



UNIVERSITÀ
DEGLI STUDI
DI TRIESTE

Deams

Dipartimento di

Scienze Economiche, Aziendali,
Matematiche e Statistiche "Bruno de Finetti"

Non parametric statistics

Introduction

Francesco Pauli

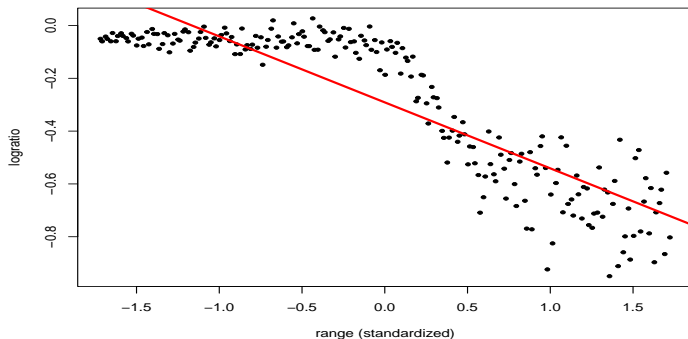
A.A. 2021/2022

Estimation of f

The linear model

$$y_i = f(x_i) + \varepsilon_i = \beta_1 + \beta_2 x_i + \varepsilon_i$$

where $\varepsilon_i \sim IID(\mathcal{N}(0, \sigma^2))$ clearly does not work

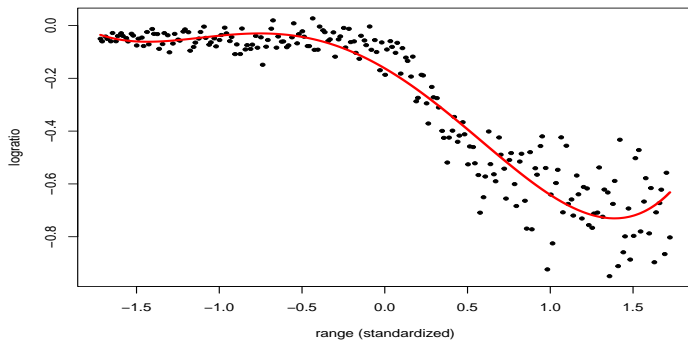


Naive solution 1: polynomial model

A first idea is using a polynomial model, effective if the degree is high enough

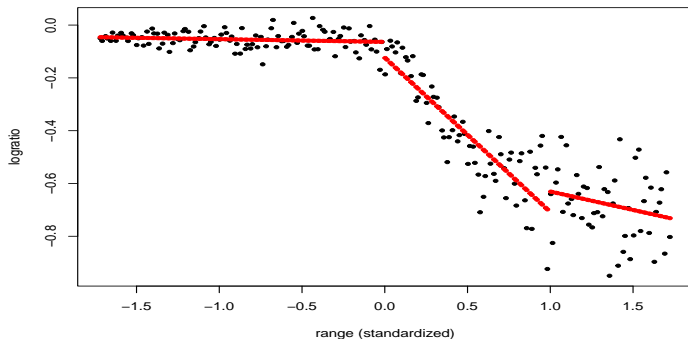
$$f(x; \beta) = \beta_1 + \beta_2 x + \beta_3 x^2 + \beta_4 x^3 + \beta_5 x^4$$

however, it has various problems.



Naive solution 2: piecewise linear model

An alternative might be using a piecewise linear model.



Indice

Spline

Penalized likelihood

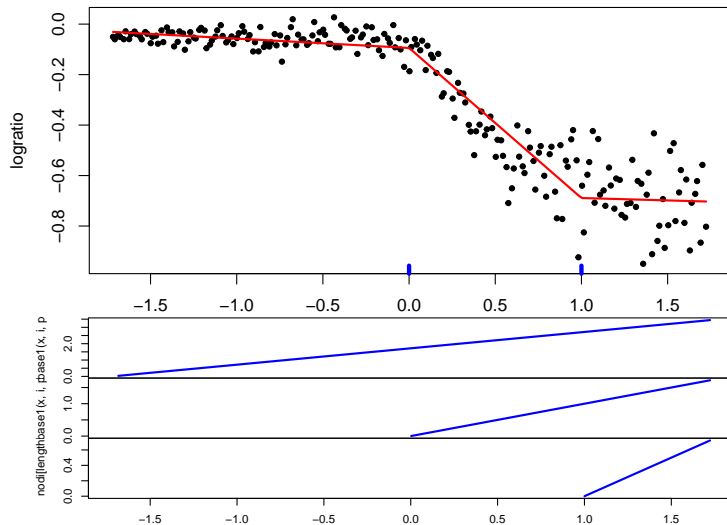
Why splines

Other basis

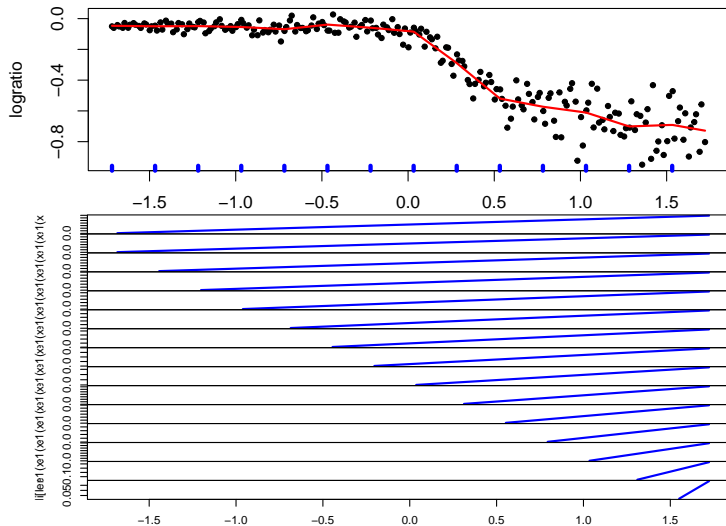
More covariates

Generalized models (non gaussian data)

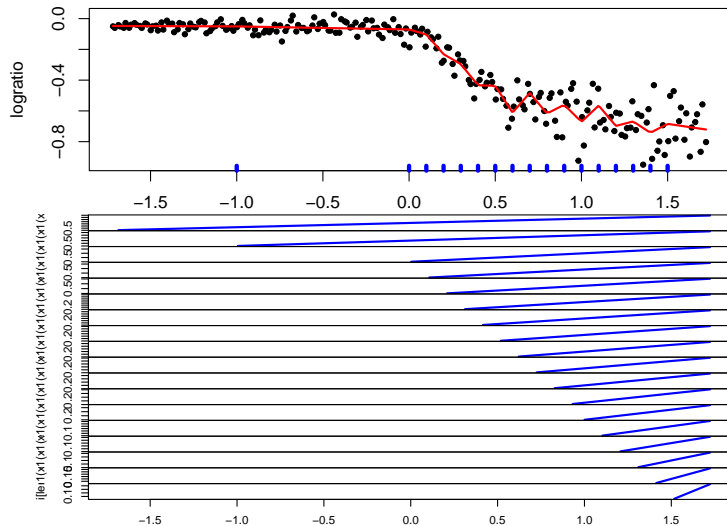
Linear spline basis: a two-knots example



Linear spline basis: general specification



Linear spline basis: general specification



Truncated power basis

An obvious extension of the linear basis is the truncated power basis

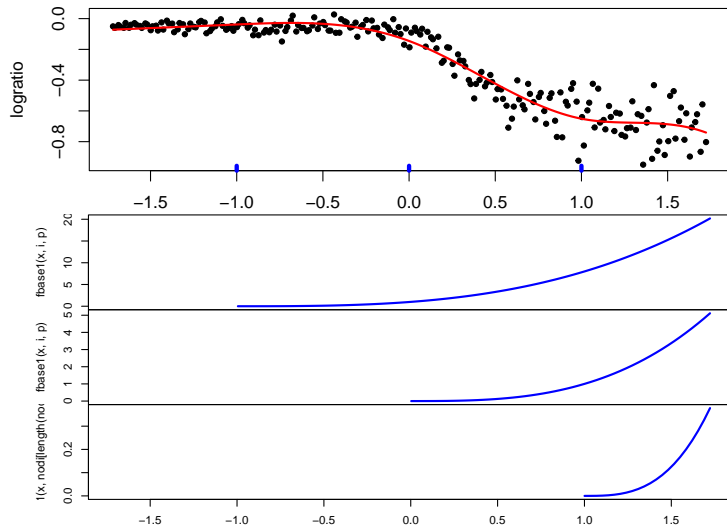
$$y_i = \beta_1 + \beta_2 x_i + \dots + \beta_{p+1} x_i^p + \sum_{k=1}^K b_k (x_i - \kappa_k)_+^p + \varepsilon_i$$

then the spline function is

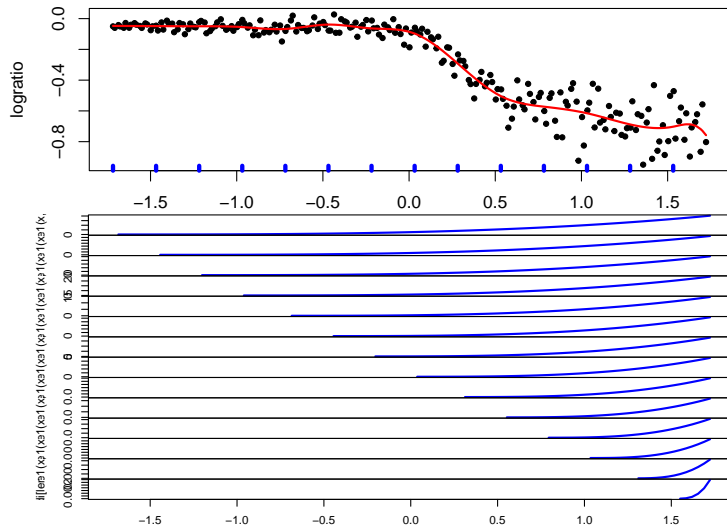
$$f(x) = \beta_1 + \beta_2 x + \dots + \beta_{p+1} x^p + \sum_{k=1}^K b_k (x - \kappa_k)_+^p$$

- ▶ A spline with degree p has $p - 1$ continuous derivatives,
- ▶ $p = 3$ is sufficient for most purposes (unless we want smooth derivatives).

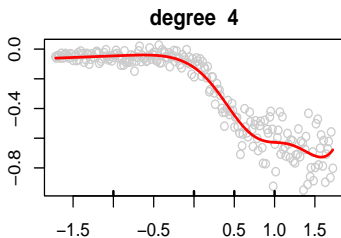
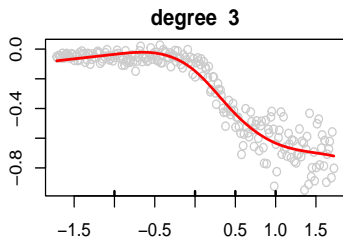
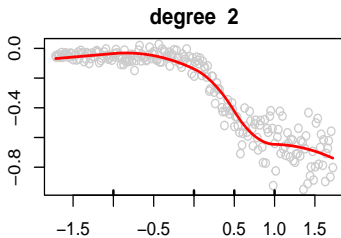
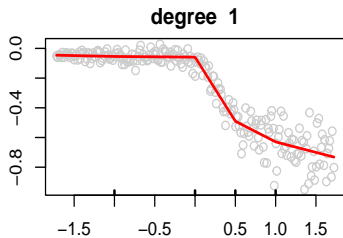
Truncated power basis



Truncated power basis



TPB: different degrees



Smoothness of the spline: bias variance trade off

Assuming that the degree is fixed (as is usually done), the choice of the number of knots is analogous to the choice of the bandwidth in kernel regression in that it implies a bias variance trade off

- ▶ more knots \leftrightarrow less smoothing \leftrightarrow less bias but more variance;
- ▶ less knots \leftrightarrow more smoothing \leftrightarrow more bias but less variance.

As for the choice of the bandwidth in kernel regression, choice of the smoothness of the spline is crucial.



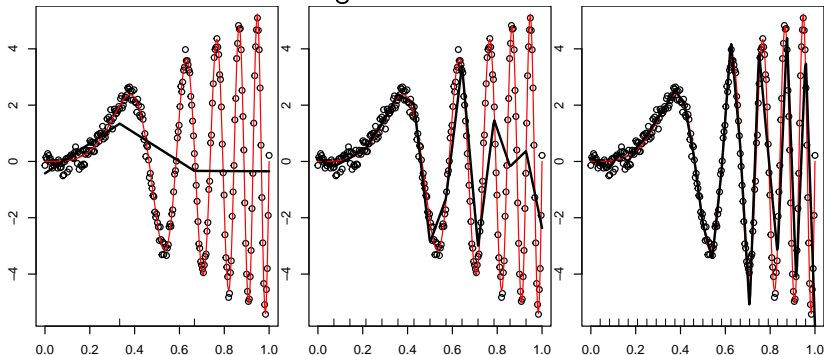
In principle, we may seek for the optimal level of smoothing by choosing the number of knots by estimating the mean square errors implied by different choices.



This means we must minimize the MSE with respect to the knots number and position, which is a difficult task (although doable, this a legitimate strategy which is actually pursued).

Number of knots and bias

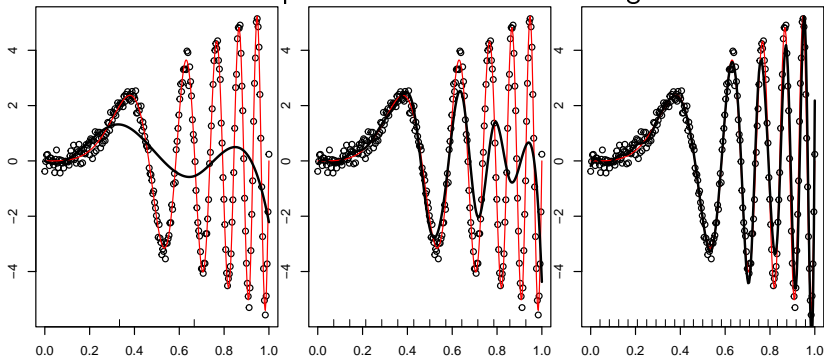
Bias is greater with curvature.



Knots represented on the x axis.

Number of knots and bias

With higher degree splines curvature is less of a problem, however if the function is complicated few knots lead to a great bias.



Knots represented on the x axis.

Smoothness of the spline: fixed knots

A different strategy may be used by first fixing the knots and then impose some restriction on the coefficients such that by changing the restriction we change the level of smoothing.

or

rather than estimating the coefficient as the minimum of the sum of squares, add a penalization to it which favours smoother functions over wigglier ones.

In this way, the smoothness can be tuned by a number (rather than a vector of unknown length), which is easier to work with.



The second alternative leads to the idea of penalized sum of squares. (Note that it is equivalent to the first for some choices of constraint/penalization.)

Indice

Spline

Penalized likelihood

Why splines

Other basis

More covariates

Generalized models (non gaussian data)

Smoothness and penalization

We consider penalized splines: instead of minimizing

$$\sum_{i=1}^n (y_i - f(x_i; \beta, \mathbf{b}))^2$$

we minimize

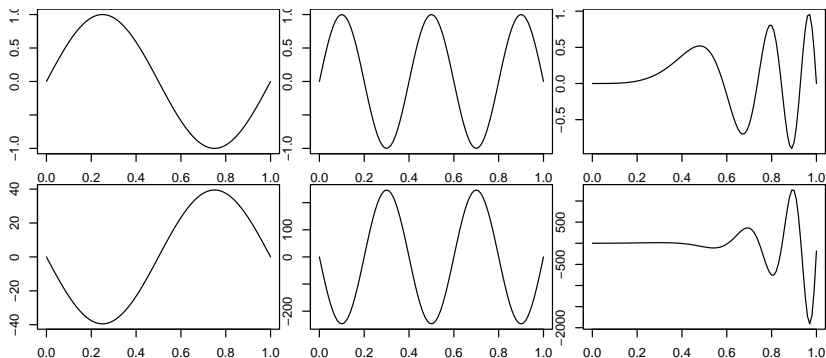
$$\sum_{i=1}^n (y_i - f(x_i, \beta, \mathbf{b}))^2 + \lambda S(f(x, \beta, \mathbf{b}))$$

where

- ▶ $S(f(x, \beta, \mathbf{b}))$ is a measure of the smoothness of f , increasing as the wiggleness of f increases
- ▶ $\lambda > 0$ is a fixed constant (for now).

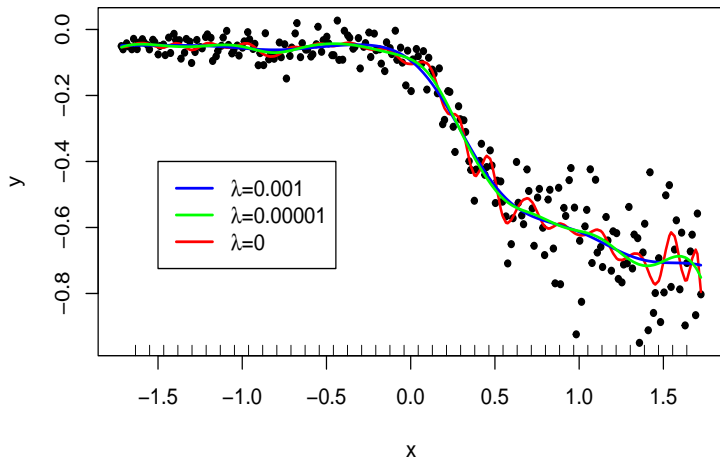
Functions and their second derivative

The second derivative, which is typically used for splines of third degree, is easily interpreted since it is a measure of the curvature (although not mathematically precise)



Penalized splines

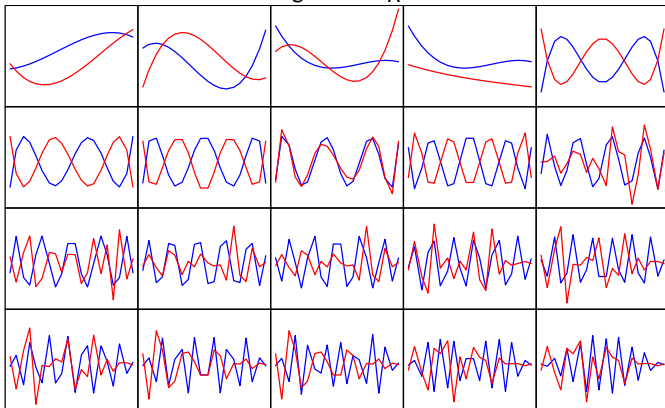
40 knots are used



Unpenalized and penalized splines

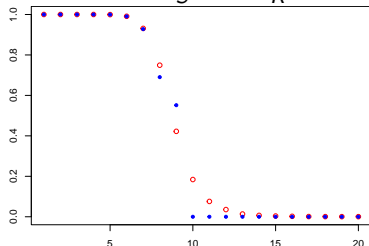
Smoothing splines v. low rank: decomposing L

To appreciate the difference more precisely consider an eigenvalue decomposition of the two matrices L_S and L_R .



Smoothing splines v. low rank: decomposing L

To appreciate the difference more precisely consider an eigenvalue decomposition of the two matrices L_S and L_R .



The eigenvalues of both the smoothing spline (red) and the regression spline (blue) show that only the first 9 eigenvectors count (in both), the following are exactly 0 for the regression spline but are very low for the smoothing spline, thus they lead to similar smooths.

Practical: fitting splines without penalization

```
## Create a sample
x=seq(0,1,length=250) #sort(runif(150,0,1)))
m=m*sin(2*pi*x^3) # (0.5+5*x)*sin(10*pi*x^3)
y=m+rnorm(length(x),0,0.4)
##
## compute the covariate matrix
base=function(x,nodi,p=1){
  X=cbind(outer(x,0:p,FUN=function(x,y) x^y),
          outer(x,nodi,FUN=function(x,y) ifelse(x-y>0,(x-y)^p,0)))
}
##
## fix knots and degree
nodi=seq(0.1,0.9,by=0.1)
p=3
X=base(x,nodi,p)
##
## If no penalization is required the fit is obtained through
fit=lm(y~X-1)
yt=X %*% fit$coef
##
## or directly
yt1 = X %*% solve(t(X) %*% X) %*% t(X) %*% y
##
## we can plot the results as
plot(x,y)
rug(nodi)
lines(x,m,col="red")
lines(x,yt,lwd=2)
lines(x,yt1,lwd=2,col="green")
```

Practical: fitting splines with penalization

```
## If a penalization is to be used we need to define the penalization matrix
## let us assume a ridge type penalization
D = diag(c(rep(0,p+1),rep(1,length(nodi))))
## then the theoretical values are found as
lambda=0.001
yt1 = X %*% solve(t(X) %*% X+lambda*D) %*% t(X) %*% y
lambda=0.00001
yt2 = X %*% solve(t(X) %*% X+lambda*D) %*% t(X) %*% y
##
## we can plot the results as
plot(x,y)
rug(nodi)
lines(x,m,col="red")
lines(x,yt1,lwd=2)
lines(x,yt2,lwd=2,col="blue")
```

Indice

Spline

Penalized likelihood

Choice of λ

Why splines

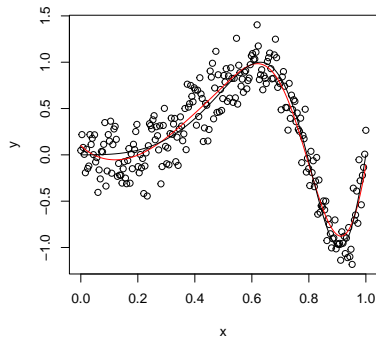
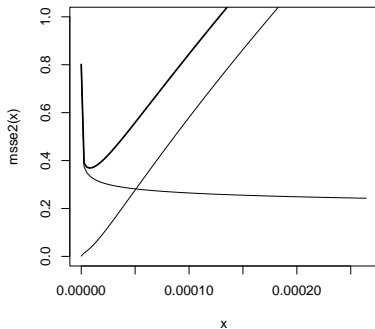
Other basis

More covariates

Generalized models (non gaussian data)

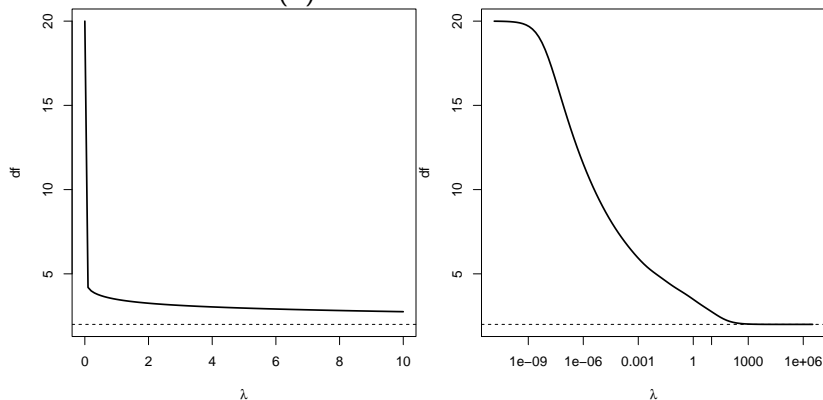
Mean summed square error

Example of MSSE decomposed into bias squared and variance



Degrees of freedom and penalization

The coefficient λ determines the smoothness of the estimated function, it is relevant to look at the relationship between λ and the degrees of freedom of the fit: $df = \text{tr}(L)$.



Fitting a spline: gam (mgcv)

There are many packages in R to estimate a spline, one of the most powerful and versatile is the package `gam` by Wood.



The functions provided have many options to tune the estimate, we will look at many of them, first, however, we use it in its simplest form.

```
fit=gam(y~s(x))
fit.s=summary(fit)
plot(fit)
```

This performs a fit with GCV choice of the smoothing parameter and default choice of knots.

Fitting a spline with gam

Arguments to the function `s()` allow to choose

`k` the basis dimension

`fx` whether a penalty should be used

```
fit2=gam(y~s(x,k=3))
fit4=gam(y~s(x,k=10))
plot(fit2)
plot(fit8)
```

Choosing the knots is also possible, but it requires discussing the basis first.

Experimenting with gam

The number of knots does not matter provided it is high enough.

```
sim=data.frame(x=seq(0,1,length=200)) #sort(runif(150,0,1)))
sim$m=sin(2*pi*sim$x^3)
sim$y=sim$m+rnorm(nrow(sim),0,0.4)
plot(sim$x,sim$y)
fit0=gam(y~s(x),data=sim)
fit1=gam(y~s(x,k=5),data=sim)
fit2=gam(y~s(x,k=12),data=sim)
fit3=gam(y~s(x,k=30),data=sim)
```


Variance estimation

The variance σ^2 of the error may be estimated, by analogy with LM, as

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{f}(x_i))^2}{n - \text{df}} = \frac{RSS}{n - \text{df}}$$

We already defined the degrees of freedom of the spline as $\text{tr}(L)$.

However, note that

$$\begin{aligned} E(RSS) &= E((y - \hat{y})^T (y - \hat{y})) \\ &= E(y^T (L - I)^T (L - I) y) \\ &= f^T (L - I)^T (L - I) f + \sigma^2 \text{tr}((L - I)^T (L - I)) \\ &= f^T (L - I)^T (L - I) f + \sigma^2 (\text{tr}(LL^T) - 2\text{tr}(L) + n) \end{aligned}$$

thus, assuming the bias is negligible, an unbiased estimator for σ^2 is

$$\tilde{\sigma}^2 = \frac{RSS}{n - 2\text{tr}(L) + \text{tr}(LL^T)}$$

Variance estimation

The variance σ^2 of the error may be estimated, by analogy with LM, as

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{f}(x_i))^2}{n - \text{df}} = \frac{RSS}{n - \text{df}}$$

We already defined the degrees of freedom of the spline as $\text{tr}(L)$.

However, note that

$$\begin{aligned} E(RSS) &= E((y - \hat{y})^T (y - \hat{y})) \\ &= f^T (L - I)^T (L - I) f + \sigma^2 (\text{tr}(LL^T) - 2\text{tr}(L) + n) \end{aligned}$$

thus, assuming the bias is negligible, an unbiased estimator for σ^2 is

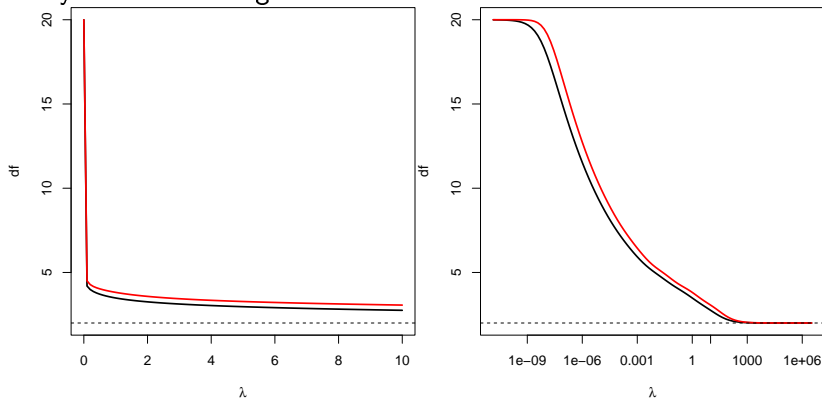
$$\tilde{\sigma}^2 = \frac{RSS}{n - 2\text{tr}(L) + \text{tr}(LL^T)}$$

Note that

- ▶ $n - 2\text{tr}(L) + \text{tr}(LL^T)$ are the residual degrees of freedom
- ▶ $2\text{tr}(L) - \text{tr}(LL^T)$ is an alternative measure of the dof of the spline

Degrees of freedom: the two versions

The two measures of the degrees of freedom of the smoother are different especially in the mid-range of λ .

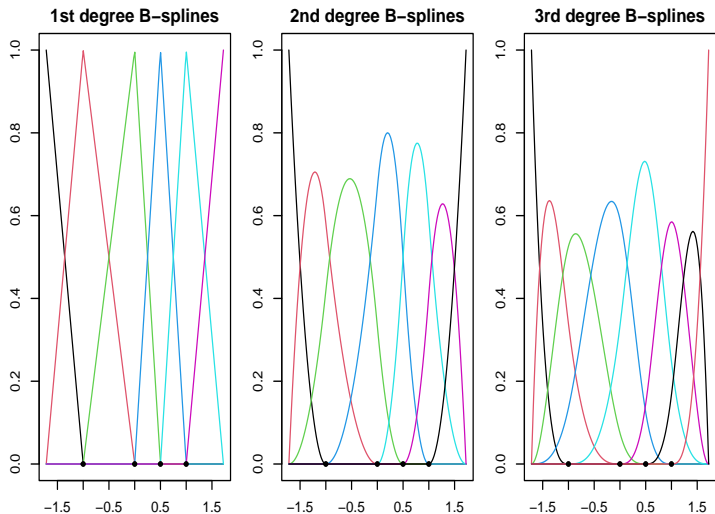


Why are the two equal for no smoothing and infinite smoothing?

Indice

Other basis

B-splines



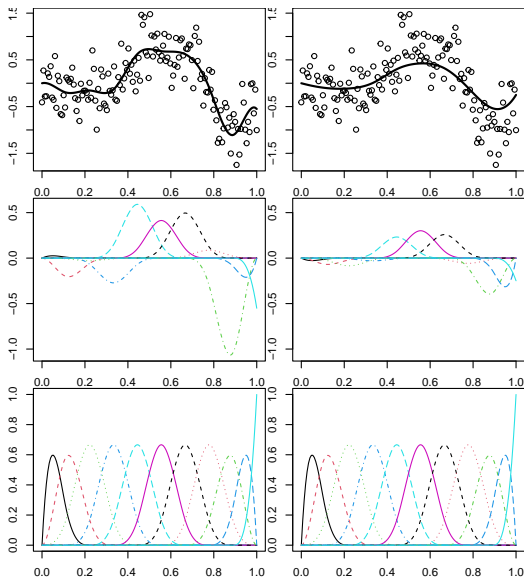
P-splines

P -splines are a low rank smoother using

- ▶ a B -splines basis
- ▶ a difference penalty (see above)

Usually they are defined on equally spaced knots (which makes the difference penalty more sensible).

Effect of difference penalty



Two spline estimates with a B -spline bases, left one has

$$\sum_{i=1}^{K-1} (b_{i+1} - b_i)^2 = 8.29$$

while for the right one

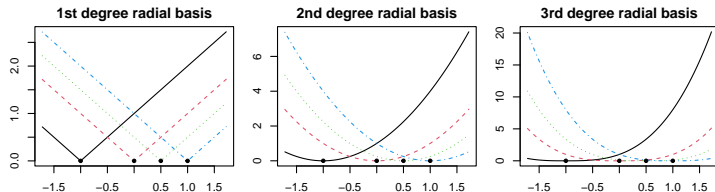
$$\sum_{i=1}^{K-1} (b_{i+1} - b_i)^2 = 0.83$$

(Middle panel: basis functions multiplied by the respective coefficients, that is, the final curve is the sum of these.)

Radial basis

A radial basis of order m with knots $\kappa_1, \dots, \kappa_K$ is defined as

$$1, x, \dots, x^m, B_k(x) = |x - \kappa_k|^m$$



The penalty matrix having elements

$$[D]_{ij} = |\kappa_i - \kappa_j|^3$$

Practical: fitting P-splines

```
## Create a sample
x=seq(0,1,length=250) #sort(runif(150,0,1)))
m=m*sin(2*pi*x^3) # (0.5+5*x)*sin(10*pi*x^3)
y=m+rnorm(length(x),0,0.4)
##
## fix knots and degree
nodi=seq(0.1,0.9,by=0.1)
p=3
library(splines)
X=bs(x,knots=nodi,degree=p)
##
## construction of the penalty
A=cbind(diag(c(rep(1,ncol(X)-1))),0)
A[col(A)==(1+row(A))]==-1
D = t(A) %*% A
##
## then the theoretical values are found as
lambda=100
yt1 = X %*% solve(t(X) %*% X+lambda*D) %*% t(X) %*% y
lambda=1
yt2 = X %*% solve(t(X) %*% X+lambda*D) %*% t(X) %*% y
##
## we can plot the results as
plot(x,y)
rug(nodi)
lines(x,m,col="red")
lines(x,yt1,lwd=2)
lines(x,yt2,lwd=2,col="blue")
```

Practical: P-splines, GCV for optimal smoothing

```
##
## we consider a sequence of values for lambda
lambdaseq=c(0.01,0.05,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1,1.5,2,2.5,3,3.5,4,4.5,5)
##
## and compute the GCV for each
gcvseq=rep(0,length(lambdaseq))
for (i in 1:length(lambdaseq)){
  L=X %>% solve(t(X) %>% X+lambdaseq[i]*D) %>% t(X)
  yt1 = L %>% y
  gcvseq[i]=mean(((y-yt1)/(1-diag(L)))^2)
}
plot(lambdaseq,gcvseq)
lambda=lambdaseq[which.min(gcvseq)]

L=X %>% solve(t(X) %>% X+lambda*D) %>% t(X)
ytopt = L %>% y
##
## we can plot the results as
plot(x,y)
rug(nodi)
lines(x,m,col="red")
lines(x,ytopt,lwd=2)
```

Practical: P-splines, bootstrap for variability bands

```
##
## Consider one lambda value
lambda=1
##
## and compute the smoothing matrix (which does not change)
L=X %*% solve(t(X) %*% X+lambda*D) %*% t(X)
yt1 = L %*% y
##
## residuals will be used for semiparametric bootstrapping
res=y-yt1
##
B=100
yt.b=matrix(NA,nrow=length(res),ncol=B)
for (i in 1:B){
  y.b=yt1+sample(res,length(res),replace=TRUE)
  yt.b[,i]=L %*% y.b
}
plot(x,y)
matlines(x,yt.b,col="black",lty=1)
```

Practical: fitting radial basis

```
## Create a sample
x=seq(0,1,length=250) #sort(runif(150,0,1)))
m=m=sin(2*pi*x^3) #(0.5+5*x)*sin(10*pi*x^3)
y=m+rnorm(length(x),0,0.4)
##
## fix knots and degree
nodi=seq(0.1,0.9,by=0.1)
p=3
X1=outer(x,0:p,FUN=function(x,y) x^y)
X2=outer(x,nodi,FUN=function(x,y) (abs(x-y))^p)
X=cbind(X1,X2)
##
## construction of the penalty
D=matrix(0,nrow=1+p+length(nodi),ncol=1+p+length(nodi))
D[-(1:(p+1)),-(1:(p+1))] = outer(nodi,nodi,FUN=function(x,y) (abs(x-y))^p)
##
## then the theoretical values are found as
lambda=1
yt1 = X %*% solve(t(X) %*% X+lambda*D) %*% t(X) %*% y
lambda=0.0001
yt2 = X %*% solve(t(X) %*% X+lambda*D) %*% t(X) %*% y
##
## we can plot the results as
plot(x,y)
rug(nodi)
lines(x,m,col="red")
lines(x,yt1,lwd=2)
lines(x,yt2,lwd=2,col="blue")
```

Indice

More covariates

Practical: multivariate spline

```

n=200
funz=function(x,z) 4*x*exp(-4*x^2 - 4*z^2)
simB=data.frame(x=runif(n,-1,1),z=runif(n,-1,1))
simB$m=funz(simB$x,simB$z)
simB$y=simB$m+rnorm(n,0,0.4*sd(simB$m))
plot(simB$x,simB$z)

nn=100
xx=zz=seq(-1,1,length=nn)
yy=outer(xx,zz,FUN=function(x,z) 4*x*exp(-4*x^2 - 4*z^2))
persp(xx,zz,yy)
contour(xx,zz,yy)

fit=gam(y~s(x,z),data=simB)
vis.gam(fit,plot.type="contour")

```

Indice

More covariates

Basis for multivariate spline

Truncated power basis

The truncated power representation (of order 1) for a univariate spline (note that there are $K + 2$ parameters)

$$f(x) = \beta_1 + \beta_2 x + \sum_{k=1}^K b_k (x - \kappa_k)_+$$

has the natural extension (with $4 + K^{(x)} + K^{(z)} + K^{(x)}K^{(z)}$ parameters)

$$f(x, z) = \beta_1^{(x)} + \beta_2^{(x)}x + \sum_{k=1}^{K^{(x)}} b_k^{(x)}(x - \kappa_k^{(x)})_+ +$$

$$\beta_2^{(z)} z + \sum_{k=1}^{K^{(z)}} b_k^{(z)} (z - \kappa_k^{(z)})_+ +$$

$$\beta_2^{(xz)} xz + \sum_{k=1}^{K(x)} \sum_{k=1}^{K(z)} b_k^{(xz)} (x - \kappa_k^{(x)})_+ (z - \kappa_k^{(z)})_+$$

Thin plate splines

Consider observations (y_i, x_i) , $x_i \in \mathbb{R}^d$ and the model

$$y_i = f(x_i) + \varepsilon_i$$

thin plate splines are defined as the function f that minimizes

$$\|y - f\| + \lambda J_{md}(f)$$

where

$$J_{md}(f) = \int \dots \int_{\mathbb{R}^d} \sum_{v_1 + \dots + v_d = m} \frac{m!}{v_1! \dots v_d!} \left(\frac{\partial^m f}{\partial x_1^{v_1} \dots \partial x_d^{v_d}} \right)^2 dx_1 \dots dx_d$$

In the $d = 2$ case

$$J_{22} = \int \int_{\mathbb{R}^2} \left(\frac{\partial^2 f}{\partial x_1^2} \right)^2 + \left(\frac{\partial^2 f}{\partial x_1 \partial x_2} \right)^2 + \left(\frac{\partial^2 f}{\partial x_2^2} \right)^2 dx_1 dx_2$$

Thin plate regression splines

The above basis has dimension n , low rank thin plate regression splines can be obtained in two ways

- ▶ knot based approximations: we let

$$f(\mathbf{x}) = \sum_{i=1}^K \delta_i \eta_{md}(\|\mathbf{x} - \boldsymbol{\kappa}_i\|) + \sum_{j=1}^M \alpha_j \phi_j(\mathbf{x})$$

and minimize

$$\|y - X\beta\| + \lambda\beta^T S\beta \quad \text{s.t. } C\beta = 0$$

where

$$\beta = \begin{bmatrix} \delta \\ \alpha \end{bmatrix} \quad X_{ij} = \begin{cases} \eta_{md}(\|x - \kappa_j\|) & j \leq K \\ \phi_{j-k}(x_i) & j > K \end{cases}$$

$$C_{ij} = \begin{cases} \phi_i(\kappa_j) & j \leq K \\ 0 & j > K \end{cases} \quad S_{ij} = \begin{cases} \eta_{md}(\|x - \kappa_j\|) & j \leq K \\ 0 & j > K \end{cases}$$

Generalized models (non gaussian data)

Generalized additive models

The generalization to non gaussian data works pretty similarly as the extension of `lm` to `glm`. Given a representation for the spline

$$f(x) = \beta_1 + \beta_2 x + \sum_{j=1}^K \beta_{1+j} B_j(x)$$

the penalized least squares criterion

$$\sum_{i=1}^n (y_i - f(x_i))^2 + \lambda S(f(x))$$

is substituted by the penalized likelihood

$$\ell(\beta, \mathbf{b}, \phi) - \lambda S(f(\mathbf{x}))$$

where

$$\ell(\beta, \mathbf{b}, \phi) = \sum_{i=1}^n \log(p(y_i; \theta_i)) = \sum_{i=1}^n (y_i \theta_i - r_i(\theta_i)) / \phi + c(\phi; y_i)$$

GCV score for GAM

The GAM fitting objective can be written in terms of the deviance

$$D(\beta) = 2(\ell(\beta_{\max}) - \ell(\beta))$$

as

$$D(\beta) + \sum_{j=1}^d \lambda_j \beta^T S_j \beta$$

whose quadratic approximation is, for a fixed λ ,

$$\left\| \sqrt{W}(\mathbf{z} - \mathbf{X}\boldsymbol{\beta}) \right\|^2 + \sum_{j=1}^d \lambda_j \boldsymbol{\beta}^T \mathbf{S}_j \boldsymbol{\beta}$$

GCV score for GAM

The GAM fitting objective can be written as

$$D(\beta) + \sum_{j=1}^d \lambda_j \beta^T S_j \beta$$

whose quadratic approximation is, for a fixed λ ,

$$\left\| \sqrt{W}(z - X\beta) \right\|^2 + \sum_{j=1}^d \lambda_j \beta^T S_j \beta$$

from which one could compute the GCV score (valid locally) and then the globally applicable GCV score

$$\frac{n \left\| \sqrt{W}(z - X\hat{\beta}) \right\|^2}{n - \text{tr}(L)} \rightarrow \frac{nD(\hat{\beta})}{n - \text{tr}(L)}$$

UBRE

Recall that the idea of CV arises from estimation of the predictive risk, $E((m(x) - \hat{m}(x))^2)$, that is

$$E(\|\mu - Ly\|^2) = \frac{1}{n}E(\|y - Ly\|^2) - \sigma^2 + 2\text{tr}(L)\frac{\sigma^2}{n}$$

when the scale parameter σ^2 is known this can be done by minimizing the UBRE (Unbiased Risk Estimator)

$$\frac{1}{n}\|y - Ly\|^2 - \sigma^2 + 2\text{tr}(L)\frac{\sigma^2}{n}$$

which is equal to Mallows' C_p .



The GCV arises as an alternative in the Gaussian case since typically σ^2 must be estimated and if this is the case then the above criterion is not suitable.



UBRE is appropriate for GAM where the scale parameter is known.

UBRE computation

Based again on

$$D(\beta) + \sum_{j=1}^d \lambda_j \beta^T S_j \beta$$

and the quadratic approximation

$$\left\| \sqrt{W}(z - X\beta) \right\|^2 + \sum_{j=1}^d \lambda_j \beta^T S_j \beta$$

we obtain the UBRE criterion

$$\frac{1}{n} \left\| \sqrt{W}(z - X\beta) \right\|^2 - \sigma^2 + \frac{2\sigma^2}{n} \text{tr}(L)$$

$$\frac{1}{n} D(\hat{\beta}) - \sigma^2 + \frac{2\sigma^2}{n} \text{tr}(L)$$