

Arthur Charpentier

arthur.charpentier@gmail.com

<http://freakonometrics.hypotheses.org/>

École Doctorale, Université de Rennes 1, March 2015

Data Science, an Overview of Classification Techniques

“An expert is a man who has made all the mistakes which can be made, in a narrow field” N. Bohr

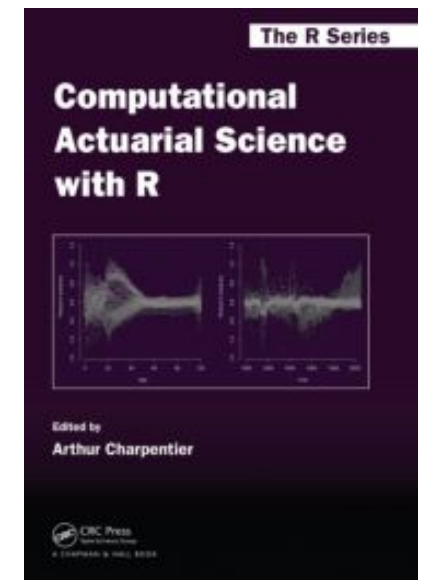
Arthur Charpentier

arthur.charpentier@gmail.com

<http://freakonometrics.hypotheses.org/>

École Doctorale, Université de Rennes 1, March 2015

Professor of Actuarial Sciences, Mathematics Department, UQàM
(previously Economics Department, Univ. Rennes 1 & ENSAE Paristech
actuary AXA General Insurance Hong Kong, IT & Stats FFSA)
PhD in Statistics (KU Leuven), Fellow Institute of Actuaries
MSc in Financial Mathematics (Paris Dauphine) & ENSAE
Editor of the freakonometrics.hypotheses.org's blog
Editor of Computational Actuarial Science, CRC



Supervised Techniques : Classification

(Linear) Discriminant Analysis

Data : $\{(\mathbf{x}_i, y_i) = (x_{1,i}, x_{2,i}, y_i), i = 1, \dots, n\}$

with $y_i \in \{0, 1\}$ or $y_i \in \{-1, +1\}$ or $y_i \in \{\bullet, \bullet\}$

$$\begin{cases} \mathbf{X}|Y = 0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \\ \mathbf{X}|Y = 1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \end{cases}$$

Fisher's linear discriminant

$$\boldsymbol{\omega} \propto [\boldsymbol{\Sigma}_0 + \boldsymbol{\Sigma}_1]^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$$

maximizes

$$\frac{\text{variance between}}{\text{variance within}} = \frac{[\boldsymbol{\omega} \cdot \boldsymbol{\mu}_1 - \boldsymbol{\omega} \cdot \boldsymbol{\mu}_0]^2}{\boldsymbol{\omega}^\top \boldsymbol{\Sigma}_1 \boldsymbol{\omega} + \boldsymbol{\omega}^\top \boldsymbol{\Sigma}_0 \boldsymbol{\omega}}$$

see Fisher (1936, wiley.com)

Supervised Techniques : Classification

Logistic Regression

Data : $\{(\mathbf{x}_i, y_i) = (x_{1,i}, x_{2,i}, y_i), i = 1, \dots, n\}$

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = \frac{\exp[\mathbf{x}^\top \boldsymbol{\beta}]}{1 + \exp[\mathbf{x}^\top \boldsymbol{\beta}]}$$

Inference using maximum likelihood techniques

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin} \left\{ \sum_{i=1}^n \log[\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}_i)] \right\}$$

and the score model is then

$$s(\mathbf{X} = \mathbf{x}) = \frac{\exp[\mathbf{x}^\top \hat{\boldsymbol{\beta}}]}{1 + \exp[\mathbf{x}^\top \hat{\boldsymbol{\beta}}]}$$

Supervised Techniques : Classification

Logistic Regression

Historically, the idea was to model the *odds ratio*,

$$\frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y \neq 1 | \mathbf{X} = \mathbf{x})} = \exp[\mathbf{x}^\top \boldsymbol{\beta}]$$

(the odds ratio is a positive number). Hence,

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = H(\mathbf{x}^\top \boldsymbol{\beta}) \text{ where } H(\cdot) = \frac{\exp[\cdot]}{1 + \exp[\cdot]}$$

is the c.d.f. of the *logistic* variable, popular in demography, see Verhulst (1845, gdz.sub.uni-goettingen.de)
cf TRISS Tauma, Boyd *et al.* (1987, journals.lww.com)

Soit p la population : représentons par dp l'accroissement infiniment petit qu'elle reçoit pendant un temps infiniment court dt . Si la population croissait en progression géométrique, nous aurions l'équation $\frac{dp}{dt} = mp$. Mais comme la vitesse d'accroissement de la population est retardée par l'augmentation même du nombre des habitants, nous devons retrancher de mp une fonction inconnue de p ; de manière que la formule à intégrer devienne

$$\frac{dp}{dt} = mp - \varphi(p).$$

L'hypothèse la plus simple que l'on puisse faire sur la forme de la fonction φ , est de supposer $\varphi(p) = np^2$. On trouve alors pour intégrale de l'équation ci-dessus

$$t = \frac{1}{m} [\log. p - \log. (m - np)] + \text{constante},$$

et il suffira de trois observations pour déterminer les deux coefficients constants m et n et la constante arbitraire.

En résolvant la dernière équation par rapport à p , il vient

$$p = \frac{mp' e^{mt}}{np' e^{mt} + m - np'} \dots \dots \dots (1)$$

en désignant par p' la population qui répond à $t = 0$, et par e la base des logarithmes népériens. Si l'on fait $t = \infty$, on voit que la valeur de p correspondante est $P = \frac{m}{n}$. Telle est donc la *limite supérieure de la population*.

Supervised Techniques : Classification

Probit Regression

Bliss (1934) sciencemag.org) suggested a model such that

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = H(\mathbf{x}^\top \boldsymbol{\beta}) \text{ where } H(\cdot) = \Phi(\cdot)$$

the c.d.f. of the $\mathcal{N}(0, 1)$ distribution. This is the *probit* model.

This yields a latent model, $y_i = \mathbf{1}(y_i^* > 0)$ where

$$y_i^* = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i \text{ is a nonobservable score.}$$

Table 3.2 Transformation of percentages to probits

%	0	1	2	3	4	5	6	7	8	9
0	—	2.67	2.95	3.12	3.25	3.36	3.45	3.52	3.59	3.66
10	3.72	3.77	3.82	3.87	3.92	3.96	4.01	4.05	4.08	4.12
20	4.16	4.19	4.23	4.26	4.29	4.33	4.36	4.39	4.42	4.45
30	4.48	4.50	4.53	4.56	4.59	4.61	4.64	4.67	4.69	4.72
40	4.75	4.77	4.80	4.82	4.85	4.87	4.90	4.92	4.95	4.97
50	5.00	5.03	5.05	5.08	5.10	5.13	5.15	5.18	5.20	5.23
60	5.25	5.28	5.31	5.33	5.36	5.39	5.41	5.44	5.47	5.50
70	5.52	5.55	5.58	5.61	5.64	5.67	5.71	5.74	5.77	5.81
80	5.84	5.88	5.92	5.95	5.99	6.04	6.08	6.13	6.18	6.23
90	6.28	6.34	6.41	6.48	6.55	6.64	6.75	6.88	7.05	7.33
—	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
99	7.33	7.37	7.41	7.46	7.51	7.58	7.65	7.75	7.88	8.09

Supervised Techniques : Classification

Logistic Regression

The classification function, from a score to a class:

if $s(\mathbf{x}) > s$, then $\hat{Y}(\mathbf{x}) = 1$ and $s(\mathbf{x}) \leq s$, then $\hat{Y}(\mathbf{x}) = 0$

Plot $TP(s) = \mathbb{P}[\hat{Y} = 1|Y = 1]$ vs. $FP(s) = \mathbb{P}[\hat{Y} = 1|Y = 0]$

Supervised Techniques : Classification

Logistic Additive Regression

Data : $(\mathbf{x}_i, y_i) = (x_{1,i}, x_{2,i}, y_i), i = 1, \dots, n$

with $y_i \in \{-0, 1\}$ or $y_i \in \{-1, +1\}$ or $y_i \in \{-\bullet, \bullet\}$

Instead of a linear function

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = \frac{\exp[\mathbf{x}^\top \boldsymbol{\beta}]}{1 + \exp[\mathbf{x}^\top \boldsymbol{\beta}]}$$

consider

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = \frac{\exp[h(\mathbf{x})]}{1 + \exp[h(\mathbf{x})]}$$

Supervised Techniques : Classification

k-Nearest Neighbors

Data : $(\mathbf{x}_i, y_i) = (x_{1,i}, x_{2,i}, y_i), i = 1, \dots, n$

with $y_i \in \{-0, 1\}$ or $y_i \in \{-1, +1\}$ or $y_i \in \{-\bullet, \bullet\}$

for each \mathbf{x} , consider the k nearest neighbors (for some distance $d(\mathbf{x}, \mathbf{x}_i)$) $V_k(\mathbf{x})$

$$s(\mathbf{x}) = \frac{1}{k} \sum_{i \in V_k(\mathbf{x})} Y_i$$

Supervised Techniques : Classification

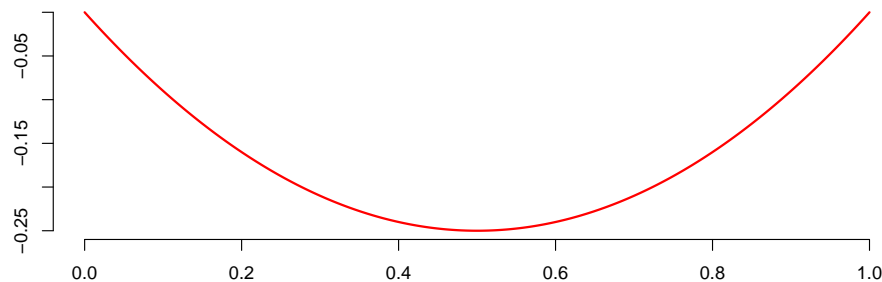
CART and Classification Tree

Data : $(\mathbf{x}_i, y_i) = (x_{1,i}, x_{2,i}, y_i), i = 1, \dots, n$

with $y_i \in \{-0, 1\}$ or $y_i \in \{-1, +1\}$ or $y_i \in \{-\bullet, \bullet\}$

Compute some impurity criteria, e.g. Gini index,

$$- \sum_{P \in \{-A, B\}} \underbrace{\mathbb{P}[\mathbf{x} \in P]}_{\text{size}} \underbrace{\mathbb{P}[Y = 1 | \mathbf{x} \in P]}_p \underbrace{\mathbb{P}[Y = 0 | \mathbf{x} \in P]}_{(1-p)}$$



Supervised Techniques : Classification

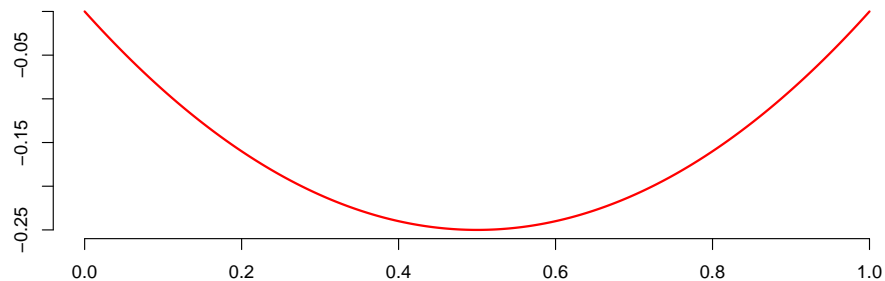
CART and Classification Tree

Data : $(\mathbf{x}_i, y_i) = (x_{1,i}, x_{2,i}, y_i), i = 1, \dots, n$

with $y_i \in \{-0, 1\}$ or $y_i \in \{-1, +1\}$ or $y_i \in \{-\bullet, \bullet\}$

Given a partition, loop,

$$- \sum_{P \in \{-A, B, C\}} \mathbb{P}[\mathbf{x} \in P] \mathbb{P}[Y = 0 | \mathbf{x} \in P] \mathbb{P}[Y = 1 | \mathbf{x} \in P]$$

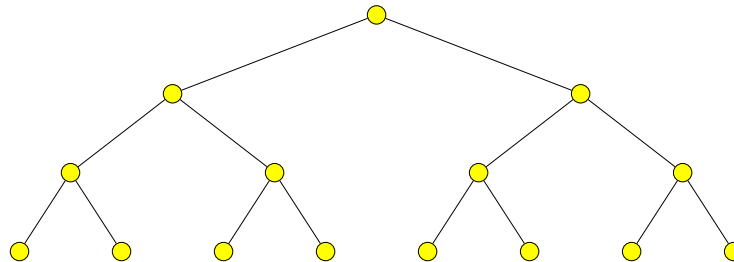


Supervised Techniques : Classification

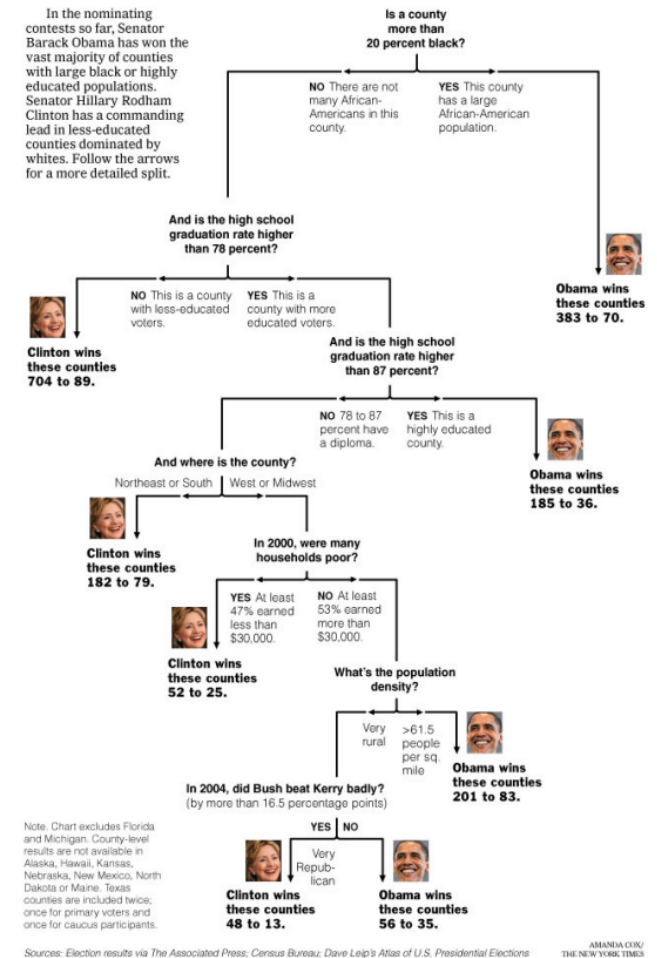
CART and Classification Tree

Breiman *et al.* (1984, stat.berkeley.edu) developed CART (Classification and Regression Trees) algorithm

One creates a complete binary tree, and then pruning starts,



Decision Tree: The Obama-Clinton Divide



Supervised Techniques : Classification

Random Forrests and Bootstraping Technique

Data : $(\mathbf{x}_i, y_i) = (x_{1,i}, x_{2,i}, y_i), i = 1, \dots, n$

with $y_i \in \{-0, 1\}$ or $y_i \in \{-1, +1\}$ or $y_i \in \{-\bullet, \bullet\}$

Estimate tree on $(\mathbf{x}_i^{*,b}, y_i^{*,b})$ on a bootstraped sample

Estimate $s(\mathbf{x})$ (or $s(\mathbf{x}_i)$ only) $\longrightarrow \hat{s}^b(\mathbf{x})$

and generate (many) other samples.

See Breiman (2001, stat.berkeley.edu)

Supervised Techniques : Classification

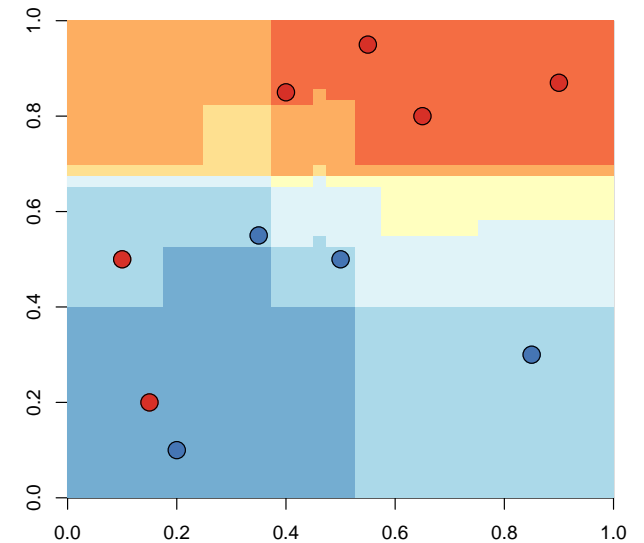
Random Forrests and Aggregation

Data : $(\mathbf{x}_i, y_i) = (x_{1,i}, x_{2,i}, y_i), i = 1, \dots, n$

with $y_i \in \{-0, 1\}$ or $y_i \in \{-1, +1\}$ or $y_i \in \{-\bullet, \bullet\}$

Define

$$\hat{s}(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B \hat{s}^b(\mathbf{x})$$



Supervised Techniques : Double Classification

Uplift Techniques

Data : $-(\mathbf{x}_i, y_i) = (x_{1,i}, x_{2,i}, y_i)$ with $y_i \in \{-\bullet, \bullet\}$

$-(\mathbf{x}_j, y_j) = (x_{1,j}, x_{2,j}, y_j)$ with $y_j \in \{-\blacksquare, \blacksquare\}$

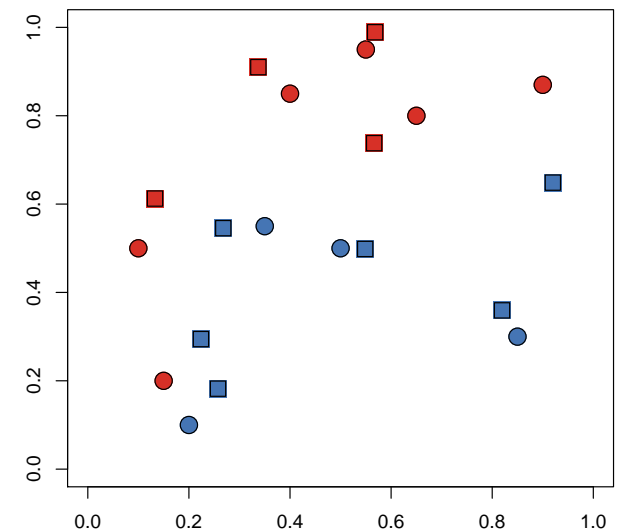
See clinical trials, treatment vs. control group

E.g. direct mail campaign in a bank

	Control ■	Promotion ●
No Purchase	85.17%	61.60%
Purchase	14.83%	38.40%

overall uplift effect +23.57%

(see Guelman *et al.*, 2014 [j.insmatheco.2014.06.009](https://doi.org/10.2139/ssrn.2571009))



Consistency of Models

In supervised models, the **errors** related to the difference between y_i 's and \hat{y}_i 's.

Consider some **loss function** $L(y_i, \hat{y}_i)$, e.g.

- in classification, $\mathbf{1}(y_i \neq \hat{y}_i)$ (misclassification)
- in regression, $(y_i - \hat{y}_i)^2$ (cf least squares, ℓ_2 norm)

Consider some statistical model, $m(\cdot)$, estimated on sample (y_i, \mathbf{x}_i) of size n . Let $\hat{m}_n(\cdot)$ denote that estimator, so that $\hat{y}_i = \hat{m}_n(\mathbf{x}_i)$. $\hat{m}_n(\cdot)$ is a regression function when y is continuous, and is a classifier when y is categorical (say $\{-0, 1\}$).

$m_n(\cdot)$ has been trained from a learning sample (y_i, \mathbf{x}_i) .

Consistency of Models

The **risk** is $R_n = \mathbb{E}(L) = \int L(y, \hat{m}_n(\mathbf{x})) d\mathbb{P}(y, \mathbf{x})$

The **empirical risk** is $\hat{R}_n = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{m}_n(\mathbf{x}_i))$ on a **training sample** $\{-y_i, \mathbf{x}_i\}$.

The **generalized risk** is $\tilde{R}_{n'} = \frac{1}{n'} \sum_{i=1}^{n'} L(\tilde{y}_i, \hat{m}_n(\tilde{\mathbf{x}}_i))$ on a **validation sample** $\{-\tilde{y}_i, \tilde{\mathbf{x}}_i\}$.

We have a **consistent** model if $\hat{R}_n \rightarrow R_n$ as $n \rightarrow \infty$.

Consistency of Models

```

1 > U <- data.frame(X1=runif(n),X2=runif(n))
2 > U$Y <- rbinom(n, size=1,prob=(U[,1]+U[,2]) /
  2)
3 > reg <- glm(Y~X1+X2, data=U, family=binomial)
4 > pd=function(x1,x2) predict(reg, newdata=
  data.frame(X1=x1,X2=x2), type="response")
  >.5
5 > MissClassU <- mean(abs(pd(U$X1,U$X2)-U$Y))

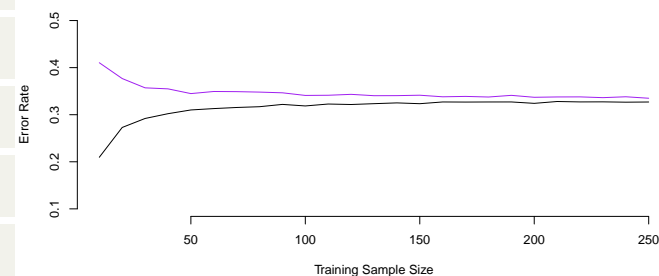
```

Consider some polynomial regression

```

1 > V <- data.frame(X1=runif(n),X2=runif(n))
2 > V$Y <- rbinom(n, size=1,prob=(V[,1]+V[,2]) /
  2)
3 > MissClassV [s <- mean(abs(pd(V$X1,V$X2)-V$Y
  ))

```



Consistency: why is it complicated ?

From the law of large numbers,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n U_i = \mathbb{E}[U]$$

when the U_i 's are i.i.d., and $U_i \sim U$. Here we look for

$$\lim_{n \rightarrow \infty} \underbrace{\frac{1}{n} \sum_{i=1}^n L(y_i, \hat{m}_n(\mathbf{x}_i))}_{R_n} = ?$$

If (y, \mathbf{x}_i) 's are independent, it is not the case for the $L(y_i, \hat{m}_n(\mathbf{x}_i))$'s, because of $\hat{m}_n(\cdot)$, also estimated on this training sample. But, on the validation sample

$$\lim_{n' \rightarrow \infty} \underbrace{\frac{1}{n'} \sum_{i=1}^{n'} L(\tilde{y}_i, \hat{m}_n(\tilde{\mathbf{x}}_i))}_{\tilde{R}_n} = \mathbb{E}[R(L(Y, m(\mathbf{X})))]$$

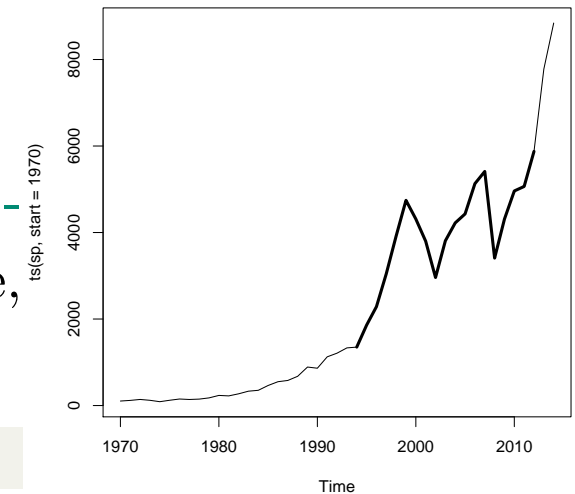
Overfitting

Inspired by Leinweber (1995, nerdsonwall-street.typepad.com). Consider S&P 500 yearly value, en.wikipedia.org

```
1 > sp <- read.table("sp.txt", header=TRUE)$SP500
2 > plot(ts(sp, start=1970))
3 > T <- 1994:2012
4 > X <- sp[T-1969]
```

Consider some polynomial regression

```
1 > df <- data.frame(X,T)
2 > reg <- lm(X~poly(T,k), data=df)
3 > tp <- seq(1990,2015,by=.1)
4 > yp <- predict(reg, newdata=data.frame(T=tp))
```



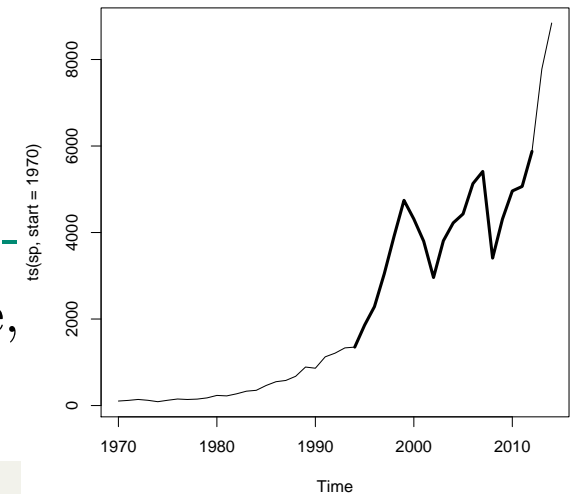
Overfitting

Inspired by Leinweber (1995, nerdsonwall-street.typepad.com). Consider S&P 500 yearly value, en.wikipedia.org

```
1 > sp <- read.table("sp.txt", header=TRUE)$SP500
2 > plot(ts(sp, start=1970))
3 > library(splines)
```

Consider some [spline regression](#)

```
1 > df <- data.frame(X,T)
2 > reg <- lm(X~bs(T,k), data=df)
3 > tp <- seq(1990,2015, by=.1)
4 > yp <- predict(reg, newdata=data.frame(T=tp))
```



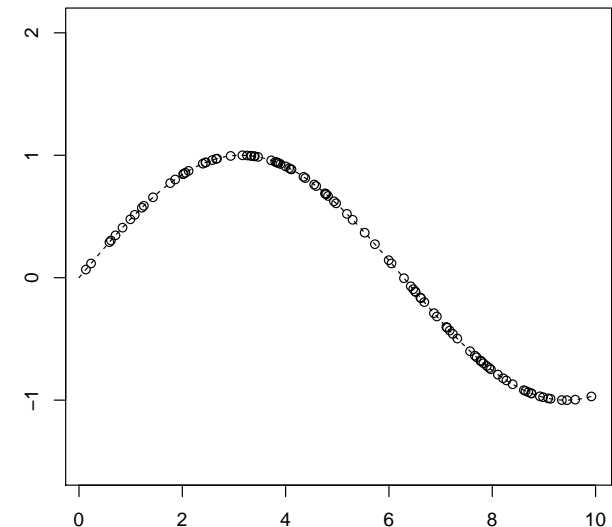
Under- and Over-fitting

Consider some simulated data

```
1 > set.seed(1)
2 > n <- 100
3 > x <- runif(n)*10
4 > y0 <- sin(x/2)
5 > y <- y0 + rnorm(length(x))/2
```

A (too) perfect sample (with no noise) would be

```
1 > plot(x, y0)
```



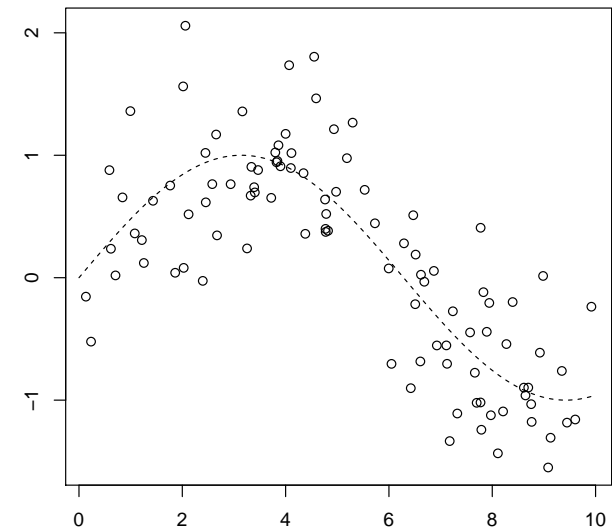
Under- and Over-fitting

Consider some simulated data

```
1 > set.seed(1)
2 > n <- 100
3 > x <- runif(n)*10
4 > y0 <- sin(x/2)
5 > y <- y0 + rnorm(length(x))/2
```

while our sample with some noise is

```
1 > plot(x, y)
```



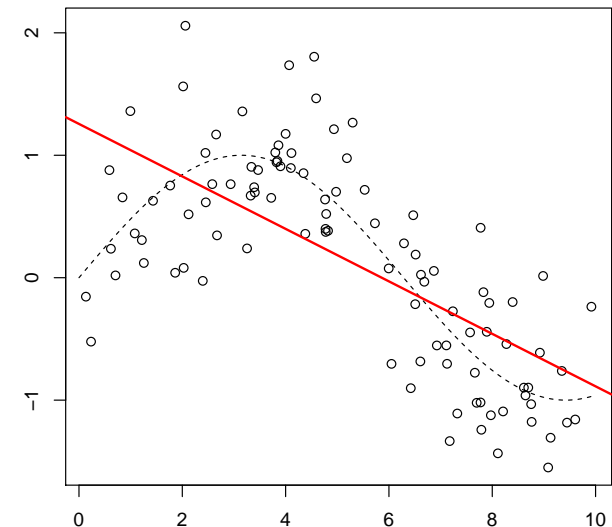
Underfitting

Consider some simulated data

```
1 > set.seed(1)
2 > n <- 100
3 > x <- runif(n)*10
4 > y0 <- sin(x/2)
5 > y <- y0 + rnorm(length(x))/2
```

An **underfitted model** is the linear model

```
1 > reg <- lm(y~x, data=db)
2 > plot(x, predict(reg))
```



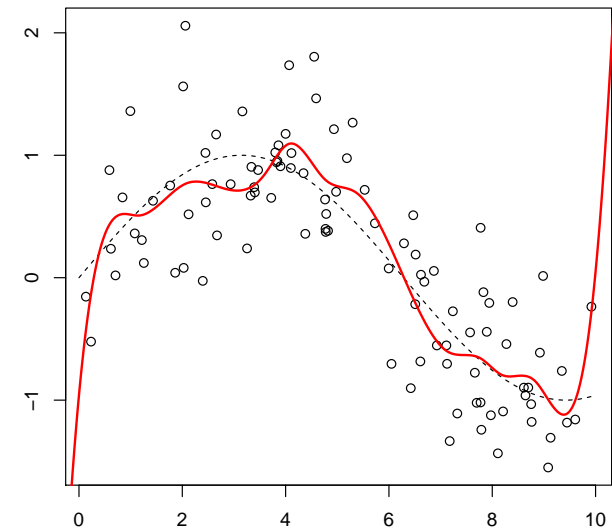
Overfitting

Consider some simulated data

```
1 > set.seed(1)
2 > n <- 100
3 > x <- runif(n)*10
4 > y0 <- sin(x/2)
5 > y <- y0 + rnorm(length(x))/2
```

while an **overfitted model** is

```
1 > library(splines)
2 > reg <- lm(y~bs(x,15), data=db)
3 > plot(x, predict(reg))
```



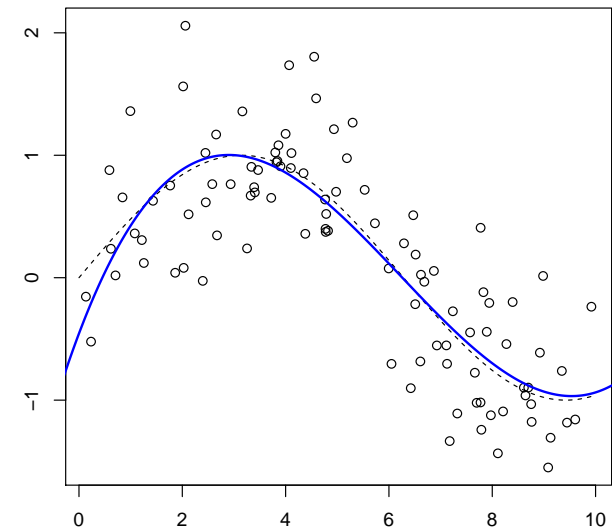
The Goldilocks Principle

Consider some simulated data

```
1 > set.seed(1)
2 > n <- 100
3 > x <- runif(n)*10
4 > y0 <- sin(x/2)
5 > y <- y0 + rnorm(length(x))/2
```

(too) perfect model

```
1 > library(splines)
2 > reg <- lm(y~bs(x), data=db)
3 > plot(x, predict(reg))
```



Overfitting and Consistency

There is a gap between R_n (on the training sample) and \hat{R}_n (on the validation sample). One can prove that

$$\hat{R}_n \leq R_n + \sqrt{\frac{\text{VC}[\log(2n/d) + 1] - \log[\alpha/4]}{n}}$$

with probability $1 - \alpha$, where **VC** denotes the Vapnik-Chervonenkis dimension.

A process is consistent if and only if $\text{VC} < \infty$.

VC can be seen as a function of complexity of some models.

Consider polynomial regressors, of degree d , with k covariates,

- if $k = 1$ then $\text{VC} = d + 1$

- if $k = 2$ then $\text{VC} = \frac{(d+1)(d+2)}{2}$ (bivariate polynomials)

if $k = 2$ then $\text{VC} = 2(d+1)$ (additive model, sum of (univariate) polynomials)

```

1 > U <- data.frame(X1=runif(n),X2=runif(n))
2 > U$Y <- rbinom(n, size=1,prob=(U[,1]+U[,2]) /
  2)
3 > reg <- glm(Y~poly(X1,s)+poly(X2,s), data=U,
4 family=binomial)
5 > pd=function(x1,x2) predict(reg, newdata=
  data.frame(X1=x1,X2=x2), type="response")
  >.5
6 > MissClassU <- mean(abs(pd(U$X1,U$X2)-U$Y))

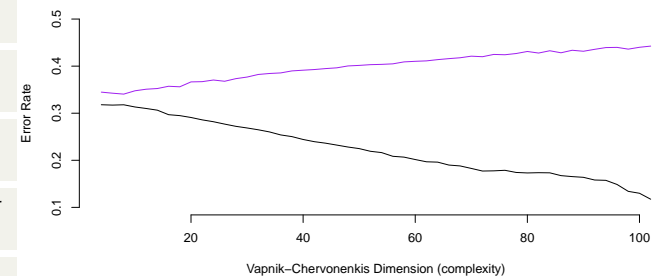
```

Consider some polynomial regression

```

1 > V <- data.frame(X1=runif(n),X2=runif(n))
2 > V$Y <- rbinom(n, size=1,prob=(V[,1]+V[,2]) /
  2)
3 > MissClassV[s <- mean(abs(pd(V$X1,V$X2)-V$Y
  ))

```



Regression vs. Statistical Learning

In a linear regression, $y = \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon$. Parameter $\boldsymbol{\beta}$ is estimated from observations (y_i, \mathbf{x}_i) 's.

The estimation is consistent if $\hat{\boldsymbol{\beta}}_n \rightarrow \boldsymbol{\beta}$, as $n \rightarrow \infty$.

In a regression, $y = m(\mathbf{x}) + \varepsilon$. Function $m(\cdot)$ is estimated from observations (y_i, \mathbf{x}_i) 's.

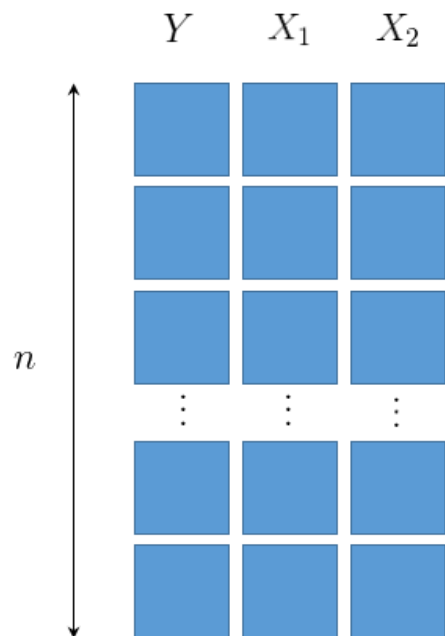
The estimation is consistent if $\hat{m}_n(\cdot) \rightarrow m(\cdot)$, as $n \rightarrow \infty$.

In the [statistical learning perspective](#), the interest is not really on $m(\cdot)$, but on the prediction error. $\hat{m}_n(\cdot)$ is said to be consistent if, given a loss function L

$$\mathbb{E}[L(\hat{m}_n(\mathbf{X}), Y)] \rightarrow \mathbb{E}[L(m(\mathbf{X}), Y)]$$

as $n \rightarrow \infty$.

Data Science, in 2015



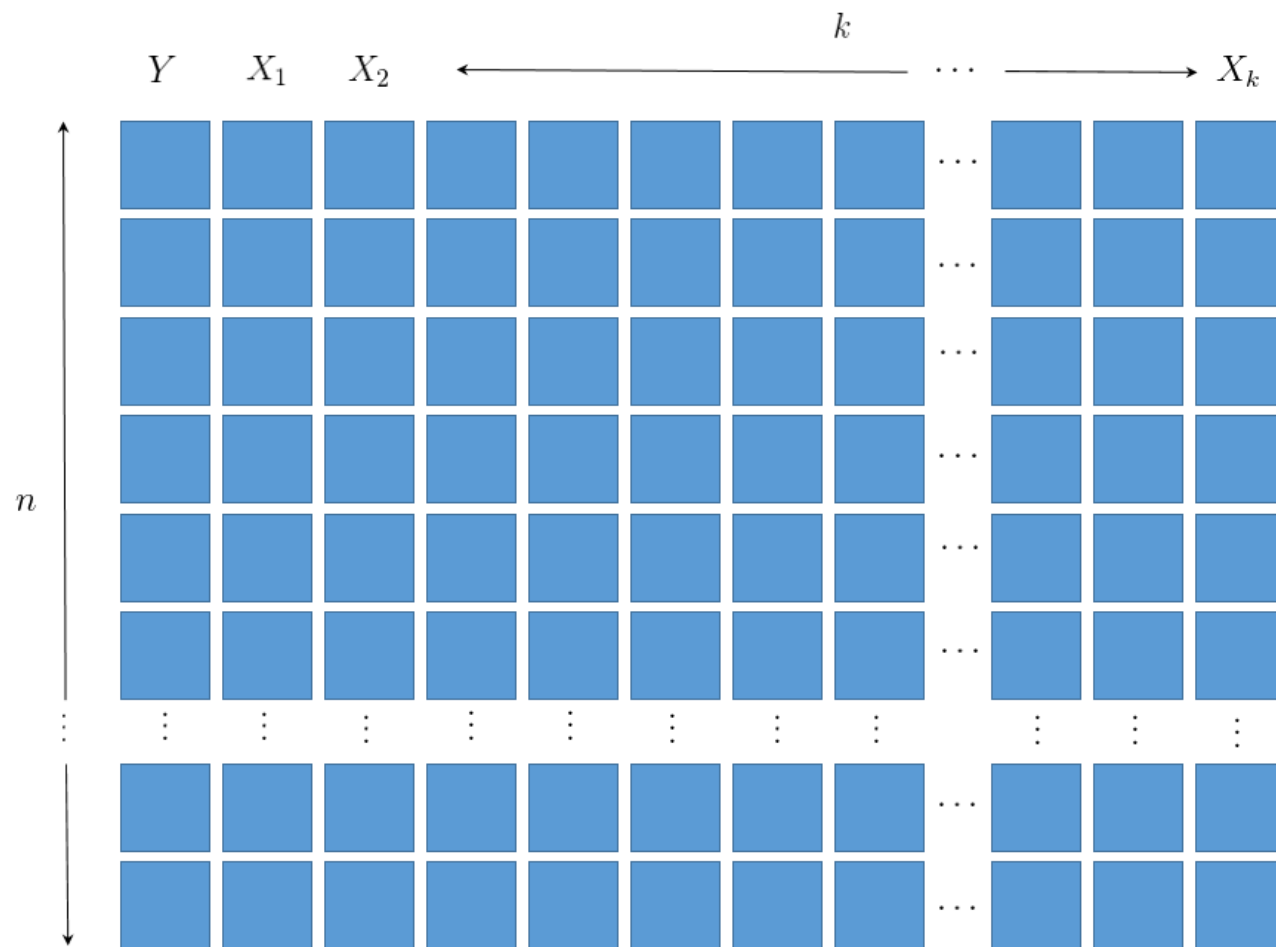
Historically, we had limited information, (Y_i, \mathbf{X}_i) , the goal was to get n as large as possible (asymptotic theory).

$\mathbb{E}(Y|X_1 = x_1) = m(x_1)$ regression **line**

$\mathbb{E}(Y|X_1 = x_1, X_2 = x_2) = m(x_1, x_2)$ regression **surface**

$\mathbb{E}(Y|\mathbf{X} = \mathbf{x}) = m(\mathbf{x})$ regression **function**

Data Science, in 2015



n can be large, but limited

(portfolio size)

large variety k

large volume nk

→ Feature Engineering

Regression?

Galton (1870, galton.org, 1886, galton.org) and Pearson & Lee (1896, jstor.org, 1903 jstor.org) studied genetic transmission of characteristics, e.g. the height.

On average the child of tall parents is taller than other children, but less than his parents.

“I have called this peculiarity by the name of regression’, Francis Galton, 1886.

REGRESSION *towards* MEDIOCRITY in HEREDITARY STATURE. By FRANCIS GALTON, F.R.S., &c.

Table 8.1. Galton's 1885 cross-tabulation of 928 adult children born of 205 midparents, by their height and their midparent's height.

Height of the mid-parent in inches	Height of the adult child											Total no. of adult children	Total no. of mid-parents	Medians				
	<61.7	62.2	63.2	64.2	65.2	66.2	67.2	68.2	69.2	70.2	71.2				72.2	73.2	>73.7	
> 73.0	—	—	—	—	—	—	—	—	—	—	—	1	3	—	4	5	—	
72.5	—	—	—	—	—	—	—	—	—	1	2	1	2	2	4	19	6	
71.5	—	—	—	—	1	3	4	5	5	10	4	9	2	2	43	111	69.8	
70.5	1	—	1	—	1	1	3	12	18	14	7	4	3	3	68	22	69.5	
69.5	—	—	1	16	4	17	27	20	33	25	20	11	4	5	183	41	68.9	
68.5	1	—	7	11	16	25	31	34	48	21	18	4	3	—	219	49	68.2	
67.5	—	3	5	13	15	36	58	28	38	19	11	4	—	—	211	33	67.6	
66.5	—	3	5	2	17	17	14	13	4	—	—	—	—	—	78	20	67.2	
65.5	1	—	9	5	7	11	11	7	7	5	2	1	—	—	66	12	66.7	
64.5	1	1	4	4	1	5	5	—	2	—	—	—	—	—	25	5	65.8	
<64.0	1	—	2	4	1	2	2	1	1	—	—	—	—	—	14	1	—	
Totals	5	7	32	59	48	117	158	120	167	99	64	41	17	14	928	205	—	
Medians	—	—	66.3	67.8	67.9	67.7	67.9	68.3	68.5	69.0	69.0	70.0	—	—	—	—	—	—

Source: Galton (1886).
 Note: All female heights were multiplied by 1.08 before tabulation. Galton added an explanatory footnote to the table: "In calculating the Medians, the entries have been taken as referring to the middle of the squares in which they stand. The reason why the headings run 62.2, 63.2, &c., instead of 62.5, 63.5, &c., is that the observations are unequally distributed between 62 and 63, 63 and 64, &c., there being a strong bias in favour of integral inches. After careful consideration, I concluded that the headings, as adopted, best satisfied the conditions. This inequality was not apparent in the case of the Mid-parents." Galton republished these data in 1889, where they are referred to as the R.F.F. Data (Record of Family Faculties); he then noted that the first row must be in error (four children cannot have five sets of parents), but he claimed that "the bottom line, which looks suspicious is correct" (p. 208).

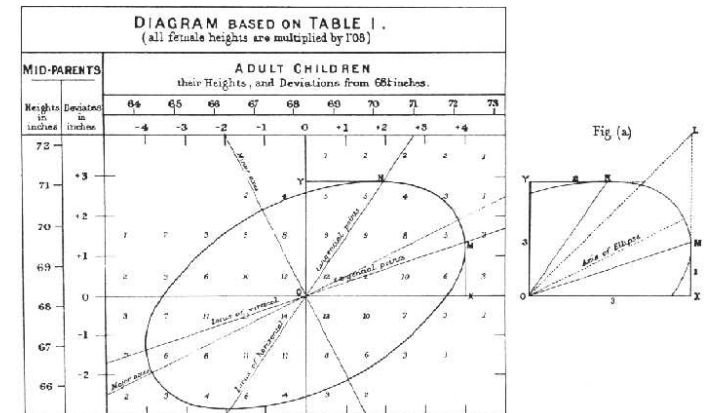
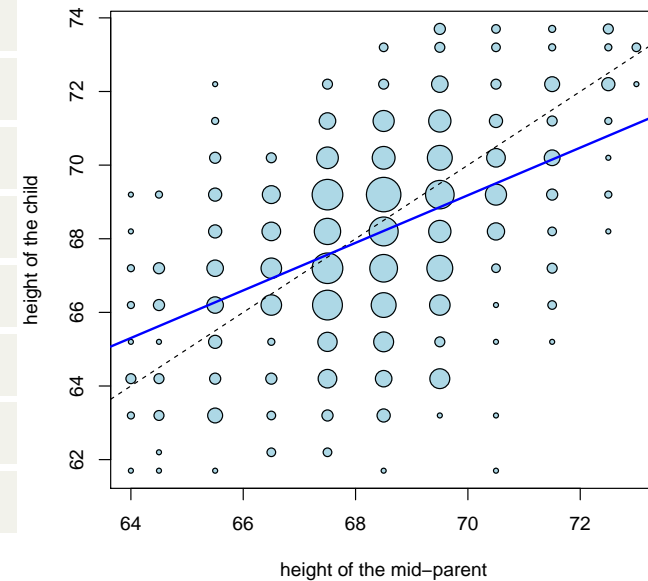


Figure 8.7. Galton's smoothed rendition of Table 8.1, with one of the "concentric and similar ellipses" drawn in. The geometric relationship of the two regression lines to the ellipse is also shown. (From Galton, 1886a.)

Regression?

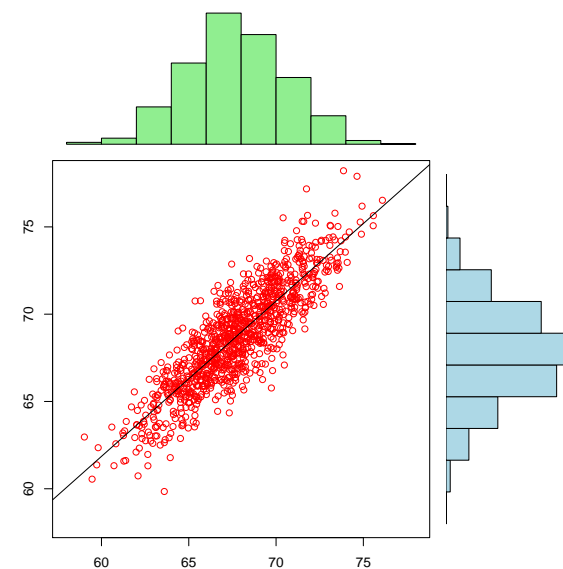
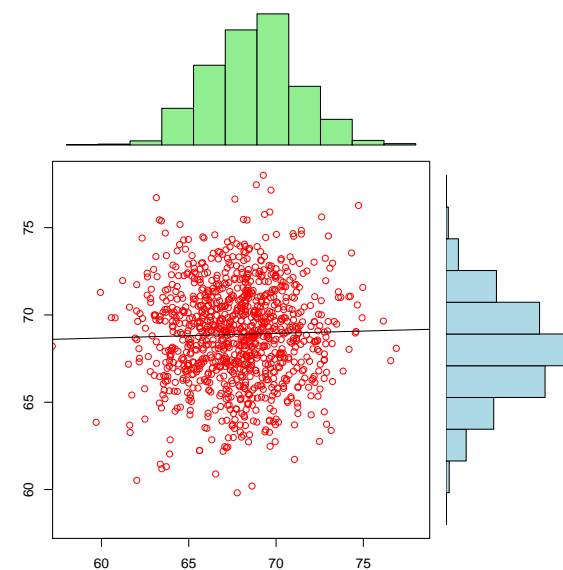
```
1 > library(HistData)
2 > attach(Galton)
3 > Galton$count <- 1
4 > df <- aggregate(Galton, by=list(parent,
  child), FUN=sum)[,c(1,2,5)]
5 > plot(df[,1:2], cex=sqrt(df[,3]/3))
6 > abline(a=0,b=1,lty=2)
7 > abline(lm(child~parent, data=Galton))
```



Regression?

Regression is a **correlation** problem

Overall, children are not smaller than parents



Least Squares?

Recall that

$$\begin{aligned} \mathbb{E}(Y) &= \operatorname{argmin}_{m \in \mathbb{R}} -\|Y - m\|_{\ell_2}^2 = \mathbb{E}([Y - m]^2) \text{ " } \\ \text{---} \\ \operatorname{Var}(Y) &= \min_{m \in \mathbb{R}} -\mathbb{E}([Y - m]^2) \text{ " } = \mathbb{E}([Y - \mathbb{E}(Y)]^2) \end{aligned}$$

The empirical version is

$$\begin{aligned} \bar{y} &= \operatorname{argmin}_{m \in \mathbb{R}} -\sum_{i=1}^n \frac{1}{n} [y_i - m]^2 \text{ " } \\ \text{---} \\ s^2 &= \min_{m \in \mathbb{R}} -\sum_{i=1}^n \frac{1}{n} [y_i - m]^2 \text{ " } = \sum_{i=1}^n \frac{1}{n} [y_i - \bar{y}]^2 \end{aligned}$$

The conditional version is

$$\begin{aligned} \mathbb{E}(Y|\mathbf{X}) &= \operatorname{argmin}_{\varphi: \mathbb{R}^k \rightarrow \mathbb{R}} -\|Y - \varphi(\mathbf{X})\|_{\ell_2}^2 = \mathbb{E}([Y - \varphi(\mathbf{X})]^2) \text{ " } \\ \text{---} \\ \operatorname{Var}(Y|\mathbf{X}) &= \min_{\varphi: \mathbb{R}^k \rightarrow \mathbb{R}} -\mathbb{E}([Y - \varphi(\mathbf{X})]^2) \text{ " } = \mathbb{E}([Y - \mathbb{E}(Y|\mathbf{X})]^2) \end{aligned}$$

Errors in Regression type Models

In predictions, there are two kinds of errors

- × error on the best estimate, \hat{Y}
- × error on the true value, Y

Recall that

$$Y = \underbrace{\mathbf{x}^\top \boldsymbol{\beta}}_{\text{model}} + \underbrace{\varepsilon}_{\text{error}} = \underbrace{\mathbf{x}^\top \hat{\boldsymbol{\beta}}}_{\text{prediction } \hat{Y}} + \hat{\varepsilon}.$$

- × error on the best estimate, $\text{Var}(\hat{Y} | \mathbf{X} = \mathbf{x}) = \mathbf{x}^\top \text{Var}(\hat{\boldsymbol{\beta}}) \mathbf{x}$ (inference error)
- × error on the true value, $\text{Var}(Y | \mathbf{X} = \mathbf{x}) = \text{Var}(\varepsilon)$ (model error)


⚠ asymptotically (under suitable conditions), $\text{Var}(\hat{\boldsymbol{\beta}}_n) \rightarrow \mathbf{0}$ as $n \rightarrow \infty$.

Errors in Regression type Models

Under the Gaussian assumption, we can derive (approximated) confidence intervals. For the best estimate,

$$\left[\hat{Y} \pm 1.96\hat{\sigma} \sqrt{\mathbf{x}^\top [\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{x}} \right]$$

```
1 > predict(lm(Y~X, data=df), newdata=data.frame(X=x), interval='
confidence')
```

 this confidence interval is a statement about estimates

For the 'true value'

$$\left[\hat{Y} \pm 1.96\hat{\sigma} \right]$$

```
1 > predict(lm(Y~X, data=df), newdata=data.frame(X=x), interval='
prediction')
```

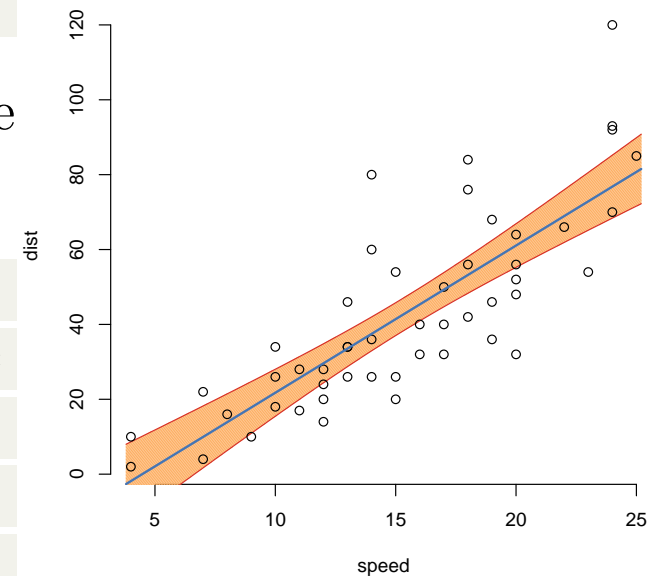
Resampling Techniques in Regression type Models

Consider some linear model

```
1 > plot(cars)
2 > reg <- lm(dist ~ speed, data=cars)
```

then if we plot the prediction with some confidence band

```
1 > u <- seq(3.8, 25.2, length=251)
2 > v <- predict(reg, newdata=data.frame(speed=
  u), interval="confidence")
3 > polygon(c(u, rev(u)), c(v[,2], rev(v[,3])),
  border=NA)
4 > lines(u, v[,1], lwd=2)
```



Resampling Techniques in Regression type Models

Resampling techniques can be used to generate a confidence interval.

1. Draw pairs from the sample

Resample from $-(\mathbf{X}_i, Y_i)$ "

```

1 > V=matrix(NA,100,251)
2 > for(i in 1:100){
3 + ind <- sample(1:n, size=n, replace=TRUE)
4 + V[i,] <- predict(lm(dist~speed, data=cars [
      ind, ]), newdata=data.frame(speed=u)) }

```

Then we can derive pointwise confidence intervals.

Resampling Techniques in Regression type Models

Confidence intervals are obtained using

```
1 > q1 <- apply(V,2,function(x) quantile(x
    ,.05))
2 > q2 <- apply(V,2,function(x) quantile(x
    ,.95))
```

Those values can be compared with the second and the third column of

```
1 > v <- predict(reg,newdata=data.frame(speed=u),interval="confidence")
```

But an alternative can be to use residuals obtained from $Y_i = \underbrace{\mathbf{X}_i^T \hat{\boldsymbol{\beta}}}_{\hat{Y}_i} + \hat{\varepsilon}_i$.

Resampling Techniques in Regression type Models

2. Draw pseudo-observations using (estimated) residuals

Resample from $-(\mathbf{X}_i, \hat{Y}_i + \varepsilon)$, where ε can be obtained by resampling from (estimated) residuals

```

1 > V=matrix(NA,100,251)
2 > for(i in 1:100){
3   + ind <- sample(1:50,replace=TRUE)
4   + cars2 <- data.frame(dist=predict(reg),
5     speed=cars$speed)
6   + cars_b <- cars2[ind,]
7   + cars_b$dist <- cars_b$dist+residuals(reg)[
8     sample(1:50,replace=TRUE)]
9   + V[i,] <- predict(lm(dist~speed,data=cars_b
10    ),newdata=data.frame(speed=u))}

```

Resampling Techniques in Regression type Models

Here again, pointwise confidence intervals can be obtained

```
1 > q1 <- apply(V,2,function(x) quantile(x  
    ,.05))
```

```
2 > q2 <- apply(V,2,function(x) quantile(x  
    ,.95))
```

Those values can be compared with the second and the third column of

```
1 > v <- predict(reg,newdata=data.frame(speed=u),interval="confidence")
```

Changing the Distance in Least-Squares?

One might consider $\hat{\beta} \in \operatorname{argmin} - \sum_{i=1}^n |Y_i - \mathbf{X}_i^\top \beta|$ ", based on the ℓ_1 -norm, and not the ℓ_2 -norm.

This is the **least-absolute deviation** estimator, related to the **median regression**, since $\operatorname{median}(X) = \operatorname{argmin} - \mathbb{E}|X - x|$ ".

💡 this regression model was introduced 50 years before Gauss and Legendre ℓ_2 -regression, by Boscovich.

⚠️ no analytical expression, since the ℓ_1 -norm is not differentiable (first order condition).

More generally, assume that, for some function $R(\cdot)$,

$$\hat{\beta} \in \operatorname{argmin} - \sum_{i=1}^n R(Y_i - \mathbf{X}_i^\top \beta) "$$

Changing the Distance in Least-Squares?


If R is differentiable, the first order condition would be

$$\sum_{i=1}^n R' \left(Y_i - \mathbf{X}_i^\top \boldsymbol{\beta} \right) \cdot \mathbf{X}_i^\top = 0.$$

i.e.

$$\sum_{i=1}^n \underbrace{\omega \left(Y_i - \mathbf{X}_i^\top \boldsymbol{\beta} \right)}_{\omega_i} \cdot \left(Y_i - \mathbf{X}_i^\top \boldsymbol{\beta} \right) \mathbf{X}_i^\top = 0 \text{ with } \omega(x) = \frac{R'(x)}{x},$$

It is the first order condition of a **weighted ℓ_2 regression**.

 weights are unknown, since they are function of residuals...

To obtain the ℓ_1 -regression, observe that $\omega = |\varepsilon|^{-1}$

Changing the Distance in Least-Squares?

💡 use iterative (weighted) least-square regressions.

Start with some standard ℓ_2 regression

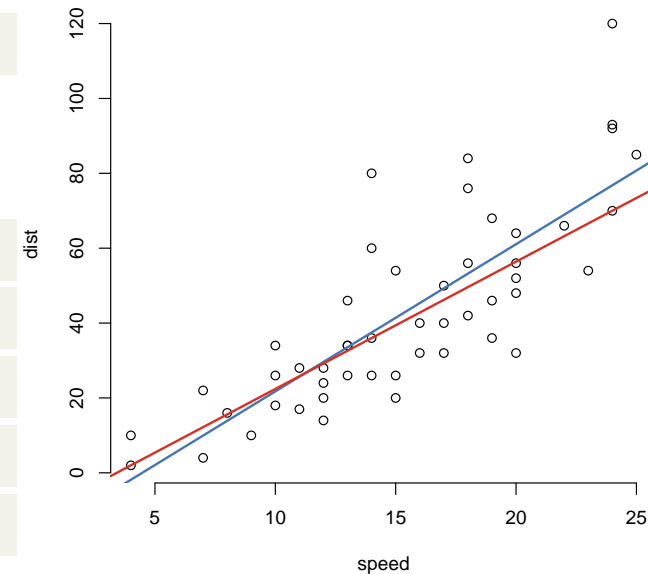
```
1 > reg_0 <- lm(Y~X, data=db)
```

For the ℓ_1 regression consider weight function

```
1 > omega <- function(e) 1/abs(e)
```

Then consider the following iterative algorithm

```
1 > resid <- residuals(reg_0)
2 > for(i in 1:100){
3 + W <- omega(e)
4 + reg <- lm(Y~X, data=db, weights=W)
5 + e <- residuals(reg)}
```



Under- and Over-Identified Models

Suppose that the true model is


$$Y_i = \mathbf{X}_{1,i}\boldsymbol{\beta}_1 + \mathbf{X}_{2,i}\boldsymbol{\beta}_2 + \varepsilon_i,$$

but we estimate the model on \mathbf{X}_1 (only)

$$Y_i = \mathbf{X}_{1,i}\mathbf{b}_1 + \eta_i.$$

$$\begin{aligned} \hat{\mathbf{b}}_1 &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top Y \\ &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top [\mathbf{X}_{1,i}\boldsymbol{\beta}_1 + \mathbf{X}_{2,i}\boldsymbol{\beta}_2 + \varepsilon] \\ &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_1 \boldsymbol{\beta}_1 + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \boldsymbol{\beta}_2 + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \varepsilon \\ &= \boldsymbol{\beta}_1 + \underbrace{(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \boldsymbol{\beta}_2}_{\boldsymbol{\beta}_{12}} + \underbrace{(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \varepsilon_i}_{\nu_i} \end{aligned}$$

i.e. $\mathbb{E}(\hat{\boldsymbol{\alpha}}_1) = \boldsymbol{\beta}_1 + \boldsymbol{\beta}_{12}$.

 if $\mathbf{X}_1^\top \mathbf{X}_2 = \mathbf{0}$ ($\mathbf{X}_1 \perp \mathbf{X}_2$), $\mathbb{E}(\hat{\mathbf{b}}_1) = \boldsymbol{\beta}_1$

Under-Identified Models, Simpson's paradox

Computer $\mathbb{E}(Y|X_1)$ where Y is the hospital

hospital	total	survivors	dead	survival rate
hospital A	1,000	800	200	80%
hospital B	1,000	900	100	90% *

For healthy people $\mathbb{E}(Y|X_1, X_2 = \text{healthy})$

hospital	total	survivors	dead	survival rate
hospital A	600	590	10	98% *
hospital B	900	870	30	97%

For sick people $\mathbb{E}(Y|X_1, X_2 = \text{sick})$

hospital	total	survivors	dead	survival rate
hospital A	400	210	190	53% *
hospital B	100	30	70	30%

Under- and Over-Identified Models

Conversely, assume that the true model is

$$Y = \mathbf{X}_1 \boldsymbol{\beta}_1 + \varepsilon$$

but we estimated on (unnecessary) variables \mathbf{X}_2

$$Y = \mathbf{X}_1 \mathbf{b}_1 + \mathbf{X}_2 \mathbf{b}_2 + \eta$$

Here estimation is unbiased $\mathbb{E}(\mathbf{b}_1) = \boldsymbol{\beta}_1$, but the estimator is not efficient...

Beyond Linear Models: Polynomial

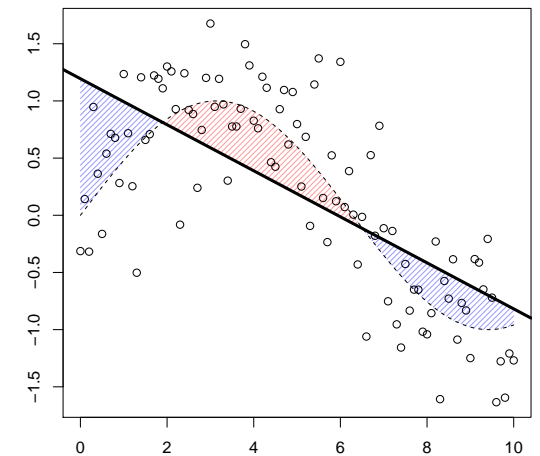
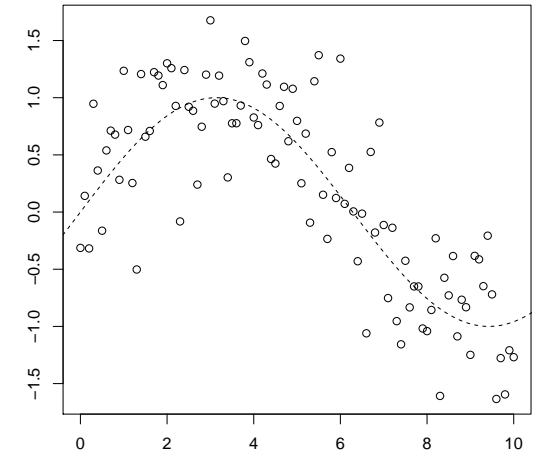
Even if $m(x) = \mathbb{E}(Y|X = x)$ is not a polynomial function, a polynomial can still be a good approximation.

From Stone-Weierstrass theorem, if $m(\cdot)$ is continuous on some interval, then there is a uniform approximation of $m(\cdot)$ by polynomial functions.

```

1 set.seed(1)
2 n=10
3 xr <- seq(0, n, by=.1)
4 yr <- sin(xr/2) + rnorm(length(xr))/2
5 db <- data.frame(x=xr, y=yr)
6 plot(db)
7 lines(xr, sin(xr/2), lty=2)

```



Beyond Linear Models: Polynomial

Assume that $m(x) = \mathbb{E}(Y|X = x) = \sum_{i=0}^k \alpha_i x^i$, where parameters $\alpha_0, \dots, \alpha_k$ will be estimated (but not k).

```
1 reg <- lm(yr~poly(xr,5),data=db)
```

```
2 reg <- lm(yr~poly(xr,25),data=db)
```

Beyond Linear Models: Local Regression

💡 (local) regression on the subset

$$-i, |X_i - x| \leq h''$$

```

1 local_reg <- function(x, h) {
2   ind <- which(abs(X-x) < h)
3   reg_loc <- lm(Y~X, data=db[ind,])
4   return(predict(reg_loc, newdata=data.frame(X=x))) }

```

OR

```

1 local_reg <- function(x, h) {
2   W <- abs(X-x) < h
3   reg_loc <- lm(Y~X, data=db, weights=W)
4   return(predict(reg_loc, newdata=data.frame(X=x))) }

```

Beyond Linear Models: Local Regression

💡 why not run a **weighted** linear regression with

$$\omega_i = K_h(X_i - x) = \frac{1}{h} K\left(\frac{X_i - x}{h}\right)$$

for some kernel $K(\cdot)$.

```

1 local_reg <- function(x, h, K) {
2   W <- K((X-x)/h)/h
3   reg_loc <- lm(Y~X, data=db, weights=W)
4   return(predict(reg_loc, newdata=data.frame(X=x))) }

```

One can also consider regression on the k -nearest neighbors

```

1 local_reg <- function(x, k) {
2   D <- rank(abs(X-x))
3   reg_loc <- lm(Y~X, data=db[which(D<=k),])
4   return(predict(reg_loc, newdata=data.frame(X=x))) }

```

Beyond Linear Models: Using Splines

💡 for linear splines, consider

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 (X_i - s)_+ + \varepsilon_i$$

```
1 positive_part <- function(x) ifelse(x>0,x,0)
2 reg <- lm(Y~X+positive_part(X-s), data=db)
```

Beyond Linear Models: Using Splines

💡 for linear splines, consider

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 (X_i - s_1)_+ + \beta_3 (X_i - s_2)_+ + \varepsilon_i$$

```
1 reg <- lm(Y~X+positive_part(X-s1)+positive_part(X-  
s2) , data=db)
```

Beyond Linear Models: Using Splines


⚠ the model estimated by R is an equivalent one, see

```
1 > library(splines)
2 > B1_1 <- bs(x, knots=4, Boundary.Knots=range(x),
3   degree=1)
4 > B2_1 <- bs(x, knots=c(2,5), Boundary.Knots=range(x),
5   degree=1)
6 > B1_2 <- bs(x, knots=4, Boundary.Knots=range(x),
7   degree=2)
8 > B2_2 <- bs(x, knots=c(2,5), Boundary.Knots=range(x),
9   degree=2)
```

Smoothing Techniques, a Unified Approach

In the linear model,

$$\hat{Y} = \mathbf{X}\hat{\beta} = \underbrace{\mathbf{X}[\mathbf{X}^T\mathbf{X}]^{-1}\mathbf{X}^T}_{\mathbf{H}}\mathbf{Y}$$

 $H_{i,i}$ is the leverage of the i th element of this hat matrix.

Write

$$\hat{Y}_i = \sum_{j=1}^n [\mathbf{X}_i^T [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T]_j Y_j = \sum_{j=1}^n [\mathcal{H}(\mathbf{X}_i)]_j Y_j$$

where

$$\mathcal{H}(\mathbf{x}) = \mathbf{x}^T [\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T$$

The prediction is

$$m(\mathbf{x}) = \mathbb{E}(Y | \mathbf{X} = \mathbf{x}) = \sum_{j=1}^n [\mathcal{H}(\mathbf{x})]_j Y_j$$

Smoothing Techniques, a Unified Approach

More generally, a predictor m is said to be linear if for all \mathbf{x} if there is $\mathcal{S}(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ such that

$$m(\mathbf{x}) = \sum_{j=1}^n \mathcal{S}(\mathbf{x})_j Y_j$$

Conversely, given $\hat{Y}_1, \dots, \hat{Y}_n$, there is a matrix \mathbf{S} $n \times n$ such that

$$\hat{\mathbf{Y}} = \mathbf{S}\mathbf{Y}$$

For the linear model, $\mathbf{S} = \mathbf{H}$.

⚠ trace(\mathbf{H}) = dim(β): degrees of freedom

⚠ $\frac{\mathbf{H}_{i,i}}{1 - \mathbf{H}_{i,i}}$ is related to Cook's distance (detect outliers)

Smoothing Techniques, a Unified Approach

For a kernel regression model, with kernel K and bandwidth h

$$S_{i,j}^{(K,h)} = \frac{K_h(X_i - X_j)}{\sum_{k=1}^n K_h(X_k - X_j)}$$

where $K_h(\cdot) = K(\cdot/h)$, while

$$S^{(K,h)}(\mathbf{x})_j = \frac{K_h(\mathbf{x} - X_j)}{\sum_{k=1}^n K_h(X_k - \mathbf{x})}$$

Smoothing Techniques, a Unified Approach

For a k -nearest neighbor,

$$S_{i,j}^{(k)} = \frac{1}{k} \mathbf{1}(j \in \mathcal{I}_{\mathbf{X}_i})$$

where $\mathcal{I}_{\mathbf{X}_i}$ are the k nearest observations to \mathbf{X}_i , while

$$\mathcal{S}^{(K,h)}(\mathbf{x})_j = \frac{1}{k} \mathbf{1}(j \in \mathcal{I}_{\mathbf{x}})$$

Testing for Linearity

Observe that $\text{trace}(S)$ is usually seen as a degree of smoothness.

Do we have to smooth? Isn't linear model sufficient?

Define

$$T = \frac{\|SY - HY\|}{\text{trace}([S - H]^T[S - H])}$$

If the model is linear, then T has a Fisher distribution.

On Optimal Smoothing

- × local regression and kernel regression: choice of bandwidth h
- × splines: choice of number and values of knots s_1, \dots, s_k

$$R_h = \mathbb{E}([Y - \hat{m}_h(\mathbf{X})]^2) + \sigma^2 + \underbrace{\mathbb{E}([m(\mathbf{x}) - \hat{m}_h(\mathbf{x})]^2)}_{\text{Mean Squared Error}}$$

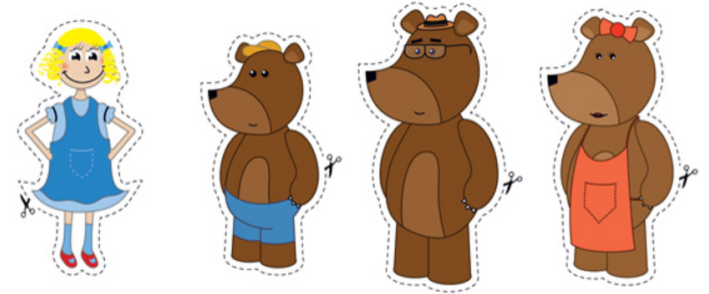
We can rewrite the Mean Squared Error as

$$MSE_h(\mathbf{x}) = \int \text{bias}(\hat{m}_h(\mathbf{x}))^2 d\mathbf{x} + \int \text{Var}(\hat{m}_h(\mathbf{x})) d\mathbf{x}$$

where $\text{bias}(\hat{m}_h(\mathbf{x})) = \mathbb{E}(\hat{m}_h(\mathbf{x}) - m(\mathbf{x}))$ and $\text{Var}(\hat{m}_h(\mathbf{x})) = \text{Var}(\hat{m}_h(\mathbf{x}))$.

Usual **bias-variance tradeoff**, or Goldilock principle:
 h should be neither too small, nor too large

- × undersmoothed: bias too large, variance too small
- × oversmoothed: variance too large, bias too small



The empirical version is

$$\hat{R}_h = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{m}_h(\mathbf{X}_i))^2$$

Use Leave-one-out Cross Validation

$$\hat{R}_h = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{m}_{h(-i)}(\mathbf{X}_i))^2$$

where $\hat{m}_{h(-i)}(\cdot)$ is the estimator obtained by omitting the i th observation.

Optimal Smoothing, Nadaraya-Watson Kernel estimator

Use a rectangular kernel $K(x) = \frac{1}{2} \mathbf{1}(|x| \leq 1)$ or a Gaussian one

$$K(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}.$$

💡 the choice of the kernel is not (that) important,
but the bandwidth is much more important.

× Understanding the Variance and Bias Tradeoff

Let us generate new samples, to understand this tradeoff

```
1 > nsim <- 500
2 > mat_sample_y <- matrix(sin(xr/2) + rnorm(n*nsim)/2, n, nsim)
```

For each x , we now have 500 $\hat{m}_s(x)$

```
1 vec_mh=function(x, h) {
2   sapply(1:nsim, function(s) {
```

```

3 db_s=data.frame(yr=mat_sample_y[,s],xr=xr)
4 reg_s = lm(yr ~ 1,weights=dnorm(xr,mean=x,sd=h),data=db_s)
5 predict(reg_s,newdata=data.frame(xr=x))})}

```

Based on those values, we can approximate $\int \text{bias}^2(\hat{m}_h(\mathbf{x}))d\mathbf{x}$

```

1 bias2_h=function(h) mean(Vectorize(function(x) mean((vec_mh(x,h)-m(x)
)) ^2))(xr))

```

and $\int \text{Var}(\hat{m}_h(\mathbf{x}))d\mathbf{x}$

```

1 var_h=function(h) mean(Vectorize(function(x) var(vec_mh(x,h)))(xr))

```

Define finally

```

1 mise_h=function(h) bias2_h(h)+var_h(h)

```

× **Cross Validation to estimate \hat{R}_h**

In real life, we cannot generate 500 other samples.


```

1 pred_i =function (i ,h) {
2 x=db$xr [ i ]
3 reg_i = lm(yr ~ 1, weights=dnorm(xr , mean=x , sd=h) , data=db , subset =(1:
  nrow(db))[-i ])
4 predict (reg_i , newdata=data . frame (xr=x) ) }

```

```

1 mse_h=function (h) {
2 mean ((db$yr - Vectorize (function (i) pred_i (i , h) ) (1:n) ) ^2) }

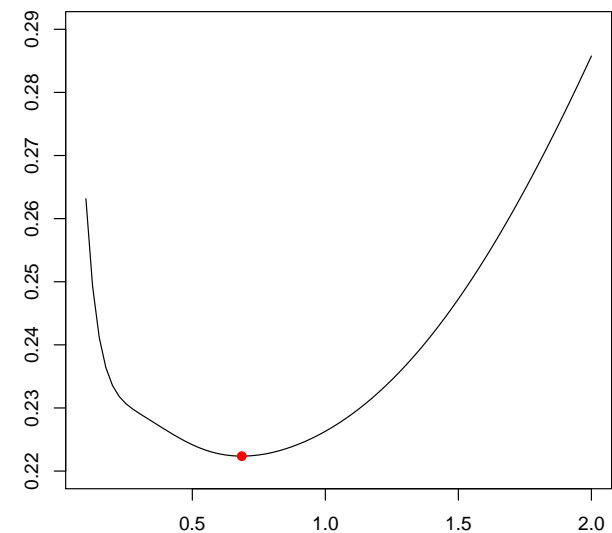
```

The optimal bandwidth parameter h^* is then

```

1 > optimize (mse_h , interval=c (.1 , 2) )
2 $minimum
3 [1] 0.6857169
4
5 $objective
6 [1] 0.2223652

```



× Optimal Bandwidth h^*

```
1 library (np)
2 h <- 1
3 nw <- npreg (yr~xr, data=db, bws=h, ckertype= "
  gaussian ")
4 u=seq (0,10,by=.025)
5 v=predict (nw,newdata=data.frame (xr=u))
6 lines (u,v,col=rgb (1,0,0,.4),lwd=2)
```

× Theoretical Results (Briefly)

One can prove that, as $n \rightarrow \infty$ and $nh \rightarrow 0$, R_h can be approximated by

$$\frac{h^4}{4} \left(\int x^2 K(\mathbf{x}) d\mathbf{x} \right)^2 \int \left(m''(\mathbf{x}) + 2m'(\mathbf{x}) \frac{f'(\mathbf{x})}{f(\mathbf{x})} \right) d\mathbf{x} + \frac{1}{nh} \sigma^2 \int K^2(x) dx \int \frac{d\mathbf{x}}{f(\mathbf{x})}$$

where f is the density of \mathbf{x} 's. Thus the optimal h is

$$h^* = n^{-\frac{1}{5}} \left(\frac{\sigma^2 \int K^2(x) dx \int \frac{d\mathbf{x}}{f(\mathbf{x})}}{\left(\int x^2 K(\mathbf{x}) d\mathbf{x} \right)^2 \int \left(\int m''(\mathbf{x}) + 2m'(\mathbf{x}) \frac{f'(\mathbf{x})}{f(\mathbf{x})} \right)^2 d\mathbf{x}} \right)^{\frac{1}{5}}$$

(hard to get a simple rule of thumb... up to a constant, $h^* \sim n^{-\frac{1}{5}}$)

One can also use (local polynomial)

```
1 > library(locfit)
```

but it is not the same smoothing technique. But here also, cross validation techniques can be used to get the optimal smoothing parameter.

Cross-Validation

💡 randomly split the sample into k chunks of (approximately) equal size, and then compute the classifier estimated with chunk κ and predict data in chunk κ .

```
1 library ( boot )  
2 model <- glm ( y~x , data=db )  
3 cv . glm ( db , model , K=10 )
```

for some 10-fold cross validation. When omitting, it is the standard leave-one-out cross validation technique.

```
1 cv . glm ( db , model )
```

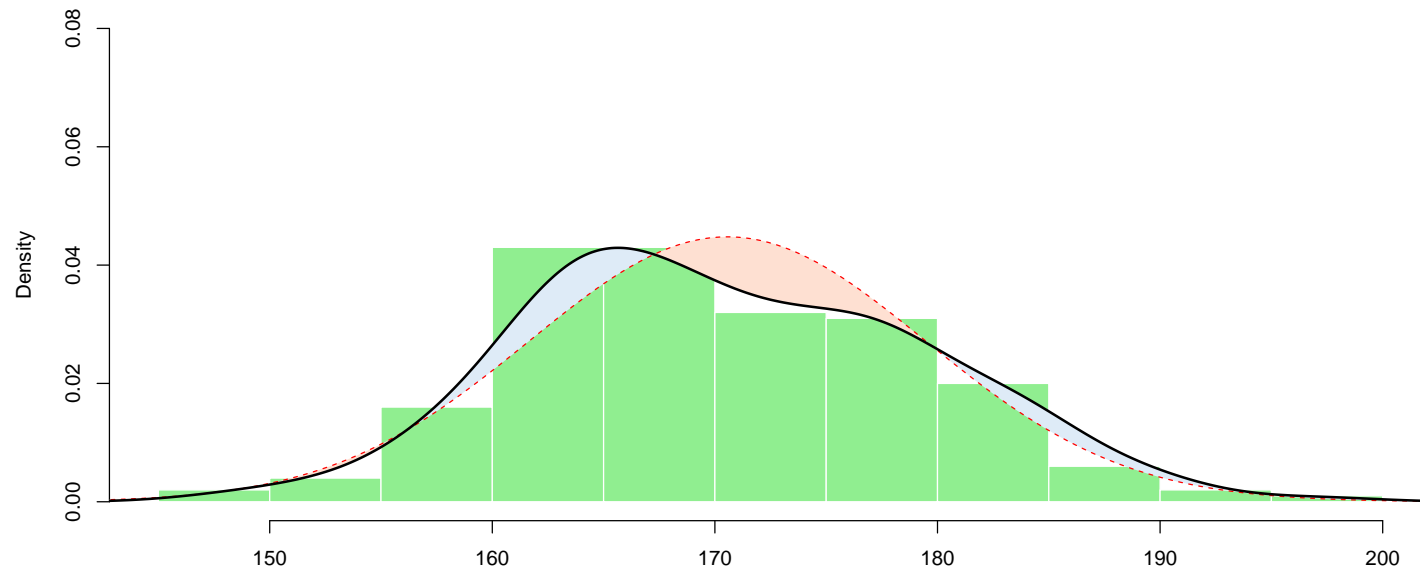
⚠ k -Fold Cross-Validation is faster to run.

Unsupervised vs. Supervised Techniques

The difference between **mixture** and **regression** models is simple

Let z_1, \dots, z_n be a sample with the height of some insured.

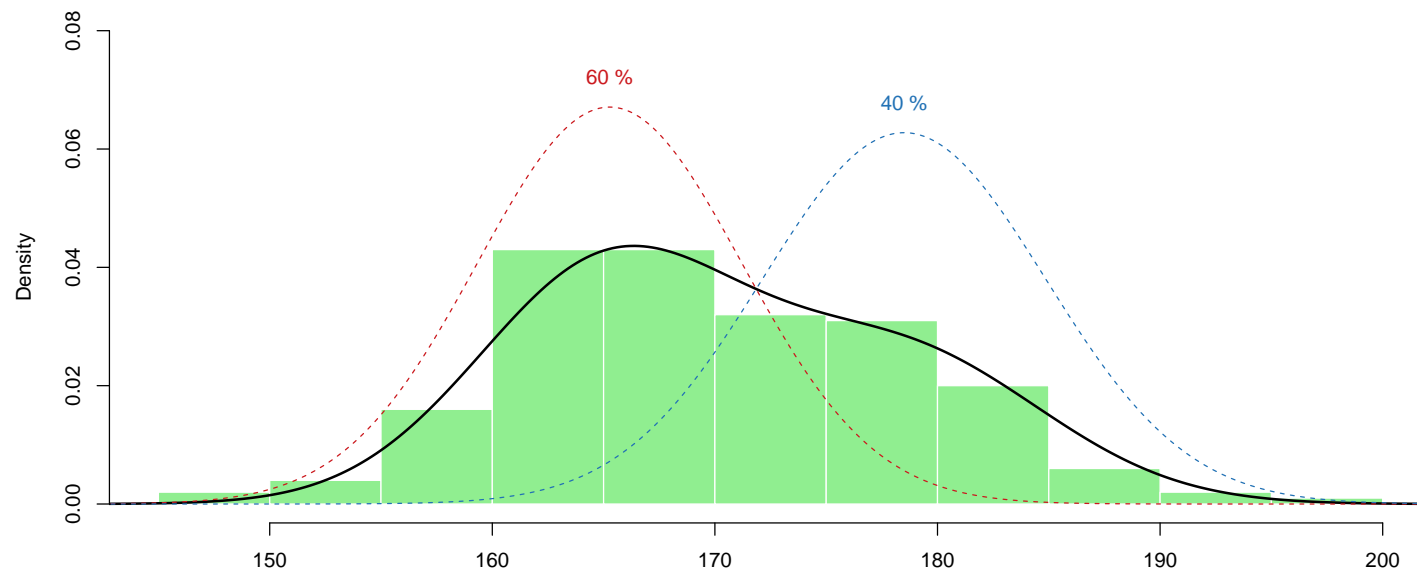
Observe that $Z \not\sim \mathcal{N}(\mu, \sigma^2)$



Unsupervised vs. Supervised Techniques

Why not a **mixture** model for $-x_i$ ''

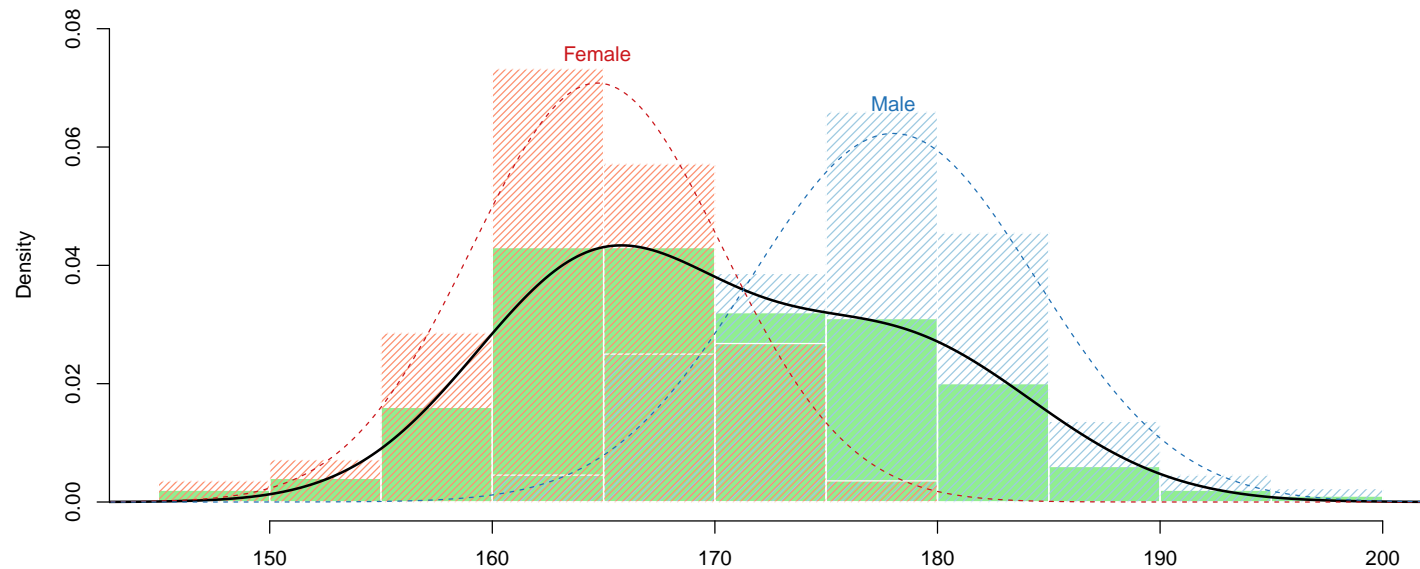
$X \sim$ – $\mathcal{N}(\mu_1, \sigma_1^2)$ with probability p_1 , i.e. with a latent variable $\Theta = 1$
 $\mathcal{N}(\mu_2, \sigma_2^2)$ with probability p_2 , i.e. with a latent variable $\Theta = 2$



Unsupervised vs. Supervised Techniques

Why not a regression model for $-y_i, x_i$ ''

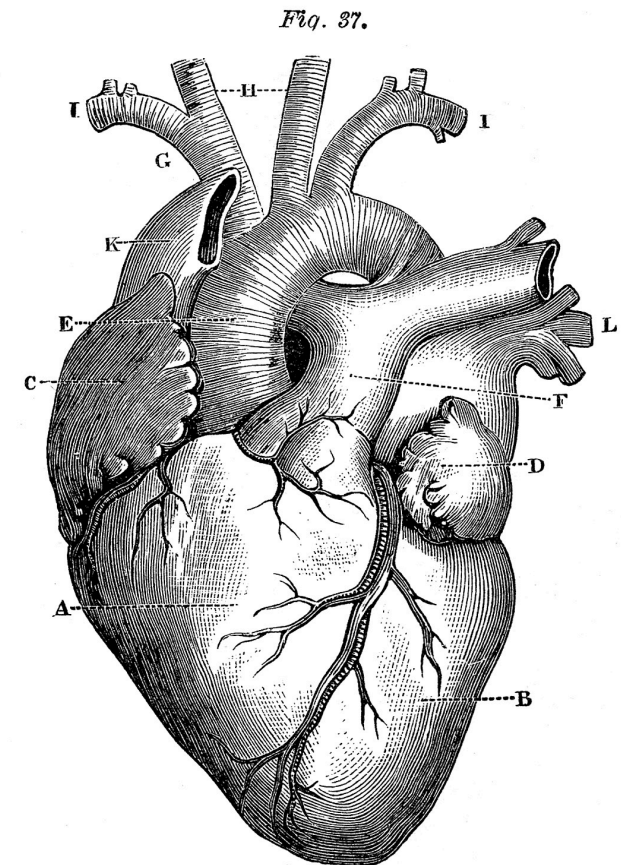
$$Y \sim \begin{cases} \mathcal{N}(\mu_M, \sigma_M^2) & \text{when } X = \text{'male'}, \text{ with probability } p_M \\ \mathcal{N}(\mu_F, \sigma_F^2) & \text{when } X = \text{'female'}, \text{ with probability } p_F \end{cases}$$



Classification: Description of the Dataset

Myocardial infarction of patients admitted in E.R.

- heart rate (FRCAR),
- heart index (INCAR)
- stroke index (INSYS)
- diastolic pressure (PRDIA)
- pulmonary arterial pressure (PAPUL)
- ventricular pressure (PVENT)
- lung resistance (REPUL)
- death or survival



Classification: Description of the Dataset

```
1 > myocarde=read.table("myocarde.csv", head=TRUE, sep=";")
```

```
2 > head(myocarde)
```

```
3   FRCAR  INCAR  INSYS  PRDIA  PAPUL  PVENT  REPUL  PRONO
4 1     90   1.71  19.0    16   19.5   16.0   912  SURVIE
5 2     90   1.68  18.7    24   31.0   14.0  1476  DECES
6 3    120   1.40  11.7    23   29.0    8.0  1657  DECES
7 4     82   1.79  21.8    14   17.5   10.0   782  SURVIE
8 5     80   1.58  19.7    21   28.0   18.5  1418  DECES
9 6     80   1.13  14.1    18   23.5    9.0  1664  DECES
```

Training and Validation Samples

⚠ in those slides, to present the technique, we fit the model on the complete dataset

To validate a model, we usually split the sample in two(or three)

- a sample to estimate the model (training)
- a sample used for cross-validation (if necessary)
- a sample to validate the model (validation)

⚠ In practice, we allocate as follows

- ~ 50%, 60%, 80% of sample to estimate the model
- ~ 50%, 40%, 20% of sample to validate the model

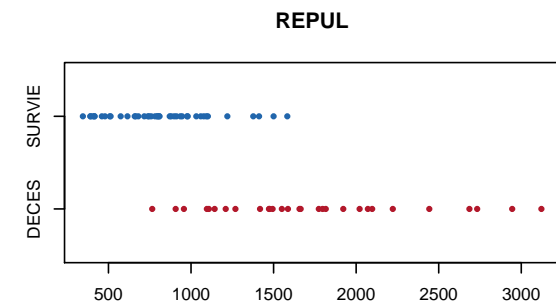
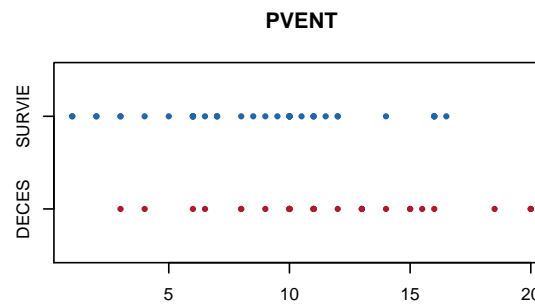
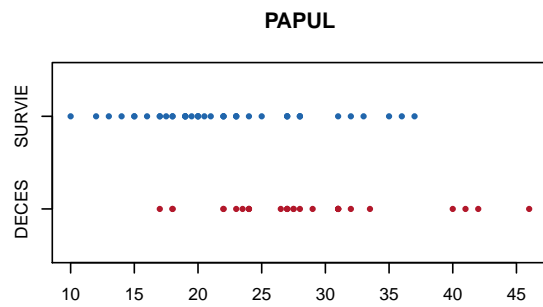
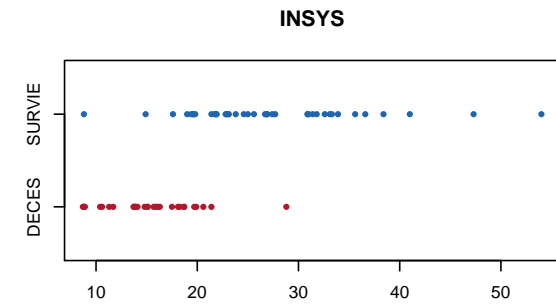
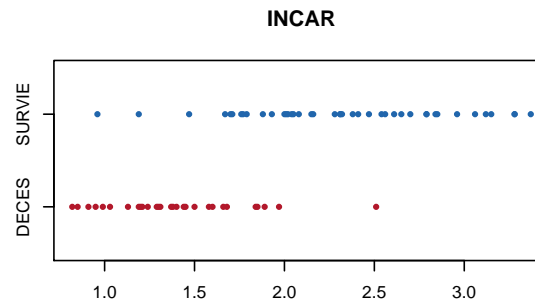
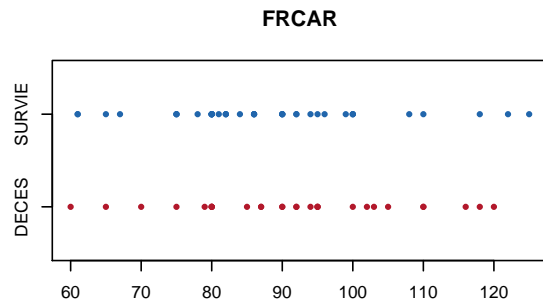
Generating Samples

```
1 splitdf <- function(dataframe, ratio=.5, seed=NULL) {  
2   if (!is.null(seed)) set.seed(seed)  
3   index <- 1:nrow(dataframe)  
4   trainindex <- sample(index, trunc(length(index)*ratio))  
5   trainset <- dataframe[trainindex, ]  
6   testset <- dataframe[-trainindex, ]  
7   return(list(trainset=trainset, testset=testset))}
```

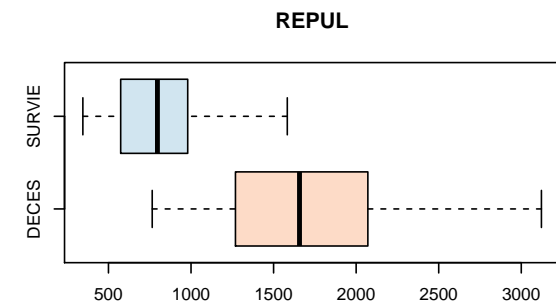
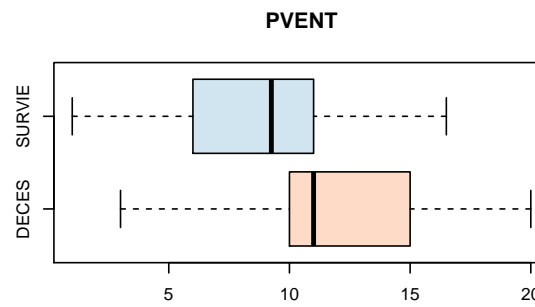
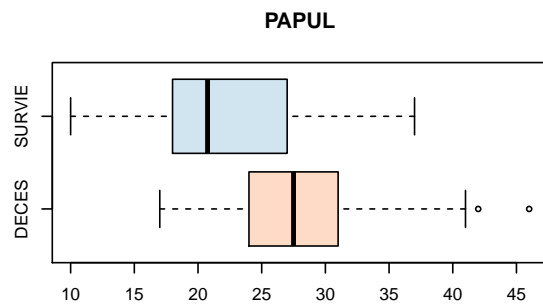
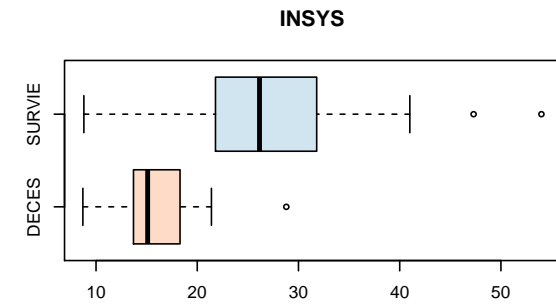
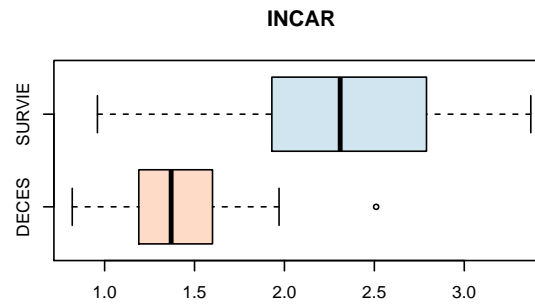
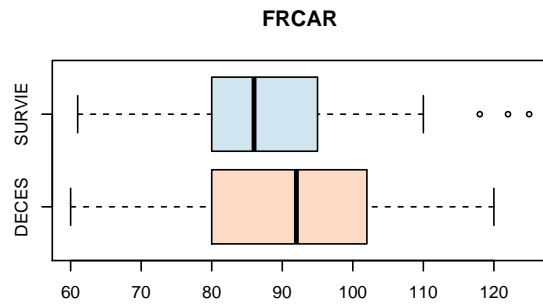
There are (at least) two generators of random numbers in SAS see [sas.com](https://www.sas.com)

- Fishman and Moore (1982, [jstor.org](https://www.jstor.org)) used for function RANUNI
- Mersenne-Twister used for the RAND function, based on Matsumoto and Nishimura (1997, math.sci.hiroshima-u.ac.jp)

Modeling a 0/1 Variable



Modeling a 0/1 Variable



Logistic Regression

```

1 > logistic <- glm(PRONO~. , data=myocarde , family=binomial)
2 > summary(logistic)
3
4 Call:
5 glm(formula = PRONO ~ ., family = binomial, data = myocarde)
6
7 Deviance Residuals:
8      Min       1Q   Median       3Q      Max
9  -2.24911  -0.41613   0.05261   0.39611   2.25781
10
11 Coefficients:
12             Estimate Std. Error z value Pr(>|z|)
13 (Intercept) -10.187642  11.895227  -0.856   0.392
14 FRCAR        0.138178   0.114112   1.211   0.226
15 INCAR       -5.862429   6.748785  -0.869   0.385
16 INSYS        0.717084   0.561445   1.277   0.202
17 PRDIA       -0.073668   0.291636  -0.253   0.801

```

```
18 PAPUL          0.016757    0.341942    0.049    0.961
19 PVENT         -0.106776    0.110550   -0.966    0.334
20 REPUL         -0.003154    0.004891   -0.645    0.519
21
22 (Dispersion parameter for binomial family taken to be 1)
23
24 Null deviance: 96.033 on 70 degrees of freedom
25 Residual deviance: 41.043 on 63 degrees of freedom
26 AIC: 57.043
27
28 Number of Fisher Scoring iterations: 7
```

Logistic Regression

```
1 > logistic <- glm(PRONO=="DECES"~. , data=myocarde , family=binomial)
2 > summary(logistic)
```

3

4 Call:

```
5 glm(formula = PRONO == "DECES" ~ ., family = binomial, data =
  myocarde)
```

6

7 Deviance Residuals:

8	Min	1Q	Median	3Q	Max
9	-2.25781	-0.39611	-0.05261	0.41613	2.24911

10

11 Coefficients:

12		Estimate	Std. Error	z value	Pr(> z)
13	(Intercept)	10.187642	11.895227	0.856	0.392
14	FRCAR	-0.138178	0.114112	-1.211	0.226
15	INCAR	5.862429	6.748785	0.869	0.385
16	INSYS	-0.717084	0.561445	-1.277	0.202


```
17 PRDIA      0.073668    0.291636    0.253      0.801
18 PAPUL     -0.016757    0.341942   -0.049     0.961
19 PVENT      0.106776    0.110550    0.966     0.334
20 REPUL      0.003154    0.004891    0.645     0.519
21
22 (Dispersion parameter for binomial family taken to be 1)
23
24 Null deviance: 96.033 on 70 degrees of freedom
25 Residual deviance: 41.043 on 63 degrees of freedom
26 AIC: 57.043
27
28 Number of Fisher Scoring iterations: 7
```

Multiple Logistic Regression

```

1 > library (VGAM)
2 > mlogistic <- vglm(PRONO~. , data=myocarde, family=multinomial)
3 > levels (myocarde$PRONO)
4 [1] "DECES" "SURVIE"
5 > summary(mlogistic)
6
7 Call:
8 vglm(formula = PRONO ~ ., family = multinomial, data = myocarde)
9
10 Pearson residuals:
11             Min           1Q       Median           3Q          Max
12 log(mu[,1]/mu[,2]) -3.4334 -0.28577 -0.03718  0.3007  3.3966
13
14 Coefficients:
15             Estimate Std. Error  z value
16 (Intercept) 10.1876411 11.8941581  0.856525
17 FRCAR       -0.1381781  0.1141056 -1.210967

```

```
18 INCAR          5.8624289   6.7484319   0.868710
19 INSYS         -0.7170840   0.5613961  -1.277323
20 PRDIA          0.0736682   0.2916276   0.252610
21 PAPUL         -0.0167565   0.3419255  -0.049006
22 PVENT          0.1067760   0.1105456   0.965901
23 REPUL          0.0031542   0.0048907   0.644939
24
25 Number of linear predictors: 1
26
27 Name of linear predictor: log(mu[,1]/mu[,2])
28
29 Dispersion Parameter for multinomial family: 1
30
31 Residual deviance: 41.04314 on 63 degrees of freedom
32
33 Log-likelihood: -20.52157 on 63 degrees of freedom
34
35 Number of iterations: 7
```

Fisher's Gradient Descent

The algorithm to estimate $\hat{\beta}$ is

1. start with some initial value β_0
2. define $\beta_k = \beta_{k-1} - H(\beta_{k-1})^{-1} \nabla \log \mathcal{L}(\beta_{k-1})$

where $\nabla \log \mathcal{L}(\beta)$ is the **gradient**, and $H(\beta)$ the **Hessian** matrix, also called Fisher's score.

The generic term of the Hessian is

$$\frac{\partial^2 \log \mathcal{L}(\beta)}{\partial \beta_k \partial \beta_\ell} = \sum_{i=1}^n X_{k,i} X_{\ell,i} [y_i - \pi_i(\beta)]$$

Define $\mathbf{\Omega} = [\omega_{i,j}] = \text{diag}(\hat{\pi}_i(1 - \hat{\pi}_i))$ so that the gradient is written

$$\nabla \log \mathcal{L}(\beta) = \frac{\partial \log \mathcal{L}(\beta)}{\partial \beta} = \mathbf{X}'(\mathbf{y} - \boldsymbol{\pi})$$

and the Hessian

$$H(\beta) = \frac{\partial^2 \log \mathcal{L}(\beta)}{\partial \beta \partial \beta'} = -\mathbf{X}'\Omega\mathbf{X}$$

The gradient descent algorithm is then

$$\beta_k = (\mathbf{X}'\Omega\mathbf{X})^{-1}\mathbf{X}'\Omega\mathbf{Z} \text{ where } \mathbf{Z} = \mathbf{X}\beta_{k-1} + \mathbf{X}'\Omega^{-1}(\mathbf{y} - \boldsymbol{\pi}),$$

```

1 > beta=as.matrix(lm(Y~0+X)$coefficients , ncol=1)
2 > for (s in 1:9) {
3 +   pi=exp(X%*%beta[,s]) / (1+exp(X%*%beta[,s]))
4 +   gradient=t(X)%*%(Y-pi)
5 +   omega=matrix(0, nrow(X), nrow(X)); diag(omega)=(pi*(1-pi))
6 +   Hessian=-t(X)%*%omega%*%X
7 +   beta=cbind(beta, beta[,s]-solve(Hessian)%*%gradient)}
8 > beta
9           [,1]      [,2]      [,3]      [,4]      [,5]
10 XConstant 0.05172274 -2.4401950 -1.2761448  2.9062584  7.9451328
11 XFRCAR    0.00316378  0.0113589 -0.0170704 -0.0723398 -0.1212657

```

12	XINCAR	-0.26677926	-0.7002255	0.3982835	2.9929899	5.2128794
13	XINSYS	0.00475193	0.0128323	-0.0992147	-0.3553974	-0.6146837
14	XPRDIA	0.00632701	0.0332371	0.0828177	0.0976728	0.0863094
15	XPAPUL	-0.00500955	-0.0486975	-0.0884098	-0.0864954	-0.0465313
16	XPVENT	0.00938577	0.0402316	0.0662682	0.0866649	0.1006702
17	XREPUL	0.00033615	0.0020090	0.0032057	0.0038279	0.0035192
18		[,6]	[,7]	[,8]	[,9]	[,10]
19	XConstant	9.9356684	10.184312	10.1876410	10.1876416	10.1876416
20	XFRCAR	-0.1364497	-0.138156	-0.1381781	-0.1381781	-0.1381781
21	XINCAR	5.8045084	5.861785	5.8624289	5.8624290	5.8624290
22	XINSYS	-0.7062320	-0.716945	-0.7170839	-0.7170840	-0.7170840
23	XPRDIA	0.0752166	0.073687	0.0736681	0.0736681	0.0736681
24	XPAPUL	-0.0206140	-0.016810	-0.0167565	-0.0167565	-0.0167565
25	XPVENT	0.1061793	0.106769	0.1067760	0.1067760	0.1067760
26	XREPUL	0.0032053	0.003154	0.0031541	0.0031541	0.0031541

After 7 iterations, convergence was reached. See

```
1 > summary(logistic)$coefficients [ ,1]
```

2	(Intercept)	FRCAR	INCAR	INSYS	PRDIA
3	10.187641685	-0.138178119	5.862429035	-0.717084018	0.073668171
4	PAPUL	PVENT	REPUL		
5	-0.016756506	0.106776012	0.003154187		
6	-0.073668171	0.016756506	-0.106776012	-0.003154187	

From maximum likelihood properties,

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, I(\beta)^{-1}).$$

From a numerical point of view, this asymptotic variance $I(\beta)^{-1}$ satisfies $I(\beta)^{-1} = -H(\beta)$,

```
1 > -solve(Hessian)
```

2		Constant	FRCAR	INCAR	INSYS
3	Constant	141.50029499	-1.1120707443	42.29308470	-5.792019e+00
4	FRCAR	-1.11207074	0.0130218132	-0.68328148	6.113030e-02
5	INCAR	42.29308470	-0.6832814845	45.54684575	-3.245723e+00
6	INSYS	-5.79201872	0.0611303045	-3.24572286	3.152285e-01

7	PRDIA	-0.25859731	-0.0026468213	0.30880800	5.050851e-05
8	PAPUL	1.15456409	0.0044720856	-0.97012941	2.087264e-03
9	PVENT	-0.03304634	-0.0002828814	0.02340121	-6.141984e-03
10	REPUL	-0.02103416	-0.0001058035	0.01811208	-1.791229e-04
11		PRDIA	PAPUL	PVENT	REPUL
12	Constant	-2.585973e-01	1.154564088	-3.304634e-02	-2.103416e-02
13	FRCAR	-2.646821e-03	0.004472086	-2.828814e-04	-1.058035e-04
14	INCAR	3.088080e-01	-0.970129411	2.340121e-02	1.811208e-02
15	INSYS	5.050851e-05	0.002087264	-6.141984e-03	-1.791229e-04
16	PRDIA	8.505231e-02	-0.076531702	-1.516899e-02	3.594489e-04
17	PAPUL	-7.653170e-02	0.116926256	1.390261e-02	-1.290423e-03
18	PVENT	-1.516899e-02	0.013902608	1.222148e-02	-4.677069e-05
19	REPUL	3.594489e-04	-0.001290423	-4.677069e-05	2.392144e-05

```

1 > sqrt(-diag(solve(Hessian)))
2      Constant      FRCAR      INCAR      INSYS      PRDIA
3 11.895389653  0.114113159  6.748840326  0.561452093  0.291637287
4      PAPUL      PVENT      REPUL
5 0.341944815  0.110550787  0.004890955
    
```



```

1 > summary(logistic)$coefficients[,2]
2 (Intercept)          FRCAR          INCAR          INSYS          PRDIA
3 11.895227082    0.114112163    6.748785486    0.561444684    0.291636009
4          PAPUL          PVENT          REPUL
5 0.341942251    0.110550096    0.004890917

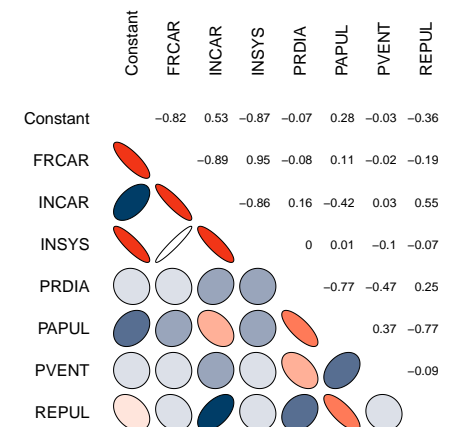
```

One can also visualise the correlation matrix of $\hat{\beta}$

```

1 > sigma=sqrt(diag(-solve(Hessian)))
2 > C <- -solve(Hessian)/(sigma %*% t(sigma))

```



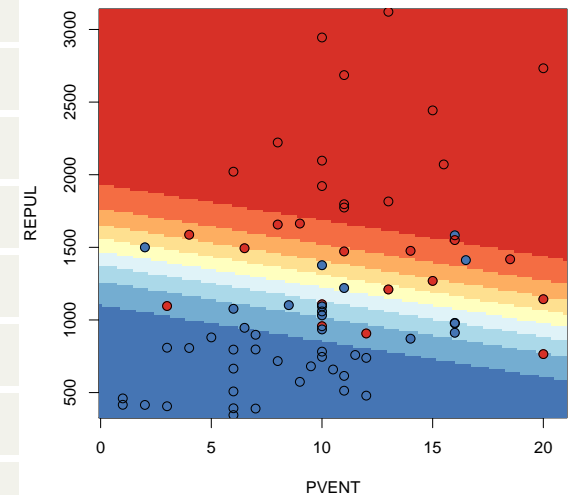
visualising a logistic regression

Consider a logistic regression with 2 covariates

```

1 > mlogistic2 <- vglm(PRONO~PVENT+REPUL , data=
  myocarde , family=multinomial)
2 > pred_logistic = function(p,r){
3   return(predict(mlogistic2 ,newdata=data.frame(
  PVENT=p,REPUL=r) ,type="response") [,1]) }
4 > vpvent=seq(min(myocarde$PVENT) -1,max(myocarde$
  PVENT) +1,length=151)
5 > vrepul=seq(min(myocarde$REPUL) -10,max(myocarde$
  REPUL) +10,length=151)
6 > image(vpvent ,vrepul ,outer(vpvent ,vrepul ,pred_
  logistic) ,col=CL2palette ,xlab="PVENT" ,ylab="
  REPUL" )

```



The logistic regression

```
1 > probabilities <- predict(logistic , myocarde , type="response")
2 > predictions <- levels(myocarde$PRONO)[(probabilities <.5)+1]
```

or for a more general classification

```
3 > probabilities <- predict(mlogistic , myocarde , type="response")
4 > probabilities [1:3 ,]
5           [,1]      [,2]
6 [1 ,] 0.3986106 0.6013894
7 [2 ,] 0.8306231 0.1693769
8 [3 ,] 0.6710440 0.3289560
9 > pred <- levels(myocarde$PRONO)[apply(probabilities ,1 ,which.max)]
```

```
1 > table(pred , myocarde$PRONO)
```

```
2
3 pred          DECES SURVIE
4          DECES      25      3
5          SURVIE      4     39
```

From $-0, 1$ to $-A, B, C$

The Bernoulli case can (easily) be extended to the multinomial distribution. In the $-0, 1$ case

$$\log \frac{\mathbb{P}(Y = 1)}{1 - \mathbb{P}(Y = 1)} = \log \frac{\mathbb{P}(Y = 1)}{\mathbb{P}(Y = 0)} = \mathbf{X}^\top \boldsymbol{\beta}$$

so that

$$\mathbb{P}(Y = 1) = \frac{e^{\mathbf{X}^\top \boldsymbol{\beta}}}{1 + e^{\mathbf{X}^\top \boldsymbol{\beta}}} = \frac{p_1}{p_0 + p_1} \quad \text{and} \quad \mathbb{P}(Y = 0) = \frac{1}{1 + e^{\mathbf{X}^\top \boldsymbol{\beta}}} = \frac{p_0}{p_0 + p_1}$$

In the case where the response variable takes values $-A, B, C$,

$$\log \frac{\mathbb{P}(Y = A)}{\mathbb{P}(Y = C)} = \mathbf{X}^\top \boldsymbol{\beta}_A \quad \text{and} \quad \log \frac{\mathbb{P}(Y = B)}{\mathbb{P}(Y = C)} = \mathbf{X}^\top \boldsymbol{\beta}_B$$

From $-0, 1$ " to $-A, B, C$ "

so that

$$\mathbb{P}(X = A) = \frac{p_A}{p_A + p_B + p_C} \quad \text{i.e.} \quad \mathbb{P}(X = A) = \frac{e^{\mathbf{X}^\top \beta_A}}{e^{\mathbf{X}^\top \beta_A} + e^{\mathbf{X}^\top \beta_B} + 1}$$

$$\mathbb{P}(X = B) = \frac{p_B}{p_A + p_B + p_C} \quad \text{i.e.} \quad \mathbb{P}(X = B) = \frac{e^{\mathbf{X}^\top \beta_B}}{e^{\mathbf{X}^\top \beta_A} + e^{\mathbf{X}^\top \beta_B} + 1}$$

and

$$\mathbb{P}(X = C) = \frac{p_C}{p_A + p_B + p_C} \quad \text{i.e.} \quad \mathbb{P}(X = C) = \frac{1}{e^{\mathbf{X}^\top \beta_A} + e^{\mathbf{X}^\top \beta_B} + 1}$$

Visualisation of a Classification

Consider here some logistic regression on k variables. Can we visualise it ?

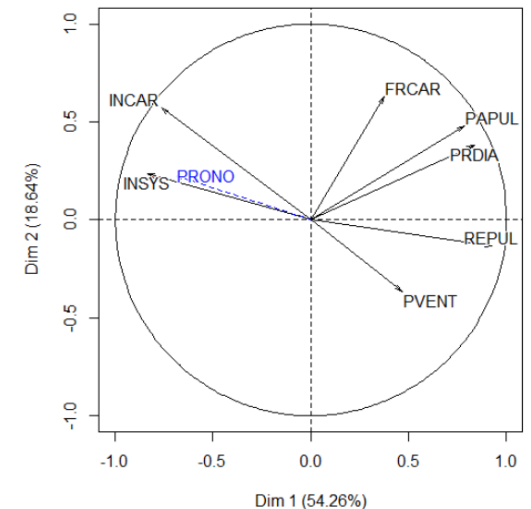
On our myocarde dataset, consider some **principal component analysis**, including the pronostic variable,

either as a continuous variable,

```

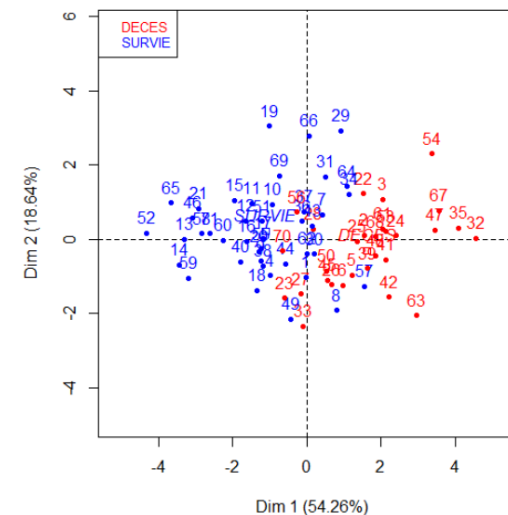
1 > library (FactoMineR)
2 > MYOCARDE2 <- MYOCARDE
3 > MYOCARDE2$PRONO <- (MYOCARDE2$PRONO == "SURVIE") * 1
4 > acp <- PCA(MYOCARDE2, quanti.sup=8, graph=TRUE)

```



or as a qualitative variable,

```
1 > acp <- PCA(MYOCARDE, quali.sup=8, graph=TRUE)
2 > plot(acp, habillage = 8, col.hab=c("red", "blue"))
```



To get from the first two component hyperplane to \mathbb{R}^k , use

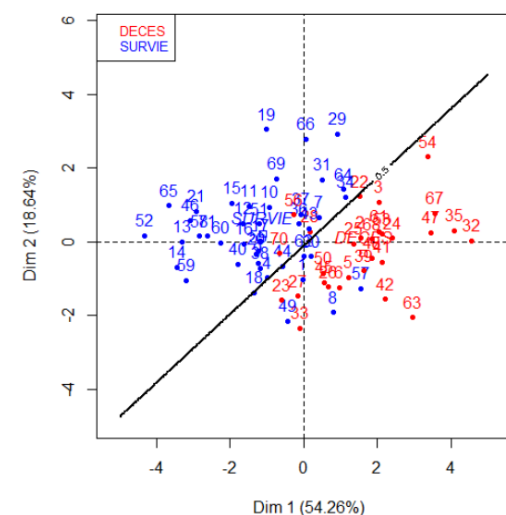
```
1 > X <- MYOCARDE[, 1:7]
2 > acp <- PCA(X, ncp=ncol(X))
3 > M <- acp$var$coord
4 > m <- apply(X, 2, mean)
5 > s <- apply(X, 2, sd)
6 > pred <- function(d1, d2, Mat, reg) {
7 +   z <- Mat %*% c(d1, d2, rep(0, ncol(X) - 2))
8 +   newd <- data.frame(t(z*s+m))
9 +   names(newd) <- names(X)
```

```
10 + predict (reg , newdata=newd , type="response" ) }
```

Hence

or as a qualitative variable,

```
1 > reg_tot <- glm (PRONO~. , data=MYOCARDE, family=
  binomial)
2 > Minv <- solve (M)
3 > p <- Vectorize (function (d1 , d2) pred (d1 , d2 , Minv ,
  reg_tot))
4 > zgrid <- outer (xgrid , ygrid , p)
5 > acp2 <- PCA (MYOCARDE, quali . sup=8, graph=TRUE)
6 > plot (acp2 , habillage = 8, col . hab=c ("red" , "blue" )
  )
7 > contour (xgrid , ygrid , zgrid , add=TRUE, levels =.5)
```



Linear Discriminant Analysis

LDA is a classification method that finds a linear combination of data attributes that best separate the data into classes.

```

1 > library(MASS)
2 > fit_lda <- lda(PRONO ~ . , data=myocarde)
3 > fit_lda
4 Call:
5 lda(PRONO ~ . , data = myocarde)
6
7 Prior probabilities of groups:
8     DECES     SURVIE
9 0.4084507 0.5915493
10
11 Group means:
12           FRCAR     INCAR     INSYS     PRDIA     PAPUL     PVENT
13           REPUL
14 DECES  91.55172  1.397931  15.53103  21.44828  28.43103  11.844828
15           1738.6897

```

```
14 SURVIE 87.69048 2.318333 27.20238 15.97619 22.20238 8.642857
    817.2143
15
16 Coefficients of linear discriminants:
17          LD1
18 FRCAR -0.012743116
19 INCAR  1.074534545
20 INSYS -0.019139867
21 PRDIA -0.025483955
22 PAPUL  0.020177505
23 PVENT -0.037804074
24 REPUL -0.001353977
25 > predictions <- predict(fit_lda, myocarde)$class
26 > table(predictions, myocarde$PRONO)
27
28 predictions  DECES SURVIE
29      DECES      25      2
30      SURVIE      4      40
```

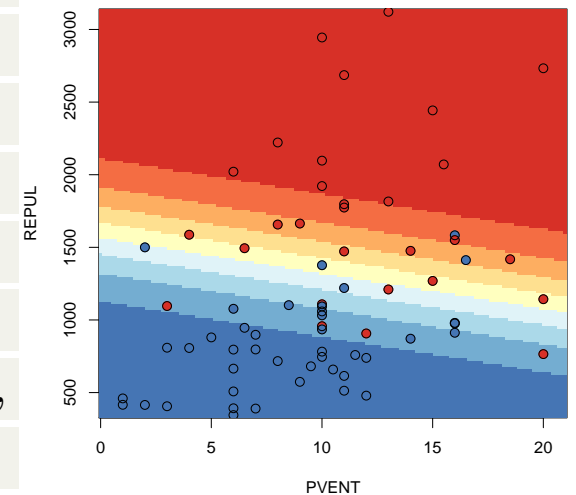
visualising a LDA

Consider a logistic regression with 2 covariates

```

1 fit_lda2 <- lda(PRONO~PVENT+REPUL , data=myocarde)
2 pred_lda = function(p,r){
3   return(predict(fit_lda2,newdata=
4     data.frame(PVENT=p,REPUL=r),
5     type="response")$posterior[, "
6     DECES" ] ) }
7 image(vpvent , vrepul , outer(vpvent , vrepul , pred_lda) ,
8   col=CL2palette ,
9   xlab="PVENT" , ylab="REPUL" )

```



Partial Least Square Discriminant Analysis

Partial Least Squares Discriminate Analysis is the application of LDA on a dimension-reducing projection of the input data (partial least squares).

```

1 > library(caret)
2 > X=as.matrix(myocarde[,1:7])
3 > Y=myocarde$PRONO
4 > fit_plsda <- plsda(X, Y, probMethod="Bayes")
5 > predictions <- predict(fit_plsda, myocarde[,1:7])
6 > table(predictions, myocarde$PRONO)
7
8 predictions  DECES  SURVIE
9      DECES      24      3
10     SURVIE      5     39

```

Assume that

$$\mathbf{X}|Y = 0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}) \text{ and } \mathbf{X}|Y = 1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma})$$


then

$$\log \frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x})} = [\mathbf{x}]^\top \boldsymbol{\Sigma}^{-1} [\boldsymbol{\mu}_y] - \frac{1}{2} [\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0]^\top \boldsymbol{\Sigma}^{-1} [\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0] + \log \frac{\mathbb{P}(Y = 1)}{\mathbb{P}(Y = 0)}$$

which is linear in \mathbf{x}

$$\log \frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x})} = \mathbf{x}^\top \boldsymbol{\beta}$$

When each groups have Gaussian distributions with identical variance matrix, then LDA and the logistic regression lead to the same classification rule.

 there is a slight difference in the way parameters are estimated.

Mixture Discriminant Analysis

```
1 > library(mda)
2 > fit_mda <- mda(PRONO ~ . , data=myocarde)
3 > fit_mda
4 Call:
5 mda(formula = PRONO ~ . , data = myocarde)
6
7 Dimension: 5
8
9 Percent Between-Group Variance Explained:
10      v1      v2      v3      v4      v5
11 82.43  97.16  99.45  99.88 100.00
12
13 Degrees of Freedom (per dimension): 8
14
15 Training Misclassification Error: 0.14085 ( N = 71 )
16
17 Deviance: 46.203
```

```
18 > predictions <- predict(fit_mda, myocarde)
```

```
19 > table(predictions, myocarde$PRONO)
```

```
20
```

```
21 predictions  DECES  SURVIE
```

```
22     DECES      24      5
```

```
23     SURVIE      5     37
```

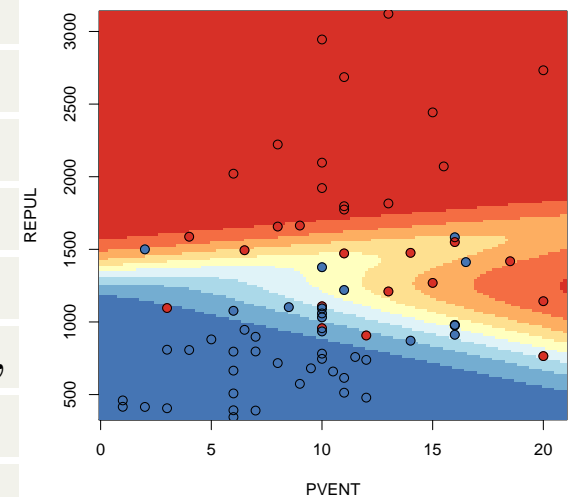
visualising a MDA

Consider a MDA with 2 covariates

```

1 fit_mda <- mda(PRONO~PVENT+REPUL , data=myocarde)
2 pred_mda = function(p,r){
3   return(predict(fit_mda,newdata=
4     data.frame(PVENT=p,REPUL=r),
5     type="posterior"))[, "DECES" ]}
6 image(vpvent , vrepul , outer(vpvent , vrepul , pred_mda) ,
7   col=CL2palette ,
8   xlab="PVENT" , ylab="REPUL" )

```



Quadratic Discriminant Analysis

QDA seeks a quadratic relationship between attributes that maximises the distance between the classes.

```

1 > library(MASS)
2 > fit_dqa <- qda(PRONO ~ . , data=myocarde)
3 > fit_dqa
4 Call:
5 qda(PRONO ~ . , data = myocarde)
6
7 Prior probabilities of groups:
8     DECES     SURVIE
9 0.4084507 0.5915493
10
11 Group means:
12           FRCAR     INCAR     INSYS     PRDIA     PAPUL     PVENT
13 DECES  91.55172  1.397931 15.53103 21.44828 28.43103 11.844828
14           1738.6897

```

```
14 SURVIE 87.69048 2.318333 27.20238 15.97619 22.20238 8.642857
    817.2143
15 > predictions <- predict(fit_dqa, myocarde)$class
16 > table(predictions, myocarde$PRONO)
17
18 predictions DECES SURVIE
19     DECES      24      5
20     SURVIE      5     37
```

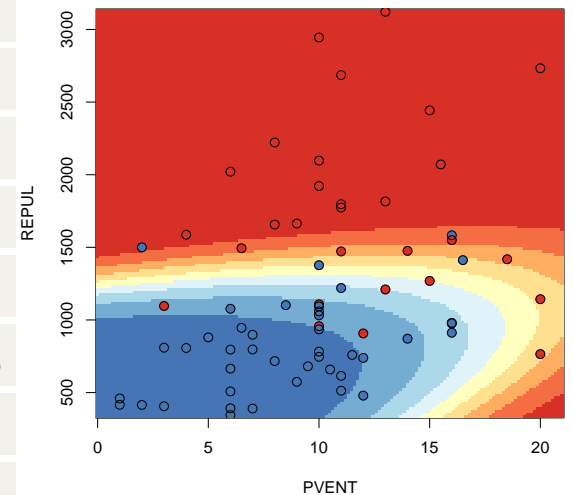
visualising a QDA

Consider a QDA with 2 covariates

```

1 fit_qda <- qda(PRONO~PVENT+REPUL , data=myocarde)
2 pred_qda = function(p,r){
3   return(predict(fit_qda,newdata=
4     data.frame(PVENT=p,REPUL=r))$
5     posterior[, "DECES" ])}
6 image(vpvent , vrepul , outer(vpvent , vrepul , pred_qda) ,
7   col=CL2palette ,
8   xlab="PVENT" , ylab="REPUL" )

```



The difference between LDA and QDA

```
1 S<-matrix(c(1,.5*2,.5*2,4),2,2)
```

```
2
```

```
3 X1<-rmnorm(300,c(-1,-1)*3,S)
```

```
4 X2<-rmnorm(300,c(1,1)*3,S)
```

```
1 S1<-matrix(c(1,.5*2,.5*2,4),2,2)
```

```
2 S2<-matrix(c(1,-.5*2,-.5*2,4),2,2)
```

```
3 X1<-rmnorm(300,c(-1,-1)*3,S1)
```

```
4 X2<-rmnorm(300,c(1,1)*3,S2)
```

The difference between LDA and QDA

```
1 fitL<-lda(y ~ x1 + x2, data=df)
```

```
2 pL<-function(u,v) predict(fitL,
```

```
3 newdata=data.frame(x1=u,x2=v)
```

```
4 )$posterior[, "1"])
```

```
1 fitQ<-qda(y ~ x1 + x2, data=df)
```

```
2 pQ<-function(u,v) predict(fitQ,
```

```
3 newdata=data.frame(x1=u,x2=v)
```

```
4 )$posterior[, "1"])
```

k-Nearest Neighbor and Classification

The *k*-Nearest Neighbor (kNN) method makes predictions by locating similar cases to a given data instance (using a similarity function) and returning the average or majority of the most similar data instances.

```
1 > library(caret)
2 > kNN <- knn3(PRONO~., data=myocarde, k=5)
3 > predictions <- predict(kNN, myocarde, type="class")
4 > table(predictions, myocarde$PRONO)
```

```
5
6 predictions  DECES  SURVIE
7      DECES      24      5
8      SURVIE      5      37
```

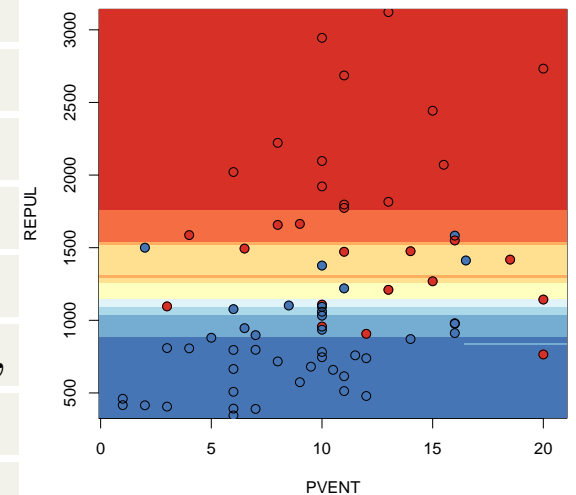
Visualising a k -Nearest Neighbor Classification

Consider a k -NN with 2 covariates

```

1 knn2 <- knn3(PRONO~PVENT+REPUL, data=myocarde, k
  =10)
2 pred_knn = function(p,r){
3   return(predict(knn2,newdata=
4   data.frame(PVENT=p,REPUL=r), type="prob")[,1])}
5 image(vpvent, vrepul, outer(vpvent, vrepul, pred_knn),
6   col=CL2palette,
  xlab="PVENT", ylab="REPUL")

```



Naive Bayes Classifier

Naive Bayes uses Bayes Theorem to model the conditional relationship of each attribute to the class variable.

```

1 > library(e1071)
2 > NB <- naiveBayes(PRONO~., data=myocarde)
3 > NB
4
5 Naive Bayes Classifier for Discrete Predictors
6
7 Call:
8 naiveBayes.default(x = X, y = Y, laplace = laplace)
9
10 A-priori probabilities:
11 Y
12     DECES     SURVIE
13 0.4084507 0.5915493
14
15 Conditional probabilities:

```



```

16      FRCAR
17 Y          [,1]      [,2]
18 DECES  91.55172  15.28441
19 SURVIE 87.69048  14.58949

```

```

21      INCAR
22 Y          [,1]      [,2]
23 DECES  1.397931  0.3808954
24 SURVIE 2.318333  0.5743880

```

etc

```

1 > predictions <- predict(NB, myocarde)
2 > table(predictions, myocarde$PRONO)
3
4 predictions DECES SURVIE
5      DECES      25      4
6      SURVIE      4      38

```

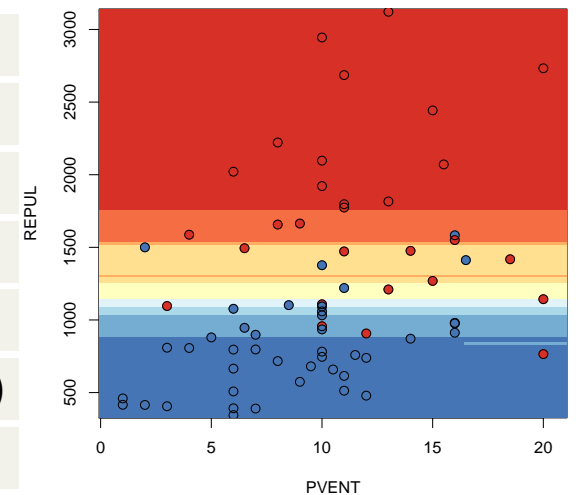
Visualising a Naive Bayes Classification

Consider a NB with 2 covariates

```

1 > NB2 <- naiveBayes(PRONO~PVENT+REPUL, data=
  myocarde)
2 > pred_NB = function(p,r){
3 + return(predict(NB2,newdata=
4 + data.frame(PVENT=p,REPUL=r), type="raw")[,1])}
5 > image(vpvent, vrepul, outer(vpvent, vrepul, pred_NB)
  , col=CL2palette, xlab="PVENT", ylab="REPUL")

```

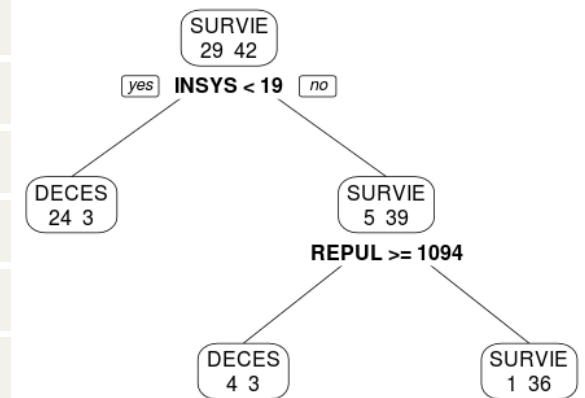


Visualising a Classification Tree

```

1 > MYOCARDE=read.table(
2 + "http://freakonometrics.free.fr/saporta.csv",
3 + head=TRUE, sep=";")
4 > library(rpart)
5 > cart<-rpart(PRONO~., data=MYOCARDE)
6 > library(rpart.plot)
7 > library(rattle)
8 > prp(cart, type=2, extra=1)

```



Visualising a Classification Tree

If X is splitted in 2 classes, denoted A and B , use either Gini index

$$\text{gini}(Y|X) = - \sum_{x \in -A, B''} \frac{n_x}{n} \sum_{y \in -0, 1''} \frac{n_{x,y}}{n_x} \left(1 - \frac{n_{x,y}}{n_x} \right)$$

```

1 > gini=function(y, classe) {
2 + T=table(y, classe)
3 + nx=apply(T, 2, sum)
4 + pxy=T/matrix(rep(nx, each=2), nrow=2)
5 + omega=matrix(rep(nx, each=2), nrow=2)/sum(T)
6 + return( -sum(omega*pxy*(1-pxy)) ) }

```

Hence,

```

1 > CLASSE=MYOCARDE[,1] <= 2.5
2 > gini(y=MYOCARDE$PRONO, classe=CLASSE)
3 [1] -0.4832375

```

Visualising a Classification Tree

or the **entropy** index,

$$\text{entropy}(Y|X) = - \sum_{x \in \{-A, B, C\}} \frac{n_x}{n} \sum_{y \in \{-0, 1\}} \frac{n_{x,y}}{n_x} \log \left(\frac{n_{x,y}}{n_x} \right)$$

```

1 > entropie=function(y, classe){
2   +   T=table(y, classe)
3   +   nx=apply(T, 2, sum)
4   +   n=sum(T)
5   +   pxy=T/matrix(rep(nx, each=2), nrow=2)
6   +   omega=matrix(rep(nx, each=2), nrow=2)/n
7   +   return( sum(omega*pxy*log(pxy)) )}

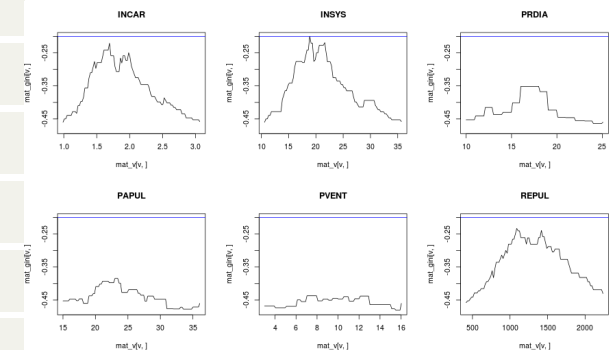
```

Visualising a Classification Tree

```

1 > mat_gini=mat_v=matrix(NA,7,101)
2 > for(v in 1:7){
3 +   variable=MYOCARDE[,v]
4 +   v_seuil=seq(quantile(MYOCARDE[,v],
5 + 6/length(MYOCARDE[,v])),
6 + quantile(MYOCARDE[,v],1-6/length(
7 + MYOCARDE[,v])),length=101)
8 +   mat_v[v,]=v_seuil
9 +   for(i in 1:101){
10 + CLASSE=variable<=v_seuil[i]
11 + mat_gini[v,i]=
12 +   gini(y=MYOCARDE$PRONO,classe=CLASSE)} }

```



Visualising a Classification Tree

```

1 > par(mfrow=c(2,3))
2 > for(v in 2:7){
3 +   plot(mat_v[v,], mat_gini[v,], type="l",
4 +   ylim=range(mat_gini),
5 +   main=names(MYOCARDE)[v])
6 +   abline(h=max(mat_gini), col="blue")
7 + }

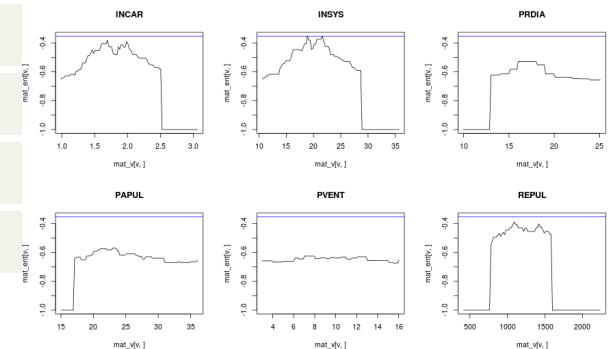
```

or we can use the entropy index.

```

1 > idx=which(MYOCARDE$INSYS>=19)

```

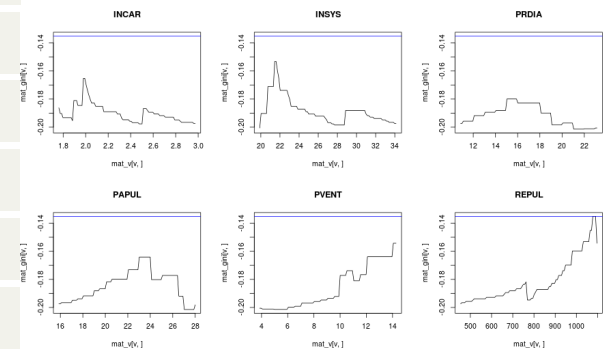


Visualising a Classification Tree

```

1 > mat_gini=mat_v=matrix(NA,7,101)
2 > for(v in 1:7){
3 +   variable=MYOCARDE[idx ,v]
4 +   v_seuil=seq(quantile(MYOCARDE[idx ,v] ,
5 + 6/length(MYOCARDE[idx ,v] ) ) ,
6 + quantile(MYOCARDE[idx ,v] ,1-6/length(
7 + MYOCARDE[idx ,v] ) ) , length=101)
8 +   mat_v[v,]=v_seuil
9 +   for(i in 1:101){
10 +     CLASSE=variable<=v_seuil[i]
11 +     mat_gini[v,i]=
12 +     gini(y=MYOCARDE$PRONO[idx] ,
13 +     classe=CLASSE) }}

```



Classification and Regression Trees (CART)

```
1 > library(rpart)
2 > cart <- rpart(PRONO~. , data=myocarde)
```

Classification and Regression Trees (CART)

```

1 > summary(cart)
2 Call:
3 rpart(formula = PRONO ~ ., data = myocarde)
4   n= 71
5
6           CP nsplit rel error      xerror      xstd
7 1 0.72413793      0 1.0000000 1.0000000 0.1428224
8 2 0.03448276      1 0.2758621 0.4827586 0.1156044
9 3 0.01000000      2 0.2413793 0.5172414 0.1186076
10
11 Variable importance
12 INSYS REPUL INCAR PAPUL PRDIA FRCAR PVENT
13    29    27    23     8     7     5     1

```

Classification and Regression Trees (CART)

```

1 Node number 1: 71 observations ,      complexity param=0.7241379
2   predicted class=SURVIE   expected loss=0.4084507   P(node) =1
3   class counts:      29      42
4   probabilities: 0.408 0.592
5   left son=2 (27 obs) right son=3 (44 obs)
6   Primary splits :
7       INSYS < 18.85   to the left ,   improve=20.112890, (0 missing)
8       REPUL < 1094.5 to the right , improve=19.021400, (0 missing)
9       INCAR < 1.69   to the left ,   improve=18.615510, (0 missing)
10      PRDIA < 17     to the right , improve= 9.361141, (0 missing)
11      PAPUL < 23.25  to the right , improve= 7.022101, (0 missing)
12      Surrogate splits :
13      REPUL < 1474   to the right , agree=0.915, adj=0.778, (0 split)
14      INCAR < 1.665  to the left , agree=0.901, adj=0.741, (0 split)
15      PAPUL < 28.5   to the right , agree=0.732, adj=0.296, (0 split)
16      PRDIA < 18.5   to the right , agree=0.718, adj=0.259, (0 split)
17      FRCAR < 99.5   to the right , agree=0.690, adj=0.185, (0 split)

```

Classification and Regression Trees (CART)

```
1 Node number 2: 27 observations
2 predicted class=DECES expected loss=0.1111111 P(node) =0.3802817
3 class counts:      24      3
4 probabilities: 0.889 0.111
```

Classification and Regression Trees (CART)

```

1 Node number 3: 44 observations ,      complexity param=0.03448276
2   predicted class=SURVIE   expected loss=0.1136364   P(node) =0.6197183
3   class counts:           5      39
4   probabilities: 0.114 0.886
5   left son=6 (7 obs) right son=7 (37 obs)
6   Primary splits :
7       REPUL < 1094.5 to the right , improve=3.489119, (0 missing)
8       INSYS < 21.55  to the left  , improve=2.122460, (0 missing)
9       PVENT < 13     to the right , improve=1.651281, (0 missing)
10      PAPUL < 23.5   to the right , improve=1.641414, (0 missing)
11      INCAR < 1.985  to the left  , improve=1.592803, (0 missing)
12      Surrogate splits :
13      INCAR < 1.685  to the left  , agree=0.886, adj=0.286, (0 split)
14      PVENT < 17.25  to the right , agree=0.864, adj=0.143, (0 split)

```

Classification and Regression Trees (CART)

1 Node number 6: 7 observations

2 predicted **class**=DECES expected loss=0.4285714 P(node)
=0.09859155

3 **class** counts: 4 3

4 probabilities: 0.571 0.429

5

6 Node number 7: 37 observations

7 predicted **class**=SURVIE expected loss=0.02702703 P(node)
=0.5211268

8 **class** counts: 1 36

9 probabilities: 0.027 0.973

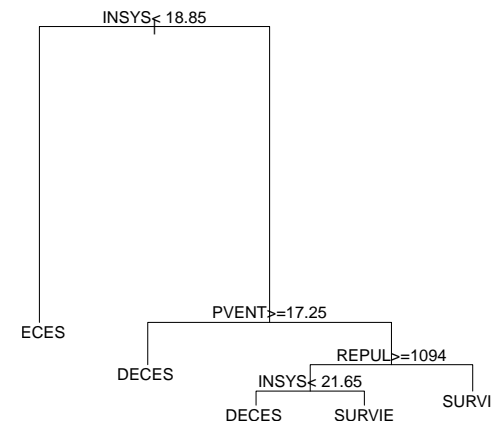
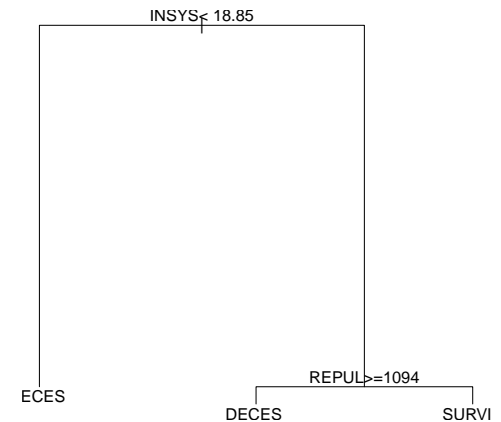
Vizualize a Trees (CART)

A basic viz' of the tree

```
1 > cart <- rpart(PRONO~. , data=myocarde)
2 > plot(cart)
3 > text(cart)
```

Each leaf contains, at least, 20 observations. But we can ask for less

```
4 > cart <- rpart(PRONO~. , data=myocarde, minsplit =
5 5)
6 > plot(cart)
7 > text(cart)
```



Vizualize a Trees (CART)

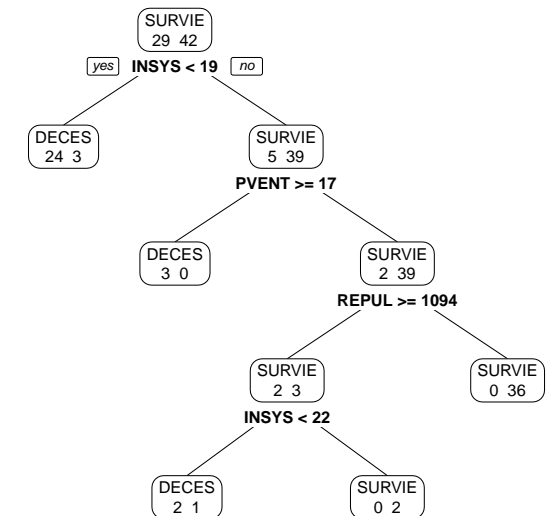
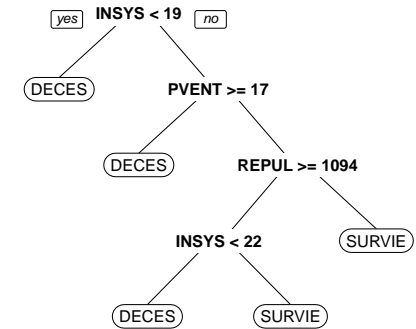
A basic viz' of the tree

```
1 > library ( rpart . plot )
```

```
2 > prp ( cart )
```

OR

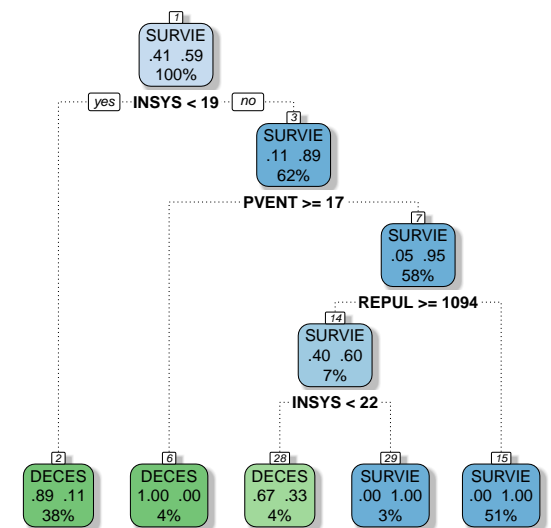
```
3 > prp ( cart , type=2, extra=1)
```



Vizualize a Trees (CART)

A basic viz' of the tree

```
1 > library(rattle)
2 > fancyRpartPlot(cart, sub=" ")
```



Vizualize a Trees (CART)

A basic viz' of the tree

```

1 cart2 <- rpart(PRONO~PVENT+REPUL , data=myocarde ,
  minsplit = 20)
2 pred_cart = function(p,r){return(predict(cart2 ,
  newdata=data.frame(PVENT=p,REPUL=r))[, "SURVIE"
  ])}
3 image(vpvent , vrepul , outer(vpvent , vrepul , pred_cart)
  , col=CL2palette ,
4   xlab="PVENT" , ylab="REPUL" )

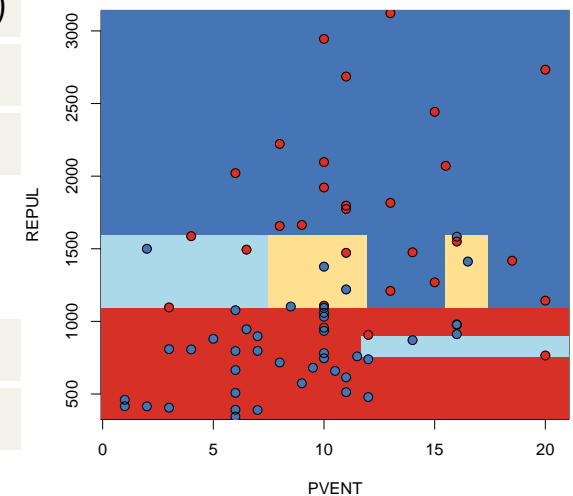
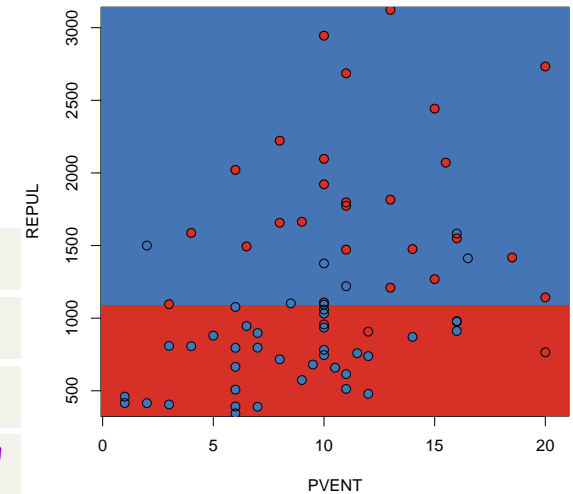
```

OR

```

1 cart2 <- rpart(PRONO~PVENT+REPUL , data=myocarde ,
  minsplit = 5)

```



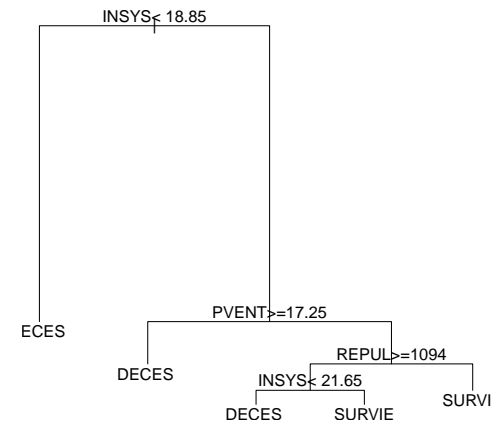
Trees (CART), Gini vs. Entropy

Tree based in **Gini** impurity index

```

1 > cart_gini <- rpart(PRONO~. , data=myocarde ,
  minsplit = 5,parms=list(split="gini"))
2 > summary(cart_gini)
3
4           CP nsplit rel error      xerror      xstd
5 1 0.72413793      0 1.0000000 1.0000000 0.1428224
6 2 0.10344828      1 0.2758621 0.5862069 0.1239921
7 3 0.01724138      2 0.1724138 0.4827586 0.1156044
8 4 0.01000000      4 0.1379310 0.4482759 0.1123721
9 > plot(cart_gini)
10 > text(cart_gini)

```



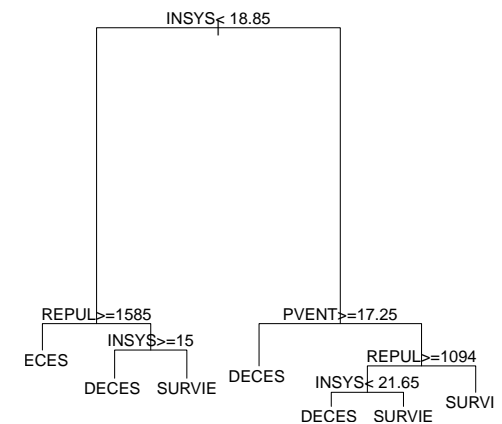
Trees (CART), Gini vs. Entropy

Tree based in the **entropy** impurity index

```

1 > cart_entropy <- rpart(PRONO~. , data=myocarde,
  minsplit = 5,parms=list(split="information"))
2 > summary(cart_entropy)
3
4           CP nsplit  rel error    xerror    xstd
5 1 0.72413793      0 1.00000000 1.0000000 0.1428224
6 2 0.10344828      1 0.27586207 0.5862069 0.1239921
7 3 0.03448276      2 0.17241379 0.4827586 0.1156044
8 4 0.01724138      4 0.10344828 0.4482759 0.1123721
9 5 0.01000000      6 0.06896552 0.4482759 0.1123721
10 > plot(cart_entropy)
11 > text(cart_entropy)

```



Trees (CART), Gini vs. Entropy

Minority class $\min\{p, 1 - p\}$ is the error rate: it measures the proportion of misclassified examples if the leaf was labelled with the majority class

Gini index $2p(1 - p)$ – this is the expected error if we label examples in the leaf randomly: positive with probability p and negative with probability $1-p$.

entropy $-p \log p - (1 - p) \log(1 - p)$ – this is the expected information

Observe that Gini index is related to the variance of a Bernoulli distribution.

With two leaves

$$\frac{n_1}{n} \underbrace{p_1(1 - p_1)}_{\text{var}_1} + \frac{n_2}{n} \underbrace{p_2(1 - p_2)}_{\text{var}_2}$$

is a weighted average variance.

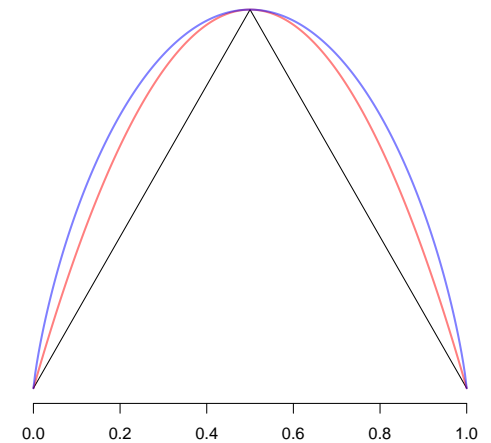
 Regression tree will be obtained by replacing the impurity measure by the variance.

Trees (CART), Gini vs. Entropy

⚠ Dietterich, Kearns and Mansour (1996) suggested to use $\sqrt{\text{Gini}}$ as a measure of impurity. Entropy and Gini index are sensitive to fluctuations in the class distribution, $\sqrt{\text{Gini}}$ isn't. See Drummond & Holte (cs.alberta.ca, 2000) for a discussion on (in)sensitivity of decision tree splitting criteria.

But standard criteria yield (usually) similar trees (rapid-i.com or quora.com)

- misclassification rate $1 - \max\{p, 1 - p\}$
- entropy $-\left[p \log p + (1 - p) \log(1 - p)\right]$
- Gini index $1 - [p^2 + (1 - p)^2] = 2p(1 - p)$



Prunning a Tree

Drop in impurity at node N is defined as

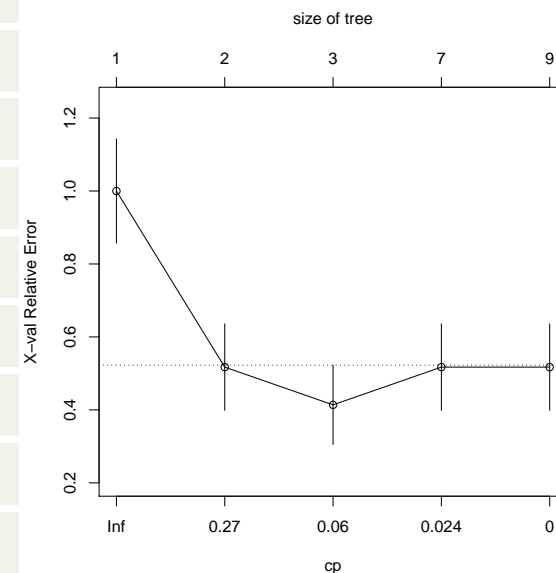
$$\Delta i(N) = i(N) - \mathbb{P}(\text{left})i(N_{\text{left}}) - \mathbb{P}(\text{right})i(N_{\text{right}})$$

If we continue to grow the tree fully until each leaf node corresponds to the lowest impurity, we will overfit. If splitting is stopped too early, error on training data is not sufficiently low and performance will suffer. Thus, stop splitting when the best candidate split at a node reduces the impurity by less than the preset amount, or when a node has a small number of observations

```

1 > cart <- rpart(PRONO~. , data=myocarde, minsplit =
      2, parms=list( split="gini" ), cp=0)
2 > printcp(cart)
3
4 Variables actually used in tree construction:
5 [1] FRCAR INCAR INSYS PVENT REPUL
6
7 Root node error: 29/71 = 0.40845
8
9 n= 71
10
11          CP nsplit rel error  xerror  xstd
12 1 0.724138      0  1.000000  1.00000  0.14282
13 2 0.103448      1  0.275862  0.51724  0.11861
14 3 0.034483      2  0.172414  0.41379  0.10889
15 4 0.017241      6  0.034483  0.51724  0.11861
16 5 0.000000      8  0.000000  0.51724  0.11861
17 > plotcp(cart)

```



Classification Pronostic

For pronostics, consider the [confusion matrix](#)

```
1 > cart <- rpart(PRONO~. , data=myocarde)
2 > predictions <- predict(cart , myocarde , type="class")
3 > table(predictions , myocarde$PRONO)
```

4

```
5 predictions  DECES  SURVIE
```

```
6     DECES      28      6
```

```
7     SURVIE      1     36
```

Classification with Categorical Variables

Consider a **spam classifier**, based on two keywords, **viagra** and **lottery**.

```
1 > load("spam.RData")
```

```
2 > head(db,4)
```

```
3      Y viagra lottery
```

```
4 27 spam      0      1
```

```
5 37 ham       0      1
```

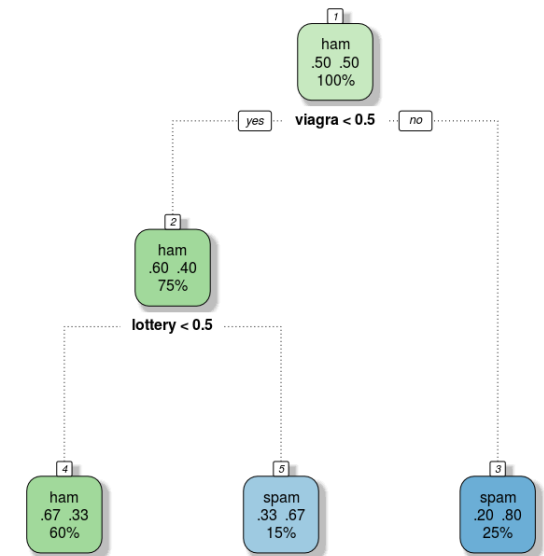
```
6 57 spam      0      0
```

```
7 89 ham       0      0
```

Classification with Categorical Variables

Consider e.g. a tree classifier

```
1 library(rpart)
2 library(rattle)
3 ctrl <- rpart.control(maxdepth=3)
4 model <- rpart(Y~viagra+lottery,
5 data=db, control=ctrl)
6 fancyRpartPlot(model)
```



Classification with C4.5 algorithm

The C4.5 algorithm is an extension of the ID3 algorithm and constructs a decision tree to maximize information gain (difference in entropy).

```

1 > library (RWeka)
2 > C45 <- J48 (PRONO~. , data=myocarde)
3 > summary (C45)
4
5 == Summary ==
6
7 Correctly Classified Instances          66           92.9577 %
8 Incorrectly Classified Instances        5            7.0423 %
9 Kappa statistic                        0.855
10 Mean absolute error                    0.1287
11 Root mean squared error                 0.2537
12 Relative absolute error                 26.6091 %
13 Root relative squared error             51.6078 %
14 Coverage of cases (0.95 level)         97.1831 %

```

```
15 Mean rel. region size (0.95 level)      69.0141 %
```

```
16 Total Number of Instances              71
```

```
17
```

```
18 == Confusion Matrix ==
```

```
19
```

```
20   a   b   <— classified as
```

```
21  27  2 |   a = DECES
```

```
22  3  39 |   b = SURVIE
```

```
1 > predictions <- predict(C45, myocarde)
```

```
2 > table(predictions, myocarde$PRONO)
```

```
3
```

```
4 predictions DECES SURVIE
```

```
5     DECES      27      3
```

```
6     SURVIE      2     39
```

Classification with C4.5 algorithm

To visualise this tree, use

```

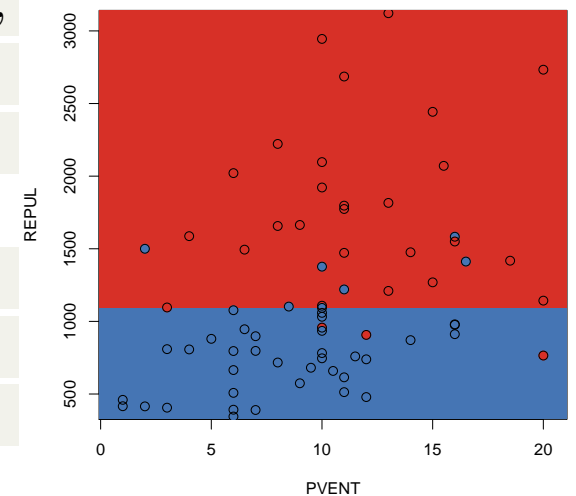
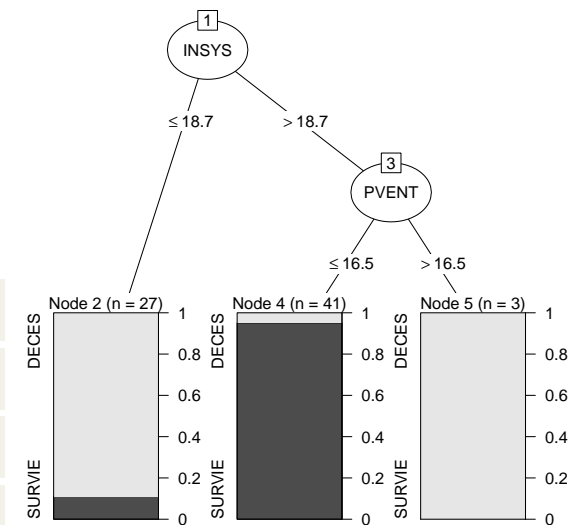
1 C452 <- J48 (PRONO~PVENT+REPUL, data=myocarde)
2 pred_C45 = function (p, r) {
3   return (predict (C452, newdata=
4     data.frame (PVENT=p, REPUL=r) ,
5     type="probability" ) [ , 1] ) }
6 image (vpvent, vrepul, outer (vpvent, vrepul, pred_C45) ,
7   col=CL2palette ,
8   xlab="PVENT" , ylab="REPUL" )

```

```

1 > library (partykit)
2 > plot (C45)
3 > plotcp (cart)

```



From C4.5 algorithm to PART

PART is a rule system that creates pruned C4.5 decision trees for the data set and extracts rules and those instances that are covered by the rules are removed from the training data. The process is repeated until all instances are covered by extracted rules.

```

1 > library (RWeka)
2 > fit_PART <- PART(PRONO~. , data=myocarde)
3 > summary (fit_PART)
4
5 == Summary ==
6
7 Correctly Classified Instances          69           97.1831 %
8 Incorrectly Classified Instances         2           2.8169 %
9 Kappa statistic                        0.9423
10 Mean absolute error                    0.0488
11 Root mean squared error                 0.1562
12 Relative absolute error                 10.0944 %

```

```
13 Root relative squared error          31.7864 %
14 Coverage of cases (0.95 level)      100      %
15 Mean rel. region size (0.95 level)  60.5634 %
16 Total Number of Instances          71
```

```
17
18 == Confusion Matrix ==
```

```
19
20 a  b  ←— classified as
```

```
21 29  0 | a = DECES
```

```
22  2 40 | b = SURVIE
```

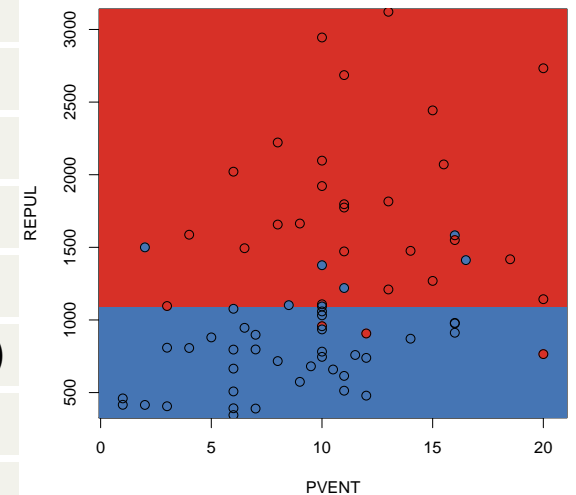

From C4.5 algorithm to PART

The code to visualise this tree is

```


1 PART2 <- PART(PRONO~PVENT+REPUL, data=myocarde)
2 pred_PART = function(p, r){
3   return(predict(PART2, newdata=
4     data.frame(PVENT=p, REPUL=r),
5     type="probability")[,1])}
6 image(vpvent, vrepul, outer(vpvent, vrepul, pred_PART)
7   , col=CL2palette,
8   xlab="PVENT", ylab="REPUL")

```



Bootstrap and Bagging

Bootstrapped Aggregation (Bagging) is an ensemble method that creates multiple models of the same type from different sub-samples of the same dataset. The predictions from each separate model are combined together to provide a superior result.

 the probability that a particular data point is not selected for a bootstrap sample of size n as $(1 - 1/n)^n \rightarrow e^{-1} \sim 36.8\%$: each bootstrap sample is likely to leave out about a third of the data points.

The Intuition of Bagging Algorithm

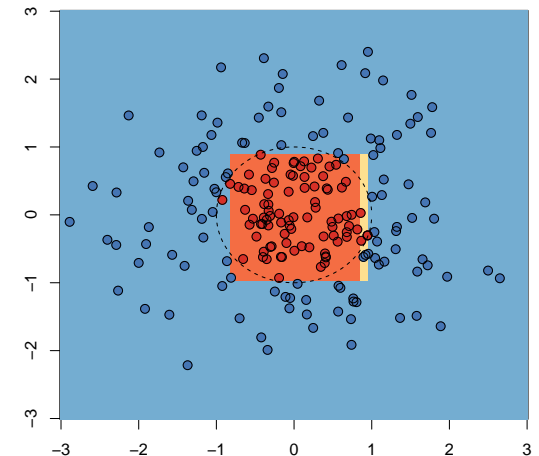
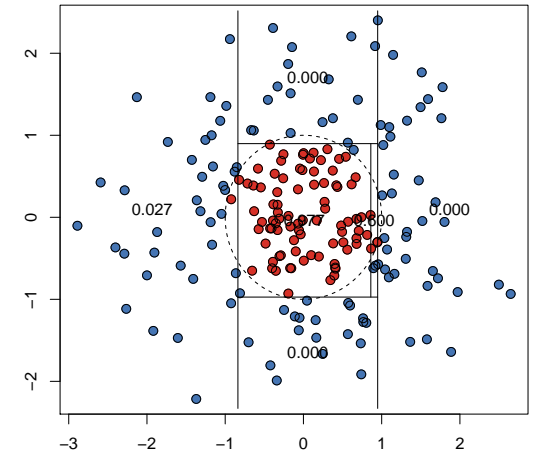
Tree cannot be used in the context of nonlinear pattern (without considering nonlinear transformation of variables)

```

1 > set.seed(1)
2 > X1 <- rnorm(200)
3 > X2 <- rnorm(200)
4 > Y <- (X1^2+X2^2)<=1
5 > df <- data.frame(Y,X1,X2)

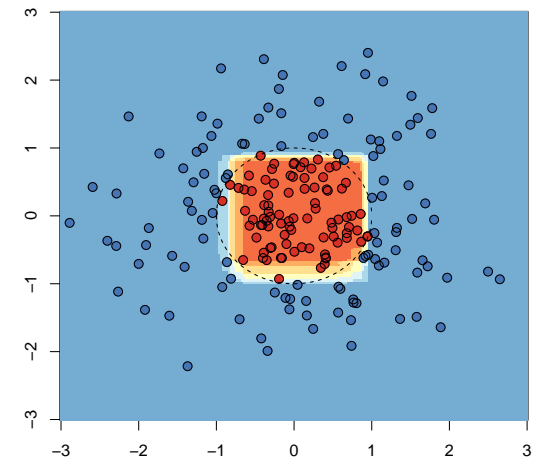
```

A strategy can be to aggregate several trees...



The Intuition of Bagging Algorithm

```
1 Y_pred <- matrix(NA, length(Y), 100)
2 for (b in 1:100) {
3   ind <- sample(1:n, size=n, replace=TRUE)
4   tree <- rpart(Y~X1+X2, data=df[ind,])
5   Y_pred[,b] <- predict(tree)}
6 Y_agg <- apply(Y_pred, 1, mean)
```



From bagging to Random Forest

Strictly speaking, when bootstrapping among observations, and aggregating, we use a bagging algorithm.

In the [random forest](#) algorithm, we combine Breiman's [bagging](#) idea and the random selection of features, introduced independently by Ho (1995, [cm.bell-labs.com](#)) and Amit & Geman (1997, [cis.jhu.edu](#))

This process is sometimes called feature bagging : typically, for a dataset with p features, \sqrt{p} features are used in each split.

Bagging for Classification

This approach has shown participtionally effective for high-variance methods such as decision trees.

```

1 > library(ipred)
2 > fit_bag <- bagging(PRONO~. , data=myocarde)
3 > summary(fit_bag)
4
5 Bagging classification trees with 25 bootstrap replications

```

Several trees are generated

```

1 $mtrees
2 $mtrees[[1]]
3 $bindx
4 [1] 13 17 44 28 33 20 2 42 43 46 28 48 27 2 23 4 2 10 14 69 33
67 42 38 52 42 16 41 54 28 32 63 3 27 49 68 12 47 52
5 [40] 1 57 20 60 15 7 71 21 53 57 29 40 14 22 68 54 57 30 55 35 31
42 4 31 27 13 48 13 42 2 44 9

```

```
6
7 $btree
8 n= 71
9
10 node), split , n, loss , yval , (yprob)
11 * denotes terminal node
12
13 1) root 71 33 SURVIE (0.46478873 0.53521127)
14   2) INSYS< 18.85 36 3 DECES (0.91666667 0.08333333)
15     4) INSYS>=8.85 32 0 DECES (1.00000000 0.00000000) *
16     5) INSYS< 8.85 4 1 SURVIE (0.25000000 0.75000000)
17       10) FRCAR< 106.5 1 0 DECES (1.00000000 0.00000000) *
18       11) FRCAR>=106.5 3 0 SURVIE (0.00000000 1.00000000) *
19     3) INSYS>=18.85 35 0 SURVIE (0.00000000 1.00000000) *
20
21 attr(,"class")
22 class
23 "sclass"
```

```

1 > predictions <- predict(fit_bag, myocarde, type="class")
2 > table(predictions, myocarde$PRONO)
3
4 predictions DECES SURVIE
5     DECES      28      0
6     SURVIE      1     42

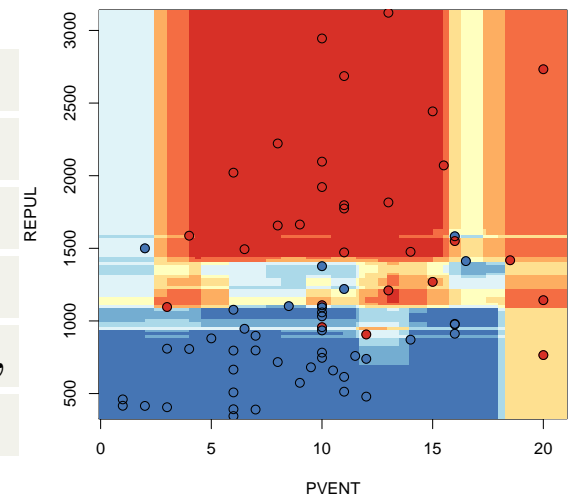
```

The output is

```

1 bag2 <- bagging(PRONO~PVENT+REPUL, data=myocarde)
2 pred_bag = function(p,r){
3   return(predict(bag2, newdata=
4     data.frame(PVENT=p, REPUL=r), type="prob")[,1]) }
5 image(vpvent, vrepul, outer(vpvent, vrepul, pred_bag),
6     col=CL2palette)

```



Random Forests for Classification

To grow forest of trees, one can consider

```

1 > library(randomForest)
2 > RF <- randomForest(PRONO~. , data=myocarde)
3 > summary(RF)
4
5          Length Class  Mode
6 call           3  -none- call
7 type           1  -none- character
8 predicted       71  factor numeric
9 err.rate       1500 -none- numeric
10 confusion       6  -none- numeric
11 votes          142  matrix numeric
12 oob.times       71  -none- numeric
13 classes         2  -none- character
14 importance       7  -none- numeric
15 importanceSD     0  -none- NULL
16 localImportance  0  -none- NULL
17 proximity         0  -none- NULL

```

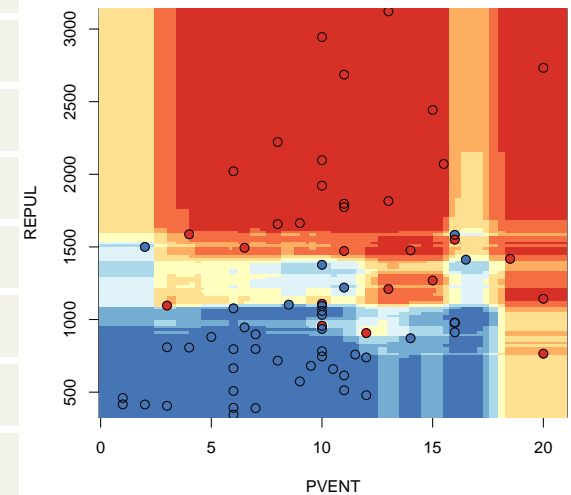
```
17 ntree          1  -none- numeric
18 mtry          1  -none- numeric
19 forest       14  -none- list
20 y            71  factor numeric
21 test         0  -none- NULL
22 inbag        0  -none- NULL
23 terms        3  terms  call
24 > predictions <- predict(RF, myocarde, type="class")
25 > table(predictions, myocarde$PRONO)
26
27 predictions DECES SURVIE
28     DECES    29     0
29     SURVIE     0    42
```

Random Forests for Classification

```

1 RF2 <- randomForest(PRONO~PVENT+REPUL, data=
  myocarde)
2 pred_RF = function(p, r) {
3   return(predict(RF2, newdata=
4     data.frame(PVENT=p, REPUL=r),
5     type="prob")[,1]) }
6 image(vpvent, vrepul, outer(vpvent, vrepul, pred_RF),
7   col=CL2palette)

```



Gradient Boosting for Classification

There is a dedicated package for Gradient Boosting

```
1 > library(gbm)
2 > fit_gbm <- gbm(PRONO~. , data=myocarde, distribution="multinomial")
3 > print(fit_gbm)
4 gbm(formula = PRONO ~ ., distribution = "multinomial", data =
  myocarde)
5 A gradient boosted model with multinomial loss function.
6 100 iterations were performed.
7 There were 7 predictors of which 3 had non-zero influence.
```

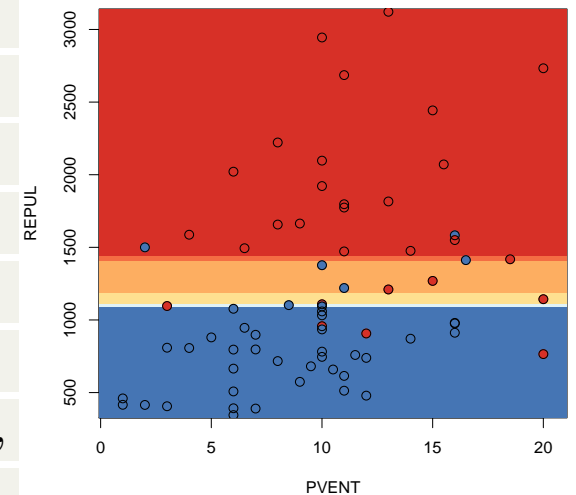
💡 This technique will be explained in [slides #4](#).

Gradient Boosting for Classification

```

1 gbm2 <- gbm(PRONO~PVENT+REPUL, data=myocarde,
  distribution="multinomial")
2 pred_gbm = function(p,r){proba=predict(gbm2,
  newdata=
3 data.frame(PVENT=p,REPUL=r),n.trees=100,type="
  response")}
4 return(matrix(proba,ncol=2)[,1])}
5 image(vpvent,vrepul,outer(vpvent,vrepul,pred_gbm),
  col=CL2palette)

```



C5.0 for Classification

```
1 > library(C50)
2 > C50 <- C5.0(PRONO~. , data=myocarde , trials=10)
3 > print(C50)
4
5 Call:
6 C5.0.formula(formula = PRONO ~ ., data = myocarde, trials = 10)
7
8 Classification Tree
9 Number of samples: 71
10 Number of predictors: 7
11
12 Number of boosting iterations: 10
13 Average tree size: 4.3
14
15 Non-standard options: attempt to group attributes
16
17 > predictions <- predict(C50, myocarde, type="class")
```

```
18 > table(predictions , myocarde$PRONO)
```

```
19
```

```
20 predictions DECES SURVIE
```

```
21     DECES      29      0
```

```
22     SURVIE      0     42
```

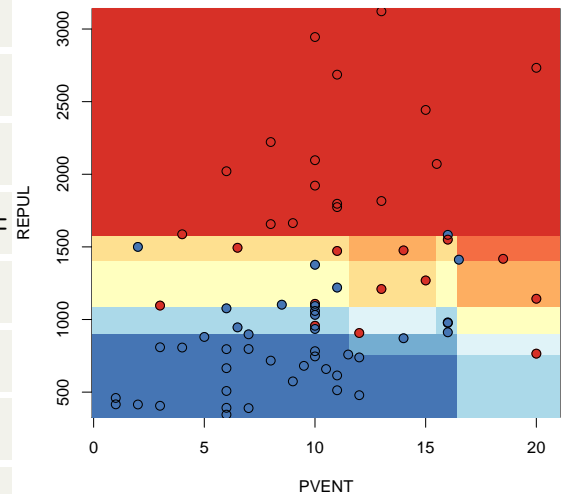
```
1 C502 <- C5.0(PRONO~PVENT+REPUL, data=myocarde ,
  trials=10)
```

```
2 pred_C50 = function(p,r){proba=predict(C502,
  newdata=
```

```
3     data.frame(PVENT=p,REPUL=
  r) , type="prob" )
```

```
4     return(proba[,1]) }
```

```
5 image(vpvent , vrepul , outer(vpvent , vrepul , pred_C50) ,
  col=CL2palette)
```



Support Vector Machine and Vapnik

SVMs were developed in the 90's based on previous work, from Vapnik & Lerner (1963, cs.iastate.edu, see Vailant (1984, people.mpi-inf.mpg.de)

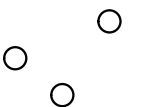
Assume that points are **linearly separable**, i.e. there is ω and b such that

$$Y = \begin{cases} +1 & \text{if } \omega^T \mathbf{x} + b > 0 \\ -1 & \text{if } \omega^T \mathbf{x} + b < 0 \end{cases}$$

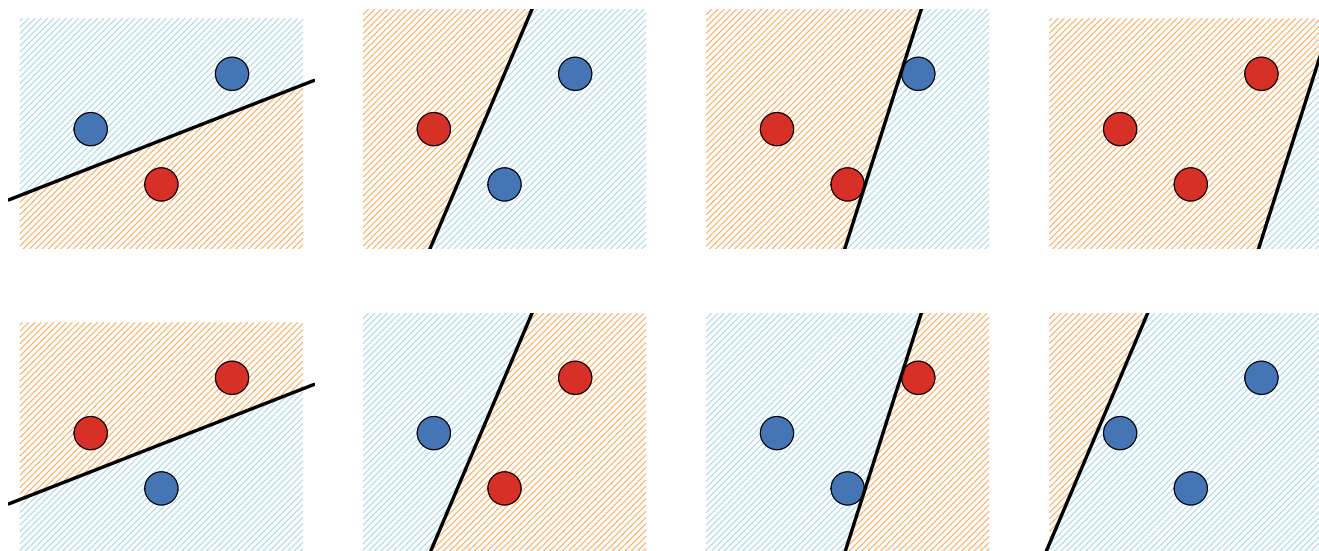
Problem: infinite number of solutions, need a **good** one, that separate the data, (somehow) far from the data.

Concept : **VC dimension**. Let $\mathcal{H} : -h : \mathbb{R}^d \mapsto \{-1, +1\}$. Then \mathcal{H} is said to **shatter** a set of points \mathbf{X} if all dichotomies can be achieved.

E.g. with those three points, all configurations can be achieved



Support Vector Machine and Vapnik



E.g. with those four points, several configurations **cannot** be achieved (with some linear separator, but they can with some quadratic one)

Support Vector Machine and Vapnik

Vapnik's (VC) dimension is the size of the largest shattered subset of \mathbf{X} .

This dimension is interesting to get an upper bound of the probability of miss-classification (with some complexity penalty, function of $VC(\mathcal{H})$).

Now, in practice, where is the optimal hyperplane ?

The distance from \mathbf{x}_0 to the hyperplane $\boldsymbol{\omega}^\top \mathbf{x} + b$ is

$$d(\mathbf{x}_0, H_{\boldsymbol{\omega}, b}) = \frac{\boldsymbol{\omega}^\top \mathbf{x}_0 + b}{\|\boldsymbol{\omega}\|}$$

and the **optimal hyperplane** (in the separable case) is

$$\operatorname{argmin} - \min_{i=1, \dots, n} d(\mathbf{x}_i, H_{\boldsymbol{\omega}, b}) "$$

Support Vector Machine and Vapnik

Define **support vectors** as observations such that

$$|\boldsymbol{\omega}^\top \mathbf{x}_i + b| = 1$$

The margin is the distance between hyperplanes defined by support vectors.

The distance from support vectors to $H_{\boldsymbol{\omega}, b}$ is $\|\boldsymbol{\omega}\|^{-1}$, and the margin is then $2\|\boldsymbol{\omega}\|^{-1}$.

💡 the algorithm is to minimize the inverse of the margins s.t. $H_{\boldsymbol{\omega}, b}$ separates ± 1 points, i.e.

$$\min -\frac{1}{2}\boldsymbol{\omega}^\top \boldsymbol{\omega} \quad \text{s.t.} \quad Y_i(\boldsymbol{\omega}^\top \mathbf{x}_i + b) \geq 1, \quad \forall i.$$

Support Vector Machine and Vapnik

Problem difficult to solve: many inequality constraints (n)

💡 solve the dual problem...

In the **primal space**, the solution was

$$\omega = \sum_{i=1} \alpha_i Y_i \mathbf{x}_i \text{ with } \sum_{i=1} \alpha_i Y_i = 0.$$

In the **dual space**, the problem becomes (hint: consider the Lagrangian)

$$\max_{\alpha} - \sum_{i=1} \alpha_i - \frac{1}{2} \sum_{i=1} \alpha_i \alpha_j Y_i Y_j \mathbf{x}_i^T \mathbf{x}_j \text{ " s.t. } \sum_{i=1} \alpha_i Y_i = 0.$$

which is usually written

$$\min_{\alpha} - \frac{1}{2} \alpha^T Q \alpha - \mathbf{1}^T \alpha \text{ " s.t. } \begin{cases} 0 \leq \alpha_i \forall i \\ \mathbf{y}^T \alpha = 0 \end{cases}$$

where $Q = [Q_{i,j}]$ and $Q_{i,j} = y_i y_j \mathbf{x}_i^T \mathbf{x}_j$.

Support Vector Machine and Vapnik

Now, what about the **non-separable case**?

Here, we **cannot** have $y_i(\omega^\top \mathbf{x}_i + b) \geq 1 \forall i$.

💡 introduce **slack variables**,

$$\begin{aligned} & \omega^\top \mathbf{x}_i + b \geq +1 - \xi_i \text{ when } y_i = +1 \\ & \omega^\top \mathbf{x}_i + b \leq -1 + \xi_i \text{ when } y_i = -1 \end{aligned}$$

where $\xi_i \geq 0 \forall i$. There is a classification error when $\xi_i > 1$.

The idea is then to solve

$$\min -\frac{1}{2}\omega^\top \omega + C \mathbf{1}^\top \mathbf{1}_{\xi > 1} \text{ ", instead of } \min -\frac{1}{2}\omega^\top \omega \text{ "}$$

Support Vector Machine and Vapnik

Here C is related to some - standard - tradeoff

- large C will penalize errors,
- small C will penalize complexity.

Note that the dual problem here is the same as the one before, with additional constraint, $0 \leq \alpha_i \leq C$.

$$\min_{\alpha} -\frac{1}{2} \alpha^T Q \alpha - \mathbf{1}^T \alpha \quad \text{s.t.} \quad \begin{array}{l} 0 \leq \alpha_i \leq C \quad \forall i \\ \mathbf{y}^T \alpha = 0 \end{array}$$

with $C \geq 0$ and $Q = [Q_{i,j}]$ where $Q_{i,j} = y_i y_j \mathbf{x}_i^T \mathbf{x}_j$.

 it is possible to consider some more general function here, instead of $\mathbf{x}_i^T \mathbf{x}_j$

Support Vector Machine and C -classification

$$\min_{\alpha} -\frac{1}{2} \alpha^T Q \alpha - \mathbf{1}^T \alpha \quad \text{s.t.} \quad \begin{cases} 0 \leq \alpha_i \leq C \quad \forall i \\ \mathbf{y}^T \alpha = 0 \end{cases}$$

where $C \geq 0$ is the upper bound, K is a Kernel, e.g.

- linear $K(\mathbf{u}, \mathbf{v}) = \mathbf{u}^T \mathbf{v}$
- polynomial $K(\mathbf{u}, \mathbf{v}) = \gamma[\mathbf{u}^T \mathbf{v} + c_0]^d$
- radial basis $K(\mathbf{u}, \mathbf{v}) = \exp(-\gamma \|\mathbf{u} - \mathbf{v}\|^2)$

and

$$Q = [Q_{i,j}] \quad \text{where} \quad Q_{i,j} = y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

- **Support Vector Machine and ν -classification**

$$\min_{\alpha} -\frac{1}{2} \alpha^T Q \alpha \quad \text{s.t.} \quad \begin{array}{l} 0 \leq \alpha_i \leq \frac{1}{d} \quad \forall i \\ \mathbf{y}^T \alpha = 0 \\ \mathbf{1}^T \alpha \geq \nu \end{array}$$

with $\nu \in (0, 1]$

- **Support Vector Machine and one class-classification**

$$\min_{\alpha} -\frac{1}{2} \alpha^T Q \alpha \quad \text{s.t.} \quad \begin{array}{l} 0 \leq \alpha_i \leq \frac{1}{\nu d} \quad \forall i \\ \mathbf{1}^T \alpha = 1 \end{array}$$

- **Support Vector Machine and ϵ -regression**

$$\min_{\alpha, \alpha^*} -\frac{1}{2} [\alpha - \alpha^*]^T Q [\alpha - \alpha^*] + \epsilon \mathbf{1}^T [\alpha + \alpha^*] + \mathbf{1}^T [\mathbf{y} \cdot (\alpha - \alpha^*)] "$$

subject to

$$\begin{aligned} & 0 \leq \alpha_i, \alpha_i^* \leq C \quad \forall i \\ & \mathbf{1}^T [\alpha - \alpha^*] = 0 \end{aligned}$$

- **Support Vector Machine and ν -regression**

$$\min_{\alpha, \alpha^*} -\frac{1}{2} [\alpha - \alpha^*]^T Q [\alpha - \alpha^*] + z^T [(\alpha - \alpha^*)] " \text{ s.t. } \begin{aligned} & 0 \leq \alpha_i, \alpha_i^* \leq C \quad \forall i \\ & \mathbf{1}^T [\alpha - \alpha^*] = 0 \\ & \mathbf{1}^T [\alpha + \alpha^*] = C\nu \end{aligned}$$

From Support Vector Machine to Perceptron

SVM's belong to the class of **linear classifiers**.

A linear classifier is define as

$$Y^*(\mathbf{x}) = \begin{cases} +1 & \text{if } B(\mathbf{x}) = \beta_0 + \mathbf{x}^\top \boldsymbol{\beta} > 0 \\ -1 & \text{otherwise} \end{cases}$$

Data as **linearly separable** if there is a hyperplane that separate (perfectly) the two classes.

- observations such that $y_i B(\mathbf{x}_i) \geq 0$ are correctly classified.
- observations such that $y_i B(\mathbf{x}_i) \leq 0$ are misclassified.

💡 Consider the separating hyperplane such that

$$B^* = \operatorname{argmin} \sum_{\text{misclassified}} y_i B(\mathbf{x}_i) "$$

From Support Vector Machine to Perceptron

The **perceptron** algorithm, introduced by Rosenblatt starts with some initial values, and then, we update

$$\begin{aligned} \beta_0 &\leftarrow \beta_0 + Y_i \\ \beta &\leftarrow \beta + Y_i \cdot \mathbf{X}_i \end{aligned}$$

 The convergence of this algorithm depends on starting values

In case of convergence, the resulting hyperplane is the maximum margin hyperplane, and points on the boundary of the margins are called support vector.

From Support Vector Machine to Perceptron

```
1 x=c(.4,.55,.65,.9,.1,.35,.5,.15,.2,.85)
2 y=c(.85,.95,.8,.87,.5,.55,.5,.2,.1,.3)
3 z=c(1,1,1,1,1,0,0,1,0,0)
4 z_sign=z*2-1
5 beta=c(0,-1,1)
6 for(i in 1:k){
7 beta=beta+c(z_sign[i],z_sign[i]*x[i],z_sign[i]*y[i]
  ])}

```

Support Vector Machines (SVM) are a method that uses points in a transformed problem space that best separate classes into two groups. Classification for multiple classes is supported by a one-vs-all method. SVM also supports regression by modeling the function with a minimum amount of allowable error.

```
1 > library(kernlab)
2 > SVM <- ksvm(PRONO~., data=myocarde)
3 Using automatic sigma estimation (sigest) for RBF or laplace kernel
4 > SVM
5 Support Vector Machine object of class "ksvm"
6
7 SV type: C-svc (classification)
8 parameter : cost C = 1
9
10 Gaussian Radial Basis kernel function.
11 Hyperparameter : sigma = 0.146414435486797
12
13 Number of Support Vectors : 41
14
```

```
15 Objective Function Value : -23.9802
16 Training error : 0.070423
17 > predictions <- predict(SVM, myocarde, type="response")
18 > table(predictions, myocarde$PRONO)
19
20 predictions DECES SURVIE
21     DECES      25      1
22     SURVIE      4     41
```

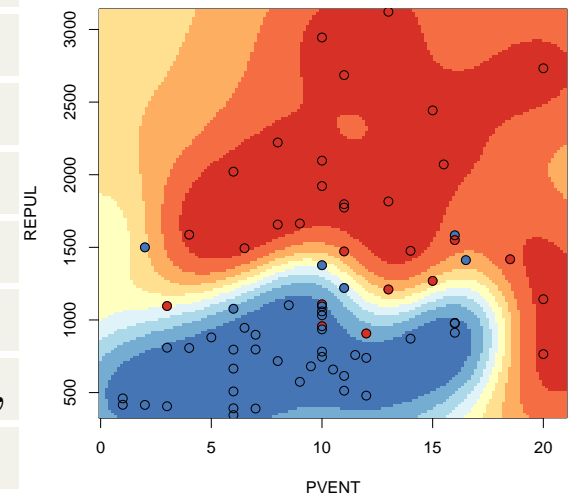
Visualising a SVM

Consider a SVM with 2 covariates

```

1 SVM2 <- ksvm(PRONO~PVENT+REPUL, data=myocarde, prob
  .model =TRUE)
2 pred_SVM = function(p, r) {
3   return(predict(SVM2, newdata=
4     data.frame(PVENT=p, REPUL=r),
5     type="probabilities")[,1]) }
6 image(vpvent, vrepul, outer(vpvent, vrepul, pred_SVM),
  col=CL2palette,
  xlab="PVENT", ylab="REPUL")

```



A short remark on kernels

A kernel is a function $k(x_i, x_j) = \varphi(x_i) \cdot \varphi(x_j)$ that calculates a dot product in some feature space, but without constructing the feature vectors $\varphi(x)$ explicitly.

Consider the Euclidean ℓ_2 distance

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_{\ell_2} = \sqrt{(\mathbf{x} - \mathbf{y}) \cdot (\mathbf{x} - \mathbf{y})} = \sqrt{\mathbf{x} \cdot \mathbf{x} - 2\mathbf{x} \cdot \mathbf{y} + \mathbf{y} \cdot \mathbf{y}}$$

More generally, define

$$d_k(\mathbf{x}, \mathbf{y}) = \sqrt{k(\mathbf{x}, \mathbf{x}) - 2k(\mathbf{x}, \mathbf{y}) + k(\mathbf{y}, \mathbf{y})}$$

Neural Network and Classification

A Neural Network (NN) is a graph of computational units that receive inputs and transfer the result into an output that is passed on. The units are ordered into layers to connect the features of an input vector to the features of an output vector. With training, such as the Back-Propagation algorithm, neural networks can be designed and trained to model the underlying relationship in data.

```
1 > library(nnet)
2 > NN <- nnet(PRONO~., data=myocarde, size=10, decay=0.0001, maxit
  =500)
3 # weights: 91
4 initial value 48.811990
5 iter 10 value 48.017252
6 iter 20 value 30.627189
7 iter 30 value 30.216558
8 iter 40 value 26.974742
```

...

```
1 iter 490 value 6.801462
2 iter 500 value 6.797179
3 final value 6.797179
4 stopped after 500 iterations
```

```
1 > predictions <- predict(NN, myocarde, type="class")
2 > table(predictions, myocarde$PRONO)
```

```
3
4 predictions DECES SURVIE
5     DECES      27      1
6     SURVIE      2     41
```

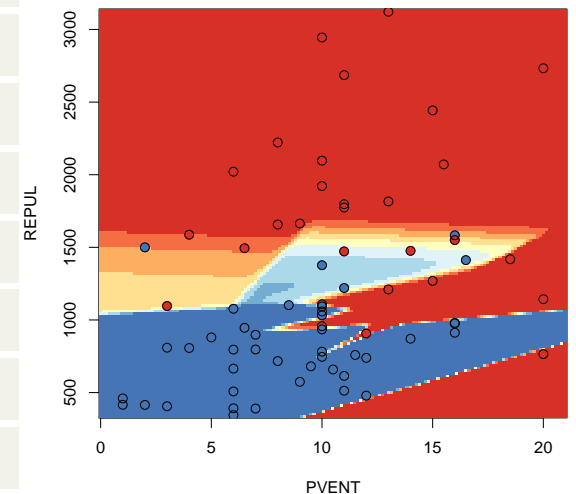
Visualising a Neural Network

Consider a NN with 2 covariates

```

1 NN2 <- nnet(PRONO~PVENT+REPUL , data=myocarde , size
  =20, decay=0.0001, maxit=500)
2 pred_NN = function(p,r) {
3   return(1 - predict(NN2, newdata=
4     data.frame(PVENT=p, REPUL=r) ,
5     type="raw" ) [ ,1] ) }
6 image(vpvent , vrepul , outer(vpvent , vrepul , pred_NN) ,
7   col=CL2palette ,
8   xlab="PVENT" , ylab="REPUL" )

```



Classification Output : terminology

	Y = 0	Y = 1	prevalence P=	
$\widehat{Y} = 0$	true negative N_{00}	false negative (type II) N_{01}	negative predictive value $NPV = \frac{N_{00}}{N_{.0}}$	false omission rate $FOR = \frac{N_{00}}{N_{.0}}$
$\widehat{Y} = 1$	false positive (type I) N_{10}	true positive N_{11}	false discovery rate $FDR = \frac{N_{00}}{N_{.0}}$	positive predictive value $PPV = \frac{N_{11}}{N_{.1}}$ (precision)
negative likelihood ratio $LR- = FNR/TNR$	true negative rate $TNR = \frac{N_{00}}{N_{.0}}$ (specificity)	false negative rate $FNR = \frac{N_{01}}{N_{.1}}$		
positive likelihood ratio $LR+ = TPR/FPR$ diagnostic odds ratio = $LR+/LR-$	false positive rate $FPR = \frac{N_{10}}{N_{.0}}$ (fall out)	true positive rate $TPR = \frac{N_{11}}{N_{.1}}$ (sensitivity)		

Measures and Curves

- sensitivity or true positive rate (TPR) eqv. with hit rate, recall

$$TPR = TP/P = TP/(TP + FN)$$

- specificity (SPC) or True Negative Rate

$$SPC = TN/N = TN/(FP + TN)$$

- precision or positive predictive value (PPV)

$$PPV = TP/(TP + FP)$$

- negative predictive value (NPV)

$$NPV = TN/(TN + FN)$$

- fall-out or false positive rate (FPR)

$$FPR = FP/N = FP/(FP + TN) = 1 - SPC$$

Measures and Curves

- false discovery rate (FDR)

$$FDR = FP / (FP + TP) = 1 - PPV$$

- miss Rate or False Negative Rate (FNR)

$$FNR = FN / P = FN / (FN + TP)$$

- accuracy (ACC)

$$ACC = (TP + TN) / (P + N)$$

A ROC Curve

Given a classifier (and a score function), the ROC curve (receiver operating characteristic) is obtained by plotting the **true positive rate** against the **false positive rate** at various threshold s settings

```

1 > roc.curve <- function(Y,S,s , print=FALSE) {
2   Ps <- (S>s) * 1
3   FP <- sum((Ps==1)*(Y==0)) / sum(Y==0)
4   TP <- sum((Ps==1)*(Y==1)) / sum(Y==1)
5   if (print==TRUE) {
6     print(table( Observed=Y, Predicted=Ps)) }
7   vect=c(FP,TP)
8   names(vect)=c("FPR" , "TPR" )
9   return(vect) }

```

 ROC = sensitivity as a function of fall-out

A ROC Curve

Construction of ROC curve points given a confusion matrix

The Area Under the Curve, AUC, can be interpreted as the probability that a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one

💡 AUC is related to Gini coefficient, since $G = 2AUC - 1$.

See Swets, Dawes & Monahan (ist-socrates.berkeley.edu, 2000)

Classification with Categorical Variables

Consider a spam classifier, based on two keywords, **viagra** and **lottery**.

```
1 > load("spam.RData")
```

```
2 > head(db,4)
```

```
3      Y viagra lottery
```

```
4 27 spam      0      1
```

```
5 37 ham      0      1
```

```
6 57 spam      0      0
```

```
7 89 ham      0      0
```

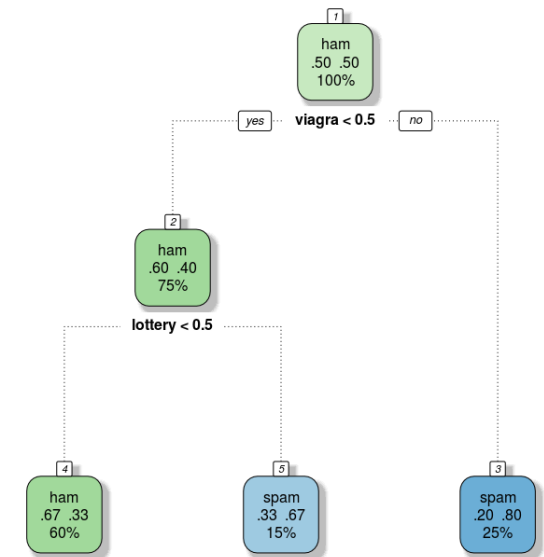
Classification with Categorical Variables

Consider e.g. a tree classifier

```

1 library(rpart)
2 library(rattle)
3 ctrl <- rpart.control(maxdepth=3)
4 model <- rpart(Y~viagra+lottery,
5 data=db, control=ctrl)
6 fancyRpartPlot(model)

```



Whatever the classifier consider ROC curve will go through 4 very specific points.

Classification with Continuous Variables

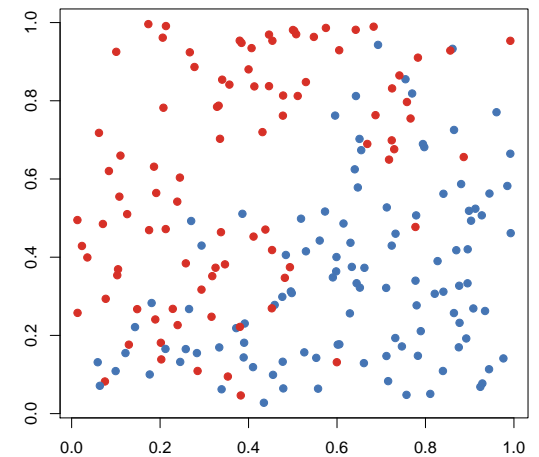
Classification with Categorical Variables

- 💡 interpolation between those points have different meanings
see convexification of ROC curves, see Flach (cs.bris.ac.uk, 2012)

Comparing ROC Curves

Consider a simple simulated dataset (just to illustrate)

```
1 > X <- runif(200)
2 > Y <- runif(200)
3 > Z <- (Y > (X + rnorm(200) / 4))
4 > df <- data.frame(X, Y, Z)
5 > plot(X, Y, col=c("red", "blue")[1+Z*1], pch=19)
```

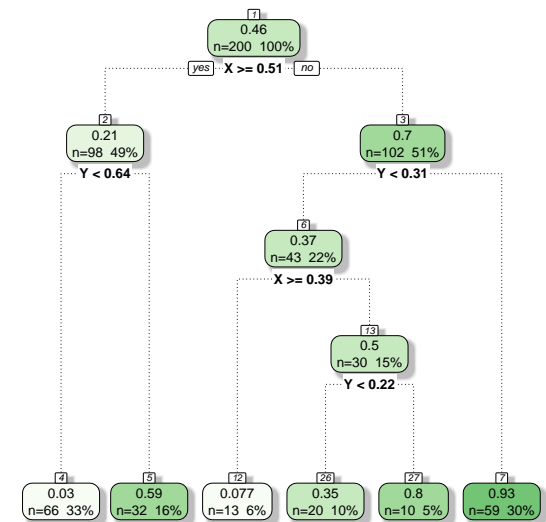
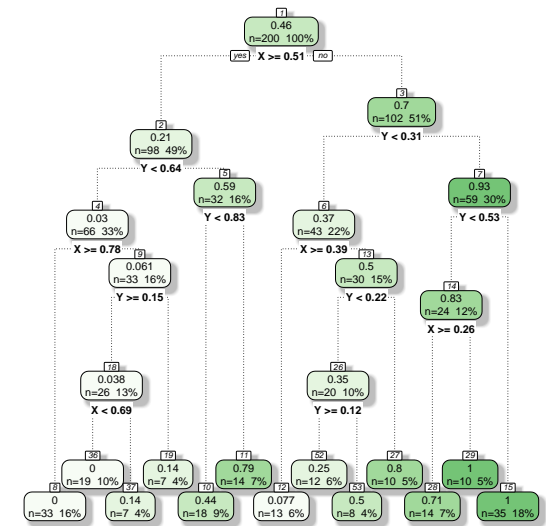


Comparing ROC Curves

On that dataset, consider two trees

```

1 > library(rpart)
2 > model1 <- rpart(Z~X+Y, data=df, control=rpart.
  control(cp = 0.00035, maxcompete = 2))
3 > model2 <- rpart(Z~X+Y, data=df, control=rpart.
  control(cp = 0.025, maxcompete = 2))
4 > library(rattle)
5 > fancyRpartPlot(model1)
6 > fancyRpartPlot(model2)
7 > df$s1 <- predict(modelprune1)
8 > df$s2 <- predict(modelprune2)
    
```



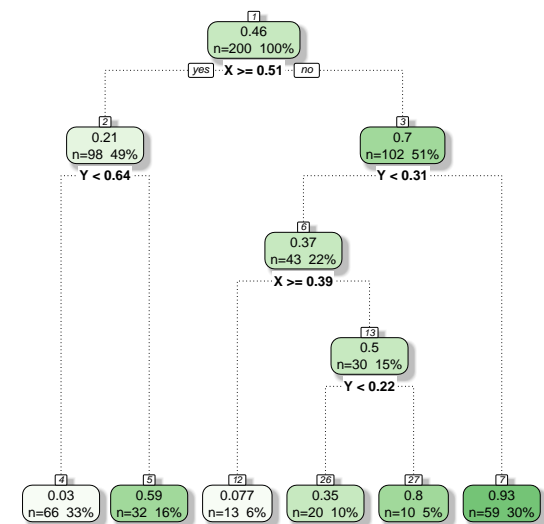
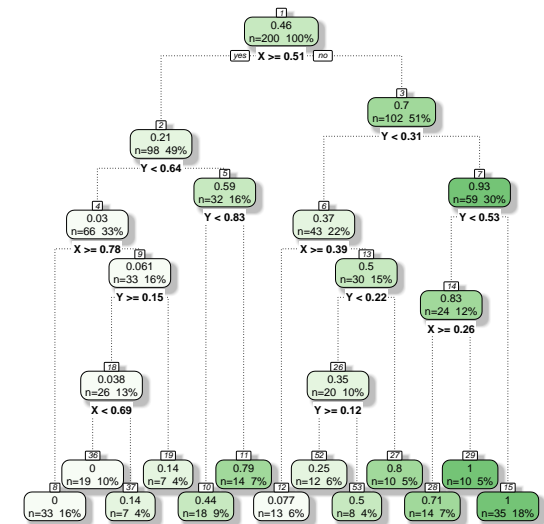
Comparing ROC Curves

On that dataset, consider two trees

```

1 > library(rpart)
2 > model1 <- rpart(Z~X+Y, data=df, control=rpart.
  control(cp = 0.00035, maxcompete = 2))
3 > model2 <- rpart(Z~X+Y, data=df, control=rpart.
  control(cp = 0.025, maxcompete = 2))
4 > library(rattle)
5 > fancyRpartPlot(model1)
6 > fancyRpartPlot(model2)
7 > df$s1 <- predict(modelprune1)
8 > df$s2 <- predict(modelprune2)

```



Comparing ROC Curves

Consider two thresholds

```

1 > s1 <- .5
2 > Ps1 <- (df$s1>s)*1
3 > Y <- df$Z
4 > FP=sum((Ps1==1)*(Y==0))/
5 sum(Y==0)
6 > TP=sum((Ps1==1)*(Y==1))/
7 sum(Y==1)
8 > table( Observed=Y, Predicted=P1)
9
10 Predicted
11 Observed 0 1
12 FALSE 99 9
TRUE 18 74

```

```

1 > s2 <- .5
2 > Ps2 <- (df$s2>s)*1
3 > Y <- df$Z
4 > FP=sum((Ps2==1)*(Y==0))/
5 sum(Y==0)
6 > TP=sum((Ps2==1)*(Y==1))/
7 sum(Y==1)
8 > table( Observed=Y, Predicted=Ps2)
9
10 Predicted
11 Observed 0 1
12 FALSE 89 19
TRUE 10 82

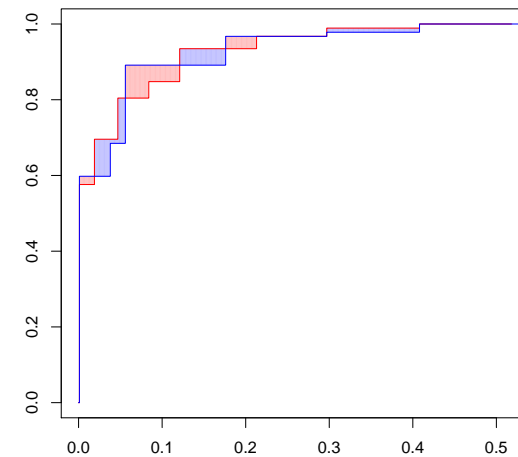
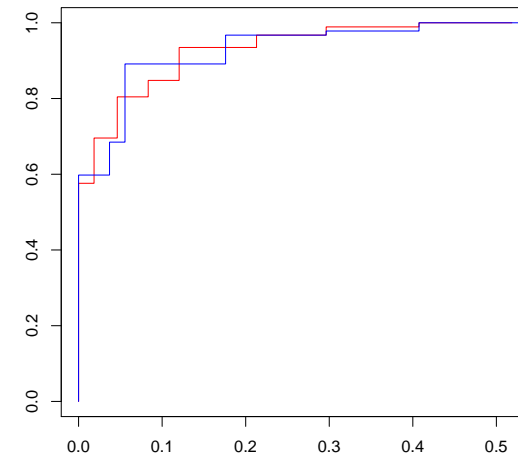
```

We have a (standard) tradeoff between type I and type II errors.

Comparing ROC Curves

```
1 > plot(roc(Z, s1, data=df), col="red")  
2 > plot(roc(Z, s2, data=df), col="blue")
```

In the case of trees, this comparison is artificial, since only (lower) corners of the curves can be reached.

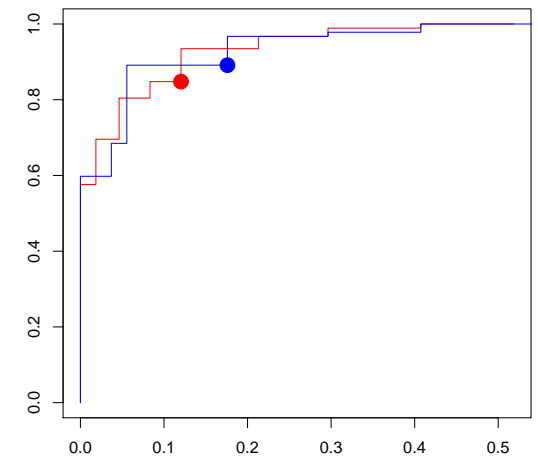
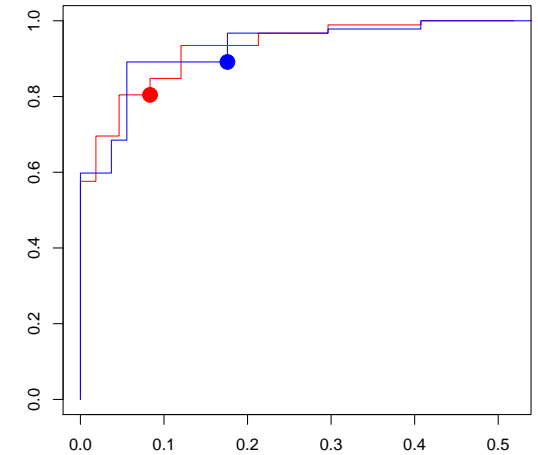


Comparing ROC Curves

```

1 > s1 <- .5
2 > s2 <- .5
3 > table( Observed=Y,
4         Predicted=Ps1)
5 Observed  0  1
6 FALSE  99  9
7 TRUE   18 74
8 > table( Observed=Y,
9         Predicted=Ps2)
10 Observed  0  1
11 FALSE  89 19
12 TRUE   10 82
1 > s1 <- .49
2 > s2 <- .49
3 > table( Observed=Y,
4         Predicted=Ps1)
5 Observed  0  1
6 FALSE  95 13
7 TRUE   14 78
8 > table( Observed=Y,
9         Predicted=Ps2)
10 Observed  0  1
11 FALSE  89 19
12 TRUE   10 82

```



R packages for ROC curves

Consider our previous logistic regression (on heart attacks)

```

1 > myocarde <- read.table("myocarde.csv", head=
  TRUE, sep=";")
2 > logistic <- glm(PRONO~. , data=myocarde ,
  family=binomial)
3 > Y <- myocarde$PRONO
4 > S <- predict(logistic , type="response")

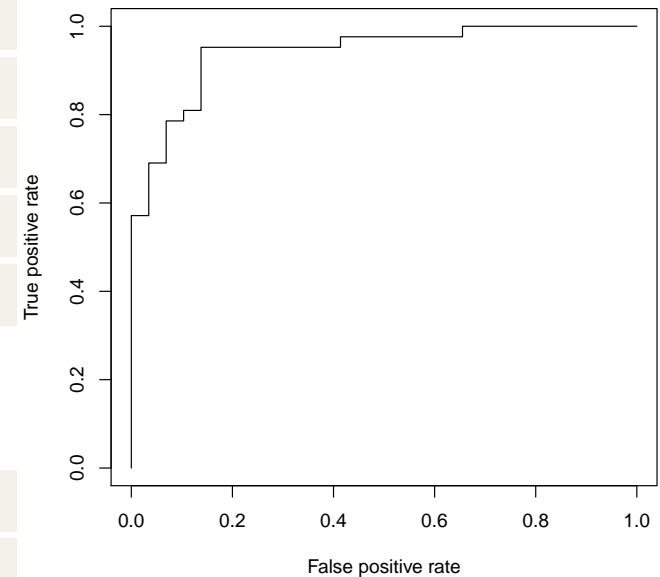
```

For a standard ROC curve

```

1 > library(ROCR)
2 > pred <- prediction(S,Y)
3 > perf <- performance(pred , "tpr" , "fpr")
4 > plot(perf)

```



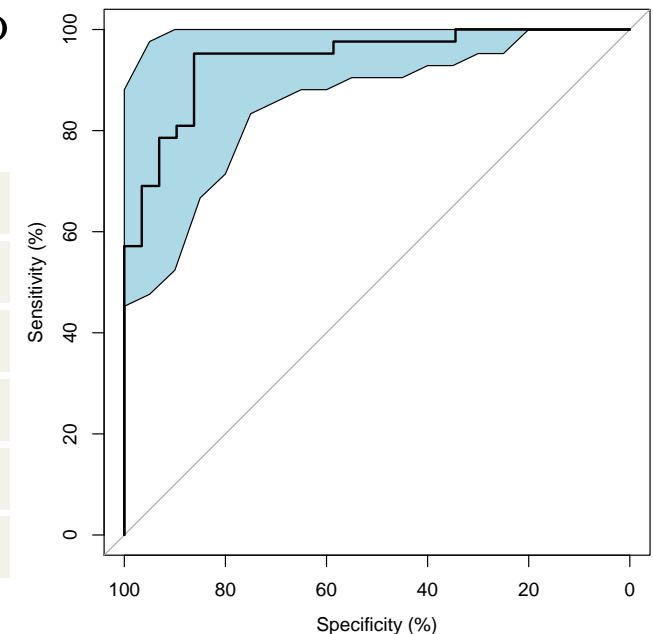
R packages for ROC curves

One can get confidence bands (obtained using bootstrap procedures)

```

1 > library(pROC)
2 > roc <- plot.roc(Y,S,main=" ", percent=TRUE, ci
   =TRUE)
3 > roc.se <- ci.se(roc,specificities=seq(0, 100,
   5))
4 > plot(roc.se, type="shape", col="light blue")

```



see also for Gains and Lift curves

```

1 > library(gains)

```

Standard Quantities Derived from a ROC curve

A standard quantity derived from the ROC curve is the **AUC**, Area Under the Curve, but many other quantities can be computed, see

```

1 > library(hmeasures)
2 > HMeasure(Y,S)$metrics[,1:5]
3 Class labels have been switched from (DECES,SURVIE) to (0,1)
4           H           Gini           AUC           AUCH           KS
5 scores 0.7323154 0.8834154 0.9417077 0.9568966 0.8144499

```

with the *H*-measure (see hmeasure.net), Gini and AUC, as well as the area under the convex hull (**AUCH**).

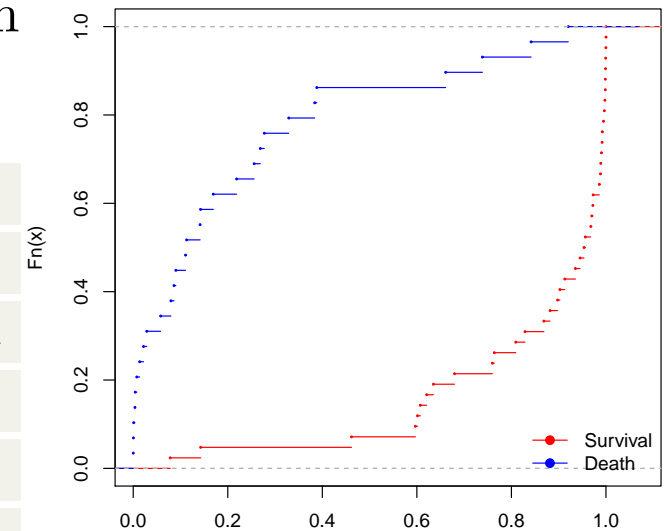
Standard Quantities Derived from a ROC curve

One can compute [Kolmogorov-Smirnov](#) statistics on the two conditional distributions of the score function (given either $Y = 1$ or $Y = 0$)

```

1 > plot(ecdf(S[Y=="SURVIE"]), main="", xlab="",
      pch=19, cex=.2, col="red")
2 > plot(ecdf(S[Y=="DECES"]), , pch=19, cex=.2, col
      ="blue", add=TRUE)
3 > max(perfy.values[[1]]-perfx.values[[1]])
4 [1] 0.8144499

```



```

1 > HMeasure(Y,S)$metrics[,6:10]
2 Class labels have been switched from (DECES,SURVIE) to (0,1)
3
4

```

	MER	MWL	Spec.Sens95	Sens.Spec95	ER
scores	0.08450704	0.08966475	0.862069	0.6904762	0.09859155

with the minimum error rate ([MER](#)), the minimum cost-weighted error rate, etc.