



## CURVE FITTING BY LINEAR SMOOTHING

multilevel regression problem :  $y_i = f(x_i) + \epsilon_i$

A) Assume  $f(x)$  is linear and zero mean

$$y_i = \beta x_i + \epsilon_i$$

and we want to estimate  $y^*$  at some new point  $x^*$

start with one predictor and linear  $f$ .  
 $\hat{y}^*$  may be expressed as a linear smoother

$$\begin{aligned} & \begin{matrix} x_i \\ \vdots \\ x_n \end{matrix} \xrightarrow{\quad \quad \quad} \frac{\sum_{i=1}^n x_i \cdot y_i}{\left(\sum_{i=1}^n x_i x_i'\right)} = \boxed{\sum_{i=1}^n w_i \cdot y_i} \\ E(y) = x \hat{\beta} &= x' (X' X)^{-1} X' y \\ E(y_i) = x_i \hat{\beta} &= x_i' (X' X)^{-1} X' y \\ & \underbrace{( \sum_{i=1}^n x_i x_i')^{-1}}_{\text{linear combination}} \underbrace{(\sum_{i=1}^n x_i y_i)}_{\text{of past } y_i} \end{aligned}$$

⇒ the smoother is a linear combination of past  $y_i$ :  
in particular it is an average of the effect coming from different observations.

comparing with alternative approach:

(K-nearest neighbor smoothing)

$$w_K(x_i, x^*) = \begin{cases} 1/K & x_i \text{ one of the closest sample points to } x^* \\ 0 & \text{otherwise} \end{cases}$$

in the first case all the other observations are playing a role, while in this case only the closest  
 $\Rightarrow$  approximation will be smoother if compared  
 to the pure case.

will have seen that more observations the  
 smoother points.

B) Kernel Function:  $K(x)$  smooth function satisfying:

- $\int_{\mathbb{R}} K(x) dx = 1$
- $\int_{\mathbb{R}} x K(x) dx = 0$
- $\int_{\mathbb{R}} x^2 K(x) dx > 0$

ex:  $\rightarrow$  uniform kernel:

$$K(x) = \frac{1}{2} \mathbb{1}(x), \quad , \quad \mathbb{1}(x) = \begin{cases} 1 & |x| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

area under the  $\mathbb{1}$  condenses some distance to zero

2) Gaussian Kernel:

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

$\Rightarrow$  used as weighting functions  
which other obs should  
be more relevant?

for taking local averages  
derive the represent.  
in that area

$$w(x_i, x^*) = \frac{1}{h} K\left(\frac{x_i - x^*}{h}\right) \Rightarrow \text{weighting function}$$

↙  
BANDWIDTH

$\Rightarrow$  we want to derive  
the weight per iteration

$\hookrightarrow$  you can use this kernel to get the  
unnormalised weights (deciding which of  
the other variables are more relevant)

( $\rightarrow$  you are using them "in a linear smoother")

$\Rightarrow$  KERNEL REGRESSION

We want to chose h well that gives the smallest possible  
prediction error.

see notebook for implementation and plots

## CROSS VALIDATION

A) write code that would split the data in two parts and chose h according to the mse criterion.

see jupyter notebook and python code for implementation and results.

B) Now repeat a similar experiment but with data simulated from the following functions:

B.1): "wiggly" function with small noise observations

B.2): "wiggly" function with high noise observations

B.3): "smooth" function with small noise observations

B.4): "smooth" function with high noise observations

see notebook and python code

D) leave one out lemma

**Note:** The process we followed to find a good value for  $\hat{h}$  has some drawbacks:

- loss of information by splitting the dataset

⇒ alternative: leave one out cross validation

$$LOOCV = \sum_{i=1}^n (y_i - \hat{y}_i^{(-i)})^2$$

predictive value of  $y_i$   
obtained by omitting the  $i$ -th  
pair  $(x_i, y_i)$  and fitting the  
model to the ~~left~~ dataset  
of size  $N-1$

→ need do it  $H_i \rightarrow N$  times

feels like we would have to compute  $\hat{h}$   
 $N$  models and then average the results



What:  $\hat{y} = H\hat{\beta} \stackrel{def}{=} X\hat{\beta} = \underbrace{X(X'X)^{-1}X'y}_H$

and considering that matrix one can write: (FOR ANY LINEAR SMOOTHER)

$$LOOCV = \sum_{i=1}^n \left( \frac{(y_i - \hat{y}_i)^2}{(1 - H_{ii})} \right)$$

## LOCAL POLYNOMIAL REGRESSION

We can see kernel regression as a "locally constant" estimator obtained from locally weighted least square.

$$\hat{f}(x) = a = \underset{R}{\operatorname{argmin}} \sum_{i=1}^n w_i (y_i - a)^2$$

$a(x) \rightarrow$  minimizing MSE  
in  $x$

↳  $w_i$  is normalized

$$\hookrightarrow \sum w_i = 1 \rightarrow \hat{f}(x) = \bar{y}$$

$$\hookrightarrow \text{if } w_i = w(x, x_i) = \underbrace{\frac{1}{h}}_{\text{unnormalized}} \kappa \left( \frac{x_i - x}{h} \right)$$

→ we get the kernel situation.

weight depends on proximity

A) Generalization of locally constant regression  $\Rightarrow$  local polynomial regression.

$$\hat{f}(x_i) = g_x(x_i | \hat{\alpha}) \rightarrow \text{each } x \text{ is } k\text{-pointable for } x_i \text{ in its neighborhood}$$

$x_i$  = target point

$u$  = points in its neighborhood

$$g_x(u | \alpha) = \alpha_0 + \sum_{k=1}^K \alpha_k (u - x)^k$$

→ quelli che vanno stimati

→ to be estimated using weighted least squares

quelli che minimizzano LSE

$$\hat{\alpha} = \underset{R \in \mathbb{R}^K}{\operatorname{argmin}} \sum w_i (y_i - g_x(x_i, \alpha))^2$$

Define  $R_x =$

$$R_x = \begin{bmatrix} 1 & (x_1 - x)^1 & (x_1 - x)^2 & \dots & (x_1 - x)^K \\ 1 & \vdots & \vdots & & \vdots \\ \vdots & (x_n - x)^1 & (x_n - x)^2 & \dots & (x_n - x)^K \end{bmatrix}$$

let's first find  $\hat{\alpha}$  that minimizes the sum of the errors:

$$\hat{\alpha} = \underset{\mathbb{R}^{D+1}}{\operatorname{arg\,min}} \sum_{i=1}^n w_i (y_i - g(x_i | \alpha))^2$$

new arrival weights

$$= \underset{\mathbb{R}^{D+1}}{\operatorname{arg\,min}} (y - Rx\alpha)' \operatorname{diag}(w_1, \dots) (y - Rx\alpha)$$

$$\Leftrightarrow y' \tilde{W} y - 2\alpha' R x' \tilde{W} y + \alpha' R x' \tilde{W} R x \alpha$$

$$\frac{\partial}{\partial \alpha} = -2 R x' \tilde{W} y + 2 \alpha' R x' \tilde{W} R x = 0$$

$$\Rightarrow \hat{\alpha} = (R x' \tilde{W} R x)^{-1} R x' \tilde{W} y$$

B) Note that when  $D=1$ :

$g(x | \alpha_0, \alpha_1) = \hat{\alpha}_0 + \hat{\alpha}_1 (x - \bar{x}) \Rightarrow$  we are fitting two terms:  $\alpha_0$  and  $\alpha_1$

so that

$$\begin{bmatrix} \hat{\alpha}_0 \\ \hat{\alpha}_1 \end{bmatrix} = (R x' \tilde{W} R x)^{-1} R x' \tilde{W} y$$

$\downarrow$

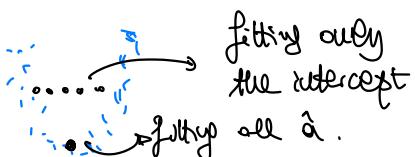
$$\begin{bmatrix} 1 & \dots & 1 \\ (x_1 - \bar{x}) & \dots & (x_n - \bar{x}) \end{bmatrix} \quad \begin{bmatrix} w_1 & \dots & 0 \\ 0 & \dots & w_n \end{bmatrix}$$

More specifically, we are asked to  
look for  $f(x)$  at the target point.

we have:  $\hat{f}(x) = g_x(u=x, \hat{a}) = \hat{a}_0$

so that:  
you end up w/ only the  
intercept, but actually  
the point of the local regression  
is to fit all the vector  $\hat{a}$ :

$$\hat{f}(x) = [1 \ 0] \begin{bmatrix} \hat{a}_0 \\ \hat{a}_1 \end{bmatrix}$$

  
fitting only  
the intercept  
fitting all  $\hat{a}$ .

$$= [1 \ 0] \left( \begin{bmatrix} 1 & \dots & 1 \\ (x_1-x) & \dots & (x_n-x) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \begin{bmatrix} 1 & (x_1-x) \\ \vdots & \vdots \\ 1 & (x_n-x) \end{bmatrix} \right)^{-1}$$

$$\begin{pmatrix} 1 & \dots & 1 \\ (x_1-x) & \dots & (x_n-x) \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$= [1 \ 0] \times \left( \begin{bmatrix} w_1 & \dots & w_n \\ w_1(x_1-x) & \dots & w_n(x_n-x) \end{bmatrix} \begin{bmatrix} 1 & (x_1-x) \\ \vdots & \vdots \\ 1 & (x_n-x) \end{bmatrix} \right)^{-1}$$

$$\begin{pmatrix} 1 & \dots & 1 \\ (x_1-x) & \dots & (x_n-x) \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$= [1 \ 0] \begin{pmatrix} \sum w_i & \sum w_i (x_i - \bar{x}) \\ \sum w_i (x_i - \bar{x}) & \sum w_i (x_i - \bar{x})^2 \end{pmatrix}^{-1} \begin{pmatrix} w_1 \dots w_n \\ w_1 (x_1 - \bar{x}) \dots w_n (x_n - \bar{x}) \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

formula of inverse for  $2 \times 2$  matrix

$$= [1 \ 0] \times \frac{1}{\sum w_i \sum_{2 \times 2}^{w_i} (x_i - \bar{x})^2 - \sum w_i (x_i - \bar{x})^2}$$

$$\times \begin{bmatrix} \sum w_i (x_i - \bar{x})^2 & \sum w_i (x_i - \bar{x}) \\ \sum w_i (x_i - \bar{x}) & \sum w_i \end{bmatrix} \times \begin{bmatrix} \sum y_i w_i \\ \sum y_i w_i (x_i - \bar{x}) \end{bmatrix}^{2 \times 1}$$

$$= \frac{\sum w_i (x_i - \bar{x})^2 \sum y_i w_i - \sum w_i (x_i - \bar{x}) (\sum y_i w_i) (x_i - \bar{x})}{\sum w_i \sum w_i (x_i - \bar{x})^2 - (\sum w_i (x_i - \bar{x}))^2}$$

$$\frac{\sum w_i y_i (\sum w_i (x_i - \bar{x})^2 - (x_i - \bar{x}) \sum w_i (x_i - \bar{x}))}{\sum w_i ((\sum w_i (x_i - \bar{x})^2 - (x_i - \bar{x}) \sum w_i (x_i - \bar{x}))^2)}$$

$$= \frac{\sum w_i ((\sum w_i (x_i - \bar{x})^2 - (x_i - \bar{x}) \sum w_i (x_i - \bar{x}))^2)}{\sum w_i ((\sum w_i (x_i - \bar{x})^2 - (x_i - \bar{x}) \sum w_i (x_i - \bar{x}))^2)}$$

$$= \frac{\sum_{i=1}^n K \left( \frac{x - x_i}{h} \right) (f_2(x) - (x_i - \bar{x}) f_1(x)) \cdot y_i}{\sum_{i=1}^n K \left( \frac{x - x_i}{h} \right) (f_2(x) - (x_i - \bar{x}) f_1(x))}$$

c) Assuming that the residuals (left out from the unknown model) have constant variance  $\sigma^2$

→ derive mean and variance of the sampling distribution for the local polynomial estimate  $\hat{f}(x)$  at some arbitrary point  $x$

Again we have:  $\hat{f}(x_i) = \hat{g}_x(x_i; \hat{\alpha}) = R_{xi} \hat{\alpha}$

$$\text{so that: } \bullet E(\hat{f}(x_i)) = E(g_x(x_i; \hat{\alpha})) = E(R_{xi} \hat{\alpha})$$

$$= R_{xi} E(\hat{\alpha})$$

$$= R_{xi} E[(R_x^T \tilde{W} R_x)^{-1} R_x^T \tilde{W} y]$$

$$= R_{xi} (R_x^T \tilde{W} R_x)^{-1} R_x^T \tilde{W} E(y)$$

$$\Rightarrow E(\hat{f}(x)) = \frac{\sum w_i(x) y_i}{\sum w_i(x)} = \tilde{w}^T y = \tilde{w}^T f(x)$$

$$\bullet V(\hat{f}(x_i)) = VAR(g_x(x_i; \hat{\alpha})) = VAR(R_{xi} \hat{\alpha})$$

$$= R_{xi} VAR[(R_x^T \tilde{W} R_x)^{-1} R_x^T \tilde{W} y] R_{xi}^T$$

$$= R_{xi} (R_x^T \tilde{W} R_x)^{-1} R_x^T \tilde{W} V(y) \tilde{W}^T R_x (R_x^T \tilde{W} R_x)^{-1} R_{xi}^T$$

$$= \sigma^2 R_{xi} (R_x^T \tilde{W} R_x)^{-1} R_x^T \tilde{W} \tilde{W}^T R_x (R_x^T \tilde{W} R_x)^{-1} R_{xi}^T$$

$$\Rightarrow V(\hat{f}(x)) = \sigma^2 \tilde{W}^T I \tilde{W} = \sigma^2 \| \tilde{W} \|_2^2 \leq \sigma^2$$

$$\text{d) } r = y - \hat{y} = y - H\hat{y} \quad \begin{array}{l} \text{projection matrix} \\ (\text{"smoothing matrix"}) \end{array}$$

$$\hat{\sigma}^2 = \frac{\|r\|_2^2}{n - 2\text{tr}(H) + \text{tr}(H^T H)}$$

$$\Rightarrow E(\hat{\sigma}^2) = E\left(\frac{\|r\|_2^2}{n - 2\text{tr}(H) + \text{tr}(H^T H)}\right)$$

$$= E\left(\frac{(y - H\hat{y})^T(y - H\hat{y})}{n - 2\text{tr}(H) + \text{tr}(H^T H)}\right)$$

$$= E\left(\frac{y'y - 2y^T H\hat{y} + y^T H^T H\hat{y}}{n - 2\text{tr}(H) + \text{tr}(H^T H)}\right)$$

$$= \frac{E(y'y) - 2E(y^T H\hat{y}) + E(y^T H^T H\hat{y})}{n - 2\text{tr}(H) + \text{tr}(H^T H)}$$

**Hint:**  $E(X'QX) = \text{tr}(Q\Sigma) + \mu' \mu \quad , \quad E(y) = \mu$

$$= \frac{\text{tr}(\Sigma) + \mu' \mu - 2\text{tr}(H\Sigma) - 2\mu^T H\mu + \text{tr}(H^T H\Sigma) + 2\mu^T H^T H\mu}{n - 2\text{tr}(H) + \text{tr}(H^T H)}$$

ausgefüllt in  $I \otimes \Sigma$

$$= \frac{\text{tr}(\delta^2 I) + \mu' \mu - 2\text{tr}(H\delta^2 I) - 2\mu^T H\mu + \text{tr}(H^T H\delta^2 I) + 3\mu^T H^T H\mu}{n - 2\text{tr}(H) + \text{tr}(H^T H)}$$

plugging in  $\hat{f}(x) = \mu = f(x)$

$$= \frac{\text{tr}(\delta^2 I) + f(x)' f(x) - 2\text{tr}(H\delta^2 I) - 2f(x)' H f(x) + \text{tr}(H' H \delta^2 I) + 2f(x)' H' H f(x)}{n - 2\text{tr}(H) + \text{tr}(H' H)}$$

$$= \delta^2 [n - 2\text{tr}(H) + \text{tr}(H' H)] + f(x)' [I - 2H + H' H] f(x)$$

$$= \delta^2 + \frac{f(x)' [(I-H)' (I-H)] f(x)}{n - 2\text{tr}(H) + \text{tr}(H' H)}$$



bias of the estimate

E) Write code to fit local linear estimator  $\Rightarrow D=1$

Apply this strategy to the dataset "utilities"

$y$  = average daily gas bill in a given month

$x$  = average temperature

with:

- gaussian kernel
- LOOCV method to chose  $h$

see python code and jupyter notebook

F) see jupyter notebook

The assumption of homoscedasticity does not look reasonable.  
The variance seem to decrease as the value on the  $x$  axis increase.

In class, we discussed that an alternative approach could be

$$\text{minimizing } y = f(x) + e$$

$$e = g^{(x)}$$

G) see jupyter notebook

## GAUSSIAN PROCESSES

gaussian process: collection of random variables

$$\{f(x) : x \in \mathcal{X}\}$$

such that, for any finite collection of indices  $x_1, \dots, x_n \in \mathcal{X}$ ,

the random vector

$[f(x_1), \dots, f(x_n)]^T$  has a multivariate normal distribution

in  $\mathbb{R}^n$  dimensional space

$f \sim GP(m, c)$  for some mean function  $m: \mathcal{X} \rightarrow \mathbb{R}$   
and some covariance function  $c: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^+$

and we have:  $E(f(x)) = m(x)$

$$\text{cov}(f(x_1), f(x_2)) = c(x_1, x_2) \quad \forall x_1, x_2$$

↳ typically covariance functions are those that are such that covariance decay when points are more far away.

↳ high covariance when  $x_1$  and  $x_2$  are close to each other



## A) NATERN CLASS of COVARIANCE FUNCTIONS:

→ class of covariance functions used in applications of multivariate statistical analysis on metric space  
( $\Leftarrow$ : spatial statistics)

According to this class of covariance functions, covariance only depends on distance between points ( $\Rightarrow$  stationary covariance).

general formula:

$$C_v(d) = \sigma^2 \frac{2^{1-v}}{\pi(v)} \left( \sqrt{2} \sqrt{\frac{d}{\rho}} \right)^v \cdot K_v \left( \sqrt{2} \sqrt{\frac{d}{\rho}} \right)$$

modified  
Bessel function  
of the second kind

↓      ↓      ↓  
 gamma function      positive parameters.

#

special case: SQUARED EXPONENTIAL

$$C_{SE}(x_1, x_2) = \gamma_1^2 \exp\left\{-\frac{1}{2} \left(\frac{d(x_1, x_2)}{b}\right)^2\right\} + \gamma_2^2 \delta(x_1, x_2)$$

$$= \underbrace{\gamma_1^2}_{\text{hyperparameters}} \exp\left\{-\frac{1}{2} \left(\frac{\|x_1 - x_2\|_2}{b}\right)^2\right\} + \underbrace{\gamma_2^2}_{\text{Kronecker delta function}} \delta(x_1, x_2)$$

⇒ note: CSE usually generates functions that are "too smooth"

EXAMPLE:  $X = [0, 1]$  ⇒ write a function that simulates a mean zero Gaussian process on  $[0, 1]$  using CSE

now: I write function that takes points in  $X$  (which now is simply the unit interval) and hyper parameters and returns the "value of the random process at each point".  
→ visualize covariance and core on R tippable f?

and "plots" the gaussian process.

Recall that has a multivariate distribution now (the vector)  
→ the marginals will not be useful and maybe that is what I can plot?

⇒ mean and variance are enough to define  $f(x)$   
But ≠ plot: I need to simulate

Correct ⇒ try out different hyperparameters (start with small)  
3) change the covariance function

Note that:

Model 1:  $x \sim GP(0, C_{\Theta})$

$\Leftarrow$  squared exponential

$$\Theta = (\beta_1, \gamma_1^2, \gamma_2^2)$$

$$\Rightarrow C_{\Theta}(x_i, x_j) = \gamma_1^2 \cdot \exp(-\frac{1}{2\beta_1} (d(x_i, x_j)^2)) + \gamma_2^2 \delta(x_i, x_j)$$

equivalent to

Model 2:  $\tilde{\theta} = (\beta_1, \gamma_1^2, 0)$

$$\tilde{x} \sim GP(0, \tilde{C}_{\tilde{\theta}})$$

$$z \sim N(0, \gamma_2^2 I)$$

$$x = \tilde{x} + z$$

Consider a different function for constructing the covariance:

$$C_{MS2}(x_1, x_2) = \gamma_1^2 \left\{ 1 + \frac{\sqrt{5}d}{b} + \frac{5d^2}{3b^2} \right\} \exp \left( -\frac{\sqrt{5}d}{b} \right) + \gamma_2^2 \delta(x_1, x_2)$$

with  $d = |x_1 - x_2|$

b) Assume you observe the value of a Gaussian process:

$$f \sim Gp(m, C) \text{ at points } x_1, \dots, x_N$$

⇒ what is the conditional distribution of the value of the process at some new point  $x^*$ ?

$$P(f(x^*) | f(x_1), \dots, f(x_N))$$

(notation:  $c(x_i, x_j) = c_{ij}$ )

Note that we have (given the definition of  $\hat{f}_t$ )

$$\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \\ f(x^*) \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \\ f(x^*) \end{bmatrix} \sim N \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} C(x, x) & & & \\ & C(x^*, x) & & \\ & & C(x^*, x^*) & \\ & & & C(x^*, x^*) \end{pmatrix}$$

variate that  
the model implies on  
the marginal of  $y_1, \dots, y_n$

$$\Rightarrow p(y_1, \dots, y_n) \sim N(0, t(x, x))$$

and we want to look for the distribution of  
 $f_t | y_1, \dots, y_n$

Now recall a general result for some multivariate:

when we have:

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right)$$

we can state:

$$x_2 | x_1 \sim N(\mu_2 + \Sigma_{21} \Sigma_{11}^{-1} (x_1 - \mu_1), \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})$$

]

Applying this result to our case, we have:

$$f_{\theta}(y_1, \dots, y_n) \sim N(\text{mean}, \text{variance})$$

$$\text{mean} = C(x_k, x) \cdot (C(x, x))^{-1} y$$

$$\text{variance} = C(x_k, x_k) - C(x_k, x) [C(x, x)]^{-1} C(x, x_k)$$

c) We want to prove the following lemma:

if :

$$\begin{aligned} \Rightarrow y | \theta &\sim N(\theta, \Sigma) \\ \Rightarrow \theta &\sim N(\mu, V) \end{aligned} \quad \Rightarrow p(\theta | y) \sim N(\mu, V)$$

$$\text{with } R, \Sigma, \mu, V \text{ constants}$$

proof. we have:

$$p(\theta | y) = p(y | \theta) \cdot p(\theta)$$

$$\propto \exp \left\{ -\frac{1}{2} \left[ (y - R\theta)' \Sigma^{-1} (y - R\theta) + (\theta - \mu)' V^{-1} (\theta - \mu) \right] \right\}$$

$$= \exp \left\{ -\frac{1}{2} \left[ y' \Sigma^{-1} y - 2\theta' R' \Sigma^{-1} y - \theta' R' \Sigma^{-1} R\theta + \theta' V^{-1} \theta - 2\mu' V^{-1} \theta + \mu' V^{-1} \mu \right] \right\}$$

Add and subtract  $2m'R'\Sigma^{-1}R\theta$

$$= \exp\left\{-\frac{1}{2}\left[y'\Sigma^{-1}y - 2\theta'R'\Sigma^{-1}y - \theta'R'\Sigma^{-1}R\theta + \theta'R^{-1}\theta - 2m'V^{-1}\theta + m'V^{-1}m - 2m'R'\Sigma^{-1}R\theta + 2m'R'\Sigma^{-1}R\theta\right]\right\}$$

$$= \exp\left\{-\frac{1}{2}\left[y'\Sigma^{-1}y - 2\theta'R'\Sigma^{-1}y - \theta'R'\Sigma^{-1}R\theta + \theta'R^{-1}\theta - 2m'V^{-1}\theta + m'V^{-1}m - 2m'R'\Sigma^{-1}R\theta + m'R'\Sigma^{-1}R\theta - \theta'R'\Sigma^{-1}Rm\right]\right\}$$

Add and subtract  $y'\Sigma^{-1}Rm$

$$= \exp\left\{-\frac{1}{2}\left[y'\Sigma^{-1}y - 2\theta'R'\Sigma^{-1}y - \theta'R'\Sigma^{-1}R\theta + \theta'R^{-1}\theta - 2m'V^{-1}\theta + m'V^{-1}m - 2m'R'\Sigma^{-1}R\theta + m'R'\Sigma^{-1}R\theta - \theta'R'\Sigma^{-1}Rm - 2m'R'\Sigma^{-1}y + y'\Sigma^{-1}Rm + m'R'\Sigma^{-1}y\right]\right\}$$

Add and subtract  $m'(R'\Sigma^{-1}R)m$

$$= \exp\left\{-\frac{1}{2}\left[\theta'V^{-1}\theta + \theta'R'\Sigma^{-1}R\theta - 2m'V^{-1}\theta - 2m'R'\Sigma^{-1}R\theta - y'\Sigma^{-1}R\theta + y'\Sigma^{-1}Rm + m'R'\Sigma^{-1}R\theta - \theta'R'\Sigma^{-1}y + \theta'R'\Sigma^{-1}Rm + m'R'\Sigma^{-1}y + y'\Sigma^{-1}y - 2m'R'\Sigma^{-1}y + m'R'\Sigma^{-1}Rm - m'R'\Sigma^{-1}Rm\right]\right\}$$

Combine piece together:

$$\propto \exp\left\{-\frac{1}{2}\left[\theta'(V^{-1} + R'\Sigma^{-1}R)\theta - 2m'(V^{-1} + R'\Sigma^{-1}R)\theta - y'\Sigma^{-1}R\theta + y'\Sigma^{-1}Rm + m'R'\Sigma^{-1}R\theta - \theta'R'\Sigma^{-1}y + \theta'R'\Sigma^{-1}Rm + m'R'\Sigma^{-1}y + y'\Sigma^{-1}y - 2m'R'\Sigma^{-1}y + m'R'\Sigma^{-1}Rm - m'R'\Sigma^{-1}Rm + m'R'\Sigma^{-1}Rm\right]\right\}$$

$$\begin{aligned}
 &= \exp\left\{-\frac{1}{2} \left[ (\boldsymbol{\theta} - \boldsymbol{m})' (\mathbf{V}' + \mathbf{R}' \boldsymbol{\Sigma}' \mathbf{R}) (\boldsymbol{\theta} - \boldsymbol{m}) + (\mathbf{y} - \mathbf{R}\boldsymbol{m})' \boldsymbol{\Sigma}' \mathbf{R} (\boldsymbol{\theta} - \boldsymbol{m}) \right. \right. \\
 &\quad \left. \left. - (\boldsymbol{\theta} - \boldsymbol{m})' \mathbf{R}' \boldsymbol{\Sigma}' (\mathbf{y} - \mathbf{R}\boldsymbol{m}) + (\mathbf{y} - \mathbf{R}\boldsymbol{m})' \boldsymbol{\Sigma}' (\mathbf{y} - \mathbf{R}\boldsymbol{m}) \right] \right\} \\
 &= \exp\left\{-\frac{1}{2} \left[ \begin{pmatrix} \boldsymbol{\theta} - \boldsymbol{m} \\ \mathbf{y} - \mathbf{R}\boldsymbol{m} \end{pmatrix}' \begin{pmatrix} \mathbf{V}' + \mathbf{R}' \boldsymbol{\Sigma}' \mathbf{R} & -\mathbf{R}' \boldsymbol{\Sigma}' \\ -\boldsymbol{\Sigma}' \mathbf{R} & \boldsymbol{\Sigma}' \end{pmatrix} \begin{pmatrix} \boldsymbol{\theta} - \boldsymbol{m} \\ \mathbf{y} - \mathbf{R}\boldsymbol{m} \end{pmatrix} \right] \right\} \\
 &\Rightarrow p(\boldsymbol{\theta}, \mathbf{y}) \propto \exp\left\{-\frac{1}{2} \left( \dots \right) \right\}
 \end{aligned}$$

$\Rightarrow$  travel of a multivariate normal with:

$$\text{mean} = \begin{bmatrix} \boldsymbol{m} \\ \mathbf{R}\boldsymbol{m} \end{bmatrix}$$

$$\text{variance} = \begin{bmatrix} \mathbf{V}' + \mathbf{R}' \boldsymbol{\Sigma}' \mathbf{R} & -\mathbf{R}' \boldsymbol{\Sigma}' \\ -\boldsymbol{\Sigma}' \mathbf{R} & \boldsymbol{\Sigma}' \end{bmatrix}^{-1}$$

Alternatively, we could also note:

$$\begin{bmatrix} \mathbf{y} \\ \boldsymbol{\theta} \end{bmatrix} = \begin{bmatrix} \mathbf{R} \\ \mathbf{I} \end{bmatrix} \boldsymbol{\theta} + \begin{bmatrix} \mathbf{I} \\ \mathbf{0} \end{bmatrix} \varepsilon$$

which is an affine transformation of a multivariate normal  
 $\Rightarrow$  from previous exercise we know this is still a multivariate normal.

## NON PARAMETRIC REGRESSION & GAUSSIAN SMOOTHING

model:  $y_i = f(x_i) + \epsilon_i$ ,  $\epsilon_i \sim N(0, \delta^2)$

$$f \sim GP(0, C)$$

A)

we observed  $x_1, \dots, x_n \rightarrow$  want to make inference  
on  $f(x_1), \dots, f(x_n)$   
 $\Rightarrow$  what is the posterior?

$$P(f, y) \propto P(y|f) \cdot P(f)$$

since  $y$  is a finite vector of observations it's just a multivariate normal

$$\propto \exp\left\{-\frac{1}{2} (y - f(x))' (\delta^2 I)^{-1} (y - f(x)) - \frac{1}{2} (f(x) - \mu)' C^{-1} (f(x) - \mu)\right\}$$

$$\propto \exp\left[-\frac{1}{2} \left[ f(x)' \delta^2 I f(x) + f'(x)' C^{-1} f(x) - f(x)' (\delta^2 I)^{-1} y \right]\right]$$

$$\Rightarrow f|y \sim MVN(\text{mean}, \text{variance})$$

$$\text{variance: } ((\delta^2 I)^{-1} + C^{-1})^{-1} = \left(\frac{I}{\delta^2} + C^{-1}\right)^{-1}$$

$$\text{mean: } ((\delta^2 I)^{-1} + C^{-1})^{-1} \cdot (\delta^2 I)^{-1} y = C(C + \delta^2 I)^{-1}$$

B) Now we are interested in the predictive distribution:

note that since  $f(x_1), \dots, f(x_n), f(x^*)$  is a GP we have:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \\ f(x^*) \end{bmatrix} \sim N \begin{pmatrix} 0 & (c(x, x) + \sigma^2 I) & c(x, x^*) \\ 0 & c(x^*, x) & c(x^*, x^*) \end{pmatrix}$$

Since because of well known results of multivariate normals, we have

$$y \sim N(0, \sigma^2 I + c(x, x)) \rightarrow \text{see ex D}$$

Now using again the result recalled in the exercise above, we have:

$$f^* | y \sim N(\text{mean}, \text{variance})$$

with:  $\text{mean} = c(x^*, x) (c(x, x) + \sigma^2 I)^{-1} y$

$\text{var} = c(x^*, x^*) - c(x^*, x) [c(x, x) + \sigma^2 I]^{-1} c(x, x^*)$

note that it can be written as  $\sum_i w_i y_i$

$\Rightarrow$  linear smoother

c) see jupyter notebook

d) Now assume  $C = \text{Motm}(\mathbb{S}/2)$

Derive the marginal of  $y_1 \dots y_n$ .

Recall well known result:

$$\begin{cases} \text{if } p(ab) = N(a|AL, \Sigma) \\ p(b) = N(b|\mu, \Sigma) \end{cases} \Rightarrow p(a) = \int p(ab) db = N(A\mu, A\Sigma A' + \Sigma)$$

1

So that in our case we have:

$$p(y) = N(0, (C + \delta^2 I))$$

Moreover, recall that the motm  $(\mathbb{S}/2)$  is the following structure:

$$\begin{array}{c} (\delta^2 I) + \\ \sim \\ \text{reduced} \\ \text{in the} \\ \text{marginal module} \end{array} \quad \left[ \begin{array}{cc} C_{11} & C_{21} \\ C_{12} & \ddots \end{array} \right] \quad \begin{array}{l} \text{Cl} = \gamma_1^2 + \gamma_2^2 \\ \Rightarrow \text{marginal variance} \\ \text{will be } \delta^2 + \gamma_1^2 + \gamma_2^2 \end{array}$$

$\Rightarrow C_{ij} = \gamma_1^2 \cdot 1 / 1 + \frac{\sqrt{5} d_{ij}}{b}, \frac{5 d_{ij}^2}{3b^2} \left( 1 \cdot \exp \left( -\frac{\sqrt{5} d_{ij}}{b} \right) \right)$

E) we are asked to evaluate the log of the marginal likelihood  
setting  $\eta_i^2 = 0$ .

in D) we derived:  $p(y) = N(\alpha_1(C + \delta^2 I))$ .

Hence we have:

$$\begin{aligned}\log(p(y)) &= \log\left(|2\pi(C + \delta^2 I)|^{-\frac{1}{2}} \cdot \exp\left\{-\frac{1}{2}y'(C + \delta^2 I)^{-1}y\right\}\right) \\ &= -\frac{1}{2}\log(|2\pi(C + \delta^2 I)|) - \frac{1}{2}y'(C + \delta^2 I)^{-1}y \\ &= -\frac{n}{2}\log(2\pi) - \frac{1}{2}\log(|C + \delta^2 I|) - \frac{1}{2}y'(C + \delta^2 I)^{-1}y\end{aligned}$$

For implementation see jupyter notebook and python code