes and the class conditional distributions. We classify to the class with the maximal product $\pi_k f_k(x)$. We have looked at LDA (multivariate normal densities with equal covariance matrices) and QDA (ditto, but each class has it's own covariance matrix).
- The **diagnostic paradigm**: We focus on *directly* estimating the posterior distribution for the classes $\Pr(Y = k \mid X = x)$. We have looked at the KNN-classifier in Part A.

Focus now is on *diagnostic paradigm* = we estimates $\Pr(Y = k \mid X = x)$ and classify a new observation based on this estimated probability.

But first, what about linear regression Y on \mathbf{x} to make a classification?

Confusion matrix, sensitivity, specificity

In a two class problem - assume the classes are labelled “-” (non disease) and “+” (disease). In a population setting we define the following event and associated number of observations.

##	[,1]	[,2]	[,3]	[,4]
## row1	"	"Predicted -"	"Predicted +"	"Total"
## row2	"True -"	"True Negative TN"	"False Positive FP"	"N"
## row3	"True +"	"False Negative FN"	"True Positive TP"	"P"
## row4	"Total"	"N*"	"P*"	" "

Sensitivity is the proportion of correctly classified positive observations: $\frac{\# \text{True Positive}}{\# \text{Condition Positive}} = \frac{\text{TP}}{\text{P}}$.

Specificity is the proportion of correctly classified negative observations: $\frac{\# \text{True Negative}}{\# \text{Condition Negative}} = \frac{\text{TN}}{\text{N}}$.

We would like that a classification rule (or a diagnostic test) have both a high sensitivity and a high specificity.

Other useful quantities:

##	[,1]	[,2]
## row1	"Name"	"Definition"
## row2	"False positive rate"	"FP/N"
## row3	"True positive rate"	"TP/P"
## row4	"Positive predictive value (PPV)"	"TP/P*"
## row5	"Negative predictive value (NPV)"	"TN/N*"

##	[,3]
## row1	"Synonyms"
## row2	"Type I error, 1-specificity"
## row3	"1-Type II error, power, sensitivity, recall"
## row4	"Precision, 1-false discovery proportion"
## row5	"

(These two tables are tables 4.6 and 4.7 in our ISL-textbook.)

Example Continued: South African heart disease

We want to evaluate our multiple logistic model for the **SAheart** data set. In order to investigate the training error and the test error, we divide the original data set, randomly, into two samples of equal size.

```
set.seed(20)
d.heart <- SAheart
train_ID = sample(1:nrow(d.heart), nrow(d.heart)/2)
train_SA = d.heart[train_ID, ]
test_SA = d.heart[-train_ID, ]
```

We now fit a logistic regression model, using the training set only:

```
glm_SA = glm(chd ~ ., data = train_SA, family = "binomial")
summary(glm_SA)
```

```
##
## Call:
## glm(formula = chd ~ ., family = "binomial", data = train_SA)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.0583  -0.7660  -0.4523   0.8679   2.2682
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -5.4999981  1.9010122  -2.893  0.00381 **
## sbp           0.0050738  0.0089677   0.566  0.57154
## tobacco      0.1228036  0.0443098   2.771  0.00558 **
## ldl           0.0315840  0.0817017   0.387  0.69907
## adiposity     0.0223541  0.0398203   0.561  0.57454
## famhistPresent 0.7763880  0.3434487   2.261  0.02379 *
## typea         0.0328267  0.0181739   1.806  0.07088 .
## obesity      -0.0487256  0.0592647  -0.822  0.41098
## alcohol      -0.0007868  0.0079569  -0.099  0.92123
## age           0.0470105  0.0170190   2.762  0.00574 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

By comparing this output with the corresponding output above, we see that the estimated coefficients slightly differ. This is because a different data set has been used to fit the model. We previously used the full data set.

We want to estimate the probability of a `chd` event for the observations in the test set. To do this we can insert the estimated coefficient into the logistic equation, remembering that `famhist` is a categorical covariate, modeled by a dummy variable:

$$x_{\text{famhist}} = \begin{cases} 1, & \text{if Present,} \\ 0, & \text{if Absent.} \end{cases}$$

The estimated probability of $Y = 1$ if `famhist` = "Present", given a vector \mathbf{X} of covariate observations is:

$$\hat{p}(\mathbf{X}) = \frac{e^{-7.43+0.01x_{\text{sbp}}+0.09x_{\text{tobacco}}+0.16x_{\text{ldh}}+0.01x_{\text{adiposity}}+1.04 \cdot 1+0.04x_{\text{typea}}-0.04x_{\text{famhist}}}}{1 + e^{-7.43+0.01x_{\text{sbp}}+0.09x_{\text{tobacco}}+0.16x_{\text{ldh}}+0.01x_{\text{adiposity}}+1.04 \cdot 1+0.04x_{\text{typea}}-0.04x_{\text{famhist}}}}$$

Whereas, if `famhist` = "Absent" the estimated probability is:

$$\hat{p}(\mathbf{X}) = \frac{e^{-7.43+0.01x_{\text{sbp}}+0.09x_{\text{tobacco}}+0.16x_{\text{ldh}}+0.01x_{\text{adiposity}}+1.04 \cdot 0+0.04x_{\text{typea}}-0.04x_{\text{famhist}}}}{1 + e^{-7.43+0.01x_{\text{sbp}}+0.09x_{\text{tobacco}}+0.16x_{\text{ldh}}+0.01x_{\text{adiposity}}+1.04 \cdot 0+0.04x_{\text{typea}}-0.04x_{\text{famhist}}}}$$

The `predict` function does these calculations for us. When specifying `type="response"` the function returns the probabilities for $Y = 1$.

```
probs_SA = predict(glm_SA, newdata = test_SA, type = "response")
```

From these probabilities we can obtain classifications, by specifying a threshold value. We have here chosen a threshold value of 0.5. By using the `ifelse` function we specify that all probabilities larger than 0.5 are to be classified as 1, while the remaining probabilities are to be classified as 0.

```
pred_SA = ifelse(probs_SA > 0.5, 1, 0)

predictions_SA = data.frame(probs_SA, pred_SA, test_SA[, 10])
colnames(predictions_SA) = c("Estim. prob. of Y=1", "Predicted class",
                             "True class")
head(predictions_SA)
```

##	Estim. prob. of Y=1	Predicted class	True class
## 1	0.7316942	1	1
## 4	0.7181876	1	1
## 7	0.2485903	0	0
## 11	0.7171302	1	1
## 12	0.7716460	1	1
## 15	0.6464965	1	0

We can now use the confusion matrix to count the number of misclassifications. The below confusion matrix is calculated using the test set and comparing the predicted classes with the true classes.

```
table(pred_SA, SAheart[-train_ID, 10])
```

```
##  
## pred_SA    0    1  
##          0 117  44  
##          1  28  42
```

The logistic model has correctly classified 130+40 times, and misclassified 24+37 times. The misclassification test error rate is thus:

$$\text{Test error} = \frac{24 + 37}{231} \approx 0.264$$

The training error can be calculated in a similar fashion, but now we use the fitted model to make prediction for the training set.

```
SA_train_prob = glm_SA$fitted.values
SA_train_pred = ifelse(SA_train_prob > 0.5, 1, 0)
conf_train = table(SA_train_pred, SAheart[train_ID, 10])
misclas_train = (231 - sum(diag(conf_train)))/231
misclas_train
```

```
## [1] 0.2380952
```

The train misclassification error rate is $\approx 25.1\%$.

ROC curves and AUC

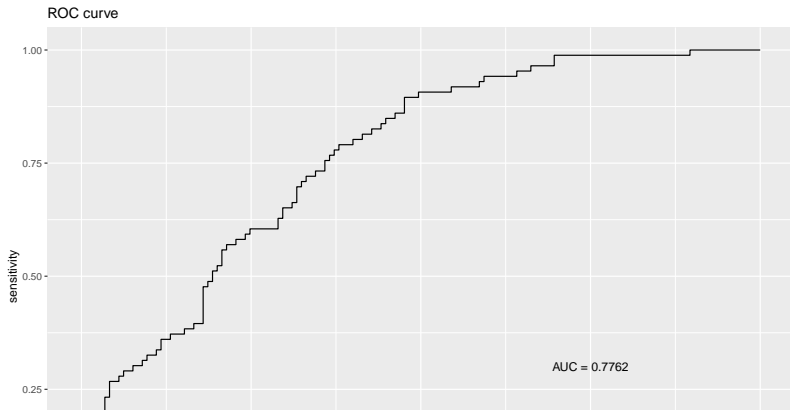
The receiver operating characteristics (ROC) curve gives a graphical display of the sensitivity against specificity, as the threshold value (cut-off on probability of success or disease) is moved over the range of all possible values. An ideal classifier will give a ROC curve which hugs the top left corner, while a straight line represents a classifier with a random guess of the outcome.

The **AUC** score is the area under the AUC curve. It ranges between the values 0 and 1, where a higher value indicates a better classifier. The AUC score is useful for comparing the performance of different classifiers, as all possible threshold values are taken into account.

Example Continued: South African heart disease

In order to see how our model performs for different threshold values, we can plot a ROC curve:

```
library(pROC)
SA_roc = roc(SAheart[-train_ID, 10], probs_SA, legacy.axes = TRUE)
ggroc(SA_roc) + ggtitle("ROC curve") + annotate("text", x = 0.25,
  label = "AUC = 0.7762")
```



To check where in the plot we find the default cut-off on 0.5, we need to calculate sensitivity and specificity for this cut-off:

```
res = table(pred_SA, SAheart[-train_ID, 10])
sens = res[2, 2]/sum(res[2, ])
spec = res[1, 1]/sum(res[1, ])
sens
```

```
## [1] 0.6
```

```
spec
```

```
## [1] 0.7267081
```

Observe that the value 0.6 (on y-axis) and 0.7267081 (on x-axis) is on our ROC curve.

The ROC-curve is made up of all possible cut-offs and their associated sensitivity and specificity.

Which classification method is the best?

Advantages of discriminant analysis

- Discriminant analysis is more stable than logistic regression when the classes are well-separated.
- Discriminant analysis is more stable than logistic regression if the number of observations n is small and the distribution of the predictors \mathbf{X} is approximately (multivariate) normal.

Linearity

Assume a binary classification problem with one covariate. Recall that logistic regression can be written:

$$\log \left(\frac{p(x)}{1 - p(x)} \right) = \beta_0 + \beta_1 x.$$

For LDA we have that $p_0(x)$ is the probability that the observation x belongs to class 0, while $p_1(x) = 1 - p_0(x)$ is the probability that it belongs to class 1.

Observe that this show that our class boundary is linear.

Compulsory Exercise 1, Problem 3a.

$$\begin{aligned}
\log \frac{\Pr(Y = 0|X = x)}{\Pr(Y = 1|X = x)} &= \log \frac{\pi_0}{\pi_1} + \log \frac{f_0(x)}{f_1(x)} \\
&= \log \frac{\pi_0}{\pi_1} - \frac{1}{2\sigma^2}(x - \mu_0)^2 + \frac{1}{2\sigma^2}(x - \mu_1)^2 \\
&= \log \frac{\pi_0}{\pi_1} - \frac{1}{2\sigma^2}(x^2 - 2x\mu_0 + \mu_0^2 - x^2 + 2x\mu_1 - \mu_1^2) \\
&= \log \frac{\pi_0}{\pi_1} - \frac{1}{2\sigma^2}(\mu_0^2 - \mu_1^2) + \frac{1}{\sigma^2}(\mu_0 - \mu_1)x \\
&= \alpha_0 + \alpha_1 x
\end{aligned}$$

The two methods can thus be written in the same form (linear in parameters and in x). The difference is in *how* the parameters $(\alpha_0, \alpha_1, \beta_0, \beta_1)$ are estimated.

LDA vs logistic regression

- Logistic regression uses the diagnostic paradigm, and models the posterior distribution $\Pr(Y = 1|\mathbf{x})$.
- Linear discriminant analysis models the class conditional densities $f_k(\mathbf{x})$.
- The results are usually quite similar, but
 - LDA is “more available” in the multi-class setting
 - if the class conditional distributions are multivariate normal then LDA (or QDA) is preferred
 - if the class conditional distributions are far from multivariate normal then logistic regression is preferred
 - in medicine for two-class problems logistic regression is often preferred (for interpretability) and (always) together with ROC and AUC (for model comparison).

and KNN?

- KNN is used when the class boundaries are non-linear.

Extensions for classifications

- Module 5: how to use cross-validation in model evaluation and model selection
- (Module 6: model selection - but mainly regression)
- Module 7: maybe a taste of nonlinear methods
- Module 8: classification trees (binary splits for the covariates)
- Module 9: support vector machines
- Module 11: neural nets

Further reading

- More on logistic regression from TMA4315 Generalized linear models H2018: [TMA4315M3: Binary regression](#)
- [Videos on YouTube by the authors of ISL, Chapter 4](#)

Want to learn more (theory) about logistic regression?

- In TMA4315 Generalized linear models there are 3 weeks with binary regression - mainly logistic regression.
- The focus there is on all parts of the regression (not classification) with a mathematical focus on estimation, inference, model fit.

R packages

packages to install before knitting this R Markdown file to kn
the Rmd

```
install.packages("knitr")
```

```
install.packages("rmarkdown")
```

nice tables in Rmd

```
install.packages("kableExtra")
```

cool layout for the Rmd plotting

```
install.packages("ggplot2") # cool plotting
```

```
install.packages("ggpubr") # for many ggplots
```

```
install.packages("GGally") # for ggpairs
```

datasets

```
install.packages("ElemStatLearn")
```

```
install.packages("ISLR")
```

data manipulations

```
install.packages("dplyr")
```

```
install.packages("reshape")
```

classificaton

```
install.packages("class")
```

```
install.packages("pROC")
```

diagnostic statistics

Acknowledgements

Thanks to Julia Debik for contributing to this module page.