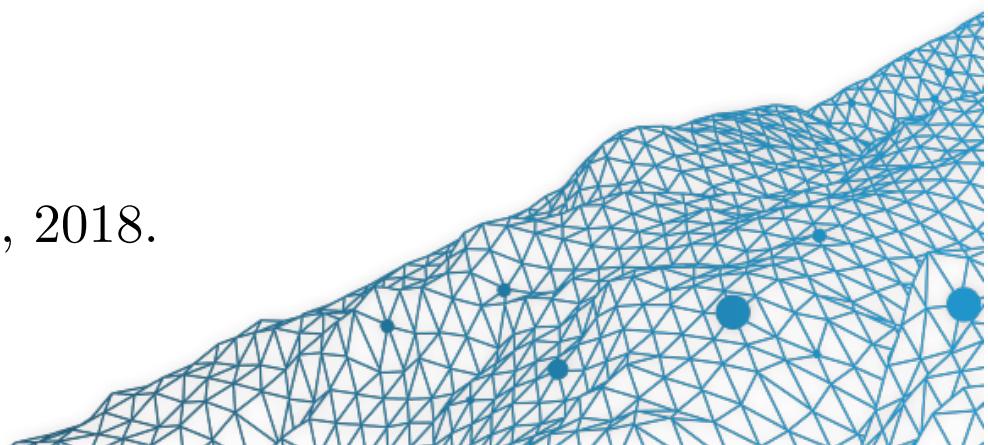


Big Data for Economics # 4

A. Charpentier (Université de Rennes 1)

<https://github.com/freakonometrics/ub>

UB School of Economics Summer School, 2018.



#4 Nonlinearities and Discontinuities

References

Motivation

Kopczuk, W. **Tax bases, tax rates and the elasticity of reported income.** JPE.

income. I experiment with 10-piece splines in logarithms of both the $t-1$ income and the “transitory” component to allow for potential nonlinear effects. Nonlinearity in the permanent component allows me to account for trends in income varying across different income classes. In principle, the transitory component can be controlled for in a linear

References

Eubank, R.L. (1999) **Nonparametric Regression and Spline Smoothing**, CRC Press.

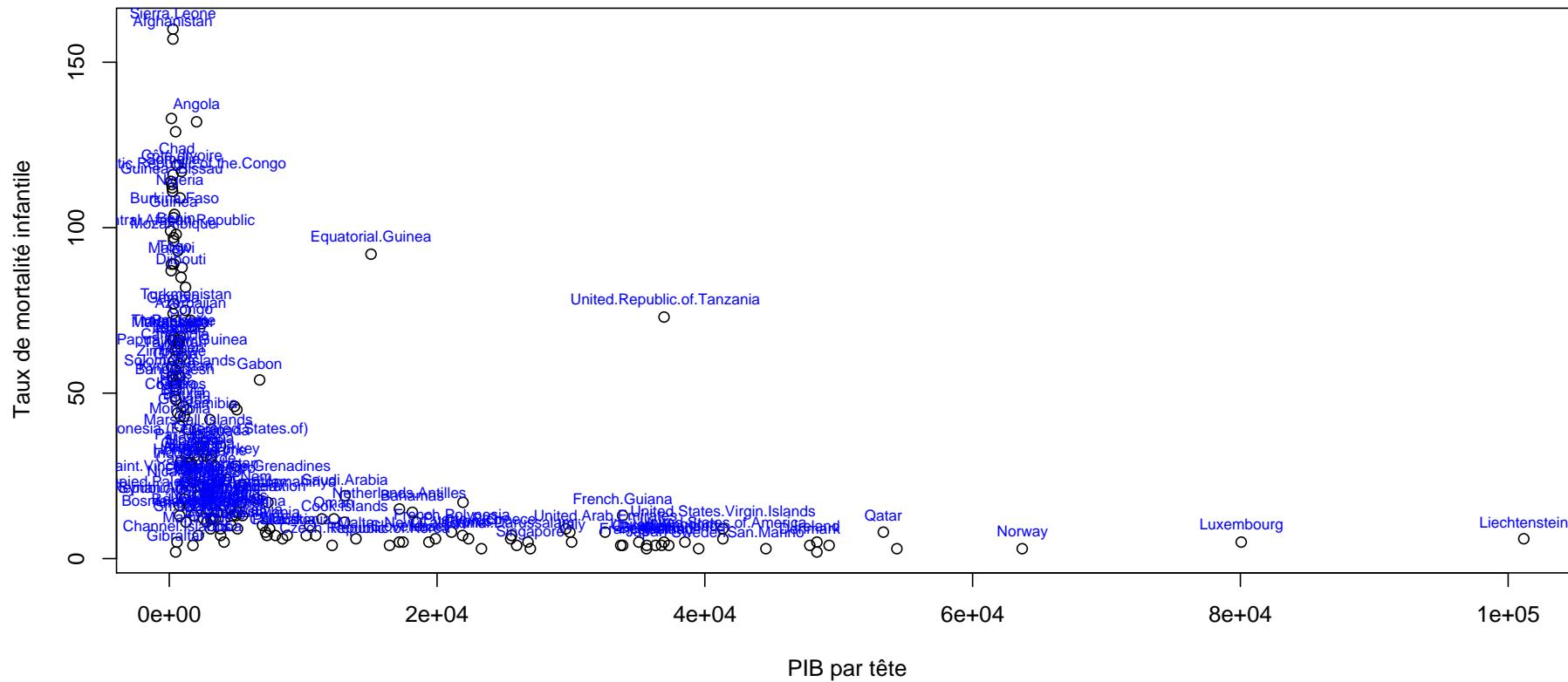
Fan, J. & Gijbels, I. (1996) **Local Polynomial Modelling and Its Applications** CRC Press.

Hastie, T.J. & Tibshirani, R.J. (1990) **Generalized Additive Models**. CRC Press

Wand, M.P & Jones, M.C. (1994) **Kernel Smoothing**. CRC Press

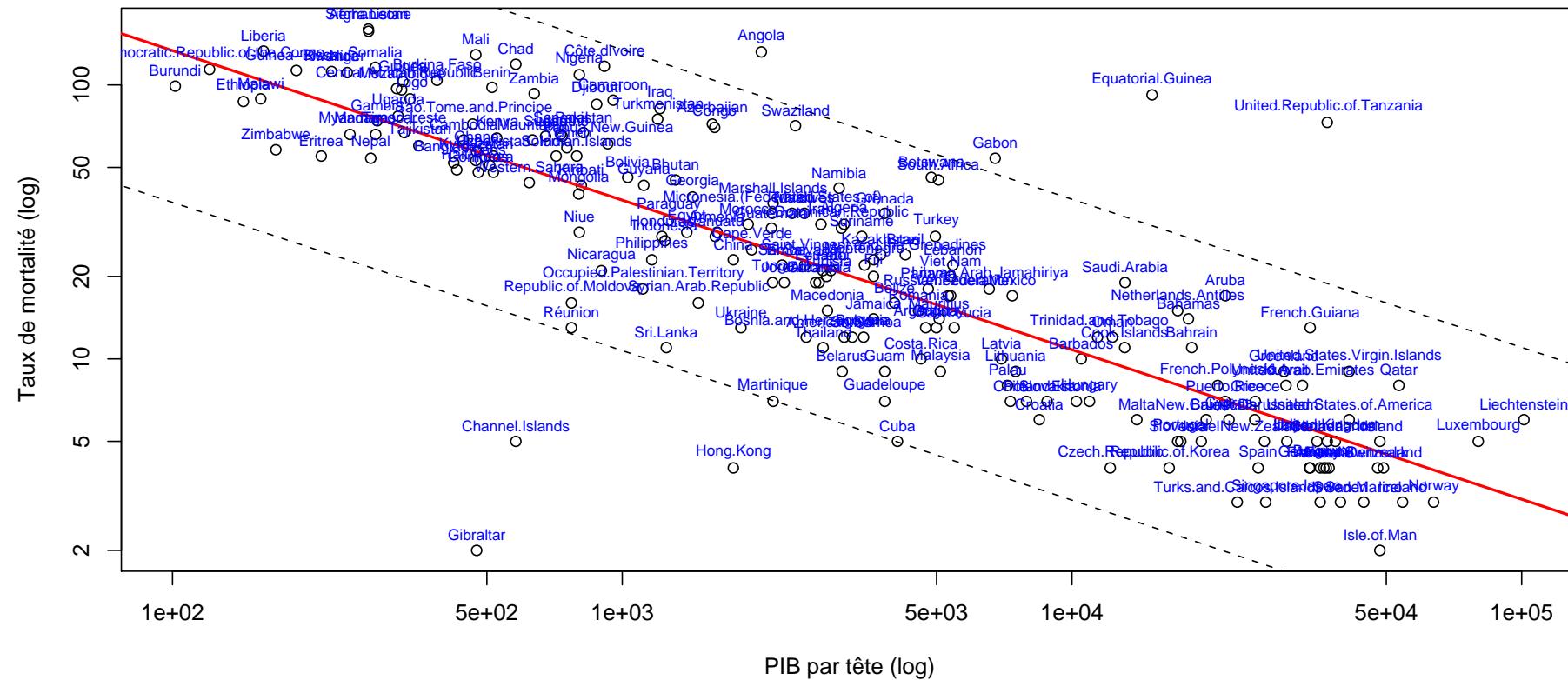
Deterministic or Parametric Transformations

Consider child mortality rate (y) as a function of GDP per capita (x).



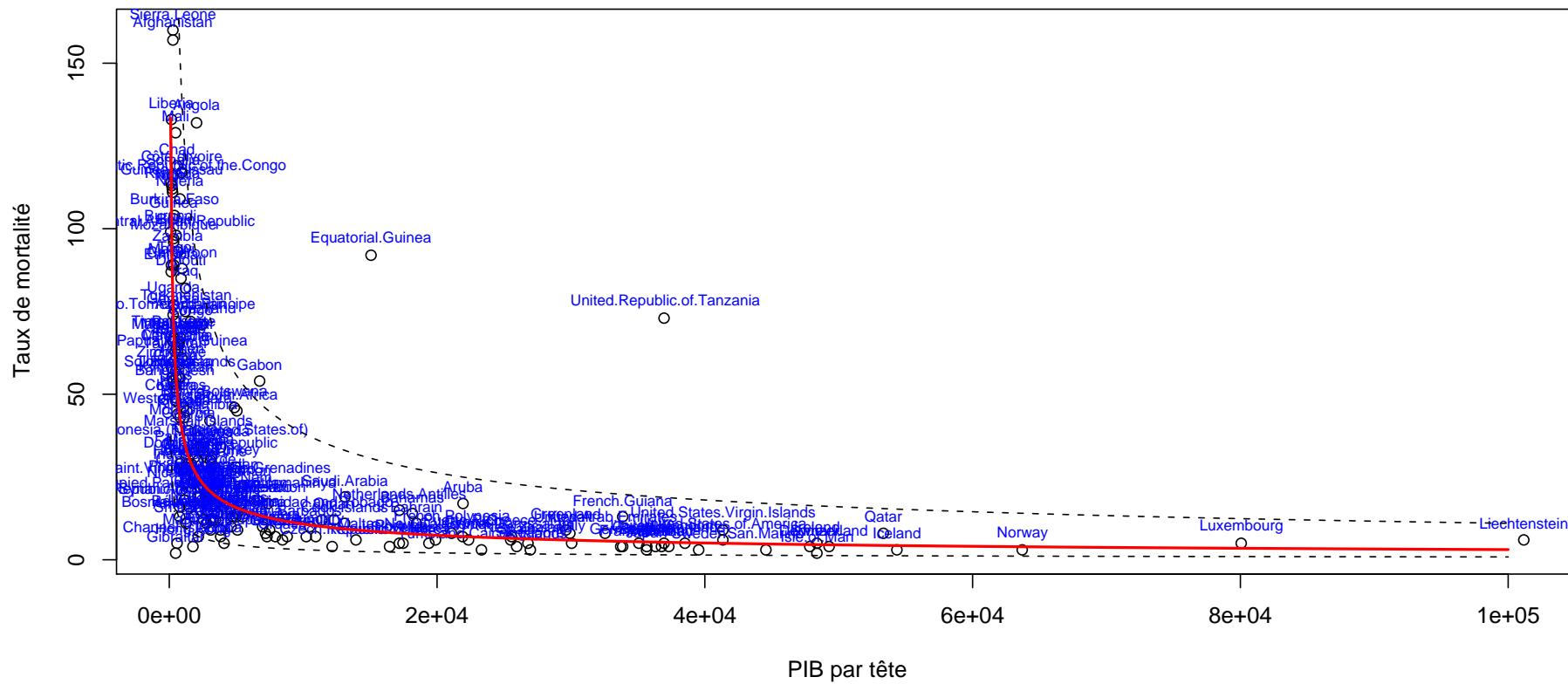
Deterministic or Parametric Transformations

Logarithmic transformation, $\log(y)$ as a function of $\log(x)$



Deterministic or Parametric Transformations

Reverse transformation



Box-Cox transformation

See Box & Cox (1964) *An Analysis of Transformations* ,

$$h(y, \lambda) = \begin{cases} \frac{y^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \log(y) & \text{if } \lambda = 0 \end{cases}$$

or

$$h(y, \lambda, \mu) = \begin{cases} \frac{[y + \mu]^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \log([y + \mu]) & \text{if } \lambda = 0 \end{cases}$$

Profile Likelihood

In a statistical context, suppose that unknown parameter can be partitioned $\boldsymbol{\theta} = (\lambda, \boldsymbol{\beta})$ where λ is the parameter of interest, and $\boldsymbol{\beta}$ is a nuisance parameter.

Consider $\{y_1, \dots, y_n\}$, a sample from distribution $F_{\boldsymbol{\theta}}$, so that the log-likelihood is

$$\log \mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^n \log f_{\boldsymbol{\theta}}(y_i)$$

$\hat{\boldsymbol{\theta}}^{MLE}$ is defined as $\hat{\boldsymbol{\theta}}^{MLE} = \operatorname{argmax} \{\log \mathcal{L}(\boldsymbol{\theta})\}$

Rewrite the log-likelihood as $\log \mathcal{L}(\boldsymbol{\theta}) = \log \mathcal{L}_{\lambda}(\boldsymbol{\beta})$. Define

$$\hat{\boldsymbol{\beta}}_{\lambda}^{pMLE} = \operatorname{argmax}_{\boldsymbol{\beta}} \{\log \mathcal{L}_{\lambda}(\boldsymbol{\beta})\}$$

and then $\hat{\lambda}^{pMLE} = \operatorname{argmax}_{\lambda} \left\{ \log \mathcal{L}_{\lambda}(\hat{\boldsymbol{\beta}}_{\lambda}^{pMLE}) \right\}$. Observe that

$$\sqrt{n}(\hat{\lambda}^{pMLE} - \lambda) \xrightarrow{\mathcal{L}} \mathcal{N}(0, [\mathbb{I}_{\lambda, \lambda} - \mathbb{I}_{\lambda, \boldsymbol{\beta}} \mathbb{I}_{\boldsymbol{\beta}, \boldsymbol{\beta}}^{-1} \mathbb{I}_{\boldsymbol{\beta}, \lambda}]^{-1})$$

Profile Likelihood and Likelihood Ratio Test

The (profile) likelihood ratio test is based on

$$2 \left(\max \{ \mathcal{L}(\lambda, \beta) \} - \max \{ \mathcal{L}(\lambda_0, \beta) \} \right)$$

If (λ_0, β_0) are the true value, this difference can be written

$$2 \left(\max \{ \mathcal{L}(\lambda, \beta) \} - \max \{ \mathcal{L}(\lambda_0, \beta_0) \} \right) - 2 \left(\max \{ \mathcal{L}(\lambda_0, \beta) \} - \max \{ \mathcal{L}(\lambda_0, \beta_0) \} \right)$$

Using Taylor's expansion

$$\frac{\partial \mathcal{L}(\lambda, \beta)}{\partial \lambda} \Big|_{(\lambda_0, \hat{\beta}_{\lambda_0})} \sim \frac{\partial \mathcal{L}(\lambda, \beta)}{\partial \lambda} \Big|_{(\lambda_0, \beta_0)} - \mathbb{I}_{\beta_0 \lambda_0} \mathbb{I}_{\beta_0 \beta_0}^{-1} \frac{\partial \mathcal{L}(\lambda_0, \beta)}{\partial \beta} \Big|_{(\lambda_0, \beta_0)}$$

Thus,

$$\frac{1}{\sqrt{n}} \frac{\partial \mathcal{L}(\lambda, \beta)}{\partial \lambda} \Big|_{(\lambda_0, \hat{\beta}_{\lambda_0})} \xrightarrow{\mathcal{L}} \mathcal{N}(0, \mathbb{I}_{\lambda_0 \lambda_0}) - \mathbb{I}_{\lambda_0 \beta_0} \mathbb{I}_{\beta_0 \beta_0}^{-1} \mathbb{I}_{\beta_0 \lambda_0}$$

and $2 \left(\mathcal{L}(\hat{\lambda}, \hat{\beta}) - \mathcal{L}(\lambda_0, \hat{\beta}_{\lambda_0}) \right) \xrightarrow{\mathcal{L}} \chi^2(\dim(\lambda))$.

Profile Likelihood and Likelihood Ratio Test

Consider some lognormal sample, and fit a Gamma distribution,

$$f(x; \alpha, \beta) = \frac{x^{\alpha-1} \beta^\alpha e^{-\beta x}}{\Gamma(\alpha)} \text{ with } x > 0 \text{ and } \boldsymbol{\theta} = (\alpha, \beta).$$

```
1 > x=exp(rnorm(100))
```

Maximum-likelihood, $\hat{\boldsymbol{\theta}} = \operatorname{argmax}\{\log \mathcal{L}(\boldsymbol{\theta})\}$.

```
1 > library(MASS)
2 > (F=fitdistr(x, "gamma"))
3     shape          rate
4     1.4214497    0.8619969
5     (0.1822570) (0.1320717)
6 > F$estimate[1]+c(-1,1)*1.96*F$sd[1]
7 [1] 1.064226 1.778673
```

Profile Likelihood and Likelihood Ratio Test

See also

```

1 > log_lik=function(theta){
2 +   a=theta[1]
3 +   b=theta[2]
4 +   logL=sum(log(dgamma(x,a,b)))
5 +   return(-logL)
6 +
7 > optim(c(1,1),log_lik)
8 $par
9 [1] 1.4214116 0.8620311

```

We can also use profile likelihood,

$$\hat{\alpha} = \operatorname{argmax}_{\beta} \left\{ \max_{\beta} \left\{ \log \mathcal{L}(\alpha, \beta) \right\} \right\} = \operatorname{argmax} \left\{ \log \mathcal{L}(\alpha, \hat{\beta}_\alpha) \right\}$$

Profile Likelihood and Likelihood Ratio Test

```

1 > prof_log_lik=function(a){
2 +   b=(optim(1,function(z) -sum(log(dgamma(x,a,z))))$par
3 +   return(-sum(log(dgamma(x,a,b))))
4 +
5
6 > vx=seq(.5,3,length=101)
7 > vl=-Vectorize(prof_log_lik)(vx)
8 > plot(vx,vl,type="l")
9 > optim(1,prof_log_lik)
10 $par
11 [1] 1.421094

```

We can use the likelihood ratio test

$$2 \left(\log \mathcal{L}_p(\hat{\alpha}) - \log \mathcal{L}_p(\alpha) \right) \sim \chi^2(1)$$

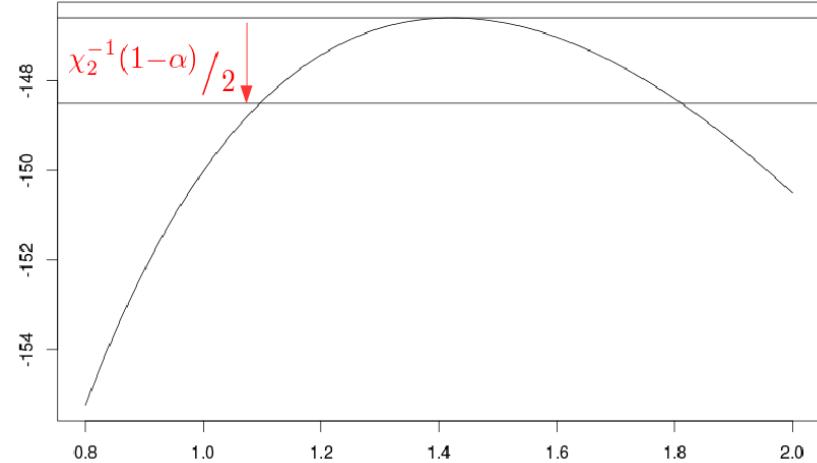
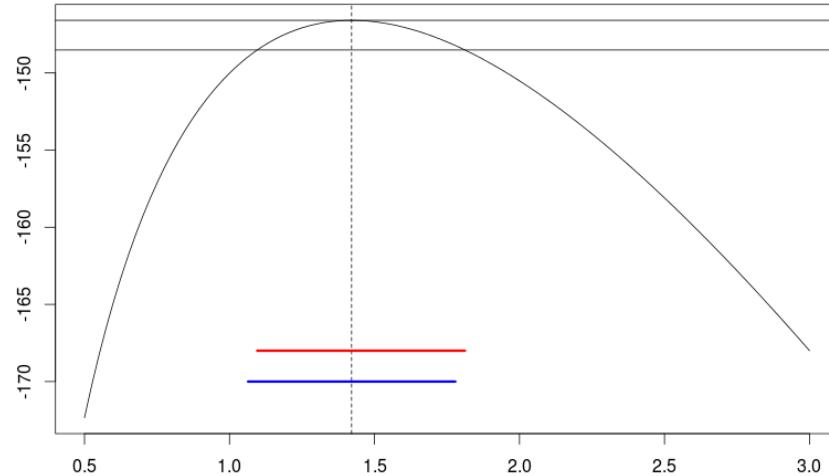
Profile Likelihood and Likelihood Ratio Test

The implied 95% confidence interval is

```

1 > (b1=uniroot(function(z) Vectorize(prof_log_lik)(z)+borne ,c(.5,1.5))
     $root)
2 [1] 1.095726
3 > (b2=uniroot(function(z) Vectorize(prof_log_lik)(z)+borne ,c
     (1.25,2.5))$root)
4 [1] 1.811809

```



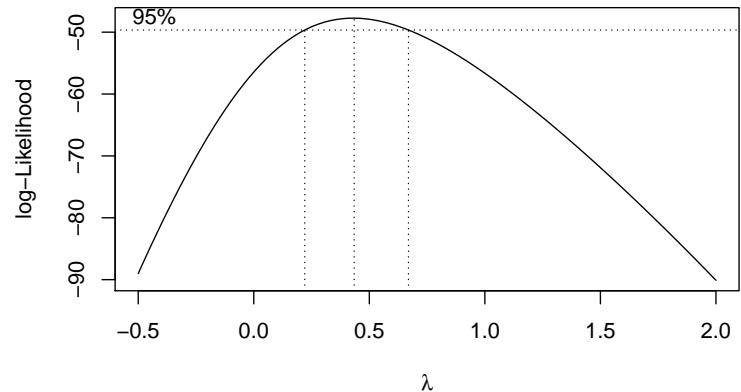
Box-Cox

```
1 > boxcox(lm(dist ~ speed, data=cars))
```

Here $h^* \sim 0.5$

Heuristically, $y_i^{1/2} \sim \beta_0 + \beta_1 x_i + \varepsilon_i$

why not consider a quadratic regression...?

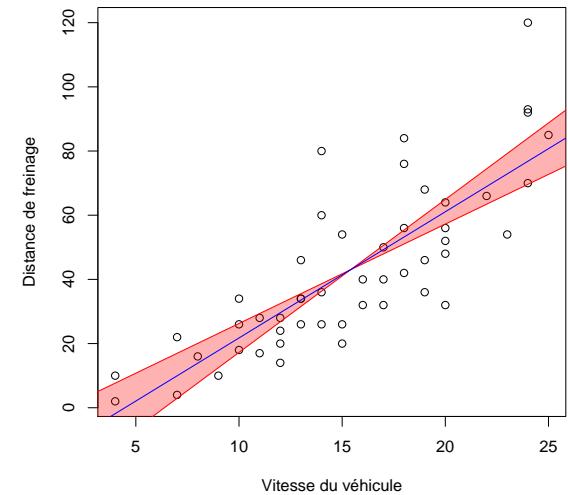
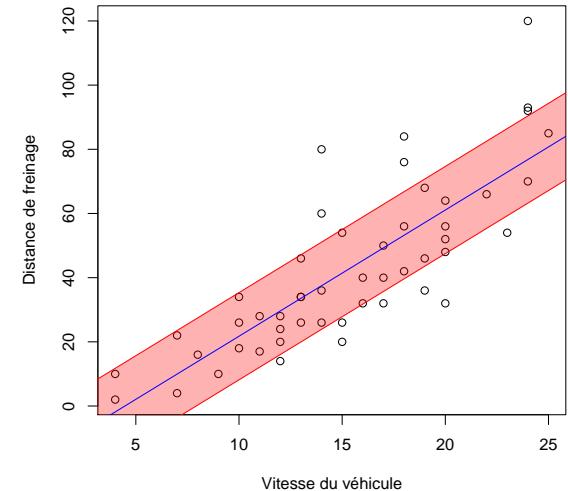


Uncertainty: Parameters vs. Prediction

Uncertainty on regression parameters (β_0, β_1)

From the output of the regression we can derive confidence intervals for β_0 and β_1 , usually

$$\beta_k \in [\hat{\beta}_k \pm u_{1-\alpha/2} \hat{s}\text{e}[\hat{\beta}_k]]$$



Uncertainty: Parameters vs. Prediction

Uncertainty on a prediction, $y = m(\mathbf{x})$. Usually

$$m(\mathbf{x}) \in [\hat{m}(\mathbf{x}) \pm u_{1-\alpha/2} \widehat{\text{se}}[m(\mathbf{x})]]$$

hence, for a linear model

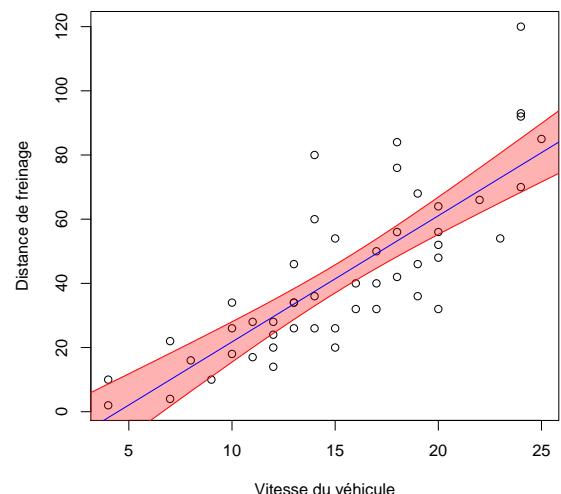
$$\left[\mathbf{x}^\top \hat{\boldsymbol{\beta}} \pm u_{1-\alpha/2} \hat{\sigma} \sqrt{\mathbf{x}^\top [\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{x}} \right]$$

i.e. (with one covariate)

$$\text{se}^2[m(x)]^2 = \text{Var}[\hat{\beta}_0 + \hat{\beta}_1 x]$$

$$\text{se}^2[\hat{\beta}_0] + \text{cov}[\hat{\beta}_0, \hat{\beta}_1]x + \text{se}^2[\hat{\beta}_1]x^2$$

```
1 > predict(lm(dist ~ speed, data=cars), newdata=data.frame(speed=x),
  interval="confidence")
```



Least Squares and Expected Value (Orthogonal Projection Theorem)

Let $\mathbf{y} \in \mathbb{R}^d$, $\bar{y} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n \frac{1}{n} \underbrace{[y_i - m]}_{\varepsilon_i}^2 \right\}$. It is the empirical version of

$$\mathbb{E}[Y] = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \int \underbrace{[y - m]}_{\varepsilon}^2 dF(y) \right\} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \mathbb{E}[(Y - m)^2] \right\}$$

where Y is a ℓ_1 random variable.

Thus, $\operatorname{argmin}_{m(\cdot): \mathbb{R}^k \rightarrow \mathbb{R}} \left\{ \sum_{i=1}^n \frac{1}{n} \underbrace{[y_i - m(\mathbf{x}_i)]}_{\varepsilon_i}^2 \right\}$ is the empirical version of $\mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$.

The Histogram and the Regressogram

Connections between the estimation of $f(y)$ and $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$. Assume that $y_i \in [a_1, a_{k+1})$, divided in k classes $[a_j, a_{j+1})$. The histogram is

$$\hat{f}_{\mathbf{a}}(y) = \sum_{j=1}^k \frac{\mathbf{1}(t \in [a_j, a_{j+1}))}{a_{j+1} - a_j} \frac{1}{n} \sum_{i=1}^n \mathbf{1}(y_i \in [a_j, a_{j+1}))$$

Assume that $a_{j+1} - a_j = h_n$ and $h_n \rightarrow 0$ as $n \rightarrow \infty$

with $nh_n \rightarrow \infty$ then

$$\mathbb{E}[(\hat{f}_{\mathbf{a}}(y) - f(y))^2] \sim O(n^{-2/3})$$

(for an optimal choice of h_n).

```
1 > hist(height)
```

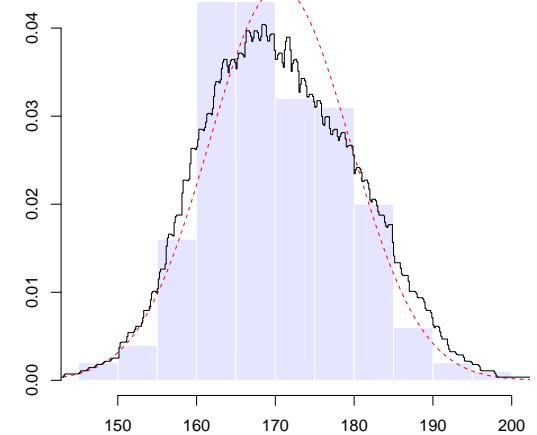
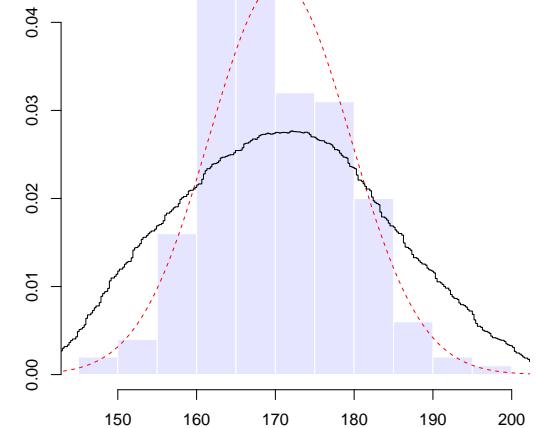
The Histogram and the Regressogram

Then a moving histogram was considered,

$$\hat{f}(y) = \frac{1}{2nh_n} \sum_{i=1}^n \mathbf{1}(y_i \in [y \pm h_n)) = \frac{1}{nh_n} \sum_{i=1}^n k\left(\frac{y_i - y}{h_n}\right)$$

with $k(x) = \frac{1}{2}\mathbf{1}(x \in [-1, 1])$, which a (flat) kernel estimator.

```
1 > density(height, kernel = "rectangular")
```



The Histogram and the Regressogram

From Tukey (1961) **Curves as parameters, and touch estimation**, the regressogram is defined as

$$\hat{m}_{\mathbf{a}}(x) = \frac{\sum_{i=1}^n \mathbf{1}(x_i \in [a_j, a_{j+1}))y_i}{\sum_{i=1}^n \mathbf{1}(x_i \in [a_j, a_{j+1}))}$$

and the moving regressogram is

$$\hat{m}(x) = \frac{\sum_{i=1}^n \mathbf{1}(x_i \in [x \pm h_n])y_i}{\sum_{i=1}^n \mathbf{1}(x_i \in [x \pm h_n])}$$

Nadaraya-Watson and Kernels

Background: Kernel Density Estimator

Consider sample $\{y_1, \dots, y_n\}$, \hat{F}_n empirical cumulative distribution function

$$\hat{F}_n(y) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(y_i \leq y)$$

The empirical measure \mathbb{P}_n consists in weights $1/n$ on each observation.

Idea: add (little) continuous noise to smooth \hat{F}_n .

Let Y_n denote a random variable with distribution \hat{F}_n and define

$$\tilde{Y} = Y_n + hU \text{ where } U \perp\!\!\!\perp Y_n, \text{ with cdf } K$$

The cumulative distribution function of \tilde{Y} is \tilde{F}

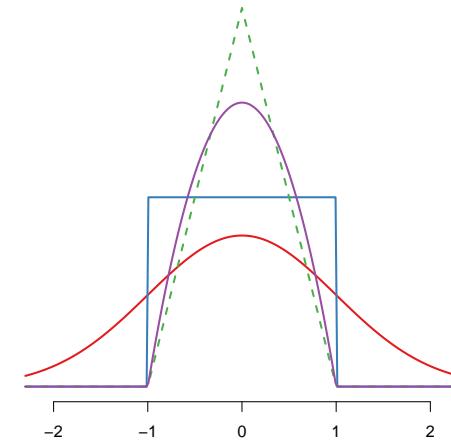
$$\tilde{F}(y) = \mathbb{P}[\tilde{Y} \leq y] = \mathbb{E}(\mathbf{1}(\tilde{Y} \leq y)) = \mathbb{E}(\mathbb{E}[\mathbf{1}(\tilde{Y} \leq y) | Y_n])$$

$$\tilde{F}(y) = \mathbb{E}\left(\mathbf{1}\left(U \leq \frac{y - Y_n}{h}\right) | Y_n\right) = \sum_{i=1}^n \frac{1}{n} K\left(\frac{y - y_i}{h}\right)$$

Nadaraya-Watson and Kernels

If we differentiate

$$\begin{aligned}\tilde{f}(y) &= \frac{1}{nh} \sum_{i=1}^n k\left(\frac{y-y_i}{h}\right) \\ &= \frac{1}{n} \sum_{i=1}^n k_h(y-y_i) \text{ with } k_h(u) = \frac{1}{h} k\left(\frac{u}{h}\right)\end{aligned}$$



\tilde{f} is the kernel density estimator of f , with kernel k and bandwidth h .

Rectangular, $k(u) = \frac{1}{2} \mathbf{1}(|u| \leq 1)$

Epanechnikov, $k(u) = \frac{3}{4} \mathbf{1}(|u| \leq 1)(1 - u^2)$

Gaussian, $k(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}}$

```
1 > density(height, kernel = "epanechnikov")
```

Kernels and Statistical Properties

Consider here an i.id. sample $\{Y_1, \dots, Y_n\}$ with density f

Given y , observe that $\mathbb{E}[\tilde{f}(y)] = \int \frac{1}{h} k\left(\frac{y-t}{h}\right) f(t) dt = \int k(u) f(y-hu) du$. Use **Taylor** expansion around $h=0$, $f(y-hu) \sim f(y) - f'(y)hu + \frac{1}{2}f''(y)h^2u^2$

$$\begin{aligned}\mathbb{E}[\tilde{f}(y)] &= \int f(y)k(u)du - \int f'(y)huk(u)du + \int \frac{1}{2}f''(y+\bar{h}u)h^2u^2k(u)du \\ &= f(y) + 0 + h^2 \frac{f''(y)}{2} \int k(u)u^2du + o(h^2)\end{aligned}$$

Thus, if f is twice continuously differentiable with bounded second derivative,

$$\int k(u)du = 1, \quad \int uk(u)du = 0 \text{ and } \int u^2k(u)du < \infty,$$

then $\mathbb{E}[\tilde{f}(y)] = f(y) + h^2 \frac{f''(y)}{2} \int k(u)u^2du + o(h^2)$

Kernels and Statistical Properties

For the heuristics on that bias, consider a flat kernel, and set

$$f_h(y) = \frac{F(y+h) - F(y-h)}{2h}$$

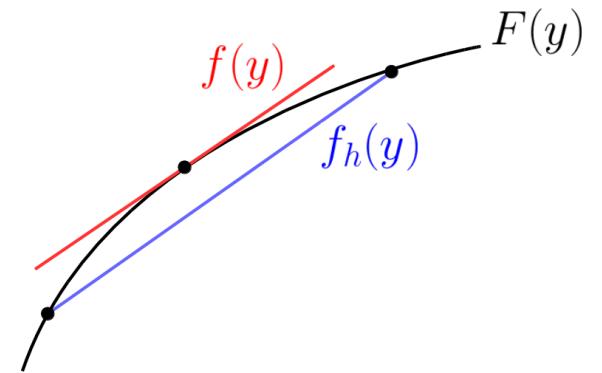
then the natural estimate is

$$\hat{f}_h(y) = \frac{\widehat{F}(y+h) - \widehat{F}(y-h)}{2h} = \frac{1}{2nh} \sum_{i=1}^n \underbrace{\mathbf{1}(y_i \in [y \pm h])}_{Z_i}$$

where Z_i 's are Bernoulli $\mathcal{B}(p_x)$ i.i.d. variables with

$p_x = \mathbb{P}[Y_i \in [x \pm h]] = 2h \cdot f_h(x)$. Thus, $\mathbb{E}(\hat{f}_h(y)) = f_h(y)$, while

$$f_h(y) \sim f(y) + \frac{h^2}{6} f''(y) \text{ as } h \sim 0.$$



Kernels and Statistical Properties

Similarly, as $h \rightarrow 0$ and $nh \rightarrow \infty$

$$\text{Var}[\tilde{f}(y)] = \frac{1}{n} \left(\mathbb{E}[k_h(z - Z)^2] - (\mathbb{E}[k_h(z - Z)])^2 \right)$$

$$\text{Var}[\tilde{f}(y)] = \frac{f(y)}{nh} \int k(u)^2 du + o\left(\frac{1}{nh}\right)$$

Hence

- if $h \rightarrow 0$ the bias goes to 0
- if $nh \rightarrow \infty$ the variance goes to 0

Kernels and Statistical Properties

Extension in Higher Dimension:

$$\tilde{f}(\mathbf{y}) = \frac{1}{n|\mathbf{H}|^{1/2}} \sum_{i=1}^n k\left(\mathbf{H}^{-1/2}(\mathbf{y} - \mathbf{y}_i)\right)$$

$$\tilde{f}(\mathbf{y}) = \frac{1}{nh^d|\boldsymbol{\Sigma}|^{1/2}} \sum_{i=1}^n k\left(\boldsymbol{\Sigma}^{-1/2} \frac{(\mathbf{y} - \mathbf{y}_i)}{h}\right)$$

Kernels and Convolution

Given f and g , set

$$(f \star g)(x) = \int_{\mathbb{R}} f(x - y)g(y)dy$$

Then $\tilde{f}_h = (\hat{f} \star k_h)$, where

$$\hat{f}(y) = \frac{\hat{F}(y)}{dy} = \sum_{i=1}^n \delta_{y_i}(y)$$

Hence, \tilde{f} is the distribution of $\hat{Y} + \varepsilon$ where

\hat{Y} is uniform over $\{y_1, \dots, y_n\}$ and $\varepsilon \sim k_h$ are independent

Nadaraya-Watson and Kernels

Here $\mathbb{E}[Y|X = x] = m(x)$. Write m as a function of densities

$$m(x) = \int y f(y|x) dy = \frac{\int y f(y, x) dy}{\int f(y, x) dy}$$

Consider some bivariate kernel k , such that

$$\int t k(t, u) dt = 0 \text{ and } \kappa(u) = \int k(t, u) dt$$

For the numerator, it can be estimated using

$$\begin{aligned} \int y \tilde{f}(y, x) dy &= \frac{1}{nh^2} \sum_{i=1}^n \int y k\left(\frac{y - y_i}{h}, \frac{x - x_i}{h}\right) dy \\ &= \frac{1}{nh} \sum_{i=1}^n \int y_i k\left(t, \frac{x - x_i}{h}\right) dt = \frac{1}{nh} \sum_{i=1}^n y_i \kappa\left(\frac{x - x_i}{h}\right) \end{aligned}$$

Nadaraya-Watson and Kernels

and for the denominator

$$\int f(y, x) dy = \frac{1}{nh^2} \sum_{i=1}^n \int k\left(\frac{y - y_i}{h}, \frac{x - x_i}{h}\right) = \frac{1}{nh} \sum_{i=1}^n \kappa\left(\frac{x - x_i}{h}\right)$$

Therefore, plugging in the expression for $g(x)$ yields

$$\tilde{m}(x) = \frac{\sum_{i=1}^n y_i \kappa_h(x - x_i)}{\sum_{i=1}^n \kappa_h(x - x_i)}$$

Observe that this regression estimator is a weighted average (see linear predictor section)

$$\tilde{m}(x) = \sum_{i=1}^n \omega_i(x) y_i \text{ with } \omega_i(x) = \frac{\kappa_h(x - x_i)}{\sum_{i=1}^n \kappa_h(x - x_i)}$$

Nadaraya-Watson and Kernels

One can prove that kernel regression bias is given by

$$\mathbb{E}[\tilde{m}(x)] \sim m(x) + h^2 \left(\frac{C_1}{2} m''(x) + C_2 m'(x) \frac{f'(x)}{f(x)} \right)$$

while $\text{Var}[\tilde{m}(x)] \sim \frac{C_3}{nh} \frac{\sigma(x)}{f(x)}$. In this univariate case, one can easily get the kernel estimator of derivatives.

Actually, \tilde{m} is a function of bandwidth h .

Note: this can be extended to multivariate x .

Nadaraya-Watson and Kernels in Higher Dimension

Here $\hat{m}_{\mathbf{H}}(\mathbf{x}) = \frac{\sum_{i=1}^n y_i k_{\mathbf{H}}(\mathbf{x}_i - \mathbf{x})}{\sum_{i=1}^n k_{\mathbf{H}}(\mathbf{x}_i - \mathbf{x})}$ for some symmetric positive definite bandwidth matrix \mathbf{H} , and $k_{\mathbf{H}}(\mathbf{x}) = \det[\mathbf{H}]^{-1}k(\mathbf{H}^{-1}\mathbf{x})$. Then

$$\mathbb{E}[\hat{m}_{\mathbf{H}}(\mathbf{x})] \sim m(\mathbf{x}) + \frac{C_1}{2} \text{trace}(\mathbf{H}^\top m''(\mathbf{x}) \mathbf{H}) + C_2 \frac{m'(\mathbf{x})^\top \mathbf{H} \mathbf{H}^\top \nabla f(\mathbf{x})}{f(\mathbf{x})}$$

while

$$\text{Var}[\hat{m}_{\mathbf{H}}(\mathbf{x})] \sim \frac{C_3}{n \det(\mathbf{H})} \frac{\sigma(\mathbf{x})}{f(\mathbf{x})}$$

Hence, if $\mathbf{H} = h\mathbb{I}$, $h^* \sim Cn^{-\frac{1}{4+\dim(\mathbf{x})}}$.

From kernels to k -nearest neighbours

An alternative is to consider

$$\tilde{m}_k(x) = \frac{1}{n} \sum_{i=1}^n \omega_{i,k}(x) y_i$$

where $\omega_{i,k}(x) = \frac{n}{k}$ if $i \in \mathcal{I}_x^k$ with

$$\mathcal{I}_x^k = \{i : x_i \text{ one of the } k \text{ nearest observations to } x\}$$

Lai (1977) Large sample properties of K-nearest neighbor procedures if $k \rightarrow \infty$ and $k/n \rightarrow 0$ as $n \rightarrow \infty$, then

$$\mathbb{E}[\tilde{m}_k(x)] \sim m(x) + \frac{1}{24f(x)^3} [(m''f + 2m'f')(x)] \left(\frac{k}{n}\right)^2$$

while $\text{Var}[\tilde{m}_k(x)] \sim \frac{\sigma^2(x)}{k}$

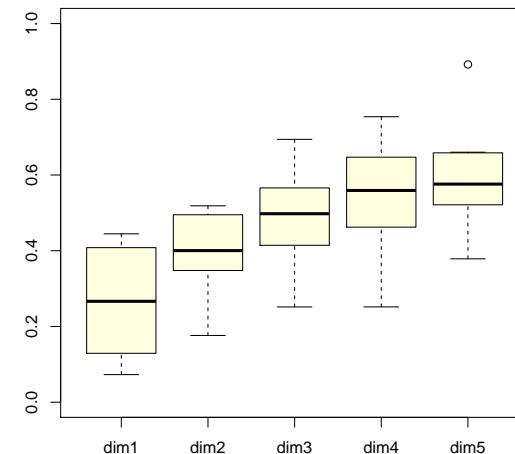
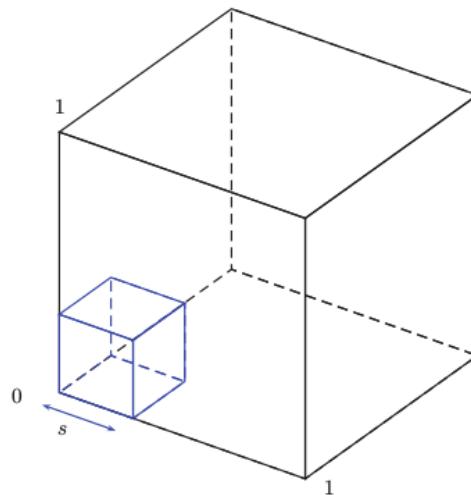
From kernels to k -nearest neighbours

Remark: Brent & John (1985) **Finding the median requires $2n$ comparisons** considered some **median smoothing** algorithm, where we consider the median over the k nearest neighbours (see section #4).

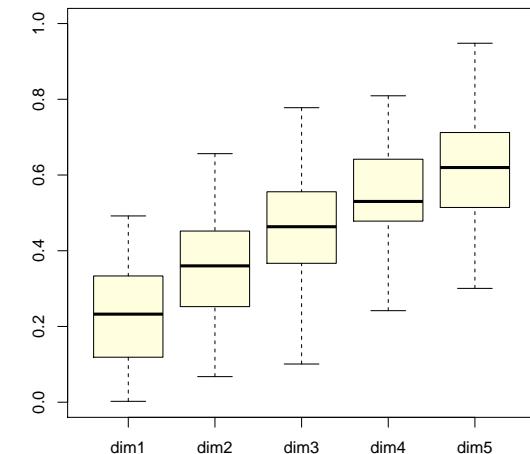
k -Nearest Neighbors and Curse of Dimensionality

The higher the dimension, the larger the distance to the closest neighbor

$$\min_{i \in \{1, \dots, n\}} \{d(\mathbf{a}, \mathbf{x}_i)\}, \mathbf{x}_i \in \mathbb{R}^{\textcolor{red}{d}}.$$



$n = 10$



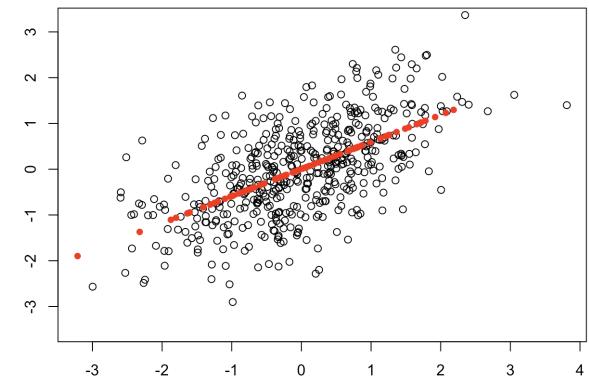
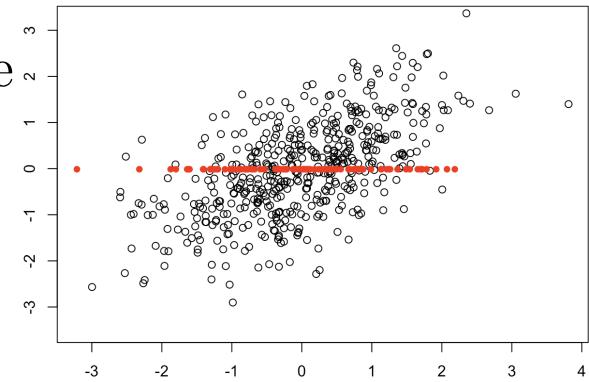
$n = 100$

k-Nearest Neighbors and Imputation

“*The best thing to do with missing values is not to have any*” (Gertrude Mary Cox)

Consider some bivariate observations \mathbf{x}_i , where some $x_{2,i}$'s can be (randomly) missing.

One can consider simple imputation $x_{2,i} = \bar{x}_2$ (computed on non-missing values).



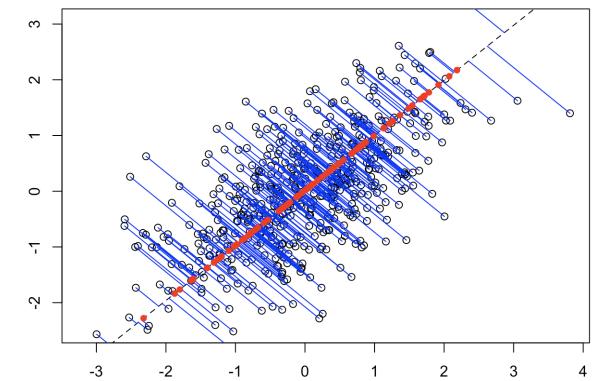
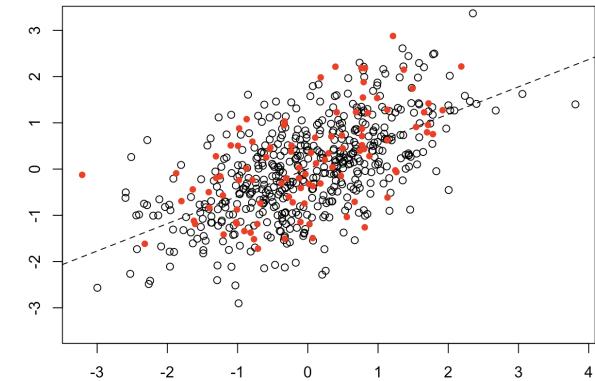
One can consider a regression of x_2 on x_1 (but that's arbitrary)

One can also draw values randomly

$$x_{2,i} = \mathbb{E}[X_2 | X_1 = x_{1,i}] + \varepsilon_i$$

where ε has a Gaussian distribution.

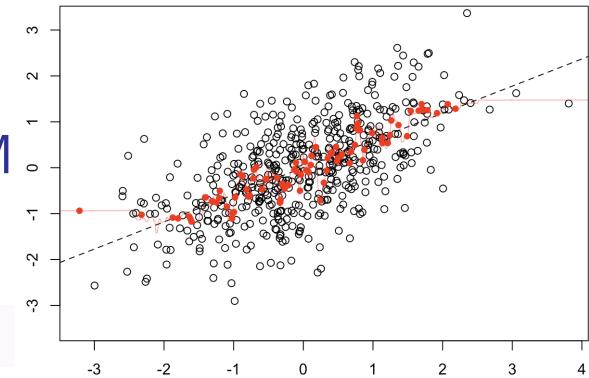
One can consider a PCA (on non-missing values)



or one can consider k nearest-neighbors.

There are packages dedicated to imputation (see **VIM** package)

```
1 > library(VIM)
```

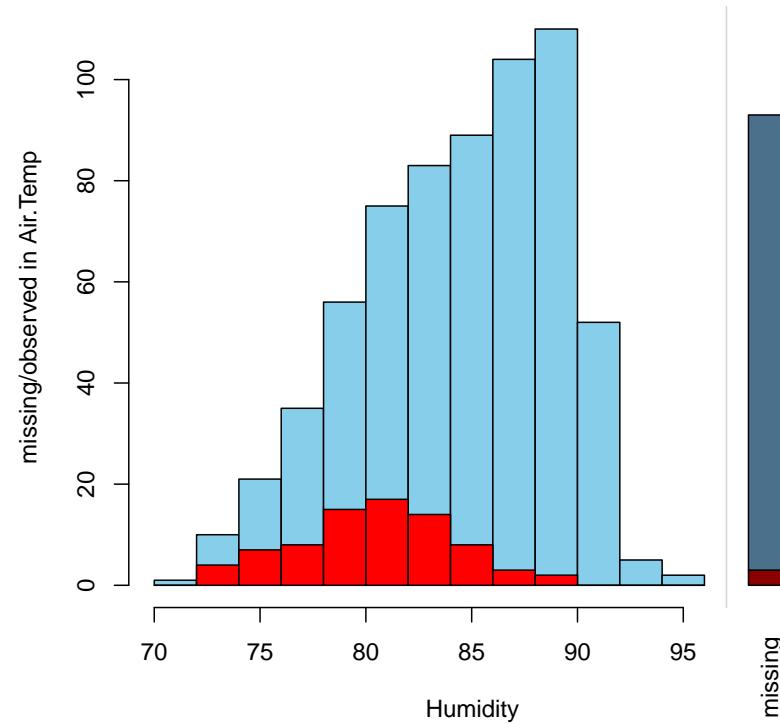
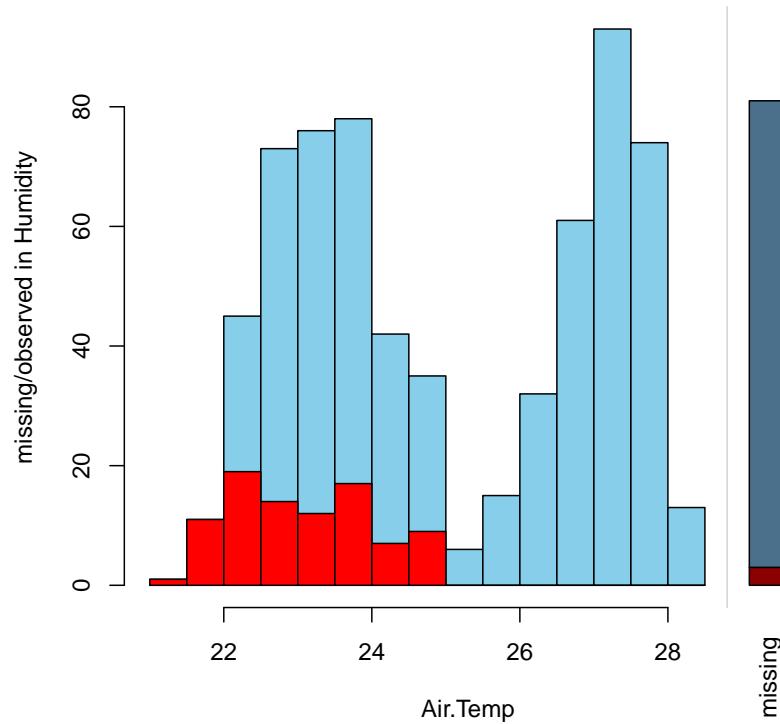


Missing humidity giving the temperature

```

1 > y= tao[,c("Air.Temp","Humidity")]
2 > histMiss(y)                                2 > histMiss(y)

```



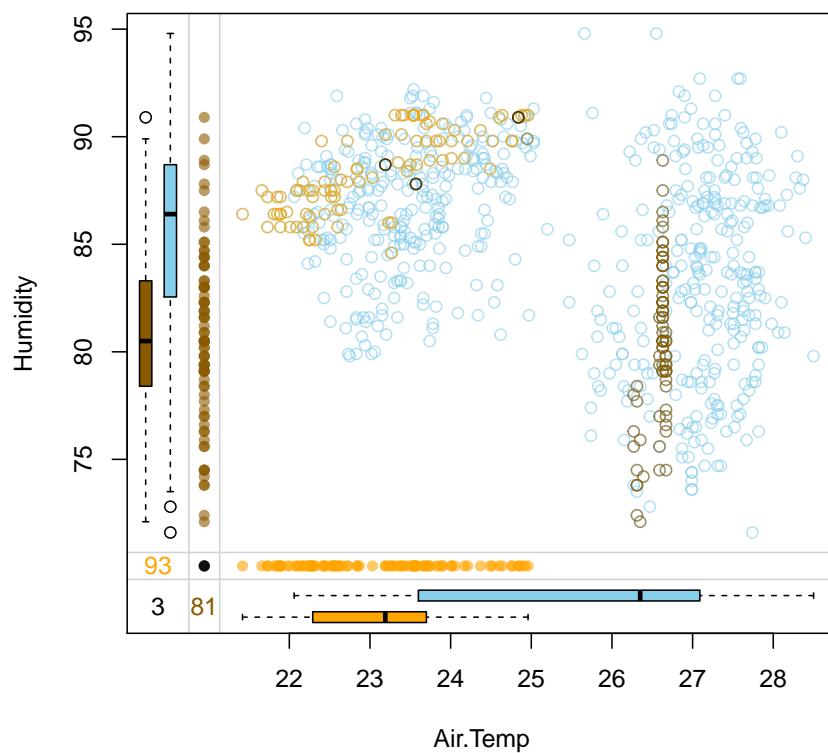
k-Nearest Neighbors and Imputation

This package contains a *k*-Nearest Neighbors algorithm for imputation

```
1 > tao_kNN <- kNN(tao, k = 5)
```

Imputation can be visualized using

```
1 vars <- c("Air.Temp", "Humidity", "  
         Air.Temp_imp", "Humidity_imp")  
2 marginplot(tao_kNN[,vars],  
            delimiter="imp", alpha=0.6)
```



Bandwidth selection : MISE for Density

$$MSE[\tilde{f}(y)] = \text{bias}[\tilde{f}(y)]^2 + \text{Var}[\tilde{f}(y)]$$

$$MSE[\tilde{f}(y)] = f(y) \frac{1}{nh} \int k(u)^2 du + h^4 \left(\frac{f''(y)}{2} \int k(u)u^2 du \right)^2 + o\left(h^4 + \frac{1}{nh}\right)$$

Bandwidth choice is based on minimization of the asymptotic integrated MSE (over y)

$$MISE(\tilde{f}) = \int MSE[\tilde{f}(y)] dy \sim \frac{1}{nh} \int k(u)^2 du + h^4 \int \left(\frac{f''(y)}{2} \int k(u)u^2 du \right)^2$$

Bandwidth selection : MISE for Density

Thus, the first-order condition yields

$$-\frac{C_1}{nh^2} + h^3 \int f''(y)^2 dy C_2 = 0$$

with $C_1 = \int k^2(u)du$ and $C_2 = \left(\int k(u)u^2 du \right)^2$, and

$$h^* = n^{-\frac{1}{5}} \left(\frac{C_1}{C_2 \int f''(y)dy} \right)^{\frac{1}{5}}$$

$h^* = 1.06n^{-\frac{1}{5}} \sqrt{\text{Var}[Y]}$ from Silverman (1986) **Density Estimation**

```

1 > bw.nrd0(cars$speed)
2 [1] 2.150016
3 > bw.nrd(cars$speed)
4 [1] 2.532241

```

with Scott correction, see Scott (1992) **Multivariate Density Estimation**

Bandwidth selection : MISE for Regression Model

One can prove that

$$\begin{aligned} MISE[\hat{m}_h] &\sim \overbrace{\frac{h^4}{4} \left(\int x^2 k(x) dx \right)^2 \int [m''(x) + 2m'(x) \frac{f'(x)}{f(x)}]^2 dx}^{\text{bias}^2} \\ &\quad + \underbrace{\frac{\sigma^2}{nh} \int k^2(x) dx \cdot \int \frac{dx}{f(x)}}_{\text{variance}} \text{ as } n \rightarrow \infty \text{ and } nh \rightarrow \infty. \end{aligned}$$

The bias is sensitive to the position of the x_i 's.

$$h^\star = n^{-\frac{1}{5}} \left(\frac{C_1 \int \frac{dx}{f(x)}}{C_2 \int [m''(x) + 2m'(x) \frac{f'(x)}{f(x)}] dx} \right)^{\frac{1}{5}}$$

Problem: depends on unknown $f(x)$ and $m(x)$.

Bandwidth Selection : Cross Validation

Consider some risk, function of some parameter h . E.g. $MISE[\hat{m}_h]$.

A first idea is to consider a validation set approach,

- Split the data in two parts
- Train the method in the first part
- Compute the error on the second part

Problem : every split yields a different estimate of the error

Bandwidth Selection : Cross Validation

Consider **leave-one-out cross validation**

For every $i \in \{1, 2, \dots, n\}$

- Train the model on every point, but i
- Compute the test error on the held out point

$$CV_{\text{LOO}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i^{(-i)})^2$$

where the **prediction** is obtained on the model based on data where observation i was removed.

It can be computationally expensive

$$CV_{\text{LOO}} = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{1 - h_{i,i}} \right)^2$$

where $h_{i,i}$ is the leverage statistic

Bandwidth Selection : Cross Validation

Consider *k*-fold cross validation

For every $j \in \{1, 2, \dots, n\}$

- Train the model on every fold, but the i th
- Compute the test error on the i th fold

As we increase k in *k*-fold cross-validation, we decrease the bias, but increase the variance. One can use **bootstrap** to estimate measures of uncertainty

Bandwidth Selection : Cross Validation

Let $R(h) = \mathbb{E}[(Y - \hat{m}_h(\mathbf{X}))^2]$.

Natural idea $\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{m}_h(\mathbf{x}_i))^2$

Instead use **leave-one-out cross validation**,

$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n \left(y_i - \hat{m}_h^{(i)}(\mathbf{x}_i) \right)^2$$

where $\hat{m}_h^{(i)}$ is the estimator obtained by omitting the i th pair (y_i, \mathbf{x}_i) or **k -fold cross validation**,

$$\hat{R}(h) = \frac{1}{n} \sum_{j=1}^k \sum_{i \in \mathcal{I}_j} \left(y_i - \hat{m}_h^{(j)}(\mathbf{x}_i) \right)^2$$

where $\hat{m}_h^{(j)}$ is the estimator obtained by omitting pairs (y_i, \mathbf{x}_i) with $i \in \mathcal{I}_j$.

Bandwidth Selection : Cross Validation

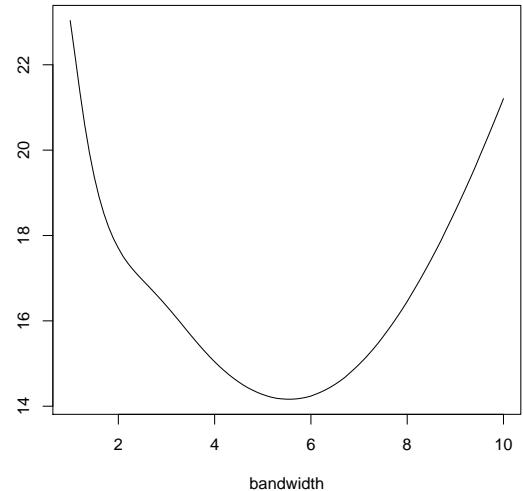
Then find (numerically)

$$h^* = \operatorname{argmin}\{\hat{R}(h)\}$$

In the context of density estimation, see Chiu (1991) **Bandwidth Selection for Kernel Density Estimation**

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



Local Linear Regression

Consider $\hat{m}(\mathbf{x})$ defined as $\hat{m}(\mathbf{x}) = \hat{\beta}_0$ where $(\hat{\beta}_0, \hat{\boldsymbol{\beta}})$ is the solution of

$$\min_{(\beta_0, \boldsymbol{\beta})} \left\{ \sum_{i=1}^n \omega_i^{(\mathbf{x})} (y_i - [\beta_0 + (\mathbf{x} - \mathbf{x}_i)^\top \boldsymbol{\beta}])^2 \right\}$$

where $\omega_i^{(\mathbf{x})} = k_h(\mathbf{x} - \mathbf{x}_i)$, e.g.

i.e. we seek the constant term in a weighted least squares regression of y_i 's on $\mathbf{x} - \mathbf{x}_i$'s. If \mathbf{X}_x is the matrix $[\mathbf{1} \ (\mathbf{x} - \mathbf{X})^\top]$, and if \mathbf{W}_x is a matrix

$$\text{diag}[k_h(\mathbf{x} - \mathbf{x}_1), \dots, k_h(\mathbf{x} - \mathbf{x}_n)]$$

$$\text{then } \hat{m}(\mathbf{x}) = \mathbf{1}^\top (\mathbf{X}_x^\top \mathbf{W}_x \mathbf{X}_x)^{-1} \mathbf{X}_x^\top \mathbf{W}_x \mathbf{y}$$

This estimator is also a linear predictor :

$$\hat{m}(\mathbf{x}) = \sum_{i=1}^n \frac{a_i(\mathbf{x})}{\sum a_i(\mathbf{x})} y_i$$

where

$$a_i(\mathbf{x}) = \frac{1}{n} k_h(\mathbf{x} - \mathbf{x}_i) \left(1 - s_1(\mathbf{x})^\top s_2(\mathbf{x})^{-1} \frac{\mathbf{x} - \mathbf{x}_i}{h} \right)$$

with

$$s_1(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n k_h(\mathbf{x} - \mathbf{x}_i) \frac{\mathbf{x} - \mathbf{x}_i}{h} \text{ and } s_2(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n k_h(\mathbf{x} - \mathbf{x}_i) \left(\frac{\mathbf{x} - \mathbf{x}_i}{h} \right) \left(\frac{\mathbf{x} - \mathbf{x}_i}{h} \right)$$

Note that Nadaraya-Watson estimator was simply the solution of

$$\min_{\beta_0} \left\{ \sum_{i=1}^n \omega_i^{(\mathbf{x})} (y_i - \beta_0)^2 \right\} \text{ where } \omega_i^{(\mathbf{x})} = k_h(\mathbf{x} - \mathbf{x}_i)$$

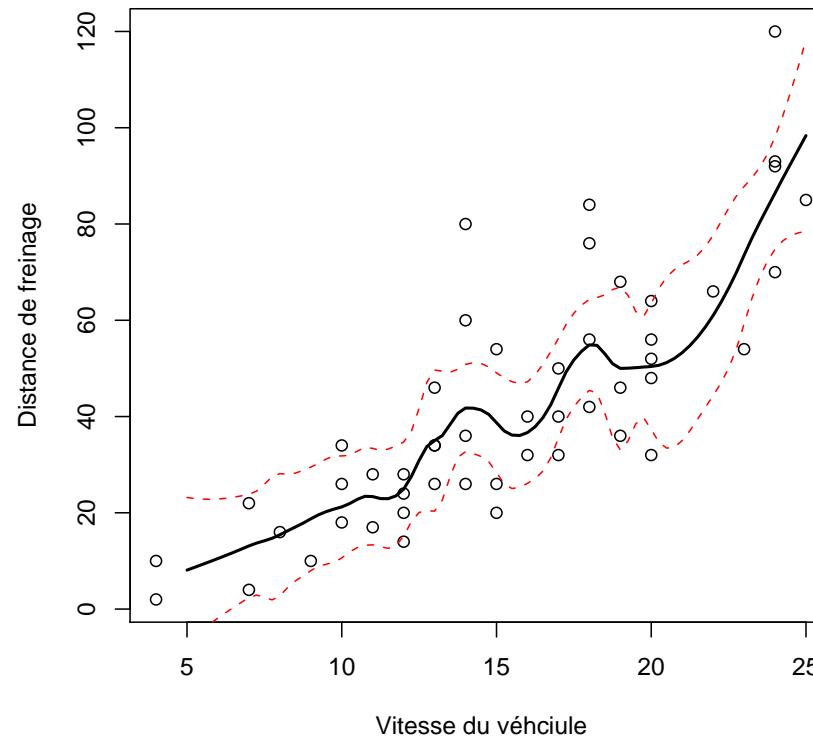
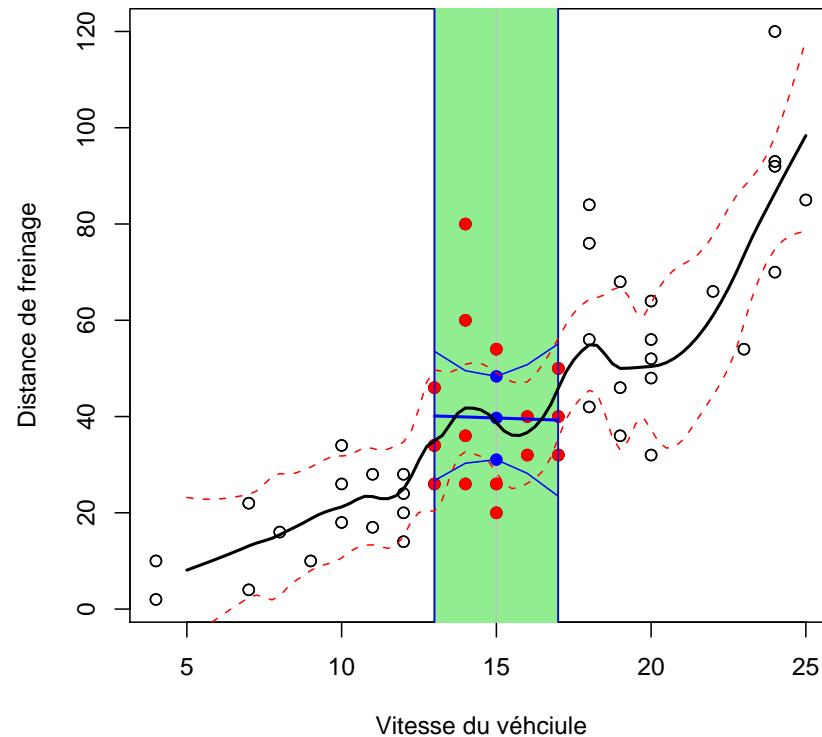
$$\mathbb{E}[\hat{m}(\mathbf{x})] \sim m(\mathbf{x}) + \frac{h^2}{2} m''(\mathbf{x}) \mu_2 \text{ where } \mu_2 = \int k(u) u^2 du.$$

$$\text{Var}[\hat{m}(\mathbf{x})] \sim \frac{1}{nh} \frac{\nu \sigma_{\mathbf{x}}^2}{f(\mathbf{x})}$$

$$\text{where } \nu = \int k(u)^2 du$$

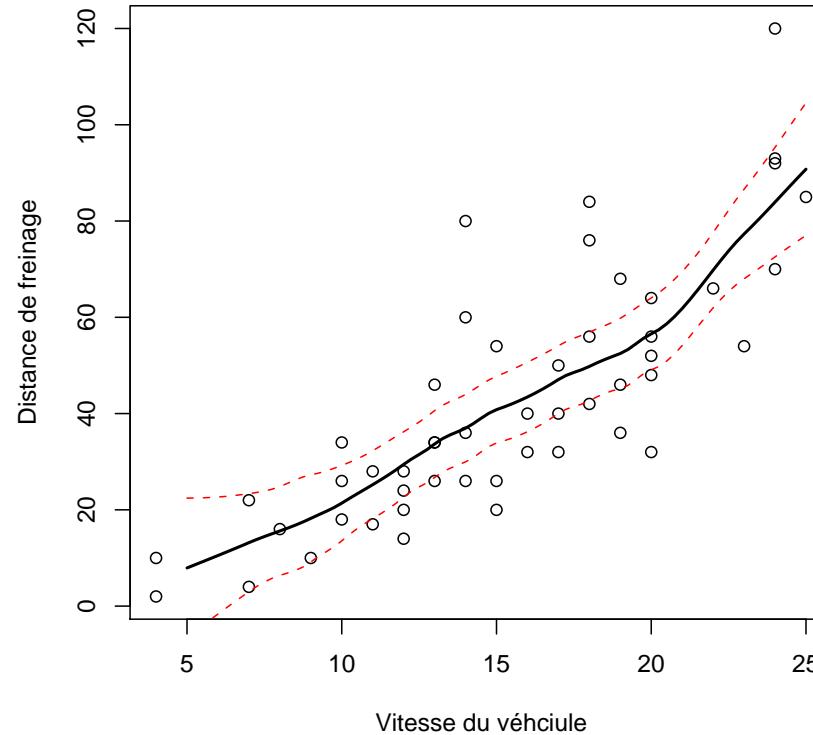
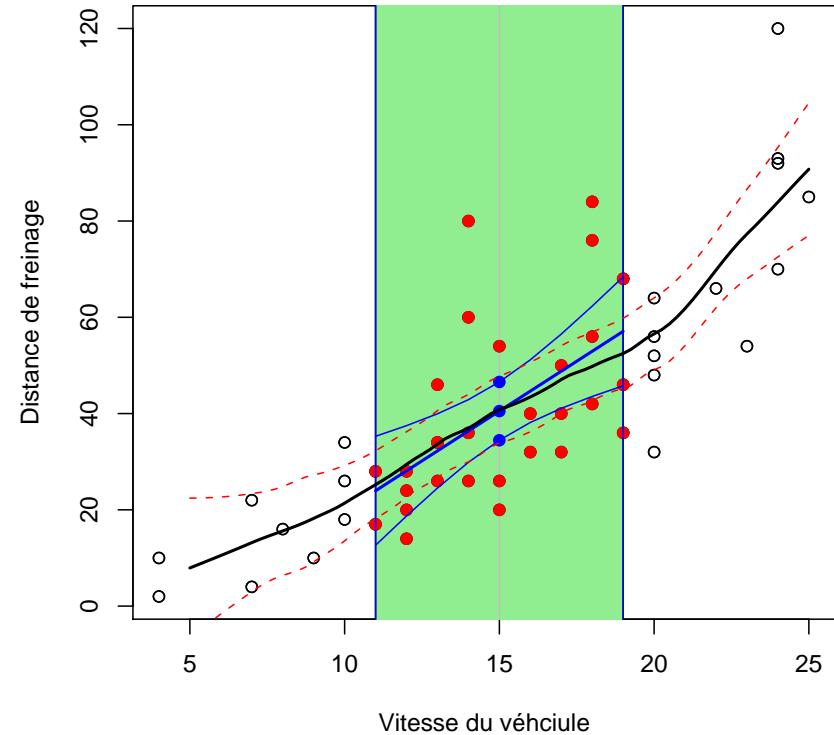
Thus, kernel regression MSE is

$$\frac{h^2}{4} \left(g''(x) + 2g'(x) \frac{f'(x)}{f(x)} \right)^2 \mu_2^2 + \frac{1}{nh} \frac{\nu \sigma_x^2}{f(x)}$$



```
1 > loess(dist ~ speed, cars, span=0.75, degree=1)
```

```
2 > predict(REG, data.frame(speed = seq(5, 25, 0.25)), se = TRUE)
```



Local polynomials

One might assume that, **locally**, $m(x) \sim \mu_x(u)$ as $u \sim 0$, with

$$\mu_x(u) = \beta_0^{(x)} + \beta_1^{(x)} + [u - x] + \beta_2^{(x)} + \frac{[u - x]^2}{2} + \beta_3^{(x)} + \frac{[u - x]^3}{2} + \dots$$

and we estimate $\beta^{(x)}$ by minimizing $\sum_{i=1}^n \omega_i^{(x)} [y_i - \mu_x(x_i)]^2$.

If \mathbf{X}_x is the design matrix $[1 \ x_i - x \ \frac{[x_i - x]^2}{2} \ \frac{[x_i - x]^3}{3} \ \dots]$, then

$$\hat{\boldsymbol{\beta}}^{(x)} = (\mathbf{X}_x^\top \mathbf{W}_x \mathbf{X}_x)^{-1} \mathbf{X}_x^\top \mathbf{W}_x \mathbf{y}$$

(weighted least squares estimators).

```
1 > library(locfit)
2 > locfit(dist ~ speed, data=cars)
```

Series Regression

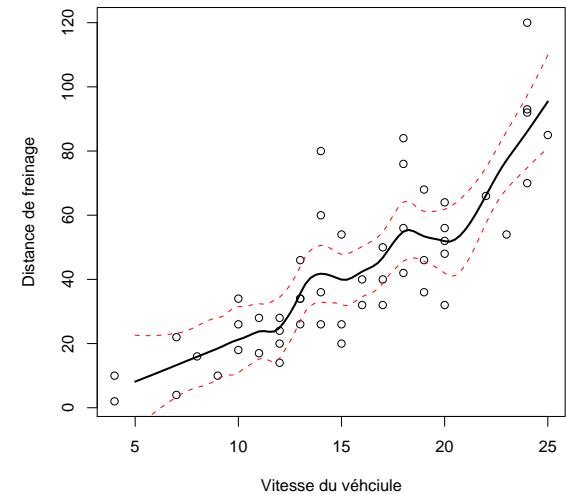
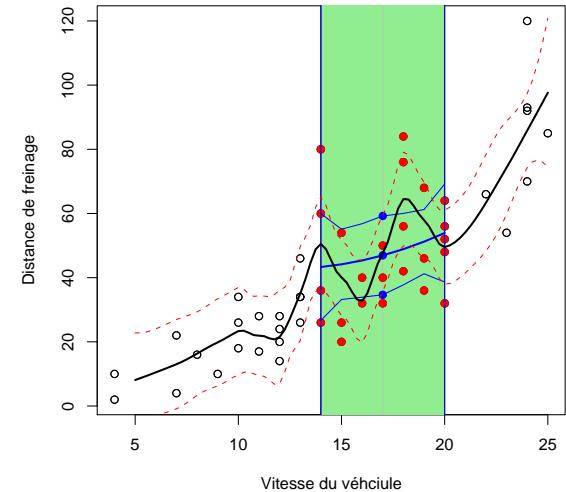
Recall that $\mathbb{E}[Y|X = x] = m(x)$.

Why not approximate m by a linear combination of approximating functions $h_1(x), \dots, h_k(x)$.

Set $\mathbf{h}(x) = (h_1(x), \dots, h_k(x))$, and consider the regression of y_i 's on $\mathbf{h}(x_i)$'s,

$$y_i = \mathbf{h}(x_i)^\top \boldsymbol{\beta} + \varepsilon_i$$

Then $\hat{\boldsymbol{\beta}} = (\mathbf{H}^\top \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{y}$

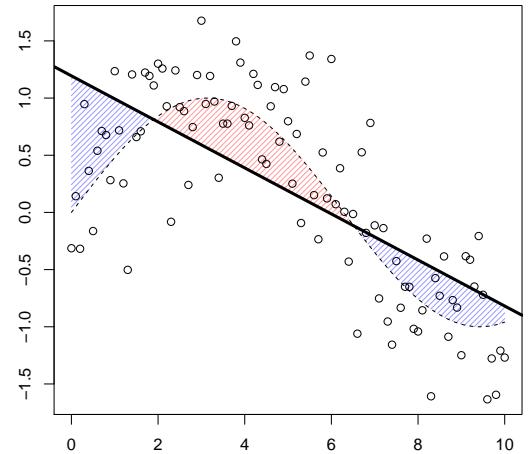
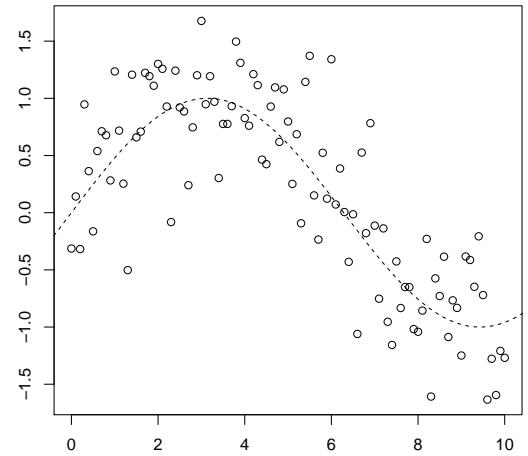


Series Regression : polynomials

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 > reg <- lm(y ~ x, data=db)
```



Series Regression : polynomials

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(y~poly(x,5),data=db)
2 > reg <- lm(y~poly(x,25),data=db)
```

Series Regression : (Linear) Splines

Consider $m + 1$ knots on \mathcal{X} , $\min\{x_i\} \leq t_0 \leq t_1 \leq \dots \leq t_m \leq \max\{x_n\}$, then define linear (degree = 1) splines positive function,

$$b_{j,1}(x) = (x - t_j)_+ = \begin{cases} x - t_j & \text{if } x > t_j \\ 0 & \text{otherwise} \end{cases}$$

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)
```

Series Regression : (Linear) 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) +
2                         positive_part(X - s2), data=db)
3 > library(bsplines)
```

A spline is a function defined by piecewise polynomials.
 b -splines are defined recursively

b-Splines (in Practice)

```

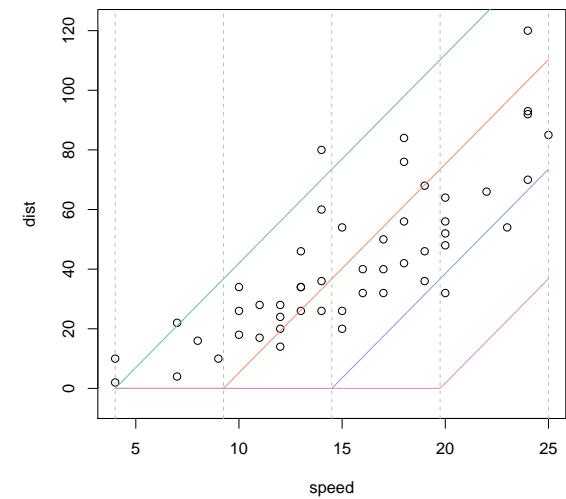
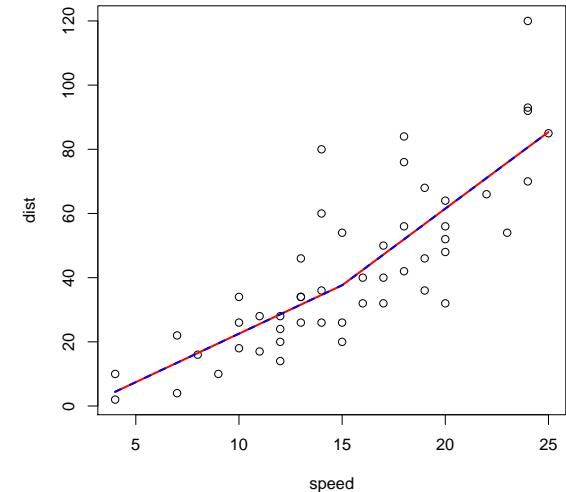
1 > reg1 <- lm(dist ~ speed + positive_part(speed - 15) ,
   data=cars)
2 > reg2 <- lm(dist ~ bs(speed, df=2, degree=1), data=
   cars)

```

Consider $m+1$ knots on $[0, 1]$, $0 \leq t_0 \leq t_1 \leq \dots \leq t_m \leq 1$,
then define recursively *b*-splines as

$$b_{j,0}(t) = \begin{cases} 1 & \text{if } t_j \leq t < t_{j+1} \\ 0 & \text{otherwise, and} \end{cases}$$

$$\begin{aligned} b_{j,n}(t) &= \frac{t - t_j}{t_{j+n} - t_j} b_{j,n-1}(t) \\ &\quad + \frac{t_{j+n+1} - t}{t_{j+n+1} - t_{j+1}} b_{j+1,n-1}(t) \end{aligned}$$



b-Splines (in Practice)

```

1 > summary(reg1)
2
3 Coefficients:
4
5             Estimate Std. Error t value Pr(>|t|) 
6 (Intercept) -7.6519   10.6254 -0.720  0.475
7 speed        3.0186    0.8627  3.499  0.001 ** 
8 (speed-15)   1.7562    1.4551  1.207  0.233
9
10
11 > summary(reg2)
12
13 Coefficients:
14
15             Estimate Std. Error t value Pr(>|t|) 
16 (Intercept) 4.423     7.343    0.602  0.5493
17 bs(speed)1 33.205    9.489    3.499  0.0012 ** 
18 bs(speed)2 80.954    8.788    9.211 4.2e-12 ***
```

b and *p*-Splines

Note that those spline function define an orthonormal basis.

O'Sullivan (1986) **A statistical perspective on ill-posed inverse problems** suggested a penalty on the second derivative of the fitted curve (see #3).

$$m(x) = \operatorname{argmin} \left\{ \sum_{i=1}^n (y_i - \mathbf{b}(x_i)^\top \boldsymbol{\beta})^2 + \lambda \int_{\mathbb{R}} \mathbf{b}''(x_i)^\top \boldsymbol{\beta} \right\}$$

Adding Constraints: Convex Regression

Assume that $y_i = m(\mathbf{x}_i) + \varepsilon_i$ where $m : \mathbb{R}^d \rightarrow \mathbb{R}$ is some convex function.

m is convex if and only if $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d, \forall t \in [0, 1]$,

$$m(t\mathbf{x}_1 + [1-t]\mathbf{x}_2) \leq tm(\mathbf{x}_1) + [1-t]m(\mathbf{x}_2)$$

Proposition (Hidreth (1954) **Point Estimates of Ordinates of Concave Functions**)

$$m^* = \underset{m \text{ convex}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - m(\mathbf{x}_i))^2 \right\}$$

Then $\boldsymbol{\theta}^* = (m^*(\mathbf{x}_1), \dots, m^*(\mathbf{x}_n))$ is unique.

Let $\mathbf{y} = \boldsymbol{\theta} + \boldsymbol{\varepsilon}$, then

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta} \in \mathcal{K}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - \theta_i)^2 \right\}$$

where $\mathcal{K} = \{\boldsymbol{\theta} \in \mathbb{R}^n : \exists m \text{ convex }, m(\mathbf{x}_i) = \theta_i\}$. I.e. $\boldsymbol{\theta}^*$ is the projection of \mathbf{y} onto the (closed) convex cone \mathcal{K} . The projection theorem gives existence and unicity.

Adding Constraints: Convex Regression

In dimension 1: $y_i = m(x_i) + \varepsilon_i$. Assume that observations are ordered $x_1 < x_2 < \dots < x_n$.

Here

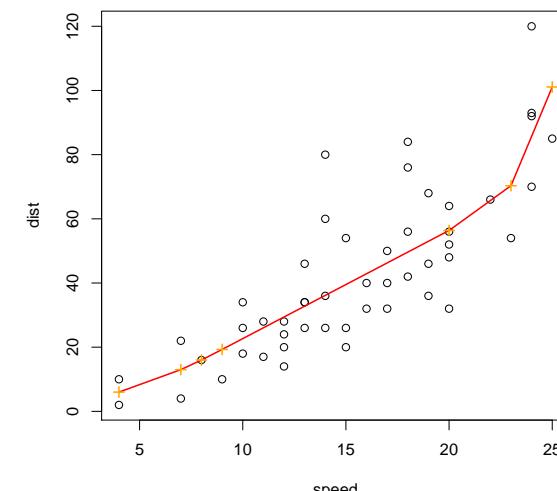
$$\mathcal{K} = \left\{ \boldsymbol{\theta} \in \mathbb{R}^n : \frac{\theta_2 - \theta_1}{x_2 - x_1} \leq \frac{\theta_3 - \theta_2}{x_3 - x_2} \leq \dots \leq \frac{\theta_n - \theta_{n-1}}{x_n - x_{n-1}} \right\}$$

Hence, quadratic program with $n - 2$ linear constraints.

m^* is a **piecewise linear function** (interpolation of consecutive pairs (x_i, θ_i^*)).

If m is differentiable, m is convex if

$$m(\mathbf{x}) + \nabla m(\mathbf{x}) \cdot [\mathbf{y} - \mathbf{x}] \leq m(\mathbf{y})$$



Adding Constraints: Convex Regression

More generally: if m is convex, then there exists $\xi_{\mathbf{x}} \in \mathbb{R}^n$ such that

$$m(\mathbf{x}) + \xi_{\mathbf{x}} \cdot [\mathbf{y} - \mathbf{x}] \leq m(\mathbf{y})$$

$\xi_{\mathbf{x}}$ is a subgradient of m at \mathbf{x} . And then

$$\partial m(\mathbf{x}) = \{m(\mathbf{x}) + \xi \cdot [\mathbf{y} - \mathbf{x}] \leq m(\mathbf{y}), \forall \mathbf{y} \in \mathbb{R}^n\}$$

Hence, $\boldsymbol{\theta}^*$ is solution of

$$\operatorname{argmin}\{\|\mathbf{y} - \boldsymbol{\theta}\|^2\}$$

subject to $\theta_i + \xi_i[\mathbf{x}_j - \mathbf{x}_i] \leq \boldsymbol{\theta}_j, \forall i, j$

and $\xi_1, \dots, \xi_n \in \mathbb{R}^n$.

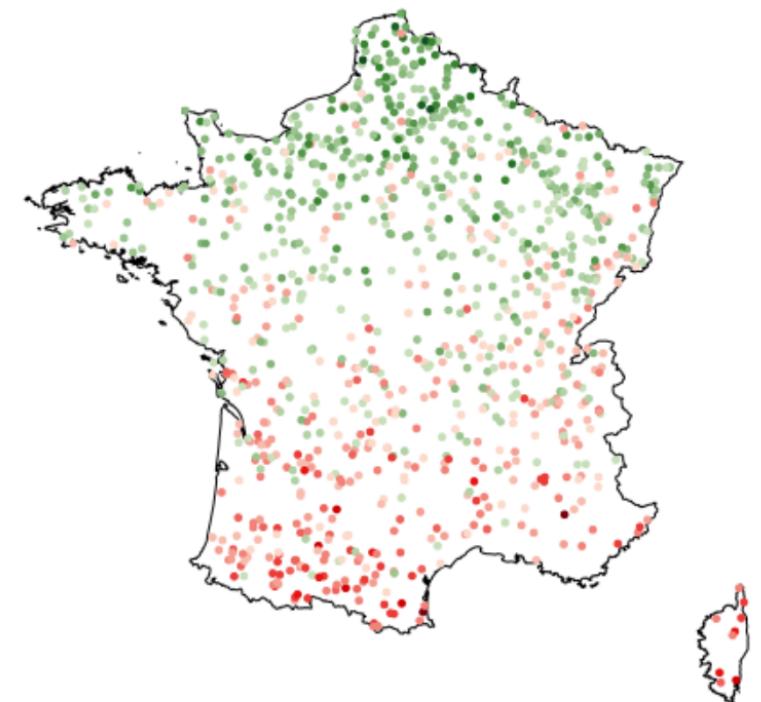
Spatial Smoothing

One can also consider some spatial smoothing, if we want to predict $\mathbb{E}[Y|X = \mathbf{x}]$ for some coordinate \mathbf{x} .

```
1 > library(rgeos)
2 > library(rgdal)
3 > library(mapproj)
4 > library(cartography)
5 > download.file("http://bit.ly/2G3KIUG","zonier.RData")
6 > load("zonier.RData")
7 > cols=rev(carto.pal(pal1="red.pal",n1=10,pal2="green.pal",n2=10))
8 > download.file("http://bit.ly/2GSvzGW","FRA_adm0.rds")
9 > download.file("http://bit.ly/2FUZ0Lz","FRA_adm2.rds")
10 > FR=readRDS("FRA_adm2.rds")
11 > donnees_carte=data.frame(FRdata)
```

Spatial Smoothing

```
1 > FR0=readRDS("FRA_adm0.rds")
2 > plot(FR0)
3 > bk = seq(-5,4.5,length=21)
4 > cuty = cut(simbase$Y,breaks=bk,labels
   =1:20)
5 > points(simbase$long,simbase$lat,col=cols
   [cuty],pch=19,cex=.5)
```



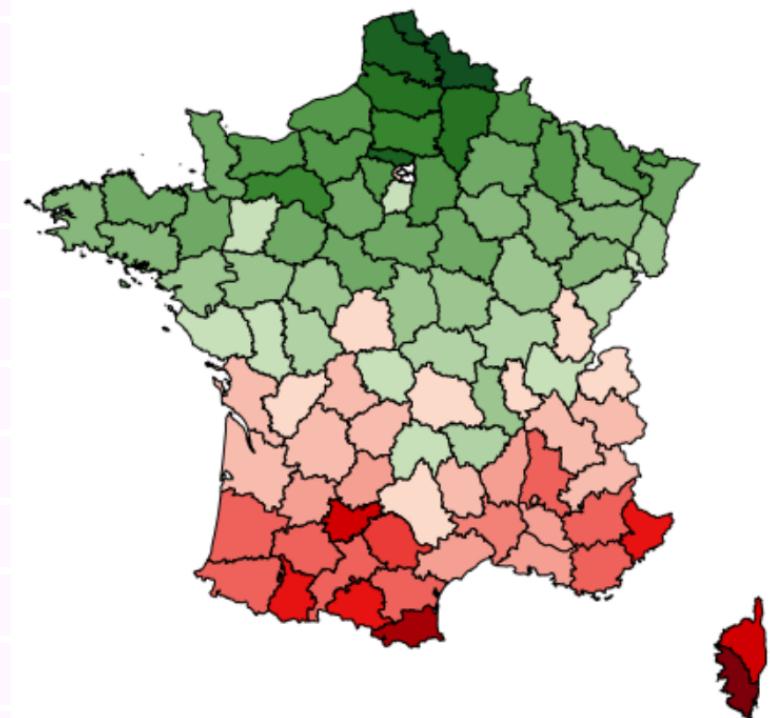
One can consider a **choropleth map** (spatial version of the histogram).

```

1 > A=aggregate(x = simbase$Y, by=list(
2   simbase$dpt), mean)
3 > names(A)=c("dpt","y")
4 > d=donnees_carte$CCA_2
5 > d[d=="2A"]="201"
6 > d[d=="2B"]="202"
7 > donnees_carte$dpt=as.numeric(as.character(d))
8 > donnees_carte=merge(donnees_carte,A,all.x=TRUE)
9 > donnees_carte=donnees_carte[order(
10   donnees_carte$OBJECTID),]
11 > bk=seq(-2.75,2.75,length=21)
12 > donnees_carte$cuty=cut(donnees_carte$y,
13   breaks=bk,labels=1:20)
14 > plot(FR, col=cols[donnees_carte$cuty],
15   xlim=c(-5.2,12))

```

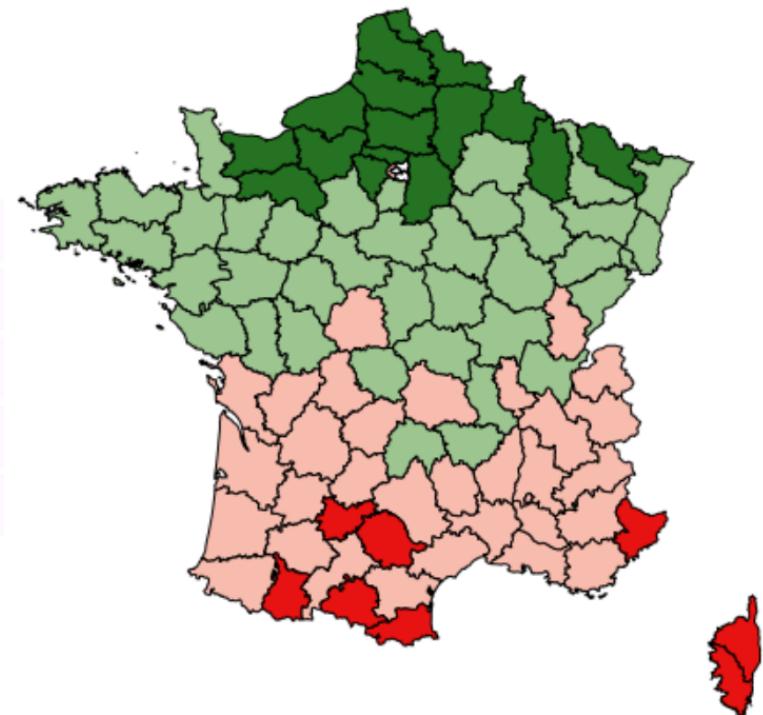
Spatial Smoothing



Spatial Smoothing

Instead of a "continuous" gradient of colors, one can consider only 4 colors (4 levels) for the prediction.

```
1 > bk=seq(-2.75,2.75,length=5)
2 > donnees_carte$cuty=cut(donnees_carte$y,
   breaks=bk,labels=1:4)
3 > plot(FR, col=cols[c(3,8,12,17)][donnees_
   carte$cuty],xlim=c(-5.2,12))
```



Spatial Smoothing

```
1 > P1 = FR0@polygons[[1]]@Polygons[[355]]  
    @coords  
2 > P2 = FR0@polygons[[1]]@Polygons[[27]]  
    @coords  
3 > plot(FR0, border=NA)  
4 > polygon(P1)  
5 > polygon(P2)  
6 > grille<-expand.grid(seq(min(simbase$long),  
    max(simbase$long), length=101), seq(min(simbase$lat),  
    max(simbase$lat), length=101))  
7 > paslong=(max(simbase$long)-min(simbase$long))/100  
8 > paslat=(max(simbase$lat)-min(simbase$lat))/100
```



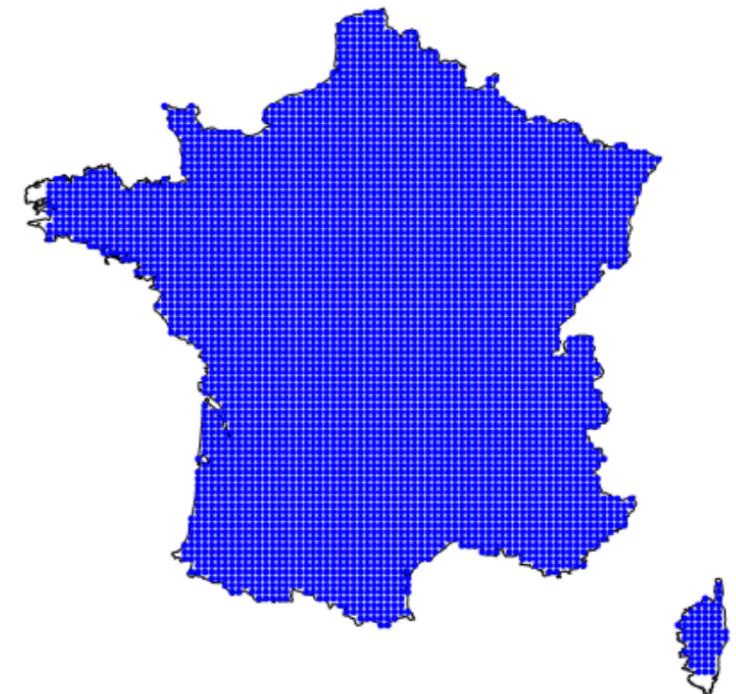
Spatial Smoothing

We need to create a grid (i.e. \mathcal{X}) on which we approximate $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$

```

1 > f=function(i){ (point.in.polygon (grille
2   [i, 1]+paslong/2 , grille[i, 2]+paslat/
3     2 , P1[,1],P1[,2])>0)+(point.in.polygon
4   (grille[i, 1]+paslong/2 , grille[i,
5     2]+paslat/2 , P2[,1],P2[,2])>0) }
6 > indic=unlist(lapply(1:nrow(grille),f))
7 > grille=grille[which(indic==1),]
8 > points(grille[,1]+paslong/2,grille[,2]-
9   paslat/2)

```



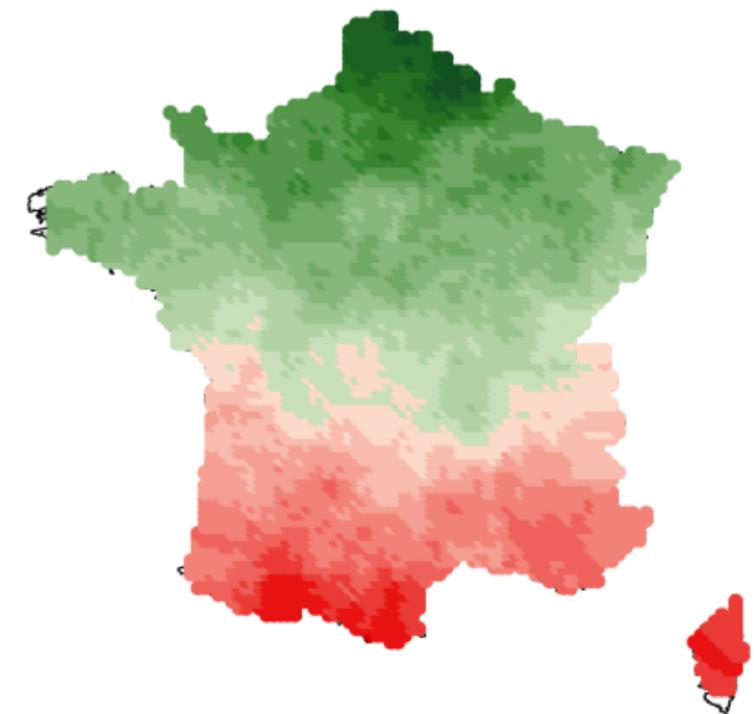
Spatial Smoothing

Consider here some k -NN, with $k = 20$

```

1 > library(geosphere)
2 > knn=function(i,k=20){
3 + d=distHaversine(grille[i,1:2],simbase[,c
+ ("long","lat")]], r=6378.137)
4 + r=rank(d)
5 + ind=which(r<=k)
6 + mean(simbase[ind,"Y"])
7 +
8 > grille$y=Vectorize(knn)(1:nrow(grille))
9 > bk=seq(-2.75,2.75,length=21)
10 > grille$cuty=cut(grille$y,breaks=bk,
+ labels=1:20)
11 > points(grille[,1]+paslong/2,grille[,2]+
+ paslat/2,col=cols[grille$cuty],pch=19)

```

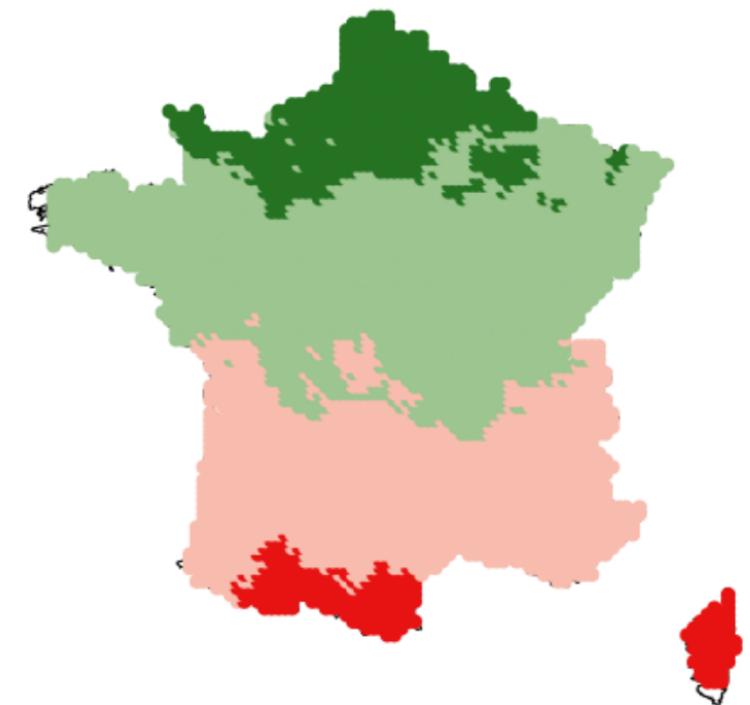


Spatial Smoothing

Again, instead of a "continuous" gradient, we can use 4 levels,

```

1 > bk=seq(-2.75,2.75,length=5)
2 > grille$cuty=cut(grille$y,breaks=bk,
  labels=1:4)
3 > plot(FR0,border=NA)
4 > polygon(P1)
5 > polygon(P2)
6 > points(grille[,1]+paslong/2,grille[,2]+
  paslat/2,col=cols[c(3,8,12,17)][grille$cuty],pch=19)
```



Testing (Non-)Linearities

In the linear model,

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \underbrace{\mathbf{X} [\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{X}^\top}_{\mathbf{H}} \mathbf{y}$$

$\mathbf{H}_{i,i}$ is the leverage of the i th element of this hat matrix.

Write

$$\hat{y}_{\textcolor{red}{i}} = \sum_{\textcolor{blue}{j}=1}^n [\mathbf{X}_{\textcolor{red}{i}}^\top [\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{X}^\top]_{\textcolor{blue}{j}} y_{\textcolor{blue}{j}} = \sum_{\textcolor{blue}{j}=1}^n [\mathcal{H}(\mathbf{X}_{\textcolor{red}{i}})]_{\textcolor{blue}{j}} y_{\textcolor{blue}{j}}$$

where

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

The prediction is

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

Testing (Non-)Linearities

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

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

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, $\mathbf{S} = \mathbf{H}$.

$\text{trace}(\mathbf{H}) = \dim(\boldsymbol{\beta})$: degrees of freedom

$\frac{\mathbf{H}_{i,i}}{1 - \mathbf{H}_{i,i}}$ is related to Cook's distance, from Cook (1977), Detection of Influential Observations in Linear Regression.

Testing (Non-)Linearities

For a **kernel regression model**, with kernel k and bandwidth h

$$S_{i,j}^{(k,h)} = \frac{k_h(x_i - x_j)}{\sum_{k=1}^n k_h(x_k - x_j)}$$

$$\text{where } k_h(\cdot) = k(\cdot/h), \text{ while } \mathcal{S}^{(k,h)}(\mathbf{x})_j = \frac{K_h(\mathbf{x} - x_j)}{\sum_{k=1}^n k_h(\mathbf{x} - x_k)}$$

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

Testing (Non-)Linearities

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

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

Define

$$T = \frac{\|\mathbf{S}\mathbf{y} - \mathbf{H}\mathbf{y}\|}{\text{trace}([\mathbf{S} - \mathbf{H}]^\top [\mathbf{S} - \mathbf{H}])}$$

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

Remark: In the case of a linear predictor, with smoothing matrix \mathbf{S}_h

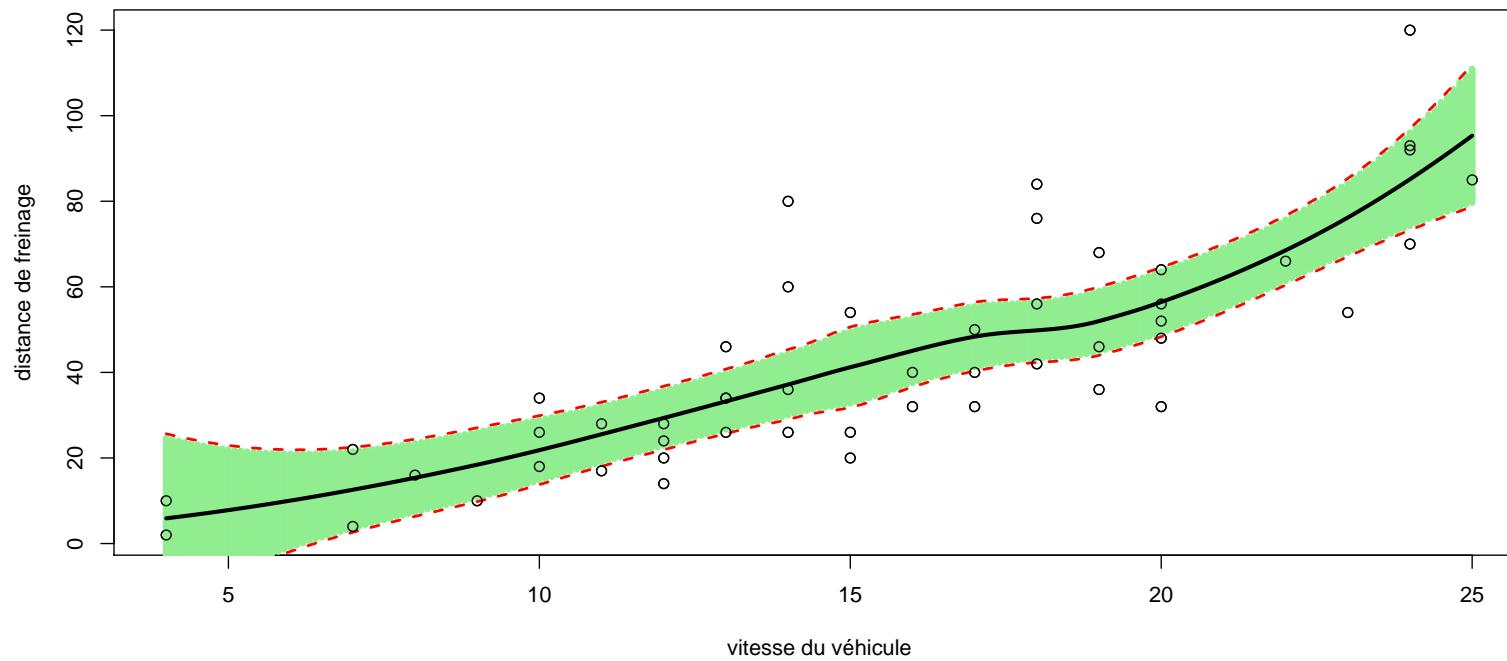
$$\widehat{R}(h) = \frac{1}{n} \sum_{i=1}^n (y_i - \widehat{m}_h^{(-i)}(\mathbf{x}_i))^2 = \frac{1}{n} \sum_{i=1}^n \left(\frac{Y_i - \widehat{m}_h(\mathbf{x}_i)}{1 - [\mathbf{S}_h]_{i,i}} \right)^2$$

We do not need to estimate n models. One can also minimize

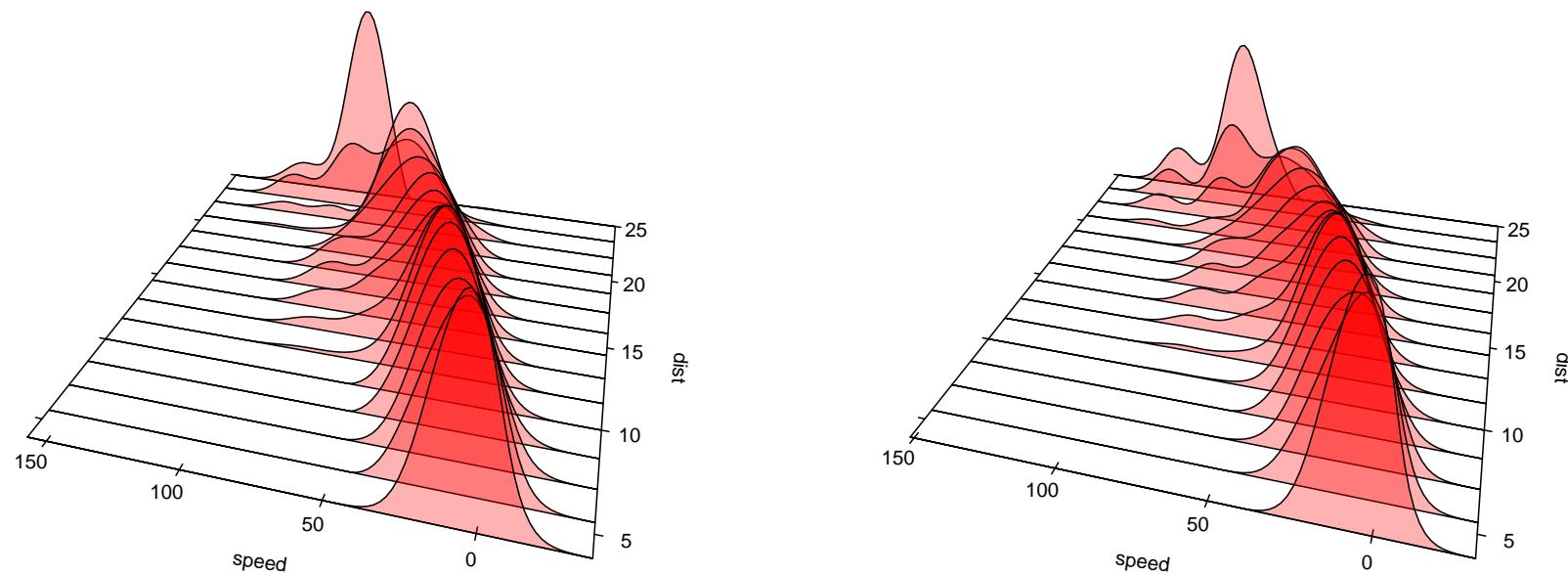
$$GCV(h) = \frac{n^2}{n^2 - \text{trace}(\mathbf{S})^2} \cdot \frac{1}{n} \sum_{i=1}^n (Y_i - \widehat{m}_h(\mathbf{x}_i))^2 \sim \text{Mallow's } C_p$$

Confidence Intervals

If $\hat{y} = \hat{m}_h(\mathbf{x}) = S_h(\mathbf{x})\mathbf{y}$, let $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{m}_h(\mathbf{x}_i))^2$ and a confidence interval is, at \mathbf{x} $\left[\hat{m}_h(\mathbf{y}) \pm t_{1-\alpha/2} \hat{\sigma} \sqrt{S_h(\mathbf{x})S_h(\mathbf{x})^\top} \right]$.



Confidence Bands



To go further see functional confidence regions [»](#)

Boosting to Capture NonLinear Effects

We want to solve

$$m^* = \operatorname{argmin} \left\{ \mathbb{E}[(Y - m(\mathbf{X}))^2] \right\}$$

The heuristics is simple: we consider an iterative process where we keep modeling the errors.

Fit model for \mathbf{y} , $h_1(\cdot)$ from \mathbf{y} and \mathbf{X} , and compute the error, $\boldsymbol{\varepsilon}_1 = \mathbf{y} - h_1(\mathbf{X})$.

Fit model for $\boldsymbol{\varepsilon}_1$, $h_2(\cdot)$ from $\boldsymbol{\varepsilon}_1$ and \mathbf{X} , and compute the error, $\boldsymbol{\varepsilon}_2 = \boldsymbol{\varepsilon}_1 - h_2(\mathbf{X})$, etc. Then set

$$m_k(\cdot) = \underbrace{h_1(\cdot)}_{\sim \mathbf{y}} + \underbrace{h_2(\cdot)}_{\sim \boldsymbol{\varepsilon}_1} + \underbrace{h_3(\cdot)}_{\sim \boldsymbol{\varepsilon}_2} + \cdots + \underbrace{h_k(\cdot)}_{\sim \boldsymbol{\varepsilon}_{k-1}}$$

Hence, we consider an iterative procedure, $m_k(\cdot) = m_{k-1}(\cdot) + h_k(\cdot)$.

Boosting

$h(\mathbf{x}) = \mathbf{y} - m_k(\mathbf{x})$, which can be interpreted as a residual. Note that this residual is the gradient of $\frac{1}{2}[\mathbf{y} - m_k(\mathbf{x})]^2$

A gradient descent is based on Taylor expansion

$$\underbrace{f(\mathbf{x}_k)}_{\langle f, \mathbf{x}_k \rangle} \sim \underbrace{f(\mathbf{x}_{k-1})}_{\langle f, \mathbf{x}_{k-1} \rangle} + \underbrace{(\mathbf{x}_k - \mathbf{x}_{k-1})}_{\alpha} \underbrace{\nabla f(\mathbf{x}_{k-1})}_{\langle \nabla f, \mathbf{x}_{k-1} \rangle}$$

But here, it is different. We claim we can write

$$\underbrace{f_k(\mathbf{x})}_{\langle f_k, \mathbf{x} \rangle} \sim \underbrace{f_{k-1}(\mathbf{x})}_{\langle f_{k-1}, \mathbf{x} \rangle} + \underbrace{(f_k - f_{k-1})}_{\beta} \underbrace{?}_{\langle f_{k-1}, \nabla \mathbf{x} \rangle}$$

where ? is interpreted as a ‘gradient’.

Boosting

Construct iteratively

$$m_k(\cdot) = m_{k-1}(\cdot) + \operatorname{argmin}_{h \in \mathcal{H}} \left\{ \sum_{i=1}^n (y_i - [m_{k-1}(\mathbf{x}_i) + h(\mathbf{x}_i)])^2 \right\}$$

$$m_k(\cdot) = m_{k-1}(\cdot) + \operatorname{argmin}_{h \in \mathcal{H}} \left\{ \sum_{i=1}^n ([y_i - m_{k-1}(\mathbf{x}_i)] - h(\mathbf{x}_i))^2 \right\}$$

where $h \in \mathcal{H}$ means that we seek in a class of **weak learner functions**.

If learner are two strong, the first loop leads to some fixed point, and there is no learning procedure, see linear regression $y = \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon$. Since $\varepsilon \perp \mathbf{x}$ we cannot learn from the residuals.

In order to make sure that we learn **weakly**, we can use some **shrinkage parameter** ν (or collection of parameters ν_j).

Boosting with Piecewise Linear Spline & Stump Functions

Instead of $\varepsilon_k = \varepsilon_{k-1} - h_k(\mathbf{x})$, set $\varepsilon_k = \varepsilon_{k-1} - \nu \cdot h_k(\mathbf{x})$

```

1 for(t in 1:100){fit = rpart(yr~x,data=df)
2                 yp = predict(fit,newdata=df)
3                 df$yr = df$yr - v*yp)}

```

Remark : bumps are related to regression trees.

Ruptures

One can use **Chow test** to test for a rupture. Note that it is simply Fisher test, with two parts,

$$\boldsymbol{\beta} = \begin{cases} \boldsymbol{\beta}_1 & \text{for } i = 1, \dots, i_0 \\ \boldsymbol{\beta}_2 & \text{for } i = i_0 + 1, \dots, n \end{cases} \quad \text{and test} \quad \begin{cases} H_0 : \boldsymbol{\beta}_1 = \boldsymbol{\beta}_2 \\ H_1 : \boldsymbol{\beta}_1 \neq \boldsymbol{\beta}_2 \end{cases}$$

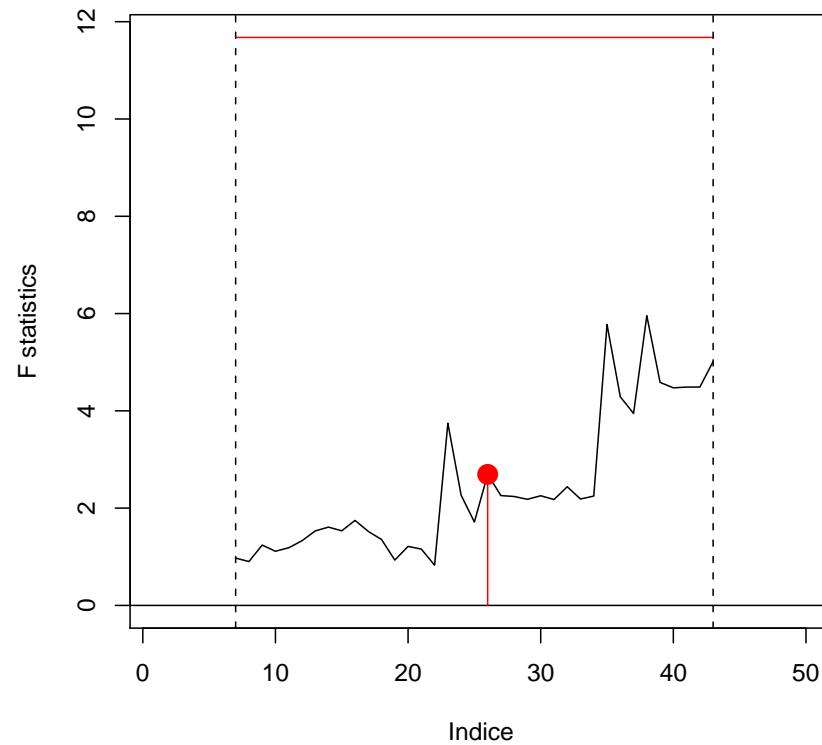
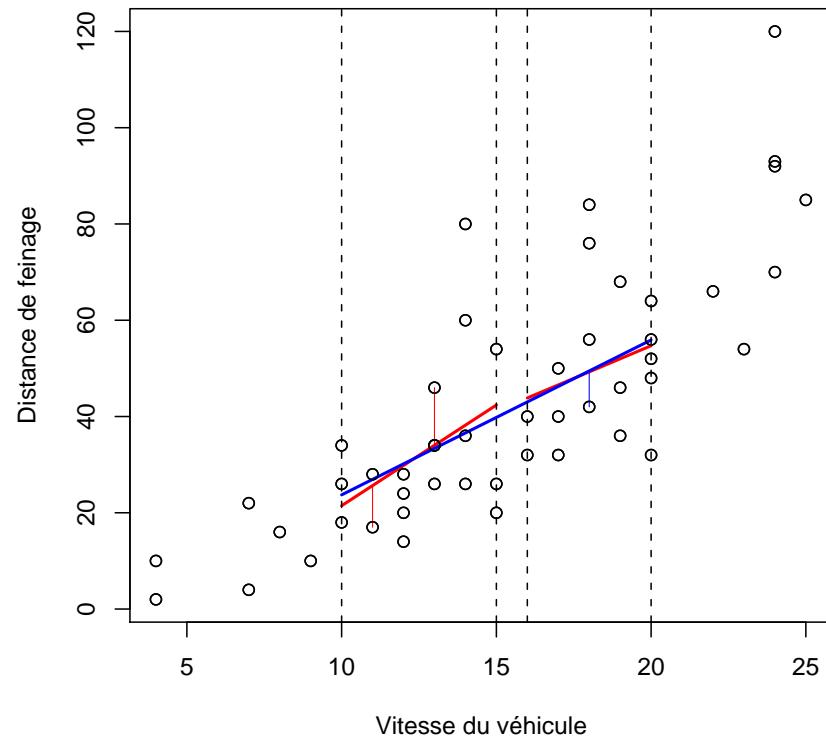
i_0 is a point between k and $n - k$ (we need enough observations). Chow (1960) **Tests of Equality Between Sets of Coefficients in Two Linear Regressions** suggested

$$F_{i_0} = \frac{\hat{\boldsymbol{\eta}}^T \hat{\boldsymbol{\eta}} - \hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}}}{\hat{\boldsymbol{\varepsilon}}^T \hat{\boldsymbol{\varepsilon}} / (n - 2k)}$$

where $\hat{\boldsymbol{\varepsilon}}_i = y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}$, and $\hat{\boldsymbol{\eta}}_i = \begin{cases} Y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_1 & \text{for } i = k, \dots, i_0 \\ Y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}_2 & \text{for } i = i_0 + 1, \dots, n - k \end{cases}$

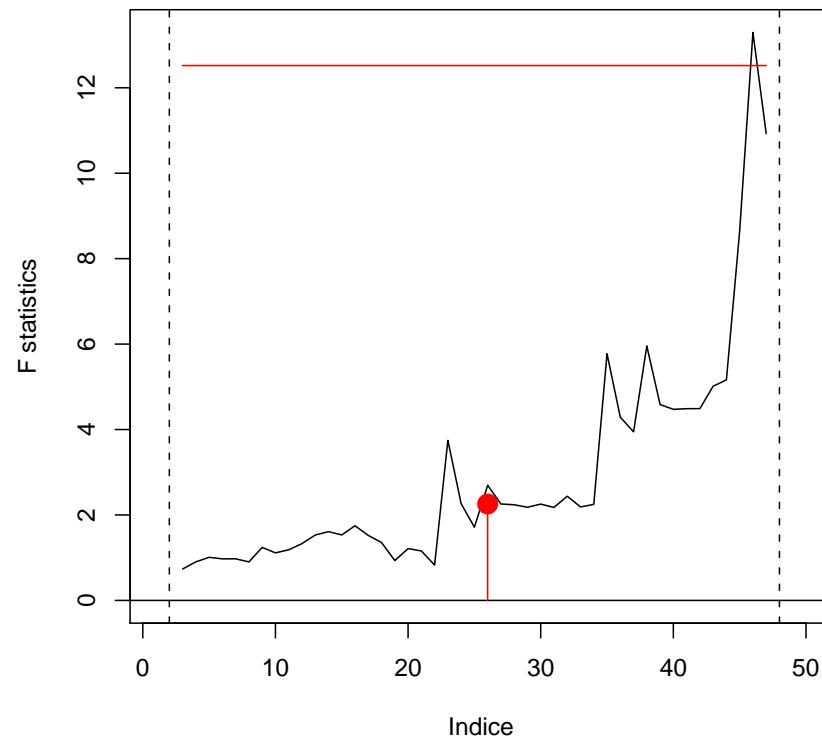
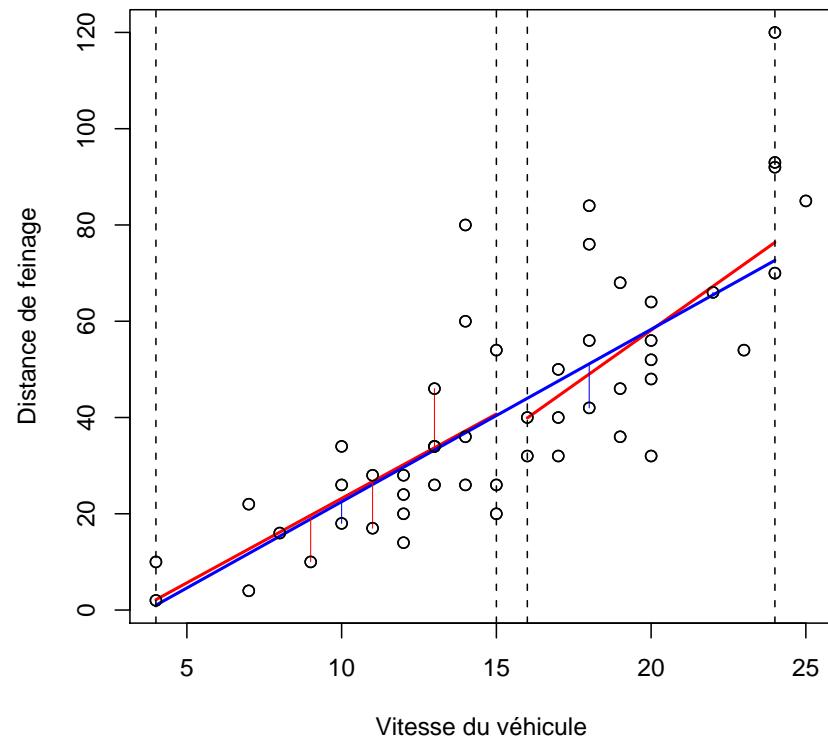
Ruptures

```
1 > Fstats(dist ~ speed, data=cars, from=7/50)
```



Ruptures

```
1 > Fstats(dist ~ speed, data=cars, from=2/50)
```



Ruptures

If i_0 is unknown, use CUSUM types of tests, see Ploberger & Krämer (1992) [The Cusum Test with OLS Residuals](#). For all $t \in [0, 1]$, set

$$W_t = \frac{1}{\hat{\sigma}\sqrt{n}} \sum_{i=1}^{\lfloor nt \rfloor} \hat{\varepsilon}_i.$$

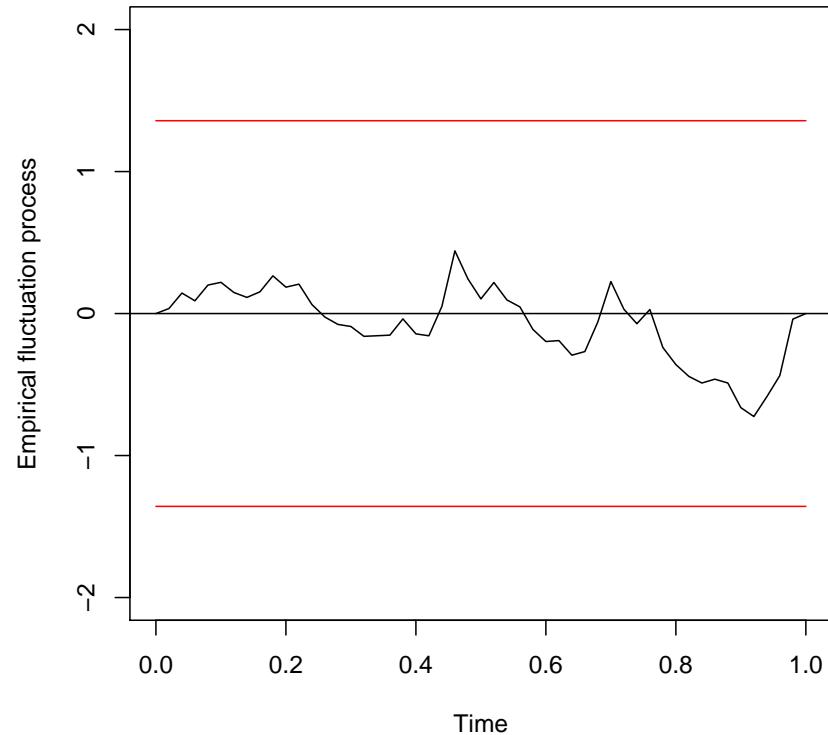
If α is the confidence level, bounds are generally $\pm\alpha$, even if theoretical bounds should be $\pm\alpha\sqrt{t(1-t)}$.

```

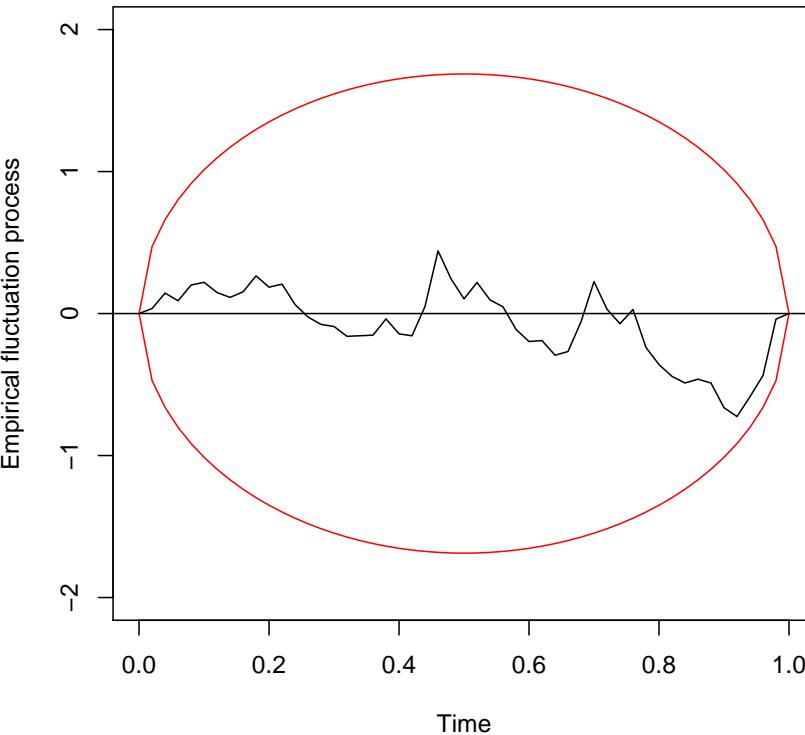
1 > cusum <- efp(dist ~ speed, type = "OLS-CUSUM", data=cars)
2 > plot(cusum, ylim=c(-2,2))
3 > plot(cusum, alpha = 0.05, alt.boundary = TRUE, ylim=c(-2,2))
```

Ruptures

OLS-based CUSUM test



OLS-based CUSUM test with alternative boundaries



From a Rupture to a Discontinuity

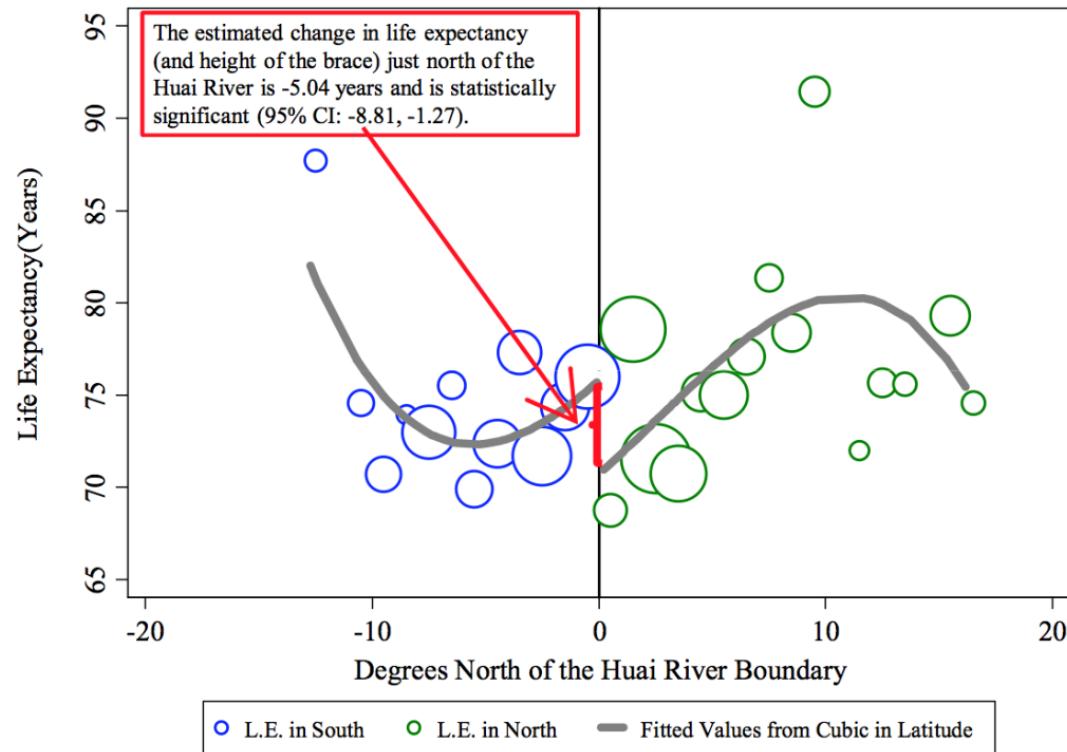


Fig. 3. The plotted line reports the fitted values from a regression of life expectancy on a cubic in latitude using the sample of DSP locations, weighted by the population at each location.

See Imbens & Lemieux (2008) **Regression Discontinuity Designs**.

From a Rupture to a Discontinuity

Consider the dataset from Lee (2008) [Randomized experiments from non-random selection in U.S. House elections](#).

```
1 > library(RDDtools)  
2 > data(Lee2008)
```

We want to test if there is a discontinuity in
0.

- with parametric tools
- with nonparametric tools

Testing for a rupture

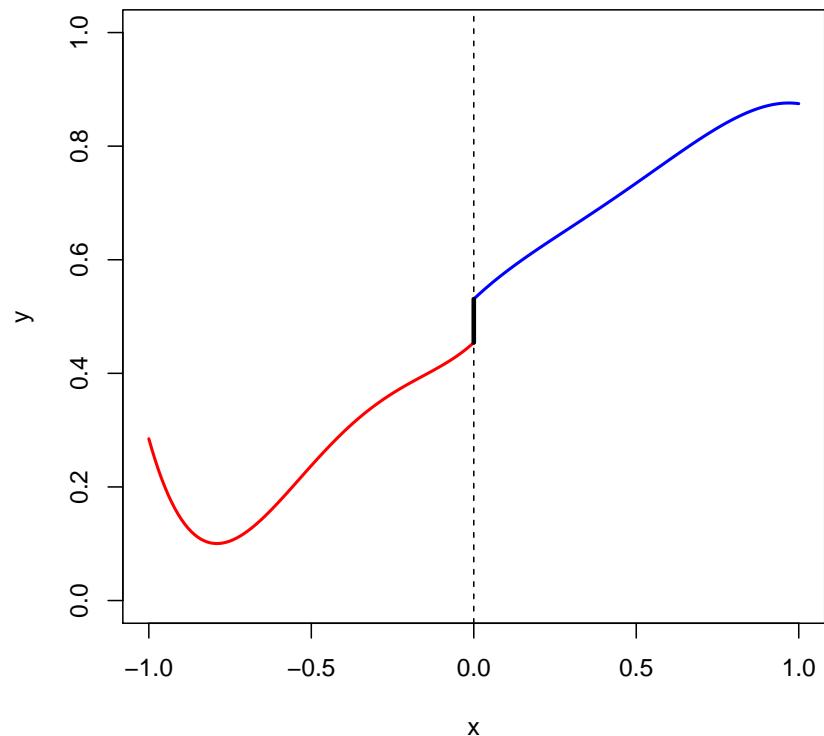
Use some 4th order polynomial, on each part

```
1 > idx1 = (Lee2008$x>0)
2 > reg1 = lm(y~poly(x,4),data=Lee2008[
   idx1,])
3 > idx2 = (Lee2008$x<0)
4 > reg2 = lm(y~poly(x,4),data=Lee2008[
   idx2,])
5 > s1=predict(reg1,newdata=data.frame(x
   =0))
6 > s2=predict(reg2,newdata=data.frame(x
   =0))
7 > abs(s1-s2)
8
9      1
9 0.07659014
```

Testing for a rupture

```

1 > reg_para <- RDDreg_lm(RDDdata(y =
2   Lee2008$y, x = Lee2008$x, cutpoint
3   = 0), order = 4)
4
5 > reg_para
6 #### RDD regression: parametric ####
7 Polynomial order: 4
8 Slopes: separate
9 Number of obs: 6558 (left: 2740,
10 right: 3818)
11
12 Coefficient:
13 Estimate Std. Error t value Pr(>|t|)
14 D 0.076590    0.013239   5.7851 7.582e-09
15 ***
```



Testing for a rupture

or use a simple local regression, see [Imbens & Kalyanaraman \(2012\)](#).

```
1 > reg1 = ksmooth(Lee2008$x[idx1],  
+                     Lee2008$y[idx1], kernel = "normal",  
+                     bandwidth = 0.1)  
2 > reg2 = ksmooth(Lee2008$x[idx2],  
+                     Lee2008$y[idx2], kernel = "normal",  
+                     bandwidth = 0.1)  
3 > s1 = reg1$y[1]  
4 > s2 = reg2$y[length(reg2$y)]  
5 > abs(s1-s2)  
6 [1] 0.09883813
```

Testing for a rupture

```

1 > reg_nonpara <- RDDreg_np(RDDobject = Lee2008_
  rdd, bw = .1)
2 > print(reg_nonpara)
3 ##> RDD regression: nonparametric local linear
4 Bandwidth: 0.1
5 Number of obs: 1209 (left: 577, right: 632)
6
7 Coefficient:
8 Estimate Std. Error z value Pr(>|z|)
9 D 0.059397 0.014119 4.207 2.588e-05 ***

```

