

SY19 – Machine Learning

Chapter 3: Linear and Quadratic Classification

Thierry Denœux

Université de technologie de Compiègne

<https://www.hds.utc.fr/~tdenoeux>

email: tdenoeux@utc.fr

Automne 2024



Classification

- In classification problems, the response variable Y is **nominal**, i.e., it takes values in a **finite and unordered set** \mathcal{C} , e.g.
 - ▶ Email is one of $\mathcal{C} = \{\text{spam}, \text{email}\}$
 - ▶ Facial expression is one of $\mathcal{C} = \{\text{sadness}, \text{joy}, \text{disgust}, \dots\}$
 - ▶ Object is one of $\mathcal{C} = \{\text{pedestrian}, \text{car}, \text{bike}, \dots\}$, etc.
- The elements in \mathcal{C} are called **classes**. They are arbitrarily numbered $1, 2, \dots, c$.
- A mapping $C : \mathbb{R}^p \rightarrow \mathcal{C}$ that predicts the class for a future predictor vector X is called a **classifier**.
- Our goals are to:
 - ▶ Learn a “good” classifier from a training set
 - ▶ Assess the **uncertainty** in the classification
 - ▶ **Understand** the effects of the different predictors on the decision



Overview

- 1 Introduction to classification
 - Basic notions
 - Bayes classifier
 - Voting K -NN rule
- 2 Logistic regression
 - Binomial logistic regression
 - Multinomial logistic regression
- 3 Linear and quadratic discriminant analysis
 - Quadratic Discriminant Analysis
 - Simplifying assumptions



Overview

- 1 Introduction to classification
 - Basic notions
 - Bayes classifier
 - Voting K -NN rule
- 2 Logistic regression
 - Binomial logistic regression
 - Multinomial logistic regression
- 3 Linear and quadratic discriminant analysis
 - Quadratic Discriminant Analysis
 - Simplifying assumptions



Formalization

- We have a **feature (predictor) vector** X , and a discrete **response variable** Y , both random.
- To represent the joint distribution of (X, Y) , we can specify:
 - 1 The marginal distribution of Y . We use the notation

$$\pi_k = \mathbb{P}(Y = k),$$

and we call π_k the **prior probability** of class k . We have

$$\sum_{k=1}^c \pi_k = 1$$

- 2 The **conditional probability density functions (pdfs)** of X given $Y = k$, for $k = 1, \dots, c$. We use the notation

$$p_k(x) = p(x | Y = k)$$



Example

- Consider a classification problem with $c = 3$ classes and $p = 1$ feature.
- Assume that

$$\pi_1 = 0.3, \quad \pi_2 = 0.5, \quad \pi_3 = 0.2$$

$$p_k(x) = \phi(x; \mu_k, \sigma_k)$$

where ϕ is the normal pdf, with

$$\mu_1 = -1, \quad \mu_2 = 0, \quad \mu_3 = 1.5$$

$$\sigma_1 = 1, \quad \sigma_2 = \sqrt{2}, \quad \sigma_3 = 0.5$$



Formalization (continued)

We can then compute

- The **marginal (mixture) pdf** of X as

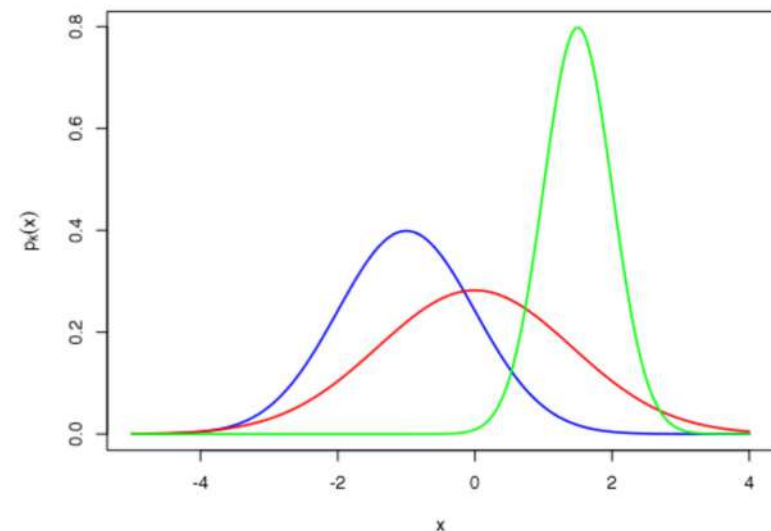
$$p(x) = \sum_{k=1}^c p_k(x) \pi_k$$

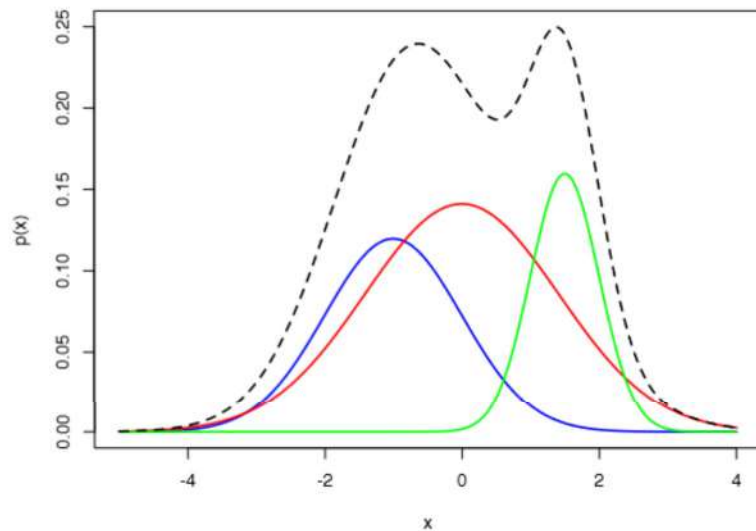
- The conditional distribution of Y given $X = x$ using **Bayes' theorem**.
Let

$$P_k(x) = \mathbb{P}(Y = k | X = x)$$

denote the **posterior (conditional) class probabilities**. We have

$$P_k(x) = \frac{p_k(x) \pi_k}{p(x)} \propto p_k(x) \pi_k, \quad k = 1, \dots, c$$

Example: conditional densities $p_k(x)$ 

Example: marginal density $p(x)$ 

Overview

1 Introduction to classification

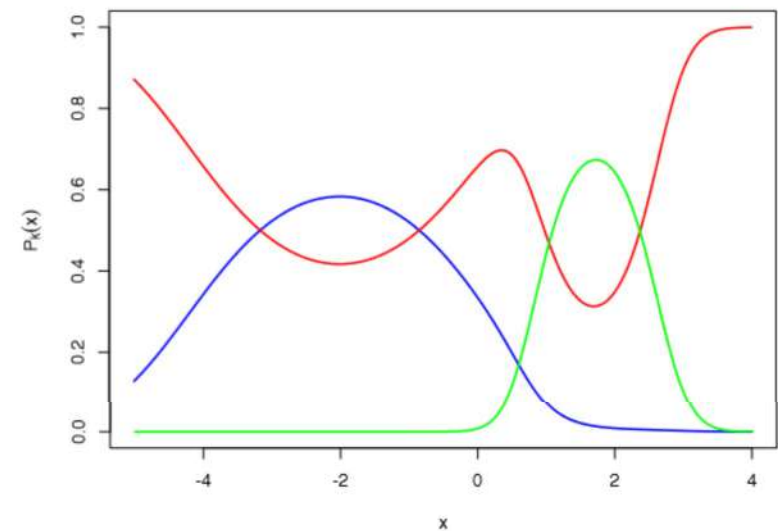
- Basic notions
- Bayes classifier
- Voting K -NN rule

2 Logistic regression

- Binomial logistic regression
- Multinomial logistic regression

3 Linear and quadratic discriminant analysis

- Quadratic Discriminant Analysis
- Simplifying assumptions

Example: posterior probabilities $P_k(x)$ 

The Bayes classifier

- The **conditional error probability** of a classifier C is

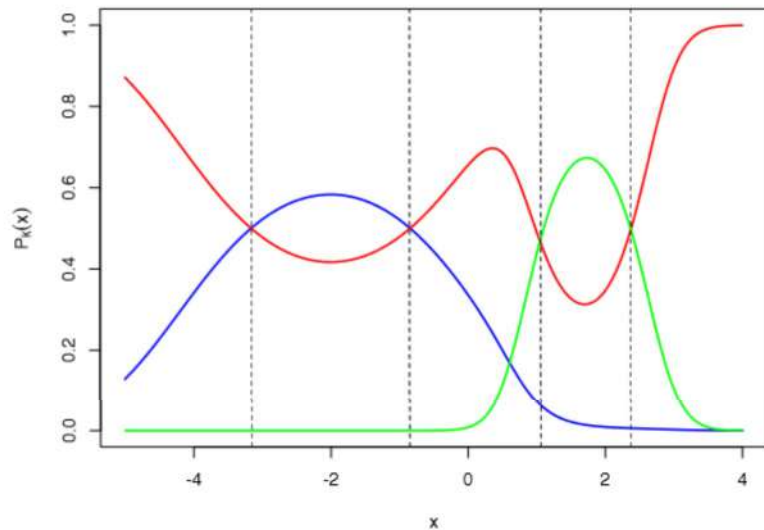
$$\begin{aligned}\mathbb{P}(\text{error} \mid X = x) &= \mathbb{P}(C(X) \neq Y \mid X = x) \\ &= 1 - \mathbb{P}(C(X) = Y \mid X = x)\end{aligned}$$

- If $C(x) = k$, then

$$\mathbb{P}(\text{error} \mid X = x) = 1 - \mathbb{P}(Y = k \mid X = x) = 1 - P_k(x)$$

- To minimize $\mathbb{P}(\text{error} \mid X = x)$, we must choose k such that $P_k(x)$ is maximum.
- The corresponding classifier $C^*(x)$ is called the **Bayes classifier**. It has the lowest error probability.

Example: decision regions of the Bayes classifier



Approximating the Bayes classifier

- The Bayes classifier is optimal but theoretical. We need practical methods to learn classifiers that will approximate the Bayes classifier.
- For this, we need to estimate the posterior probabilities $P_k(x)$.
- As in regression, we distinguish between
 - ▶ **Parametric** methods that postulate a model (of the densities $p_k(x)$, the posterior probabilities $P_k(x)$ or the decisions $C(x)$) depending on a limited number of parameters
 - ▶ **Nonparametric** methods, which make minimal assumptions about the distribution of the data.
- A widely used nonparametric method is the **voting K nearest neighbor (K -NN)** method.

Bayes error rate

- For $X = x$, the Bayes classifier predicts the class k^* such that $P_{k^*}(x) = \max_k P_k(x)$, and the conditional error probability is

$$1 - P_{k^*}(x) = 1 - \max_k P_k(x)$$

- The error probability of the Bayes classifier (averaged over all values of X) is

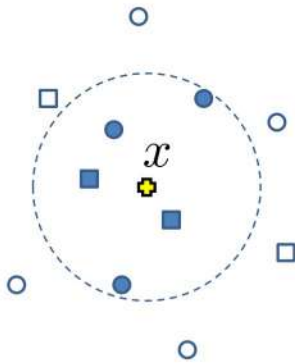
$$\text{Err}_B = \mathbb{E}_X \left[1 - \max_k P_k(X) \right] = \int \left[1 - \max_k P_k(x) \right] p(x) dx$$

- This probability is called the **Bayes error rate**. It is the lowest error probability that can be achieved by a classifier. It characterizes the difficulty of the classification task.

Overview

- 1 Introduction to classification
 - Basic notions
 - Bayes classifier
 - Voting K -NN rule
- 2 Logistic regression
 - Binomial logistic regression
 - Multinomial logistic regression
- 3 Linear and quadratic discriminant analysis
 - Quadratic Discriminant Analysis
 - Simplifying assumptions

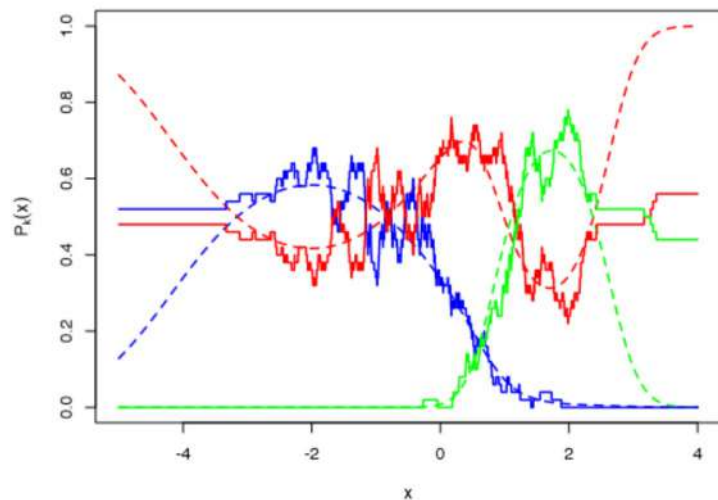
K nearest neighbors



- Nearest-neighbor averaging can be used as in regression.
- Let $x_{(1)}, \dots, x_{(K)}$ denote the K nearest neighbors of x in the learning set, and $y_{(1)}, \dots, y_{(K)}$ the corresponding class labels.



Example: voting K-NN rule with $n = 1000$ and $K = 50$



Voting K-nearest-neighbor rule

- The posterior probability of class k can be estimated by the proportion of observations from that class among the K nearest neighbors of x :

$$\hat{P}_k(x) = \frac{1}{K} \# \{i \in \{1, \dots, K\} : y_{(i)} = k\}$$

- Voting K-nearest neighbor (K-NN) rule: select the majority class among the K nearest neighbors:

$$C_K(x) = \arg \max_k \hat{P}_k(x).$$

- As in regression, the K-NN rule breaks down as dimension grows. However, the impact on $C_K(x)$ is less than that on the probability estimates $\hat{P}_k(x)$.



Error probability estimation

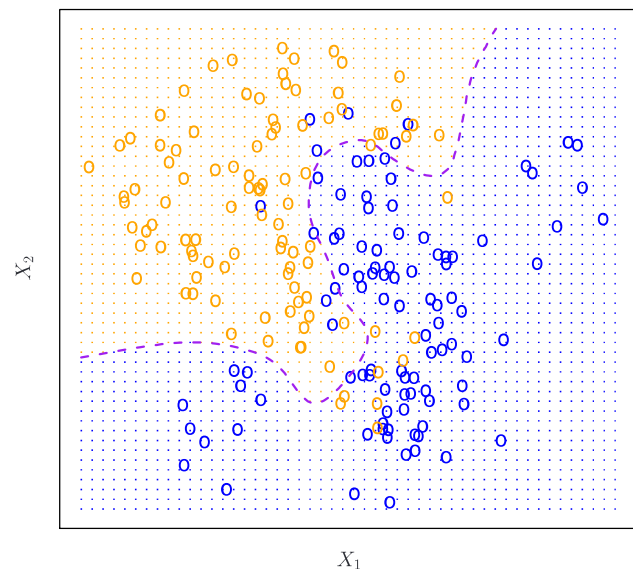
- Typically, we estimate the error probability of a classifier C by its error rate on a test set $\mathcal{T} = \{(x'_i, y'_i)\}_{i=1}^m$:

$$\text{Err}_{\mathcal{T}} = \frac{1}{m} \# \{i \in \{1, \dots, m\} : y'_i \neq C(x'_i)\}$$

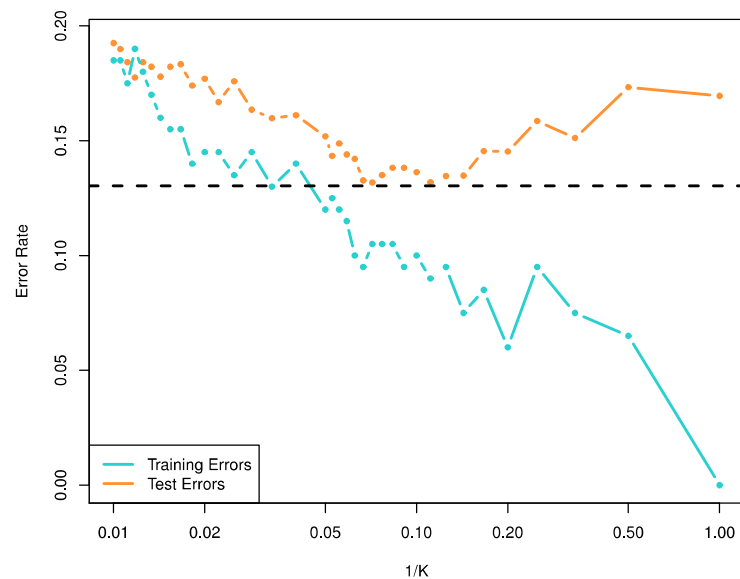
- The test error rate can be used as a criterion for selecting the best model in a set of candidate models (more on this later).



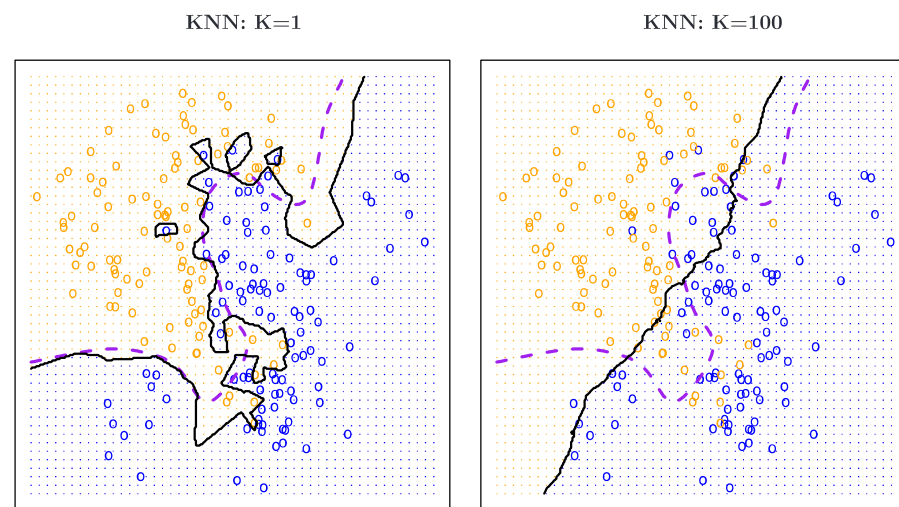
Example: simulated data and Bayes decision boundary



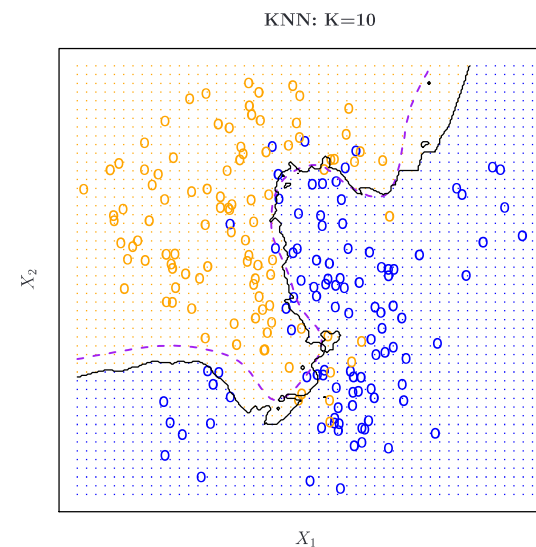
Thierry Denœux SY19 – Linear/Quadratic Classification A24 21 / 81

Training and test error rates vs. $1/K$ 

Thierry Denœux SY19 – Linear/Quadratic Classification A24 23 / 81

Decision boundaries for $K = 1$ and $K = 100$ 

Thierry Denœux SY19 – Linear/Quadratic Classification A24 22 / 81

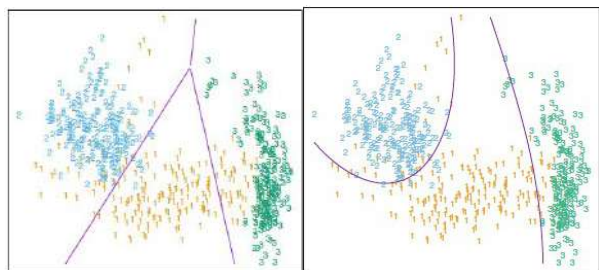
Decision boundary for the best value of K 

Thierry Denœux SY19 – Linear/Quadratic Classification A24 24 / 81

Linear/quadratic classification

In this chapter we focus mostly on **linear** and **quadratic** methods for classification, i.e., we assume the decision boundaries to have equations of the form

$$\beta^T x + \beta_0 = 0 \quad (\text{linear}) \quad \text{or} \\ x^T Q x + \beta^T x + \beta_0 = 0 \quad (\text{quadratic})$$



Overview

- 1 Introduction to classification
 - Basic notions
 - Bayes classifier
 - Voting K-NN rule
- 2 Logistic regression
 - Binomial logistic regression
 - Multinomial logistic regression
- 3 Linear and quadratic discriminant analysis
 - Quadratic Discriminant Analysis
 - Simplifying assumptions



Generative vs. discriminative models

- To approximate Bayes' rule, we need to estimate the posterior probabilities $P_k(x) = P(Y = k | X = x)$.
- We can distinguish two kinds of models for classification:
 - Discriminative models** represent the conditional probabilities $P_k(x)$ directly, or a direct map from inputs x to \mathcal{C} .
 - Generative models** represent the conditional pdfs $p_k(x)$ and the prior probabilities π_k . Using Bayes' theorem, we then get the posterior probabilities $P_k(x)$.
- In this chapter, we will focus on two families of classifiers:
 - 1 A linear classifier based on a discriminative model: **Logistic regression**.
 - 2 Linear and quadratic classifiers based on a generative model: **Linear Discriminant Analysis (LDA)** and **Quadratic Discriminant Analysis (QDA)**.



Overview

- 1 Introduction to classification
 - Basic notions
 - Bayes classifier
 - Voting K-NN rule
- 2 Logistic regression
 - Binomial logistic regression
 - Multinomial logistic regression
- 3 Linear and quadratic discriminant analysis
 - Quadratic Discriminant Analysis
 - Simplifying assumptions



Searching for a linear model

- Consider a **binary classification** problem with $c = 2$ classes, $Y \in \{0, 1\}$.
- Let $P(x) = \mathbb{P}(Y = 1 \mid X = x)$ be the conditional probability of class $Y = 1$.
- We want to find a simple model for $P(x)$. An idea could be to use a linear model of the form

$$P(x) = \beta_0 x_0 + \beta_1 x_1 + \dots + \beta_p x_p = \beta^T x,$$

where x is the augmented feature vector with $x_0 = 1$. However, this is not suitable because $\beta^T x$ can take any value in \mathbb{R} , whereas $P(x) \in [0, 1]$.



Conditional likelihood function

- The Logit model is usually fit by maximizing the **conditional likelihood**, which is the likelihood function, assuming the x_i are fixed.
- Assuming Y_1, \dots, Y_n to be independent conditionally on $X_1 = x_1, \dots, X_n = x_n$, the conditional likelihood is

$$\begin{aligned} L(\beta) &= \mathbb{P}(Y_1 = y_1, \dots, Y_n = y_n \mid X_1 = x_1, \dots, X_n = x_n) \\ &= \prod_{i=1}^n \mathbb{P}(Y_i = y_i \mid X_i = x_i; \beta) \\ &= \prod_{i=1}^n P(x_i; \beta)^{y_i} [1 - P(x_i; \beta)]^{1-y_i} \end{aligned}$$

where $y_i \in \{0, 1\}$ and $P(x_i; \beta) = \mathbb{P}(Y = 1 \mid X = x_i; \beta)$.



Logit model

- A better solution is to assume that **log-odds ratio** is linear in x :

$$\log \frac{P(x)}{1 - P(x)} = \beta^T x$$

- Solving for $P(x)$, we get

$$P(x) = \frac{1}{1 + \exp(-\beta^T x)} = \frac{\exp(\beta^T x)}{1 + \exp(\beta^T x)}$$

- The mapping $u \mapsto \frac{1}{1 + \exp(-u)}$ is the cdf of the standard **logistic distribution**.



Conditional log-likelihood

- The conditional log-likelihood is

$$\begin{aligned} \ell(\beta) &= \sum_{i=1}^n \{y_i \log P(x_i; \beta) + (1 - y_i) \log(1 - P(x_i; \beta))\} \\ &= \sum_{i=1}^n \left\{ y_i \beta^T x_i - \log(1 + \exp(\beta^T x_i)) \right\}, \end{aligned}$$

Proof

- This function is non linear and the score equation $\frac{\partial \ell}{\partial \beta} = 0$ does not have a closed-form solution: we need to use an iterative nonlinear optimization procedure such as the **Newton-Raphson algorithm**.
- As the log-likelihood function is **concave**, it has only one maximum and the convergence of the Newton-Raphson algorithm is guaranteed.



Update equation

- Let \mathbf{y} denote the vector of y_i values, \mathbf{X} the $n \times (p+1)$ matrix of x_i values, \mathbf{p} the vector of fitted probabilities with i -th element $P(x_i; \beta)$.
- The gradient and Hessian of $\ell(\beta)$ can be written as

$$\frac{\partial \ell}{\partial \beta} = \mathbf{X}^T(\mathbf{y} - \mathbf{p}) \quad \text{and} \quad \frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} = -\mathbf{X}^T \mathbf{W} \mathbf{X},$$

where \mathbf{W} an $n \times n$ diagonal matrix of weights with i -th diagonal element $P(x_i; \beta) \{1 - P(x_i; \beta)\}$. Proof

- The update equation is, thus,

$$\beta^{(t+1)} = \beta^{(t)} + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p})$$



Example: Pima diabetes dataset

- Data about diabetes in the population of Pima Indians leaving near Phoenix, Arizona, USA.
- All 768 patients were females and at least 21 years old.
- Variables:
 - ① Number of times pregnant
 - ② Plasma glucose concentration a 2 hours in an oral glucose tolerance test
 - ③ Diastolic blood pressure (mm Hg)
 - ④ Triceps skin fold thickness (mm)
 - ⑤ 2-Hour serum insulin (mu U/ml)
 - ⑥ Body mass index (weight in kg/(height in m)²)
 - ⑦ Diabetes pedigree function
 - ⑧ Age (years)
 - ⑨ Tested positive (1) or negative (0) for diabetes
- Problem: predict the test result for the 8 predictors.

Asymptotic distribution of $\hat{\beta}$

- From the properties of the MLE, we have, asymptotically (as $n \rightarrow +\infty$),

$$\hat{\beta} \stackrel{a}{\sim} \mathcal{N}(\beta_0, \mathcal{I}^{-1})$$

where β_0 is the true value of the parameter and \mathcal{I} is the **Fisher information matrix**:

$$\mathcal{I} = \mathbb{E} \left[- \frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} \Big|_{\beta=\beta_0} \right]$$

which can be estimated by the **observed information matrix**

$$\hat{\mathcal{I}} = - \frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} \Big|_{\beta=\hat{\beta}} = \mathbf{X}^T \hat{\mathbf{W}} \mathbf{X}$$

- This result makes it possible to compute **confidence intervals** and to **test the significance** of the coefficients β_j .



Binomial logistic regression in R

```
glm.fit <- glm(class ~ ., data=pima.train, family=binomial)
```

```
> summary(glm.fit)
```

```
Call:
glm(formula = class ~ ., family = binomial, data = pima.train)
```

```
Deviance Residuals:
```

```
      Min       1Q   Median       3Q      Max
-2.5516  -0.7228  -0.3801   0.7290   2.8209
```

```
Coefficients:
```

```
             Estimate Std. Error z value Pr(>|z|)
(Intercept) -8.733508   0.923597  -9.456 < 2e-16 ***
pregnant     0.110968   0.042667   2.601  0.0093 **
glucose      0.034018   0.004535   7.501 6.35e-14 ***
BP           -0.017566   0.007367  -2.384  0.0171 *
skin         0.004783   0.008659   0.552  0.5807
insulin     -0.001856   0.001197  -1.551  0.1209
bmi          0.109417   0.019992   5.473 4.43e-08 ***
diabetes     1.006970   0.379084   2.656  0.0079 **
age          0.017200   0.012145   1.416  0.1567
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
(Dispersion parameter for binomial family taken to be 1)
```

```
Null deviance: 643.65  on 499  degrees of freedom
Residual deviance: 462.01  on 491  degrees of freedom
AIC: 480.01
```



Prediction

```
pred.pima.glm<-predict(glm.fit,newdata=pima.test,type='response')
```

```
table(pima.test$class,pred.pima.glm>0.5)
```

	FALSE	TRUE
0	158	14
1	41	55

The error rate is $(14 + 41)/268 \approx 0.21$.



True and false positive rates

- Call one class “positive” and the other one “negative”.
- For a given threshold s , we get a **confusion matrix** such as

true	predicted	
	N	P
N	true negative (TN)	false positive (FP)
P	false negative (FN)	true positive (TP)

- The **true positive rate** (sensitivity) and **false positive rate** (1-specificity) are defined, respectively, as

$$TPR = \frac{TP}{FN + TP}, \quad FPR = \frac{FP}{TN + FP}$$

- If we decrease s , we increase both the TPR and the FPR.
- The ROC curve is a plot of the TPR as a function of the FPR, for different values of s .



Changing the decision threshold

- In the previous example, a decision of assignment to the positive class was made if

$$\log \frac{\hat{P}(x)}{1 - \hat{P}(x)} > 0$$

- If the class proportions in the training set do not reflect the proportions in the population, or if the model assumptions are not verified, a different threshold may give better results.
- In general, a decision is made by comparing the log-odds ratio to some threshold s .
- The **Receiver Operating Characteristic (ROC)** curve describes the performance of the classifier for any value of s .



Example

- For the Pima dataset, we had the confusion matrix

	FALSE	TRUE
0	158	14
1	41	55

- Here, the TPR is

$$55/(41 + 55) = 0.57$$

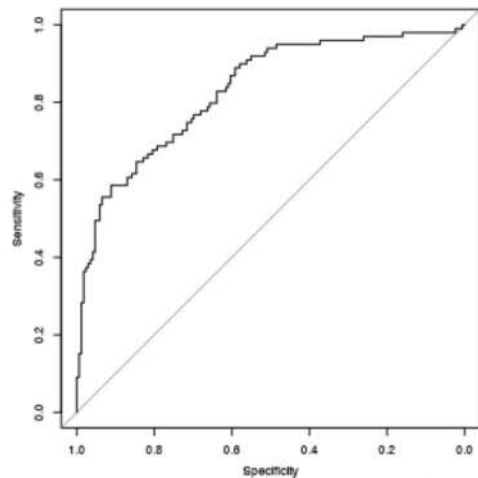
and the FPR is

$$14/(158 + 14) = 0.081$$



ROC curve for the Pima dataset

```
library(pROC)
roc_curve <- roc(pima.test$class, pred.pima.glm)
plot(roc_curve)
```



Model

- Multinomial logistic regression extends binomial logistic regression to $c > 2$ by choosing one class as the baseline (e.g., $Y = 1$) and assuming the log-odds ratios for each class $k \neq 1$ vs. class 1 to be linear in x :

$$\log \frac{P_k(x)}{P_1(x)} = \beta_k^T x, \quad k = 2, \dots, c$$

- Solving for $P_k(x)$, $k = 1, \dots, c$, we get

$$P_1(x) = \frac{1}{1 + \sum_{l=2}^c \exp(\beta_l^T x)}$$

and

$$P_k(x) = \frac{\exp(\beta_k^T x)}{1 + \sum_{l=2}^c \exp(\beta_l^T x)}, \quad k = 2, \dots, c.$$

Overview

- 1 Introduction to classification
 - Basic notions
 - Bayes classifier
 - Voting K-NN rule
- 2 Logistic regression
 - Binomial logistic regression
 - Multinomial logistic regression
- 3 Linear and quadratic discriminant analysis
 - Quadratic Discriminant Analysis
 - Simplifying assumptions

Learning

- The conditional likelihood for the multinomial model is now

$$\begin{aligned} L(\beta) &= \prod_{i=1}^n \mathbb{P}(Y_i = y_i \mid X_i = x_i; \beta) \\ &= \prod_{i=1}^n \prod_{k=1}^c [P_k(x_i; \beta)]^{y_{ik}} \end{aligned}$$

where $y_{ik} = I(y_i = k)$.

- The conditional log-likelihood is

$$\ell(\beta) = \sum_{i=1}^n \sum_{k=1}^c y_{ik} \log P_k(x_i; \beta),$$

- It can be maximized by the Newton-Raphson algorithm as in the binary case.

Example: Letter recognition dataset

- Source: P. W. Frey and D. J. Slate, *Machine Learning*, Vol 6 #2, March 91.
- Objective: identify black-and-white rectangular pixel displays as one of the 26 capital letters in the English alphabet.
- The character images were based on 20 different fonts and each letter within these 20 fonts was randomly distorted to produce a file of 20,000 instances.
- Each instance was converted into 16 primitive numerical attributes (statistical moments and edge counts) which were scaled to fit into a range of integer values from 0 through 15.



Overview

- 1 Introduction to classification
 - Basic notions
 - Bayes classifier
 - Voting K -NN rule
- 2 Logistic regression
 - Binomial logistic regression
 - Multinomial logistic regression
- 3 Linear and quadratic discriminant analysis
 - Quadratic Discriminant Analysis
 - Simplifying assumptions



Multinomial logistic regression in R

```
letter <- read.table("letter-recognition.data",header=FALSE)
n<-nrow(letter)
napp=15000; ntst=n-napp
train<-sample(1:n,napp)
letter.test<-letter[-train,]
letter.train<-letter[train,]
```

```
library(nnet)
fit<-multinom(V1~.,data=letter.train)
pred.letters<-predict(fit,newdata=letter.test)
```

```
perf <-table(letter.test$V1,pred.letters)
1-sum(diag(perf))/ntst
```

0.285



Overview

- 1 Introduction to classification
 - Basic notions
 - Bayes classifier
 - Voting K -NN rule
- 2 Logistic regression
 - Binomial logistic regression
 - Multinomial logistic regression
- 3 Linear and quadratic discriminant analysis
 - Quadratic Discriminant Analysis
 - Simplifying assumptions



Basic assumption

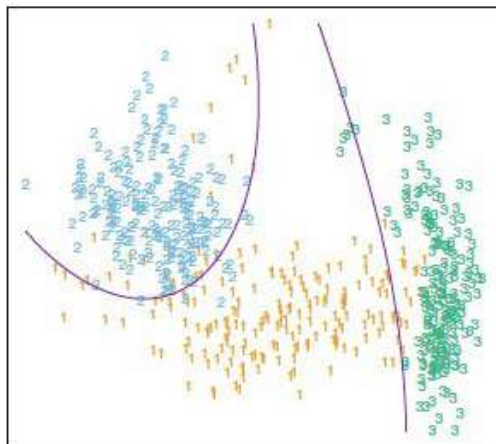
- Quadratic Discriminant Analysis (QDA) is based on the assumption that the class-conditional densities $p_k(x)$ are **multivariate normal**:

$$p_k(x) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) \right\}$$

where $\mu_k = \mathbb{E}(X | Y = k)$ and $\Sigma_k = \text{Var}(X | Y = k)$.

- The parameters of the model are the class-conditional means μ_k and covariance matrices Σ_k , as well as the prior probabilities π_k , $k = 1, \dots, c$.

Example



Optimal decision boundaries

- The boundary between optimal decision regions \mathcal{R}_k and \mathcal{R}_ℓ is defined by the equation

$$P_k(x) = P_\ell(x)$$

- Applying the logarithm to both sides and using the Bayes' theorem

$$P_k(x) \propto p_k(x) \pi_k,$$

we get

$$\log p_k(x) + \log \pi_k = \log p_\ell(x) + \log \pi_\ell \quad (1)$$

- Now

$$\log p_k(x) = -\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) - \frac{1}{2} \log |\Sigma_k| + \text{cst} \quad (2)$$

- From (1) and (2), the boundary equation can be put in the form $x^T Qx + \beta^T x + \beta_0 = 0$: **the decision boundary is a quadric**.

Estimation of parameters

- Let θ be the vector of all parameters.
- Assumption: the sample $(X_1, Y_1), \dots, (X_n, Y_n)$ is independent and identically distributed (iid).
- The **likelihood function** is

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(x_i, y_i) = \prod_{i=1}^n \underbrace{p(x_i | y_i)}_{\prod_{k=1}^c p_k(x_i)^{y_{ik}}} \underbrace{p(y_i)}_{\prod_{k=1}^c \pi_k^{y_{ik}}} \\ &= \prod_{i=1}^n \prod_{k=1}^c \phi(x_i; \mu_k, \Sigma_k)^{y_{ik}} \pi_k^{y_{ik}} \end{aligned}$$

where $\phi(x; \mu_k, \Sigma_k)$ is the normal density with mean μ_k and variance Σ_k .

Maximum likelihood estimates

- The MLEs are

$$\hat{\pi}_k = \frac{n_k}{n}, \quad \hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^n y_{ik} x_i, \quad \text{and}$$

$$\hat{\Sigma}_k = \frac{1}{n_k} \sum_{i=1}^n y_{ik} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T,$$

with $n_k = \sum_{i=1}^n y_{ik}$.

- These estimators are consistent (they converge to the true parameter values when the sample size tends to infinity).



Example (continued)

```
library(MASS)
qda.letter<- qda(V1~.,data=letter.train)
pred.letters.qda<-predict(qda.letter,newdata=letter.test)
```

```
perf <-table(letter.test$V1,pred.letters.qda$class)
1-sum(diag(perf))/ntst
```

0.1166



Implementation of QDA

- To implement QDA, we plug-in the parameter estimates into the expressions of the posterior probabilities. The decision rule is then

$$C(x) = \arg \max_k \hat{P}_k(x)$$

- We actually only need to compute monotonic transformations of the estimated posterior probabilities $\hat{P}_k(x)$, called discriminant functions (DFs).
- We get quadratic DFs by applying a logarithmic transformation:
 - Case $c = 2$: we only need one DF $\delta(x) = \log \hat{P}_1(x) - \log \hat{P}_2(x)$ and the decision rule is

$$C(x) = \begin{cases} 1 & \text{if } \delta(x) > 0 \\ 2 & \text{otherwise.} \end{cases}$$

- Case $c > 2$: we need c DFs $\delta_k(x) = \log \hat{P}_k(x)$ and the decision rule is

$$C(x) = \arg \max_k \delta_k(x)$$



Confusion matrix

```
> print(perf)
```

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	
A	182	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	2	1	0	0	0	3	0
B	0	165	0	3	2	1	0	4	1	0	0	0	0	0	2	0	0	8	4	0	0	0	0	0	0	0	0
C	0	0	164	0	2	1	6	0	0	0	5	1	0	0	3	0	0	0	0	0	1	0	0	0	0	0	0
D	0	3	0	196	0	0	0	0	0	0	1	0	3	0	3	1	1	4	3	3	1	1	0	0	0	0	0
E	0	1	3	0	163	1	7	1	0	0	1	3	0	0	0	0	1	1	0	0	0	0	0	2	0	3	
F	0	1	0	1	1	163	2	1	2	0	0	0	0	4	0	2	0	1	0	6	0	0	1	0	2	0	
G	0	3	4	1	0	3	145	0	0	0	3	3	0	0	7	0	1	3	7	0	0	0	0	0	0	0	
H	0	3	1	10	2	4	3	108	0	0	14	0	1	0	3	1	0	6	0	0	1	1	0	4	3	0	
I	0	1	0	4	1	5	0	0	157	7	1	0	0	1	0	0	0	1	3	0	0	0	0	2	0	1	
J	0	0	0	1	0	1	0	3	3	145	0	0	0	1	1	0	2	0	2	0	0	0	0	0	0	2	
K	0	0	5	1	1	0	0	5	0	0	161	0	0	0	0	0	0	12	1	0	0	1	0	5	0	0	
L	0	0	0	0	1	0	0	3	0	0	0	189	0	0	0	0	1	2	10	0	0	0	0	2	1	0	
M	2	1	0	0	0	0	0	0	0	0	0	0	181	0	0	0	0	1	0	0	0	1	3	0	0	0	
N	0	0	0	4	0	1	0	2	0	0	1	0	0	163	2	1	0	2	0	0	0	0	1	0	1	0	
O	1	0	0	2	0	0	1	3	0	0	0	0	2	0	173	0	9	1	0	0	0	0	1	0	0	0	
P	0	0	0	0	0	11	2	1	0	0	1	1	0	0	0	186	2	2	0	0	0	0	0	0	1	0	
Q	2	0	0	0	0	0	0	0	0	0	0	1	0	0	9	0	179	1	5	0	0	0	0	0	0	0	
R	0	5	1	3	0	0	1	2	0	0	2	0	0	0	0	0	0	1	180	1	0	1	0	0	0	0	
S	0	5	0	1	7	4	0	0	1	1	1	0	0	0	0	0	0	0	178	0	0	0	0	0	0	6	
T	1	0	0	0	5	3	2	2	0	0	1	1	0	0	0	0	0	3	1	168	5	0	0	1	2	0	
U	1	0	0	0	0	0	1	0	0	1	1	0	3	0	3	0	0	0	0	0	210	1	5	0	0	0	
V	0	2	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	181	3	0	3	0	
W	0	0	0	0	0	0	0	1	0	0	1	0	4	1	2	0	0	0	0	0	2	0	169	0	1	0	
X	0	0	0	0	2	0	0	5	3	1	5	1	0	0	0	0	1	1	1	1	3	0	154	1	1	1	
Y	0	0	0	0	0	3	0	0	0	0	0	0	0	0	2	0	0	0	1	4	1	29	0	162	0		
Z	2	0	0	0	2	0	0	0	0	0	1	0	0	0	0	0	2	0	4	2	1	0	0	1	0	195	



Overview

- 1 Introduction to classification
 - Basic notions
 - Bayes classifier
 - Voting K -NN rule
- 2 Logistic regression
 - Binomial logistic regression
 - Multinomial logistic regression
- 3 Linear and quadratic discriminant analysis
 - Quadratic Discriminant Analysis
 - Simplifying assumptions



Linear Discriminant Analysis (LDA)

- LDA is based on the assumption that the class-conditional covariance matrices are equal:

$$\Sigma_k = \Sigma, \quad \text{for all } k.$$

- In that case, the equation of the optimal boundary between regions \mathcal{R}_k and \mathcal{R}_ℓ ,

$$\log p_k(x) + \log \pi_k = \log p_\ell(x) + \log \pi_\ell \quad (3)$$

is **linear**. Indeed, we now have

$$\log p_k(x) = -\frac{1}{2}(x - \mu_k)^T \Sigma^{-1}(x - \mu_k) + \text{cst} \quad (4)$$

$$= \mu_k^T \Sigma^{-1} x - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \text{cst} \quad (5)$$

- The DF and the decision boundaries are now **linear** in x .

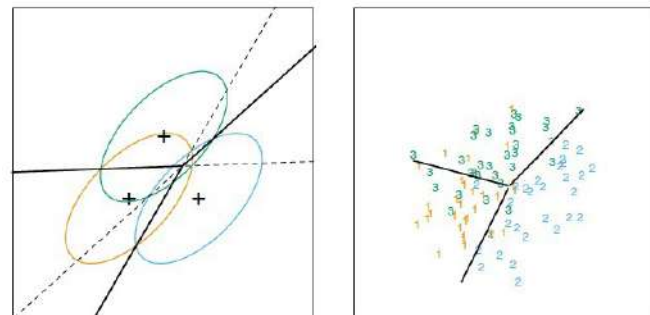


Number of parameters for QDA

- The number of parameters for QDA is $c[p + p(p + 1)/2] + c - 1$ (c means, c covariance matrices and $c - 1$ prior probabilities).
- This number is quadratic in p : the method becomes impractical when p is large, and QDA may perform badly when p is large and n is small.
- We can decrease the number of parameters to estimate by making **simplifying assumptions**. We will consider two such assumptions:
 - 1 Equality of covariance matrices (homoscedasticity) \rightarrow **Linear Discriminant Analysis (LDA)**
 - 2 Conditional independence of the predictors X_j given the class variable $Y \rightarrow$ **Naive Bayes classifiers**



Example



Left: contours of constant density enclosing 95% of the probability in each case. The Bayes decision boundaries between each pair of classes are shown (broken straight lines), and the Bayes decision boundaries separating all three classes are the thicker solid lines. Right: a sample of size 30 drawn from each distribution, and the fitted LDA decision boundaries.



Estimation of parameters

- The MLEs are

$$\hat{\pi}_k = \frac{n_k}{n}, \quad \hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^n y_{ik} x_i, \quad \text{and} \quad \hat{\Sigma} = \frac{1}{n} \sum_{k=1}^c n_k \hat{\Sigma}_k$$

where, as before $\hat{\Sigma}_k$ is the sample covariance matrix in class k :

$$\hat{\Sigma}_k = \frac{1}{n_k} \sum_{i=1}^n y_{ik} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T,$$

and $n_k = \sum_{i=1}^n y_{ik}$.

- It can be shown that $\hat{\Sigma}$ is biased. An unbiased estimator of Σ is

$$S = \frac{n}{n-c} \hat{\Sigma}.$$



Comparing classifiers: McNemar's test

- The test error rate was 0.1166 for QDA, and it is 0.2996 for LDA. Is this difference statistically significant?
- To answer such a question, we typically use McNemar's test for 2×2 contingency tables.
- We consider the following table:

	classifier 2 wrong	classifier 2 correct
classifier 1 wrong	n_{00}	n_{01}
classifier 1 correct	n_{10}	n_{11}



LDA in R

```
lda.letter <- lda(V1~., data=letter.train)
```

```
pred.letters.lda <- predict(lda.letter, newdata=letter.test)
```

```
perf <- table(letter.test$V1, pred.letters.lda$class)
1 - sum(diag(perf)) / ntst
0.2996
```



Comparing classifiers: McNemar's test (continued)

- Under the null hypothesis that the error probabilities of the two classifiers are equal, the statistics

$$D^2 = \frac{(|n_{01} - n_{10}| - 1)^2}{n_{01} + n_{10}}$$

is distributed approximately as χ^2 with 1 degree of freedom.

- The p-values is $p = \mathbb{P}_{H_0}(\chi_1^2 \geq d^2)$.
- Remark: when comparing more than two classifiers, we have more chance of rejecting the null hypothesis for at least one pair of classifiers. To address this problem, we can use the Bonferroni correction: we reject the null hypothesis at level α for any two classifiers if $p \leq \alpha/m$, where m is the number of classifier pairs.



McNemar's test in R

```
correct.lda<-letter.test$V1==pred.letters.lda$class
correct.qda<-letter.test$V1==pred.letters.qda$class
mcnemar.test(correct.lda,correct.qda)
```

McNemar's Chi-squared test with continuity correction

data: correct.lda and correct.qda

McNemar's chi-squared = 767.12, df = 1, p-value < 2.2e-16

Naive Bayes classifiers

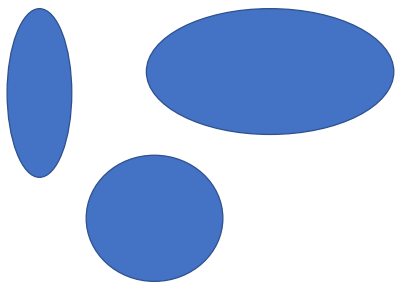
- In LDA and QDA, we need to estimate covariance matrices with $p(p+1)/2$ parameters, which can yield poor results (or can even be unfeasible) when p is very large.
- Starting from the QDA model, we get a simpler model by assuming that the covariance matrices Σ_k are diagonal:

$$\Sigma_k = \text{diag}(\sigma_{k1}^2, \dots, \sigma_{kp}^2),$$

where $\sigma_{kj}^2 = \text{Var}(X_j | Y = k)$ (see next slide).

- We get a naive QDA classifier, a special kind of naive Bayes classifier.

Naive QDA model



Conditional independence assumption

- The assumption that the covariance matrices are diagonal means that the predictors are conditionally independent given the class variable Y , i.e., for all $k \in \{1, \dots, c\}$,

$$p_k(x_1, \dots, x_p) = \prod_{j=1}^p p_{kj}(x_j)$$

- Remark: conditional independence does not imply independence. (Example: Height and vocabulary of kids are not independent; but they are conditionally independent given age).

Naive Bayes classifiers (continued)

- To estimate Σ_k under the conditional independence assumption, we simply set the off-diagonal terms in $\hat{\Sigma}_k$ to 0. The variance σ_{kj}^2 of X_j conditionally on $Y = k$ is estimated by

$$\hat{\sigma}_{kj}^2 = \frac{1}{n_k} \sum_{i=1}^n y_{ik} (x_{ij} - \hat{\mu}_{kj})^2.$$

- A further simplification is achieved by assuming that the covariance matrices are diagonal and equal:

$$\Sigma_1 = \dots = \Sigma_c = \Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_p^2).$$

This model can be called "Naive LDA" (see next slide).

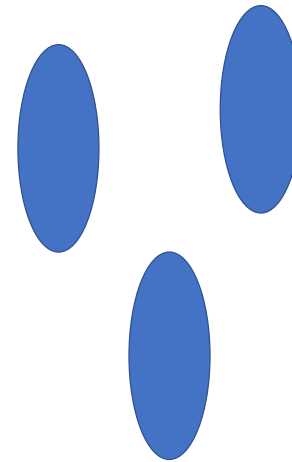


Advantages of Naive Bayes classifiers

- In spite of their simplicity, naive Bayes classifiers often (but not always) have very good performances, especially when the number p of predictors is large.
- They can accommodate **mixed feature vectors** (qualitative and quantitative). If X_j is qualitative, we can estimate the probability mass functions $p_{kj}(x_j)$ using histograms over discrete categories.



Naive LDA model



Naive Bayes classifier in R

```
library(naivebayes)
naive.letter<- naive_bayes(V1~.,data=letter.train)
pred.letters.naive<-predict(naive.letter,newdata=letter.test)

perf.naive <-table(letter.test$V1,pred.letters.naive)
1-sum(diag(perf.naive))/ntst

0.3554
# Comparison with LDA
correct.naive<-letter.test$V1==pred.letters.naive
mcnemar.test(correct.lda,correct.naive)
```

McNemar's Chi-squared test with continuity correction

```
data: correct.lda and correct.naive
McNemar's chi-squared = 83.731, df = 1, p-value < 2.2e-16
```



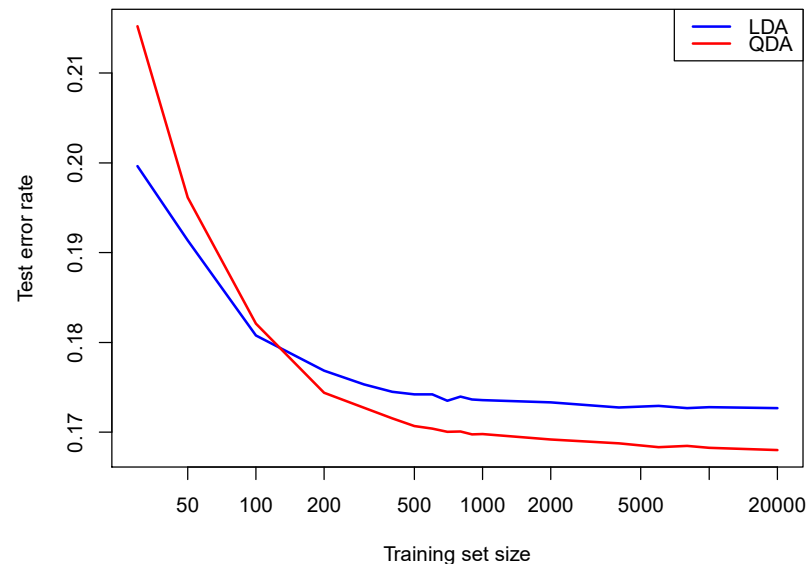
Comparison of the different models

Model	Number of parameters
QDA	$c \left(p + \frac{p(p+1)}{2} \right) + c - 1$
naive QDA	$2cp + c - 1$
LDA	$cp + \frac{p(p+1)}{2} + c - 1$
naive LDA	$cp + p + c - 1$

- QDA is the most general model. However, it does not always yield the best performances, because it has the largest number of parameters.
- Although LDA also has a number of parameters proportional to p^2 , it is usually **much more stable** than QDA. This method is recommended when n is small.
- Naive Bayes classifiers have a number of parameters proportional to p . They often outperform other methods when p is very large.



Result



Example

- We consider $c = 2$ classes with $p = 3$ normally distributed input variables, with the following parameters

$$\pi_1 = \pi_2 = 0.5$$

$$\mu_1 = (0, 0, 0)^T, \quad \mu_2 = (1, 1, 1)^T$$

$$\Sigma_1 = \mathbf{I}_3, \quad \Sigma_2 = 0.7\mathbf{I}_3.$$

- LDA and QDA classifiers were trained using training sets of different sizes between 30 and 20,000, and their error probability was estimated using a test set of size 20,000.
- For each training set size, the experiment was repeated 20 times. The next figure shows error rates over the 20 replications.



Log-likelihood of binary logistic regression

From

$$P(x_i) = \frac{1}{1 + \exp(-\beta^T x_i)} \quad \text{and} \quad 1 - P(x_i) = \frac{\exp(-\beta^T x_i)}{1 + \exp(-\beta^T x_i)}$$

we get

$$\ell(\beta) = \sum_{i=1}^n \left\{ -y_i \log[1 + \exp(-\beta^T x_i)] - \underbrace{\beta^T x_i - \log[1 + \exp(-\beta^T x_i)]}_{-\log[1 + \exp(\beta^T x_i)]} + y_i \beta^T x_i + y_i \log[1 + \exp(-\beta^T x_i)] \right\}$$

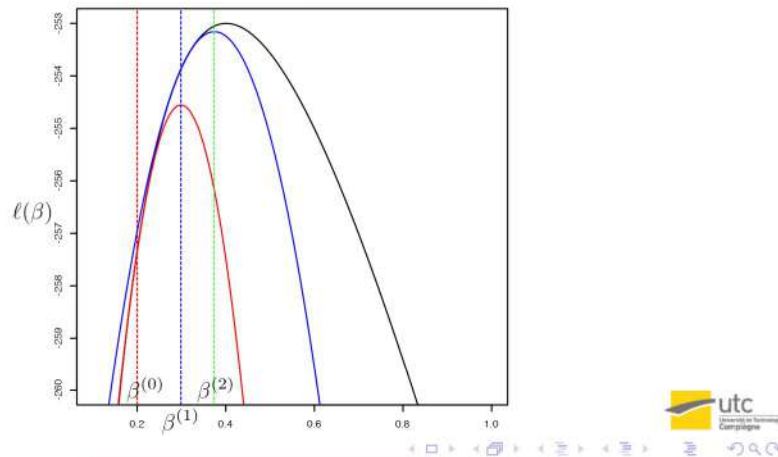
$$\ell(\beta) = \sum_{i=1}^n \left\{ y_i \beta^T x_i - \log[1 + \exp(\beta^T x_i)] \right\}$$



The Newton-Raphson algorithm

Main ideas

- An iterative optimization algorithm.
- Basic idea: at each time step, approximate $\ell(\beta)$ around the current estimate $\beta^{(t)}$ by the **second-order Taylor series expansion**.



Gradient of $\ell(\beta)$

From

$$\ell(\beta) = \sum_{i=1}^n \left\{ y_i \beta^T x_i - \log(1 + \exp(\beta^T x_i)) \right\}$$

the gradient is

$$\begin{aligned} \frac{\partial \ell}{\partial \beta} &= \sum_{i=1}^n y_i x_i - \underbrace{\frac{\exp(\beta^T x_i)}{1 + \exp(\beta^T x_i)}}_{P(x_i; \beta)} x_i \\ &= \sum_{i=1}^n x_i (y_i - P(x_i; \beta)) = \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \end{aligned}$$

where \mathbf{y} denote the vector of y_i values, \mathbf{X} the $n \times (p+1)$ matrix of x_i values, \mathbf{p} the vector of fitted probabilities with i -th element $P(x_i; \beta)$.

The Newton-Raphson algorithm

- We have

$$\begin{aligned} \ell(\beta) &\approx \ell(\beta^{(t)}) + (\beta - \beta^{(t)})^T \frac{\partial \ell(\beta^{(t)})}{\partial \beta} + \\ &\quad \frac{1}{2} (\beta - \beta^{(t)})^T \frac{\partial^2 \ell(\beta^{(t)})}{\partial \beta \partial \beta^T} (\beta - \beta^{(t)}). \end{aligned}$$

- Differentiating both sides w.r.t. β , we get

$$\frac{\partial \ell(\beta)}{\partial \beta} \approx \frac{\partial \ell(\beta^{(t)})}{\partial \beta} + \frac{\partial^2 \ell(\beta^{(t)})}{\partial \beta \partial \beta^T} (\beta - \beta^{(t)}).$$

- Setting $\frac{\partial \ell}{\partial \beta}(\beta) = 0$, we get the update equation

$$\beta^{(t+1)} = \beta^{(t)} - \left(\frac{\partial^2 \ell(\beta^{(t)})}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial \ell(\beta^{(t)})}{\partial \beta}$$

Back

Hessian of $\ell(\beta)$

- From

$$\frac{\partial \ell}{\partial \beta_j} = \sum_{i=1}^n x_{ij} \left(y_i - \underbrace{P(x_i; \beta)}_{\Lambda(\beta^T x)} \right)$$

and $\Lambda'(u) = \Lambda(u)[1 - \Lambda(u)]$, we have

$$\frac{\partial^2 \ell}{\partial \beta_j \partial \beta_k} = - \sum_{i=1}^n x_{ij} x_{ik} P(x_i; \beta) [1 - P(x_i; \beta)]$$

Hessian of $\ell(\beta)$ II

- The Hessian matrix can, thus, be written as

$$\begin{aligned}\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} &= - \sum_{i=1}^n x_i x_i^T P(x_i; \beta) [1 - P(x_i; \beta)] \\ &= -\mathbf{X}^T \mathbf{W} \mathbf{X},\end{aligned}$$

where \mathbf{W} an $n \times n$ diagonal matrix of weights with i -th diagonal element $P(x_i; \beta) [1 - P(x_i; \beta)]$.

[Back](#)