

Gaussian process regression in the flat limit

Pierre-Olivier Amblard

from a work with Simon Barthelmé, Nico Tremblay and Kostya
Usevitch

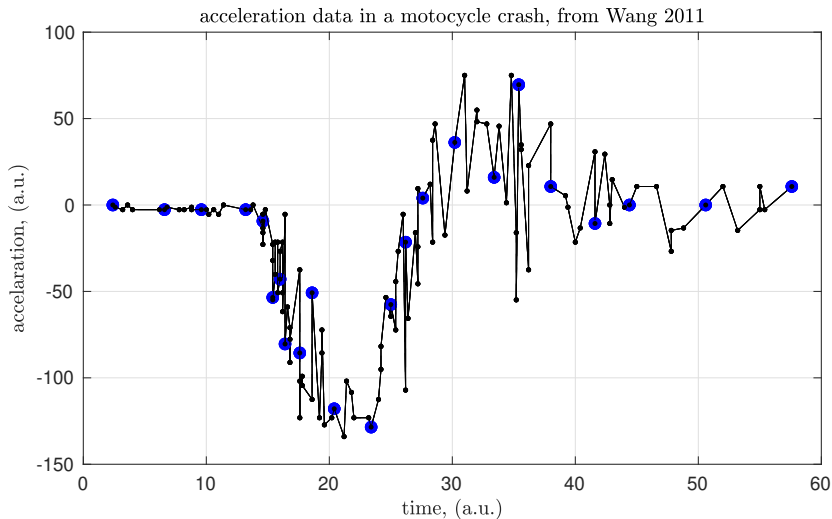
CNRS, GIPSA-lab, U. Grenoble-Alpes, France

W'p "Kernel and sampling methods for design and quantization",
Gdr MascoNum 2021



Gaussian Process Regression

↪ Find f such that $y_i = f(x_i) + \sigma w_i$ for $i = 1, \dots, N$



Gaussian Process Regression

↪ Find f such that $y_i = f(x_i) + \sigma w_i$ for $i = 1, \dots, N$

↪ In our works, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$, and w_i a sequence of i.i.d. $\mathcal{N}(0, 1)$.
BUT here, $d = 1$!

↪ In a Bayesian approach, $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is given a Gaussian process prior, with zero mean and covariance $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$.

↪ The *a posteriori* $P(y | \mathbf{y}, f, \mathbf{x}, x)$ at x is Gaussian with

$$y^*(x) = E[y | \mathbf{y}, f] = \sum_i \alpha_i k(x, x_i) = \boldsymbol{\alpha}^\top \mathbf{k}_x \text{ with } \boldsymbol{\alpha} = (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

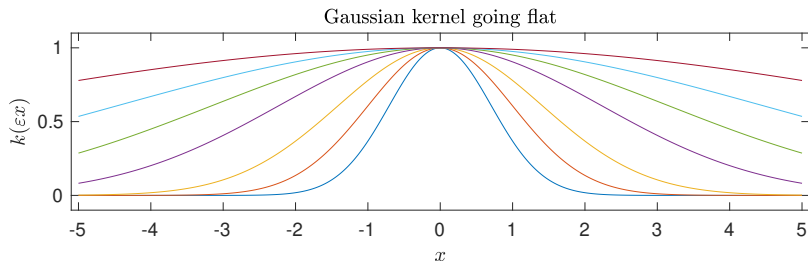
$$\text{Var}[y | \mathbf{y}, f] = k(x, x) - \mathbf{k}_x^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_x$$

$$\mathbf{K}_{i,j} = k(x_i, x_j) \quad \mathbf{k}_{x,i} = k(x, x_i)$$

↪ $\sigma^2 \rightarrow 0$: recovering interpolation

Radial Basis Interpolation in the flat limit

- Radial Basis \sim stat./isotropic Gauss. proc. : $k(\mathbf{x}, \mathbf{y}) = k(\|\mathbf{x} - \mathbf{y}\|)$.



\hookrightarrow Let ϵ be the scale parameter, $k(\mathbf{x} - \mathbf{y}) = k(\epsilon\|\mathbf{x} - \mathbf{y}\|)$.

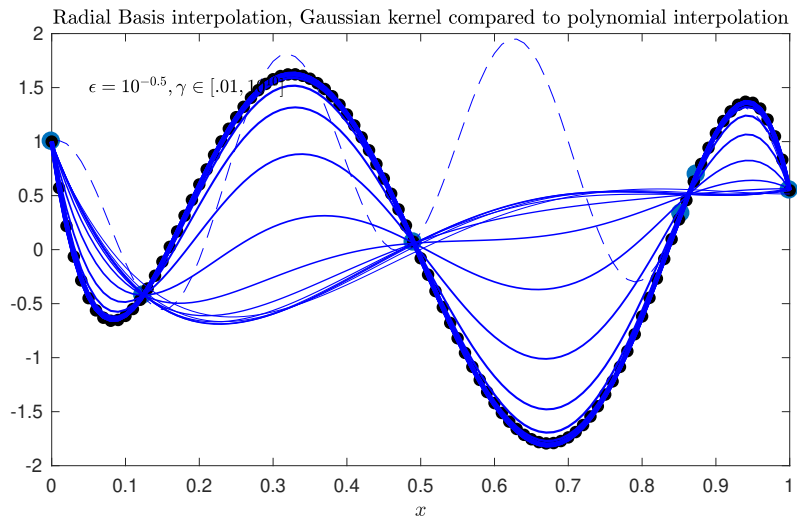
$$y^*(x) = E[y|\mathbf{y}, f] = \sum_i \alpha_i k(\epsilon(\mathbf{x} - \mathbf{x}_i)).$$

Driscoll & Fornberg 2002 proved for $k(\mathbf{x} - \mathbf{y}) = \exp(-\epsilon^2\|\mathbf{x} - \mathbf{y}\|^2)$

Gaussian RBF interpolation $\xrightarrow{\epsilon \rightarrow 0}$ polynomial interpolation

Let's have a look !

RBF interpolation in the flat limit



Radial Basis Interpolation in the flat limit

- Radial Basis \sim stat./isotropic Gauss. proc. : $k(\mathbf{x}, \mathbf{y}) = k(\|\mathbf{x} - \mathbf{y}\|)$.

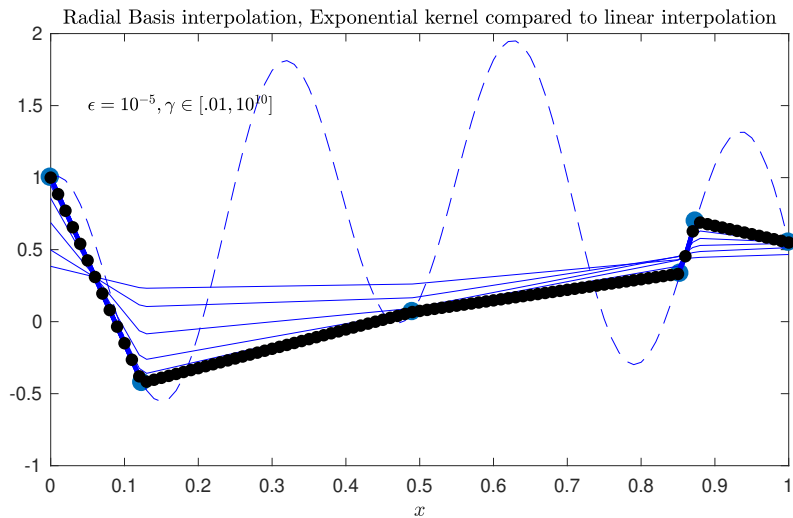
\hookrightarrow Let ε be the scale parameter, $k(\mathbf{x} - \mathbf{y}) = k(\varepsilon\|\mathbf{x} - \mathbf{y}\|)$.

$$y^*(x) = E[y|\mathbf{y}, f] = \sum_i \alpha_i k(\varepsilon(\mathbf{x} - \mathbf{x}_i)).$$

Driscoll & Fornberg 2002 proved for $k(\mathbf{x} - \mathbf{y}) = \exp(-\varepsilon^2\|\mathbf{x} - \mathbf{y}\|^2)$
Gaussian RBF interpolation $\xrightarrow{\varepsilon \rightarrow 0}$ polynomial interpolation

What's up with a harder kernel, e.g. $k(\mathbf{x} - \mathbf{y}) = \exp(-\varepsilon\|\mathbf{x} - \mathbf{y}\|)$

RBF interpolation in the flat limit



Radial Basis Interpolation in the flat limit

- Radial Basis \sim stat./isotropic Gauss. proc. : $k(\mathbf{x}, \mathbf{y}) = k(\|\mathbf{x} - \mathbf{y}\|)$.

\hookrightarrow Let ε be the scale parameter, $k(\mathbf{x} - \mathbf{y}) = k(\varepsilon\|\mathbf{x} - \mathbf{y}\|)$.

$$y^*(x) = E[y|\mathbf{y}, f] = \sum_i \alpha_i k(\varepsilon(\mathbf{x} - \mathbf{x}_i)).$$

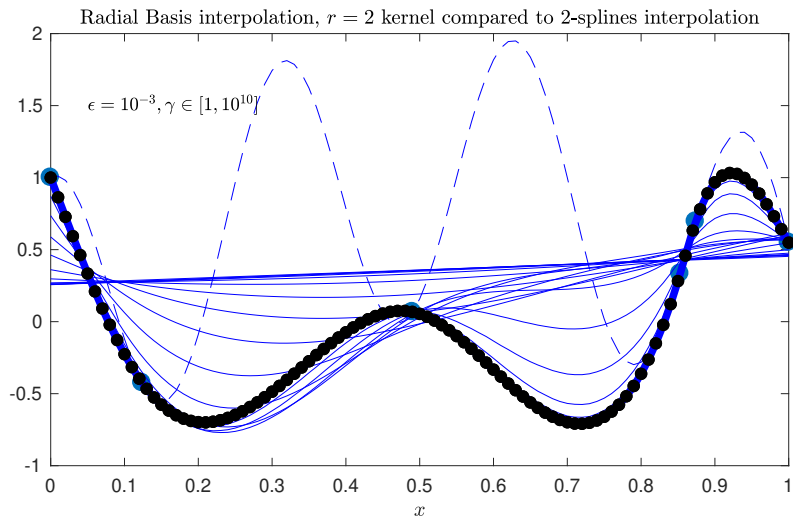
Driscoll & Fornberg 2002 proved for $k(\mathbf{x} - \mathbf{y}) = \exp(-\varepsilon^2\|\mathbf{x} - \mathbf{y}\|^2)$

Gaussian RBF interpolation $\xrightarrow{\varepsilon \rightarrow 0}$ polynomial interpolation

Nice implementation of linear interpolation!!

Let us smooth it a bit, e.g. $k(\mathbf{x} - \mathbf{y}) = (1 + \varepsilon\|\mathbf{x} - \mathbf{y}\|) \exp(-\varepsilon\|\mathbf{x} - \mathbf{y}\|)$

RBF interpolation in the flat limit



What is hidden ?

Recall : $E[y|\mathbf{y}, f] = \mathbf{k}_x^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$

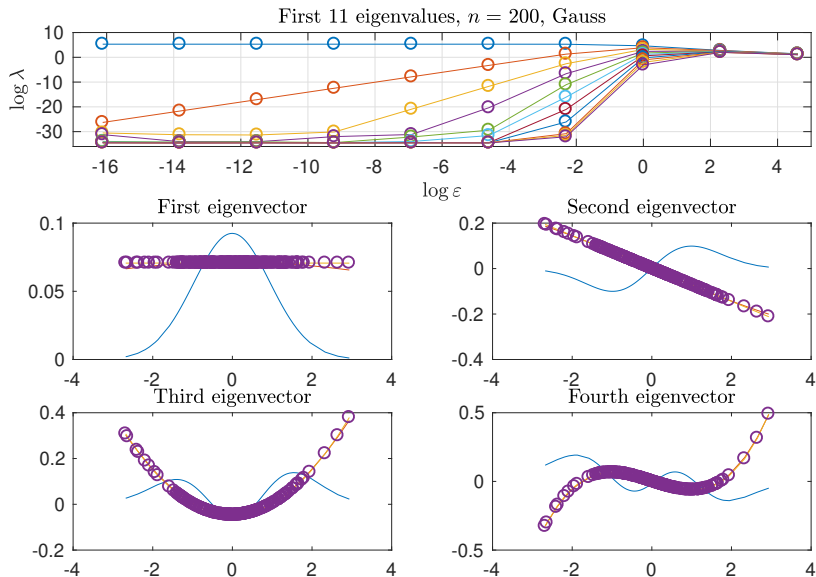
- The kernel is going flat, and thus the Gram matrix goes to a rank 1 matrix !

↪ All but one eigenvalues go to zero ! But wait, let us look at the rate !

↪ Eigen analysis of the Gram matrix $\mathbf{K} = (k(\varepsilon \|x_k - x_l\|))_{k,l}$

$$\mathbf{K} \mathbf{v}_i(\varepsilon) = \lambda_i(\varepsilon) \mathbf{v}_i(\varepsilon)$$

Eigen analysis of the Gram matrix w.r.t. ε , Gauss. k



Smooth kernel : no singularity at zero at any order

- typical example : the Gaussian kernel

$$\exp(-|x|^2) = f_0 + f_2|x|^2 + f_4|x|^4 + \dots$$

↪ for scalar inputs, as $\varepsilon \rightarrow 0$,

$$\lambda_i(\varepsilon) = \varepsilon^{2i}(\tilde{\lambda}_i + O(\varepsilon)), \forall i = 0, \dots, n-1$$

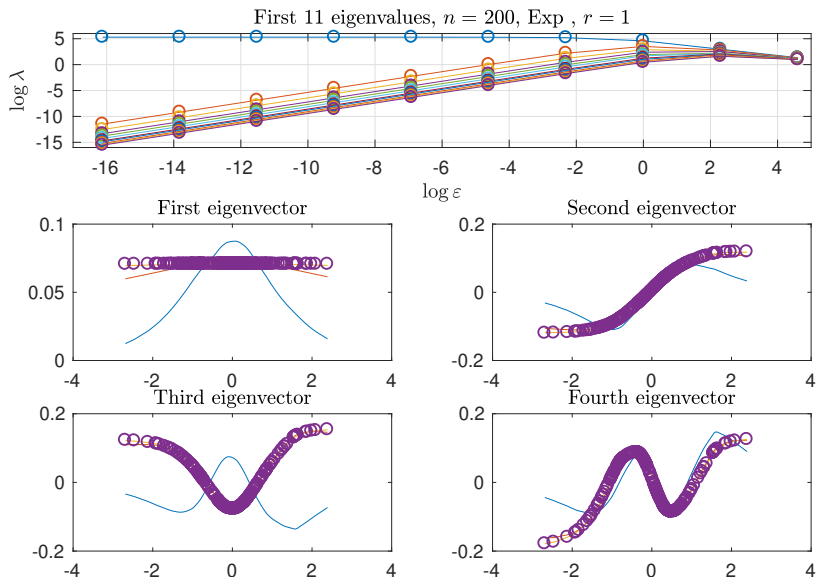
$$\mathbf{v}_i = \mathbf{e} \cdot \vec{v}_i \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{n-1} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^{n-1} \end{pmatrix}$$

i.e. orthogonal polynomials.

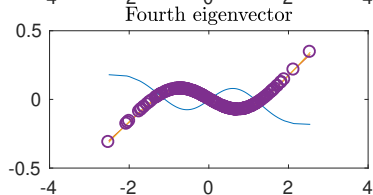
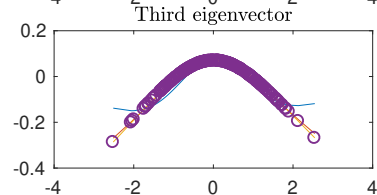
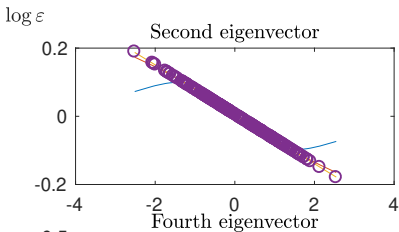
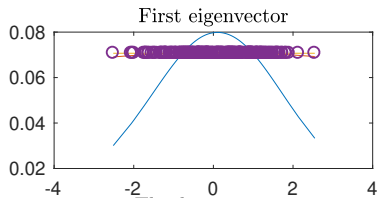
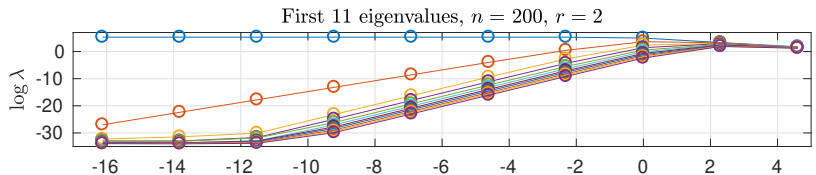
But what if :

$$k(|x|) = f_0 + f_2|x|^2 + f_4|x|^4 + \dots + f_{2r-1}|x|^{2r-1} + f_{2r}|x|^{2r} + \dots$$

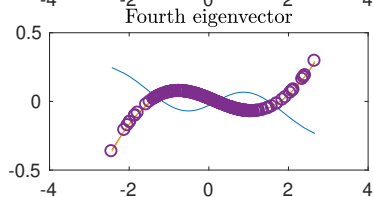
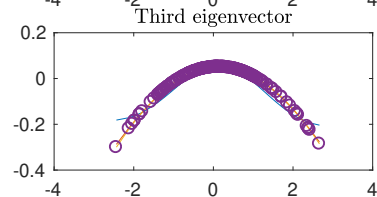
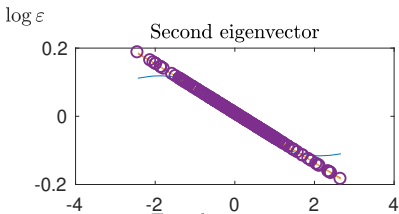
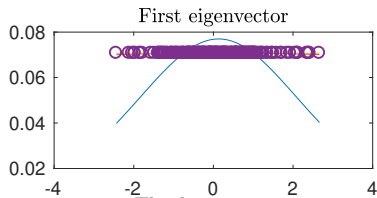
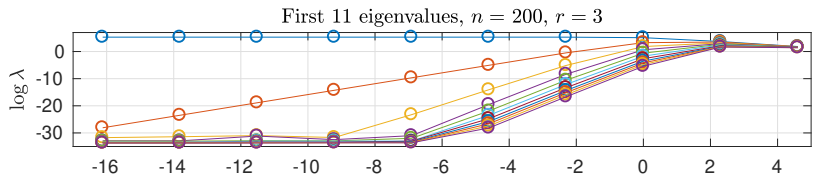
Eigen analysis of the Gram matrix w.r.t. ε , $r = 1$



Eigen analysis of the Gram matrix w.r.t. ε , $r = 2$



Eigen analysis of the Gram matrix w.r.t. ε , $r = 3$



Non smooth kernel : singularity at zero somewhere

$$k(|x|) = f_0 + f_2|x|^2 + f_4|x|^4 + \dots + f_{2r-1}|x|^{2r-1} + f_{2r}|x|^{2r} + \dots$$

↪ for scalar inputs, as $\varepsilon \rightarrow 0$,

$$\lambda_i(\varepsilon) = \varepsilon^{2i}(\tilde{\lambda}_i + \mathcal{O}(\varepsilon)), \forall i = 0, \dots, r-1$$

$$\lambda_i(\varepsilon) = \varepsilon^{2r-1}(\tilde{\lambda}_i + \mathcal{O}(\varepsilon)), \forall i = r, \dots, n$$

$$\mathbf{v}_i = \mathbf{e} \cdot \vec{\mathbf{v}} \cdot \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^{r-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{r-1} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^{r-1} \end{pmatrix} \forall i = 0, \dots, r-1$$

$$\mathbf{v}_i = \text{linked to } (\|x_k - x_l\|^{2r-1})_{k,l}, \quad \forall i = r, \dots, n$$

Back to GPR, degrees of freedom

$$y^*(x) = E[y|\mathbf{y}, f] = \sum_i \alpha_i k(x, x_i) = \boldsymbol{\alpha}^\top \mathbf{k}_x^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$$

$$\text{Var}[y|\mathbf{y}, f] = k(x, x) - \mathbf{k}_x^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_x$$

Prediction mean and variance at x given \mathbf{x} can be computed from the smoother matrix for (\mathbf{x}, x)

- For linear smoothers $y^*(\mathbf{x}) = \mathbf{M}\mathbf{y}$, d.o.f. : $\text{dof} = \text{Tr}(\mathbf{M})$ (Hastie/Tibs./Fried, TESL)

all info. in the smoothing matrix, here : $\mathbf{M}(\varepsilon) = \mathbf{K}(\mathbf{K} + \sigma^2 \mathbf{I})^{-1}$!

↪ For GPR

$$\text{dof}(\varepsilon) = \text{Tr}(\mathbf{M}) = \sum_{i=1}^n \frac{\lambda_i(\varepsilon)}{\lambda_i(\varepsilon) + \sigma^2} \approx \sum_{i=1}^n \frac{\tilde{\lambda}_i \varepsilon^{\mathbf{e}_i}}{\tilde{\lambda}_i \varepsilon^{\mathbf{e}_i} + \sigma^2} \xrightarrow{\varepsilon \rightarrow 0} \frac{\tilde{\lambda}_1}{\tilde{\lambda}_1 + \sigma^2} \leq 1!$$

A proj. (e.g. poly. reg.) onto any $q \leq n$ dimensional space : $\text{dof} = q!!$

↪ We must add another parameter to control the d.o.f.

Rescaling

We didn't control the amplitude !

↪ Let the kernel be $\gamma(\varepsilon)k(\varepsilon\|\mathbf{x} - \mathbf{y}\|)$.

The d.o.f. writes

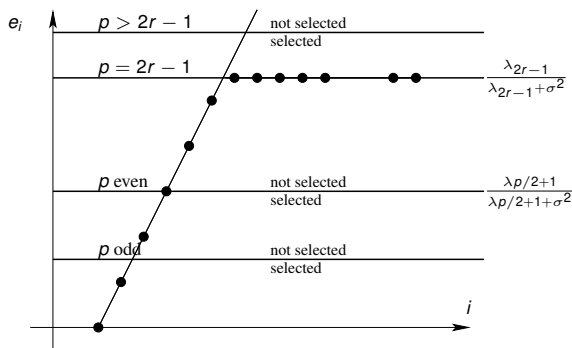
$$\text{Tr}(\mathbf{M}) = \sum_{i=1}^n \frac{\gamma(\varepsilon)\lambda_i(\varepsilon)}{\gamma(\varepsilon)\lambda_i(\varepsilon) + \sigma^2} \approx \sum_{i=1}^n \frac{\tilde{\lambda}_i}{\tilde{\lambda}_i + \gamma(\varepsilon)^{-1}\varepsilon^{-e_i}\sigma^2}$$

Natural choice $\gamma(\varepsilon) \propto \varepsilon^{-p}$: p controls the d.o.f. through ε^{p-e_i}

$$\mathbf{y}^*(\mathbf{x}) = \sum_{i=1}^n \frac{\tilde{\lambda}_i}{\tilde{\lambda}_i + \gamma(\varepsilon)^{-1}\varepsilon^{-e_i}\sigma^2} \mathbf{v}_i(\varepsilon)\mathbf{v}_i^\top(\varepsilon)\mathbf{y}$$

$$\frac{\tilde{\lambda}_i}{\tilde{\lambda}_i + \varepsilon^{p-e_i}\sigma^2} \rightarrow \begin{cases} 1 : \text{eigen vector } i \text{ is kept} \\ \text{a constant} : \text{eigen vector } i \text{ is penalized} \\ 0 : \text{eigen vector } i \text{ is left} \end{cases}$$

$\lambda/(\lambda + \sigma^2)$ as a filter



$$\frac{\tilde{\lambda}_i}{\tilde{\lambda}_i + \varepsilon^{p-2(i-1)}\sigma^2} \rightarrow \begin{cases} 1 : \text{the eigen vector is kept} \\ \text{a constant} : \text{the eigen vector is penalized} \\ 0 : \text{the eigen vector is left} \end{cases}$$

$$\frac{\tilde{\lambda}_i}{\tilde{\lambda}_i + \varepsilon^{p-(2r-1)}\sigma^2} \rightarrow \begin{cases} 1 : \text{the eigen vectors are kept} \\ \text{a constant} : \text{the eigen vectors are penalized} \\ 0 : \text{the eigen vectors are left} \end{cases}$$

Theorem

For the kernel $\varepsilon^{-p}k(\varepsilon\|\mathbf{x} - \mathbf{y}\|)$ with regularity r , the associated Gaussian process regression converges as $\varepsilon \rightarrow 0$ toward :

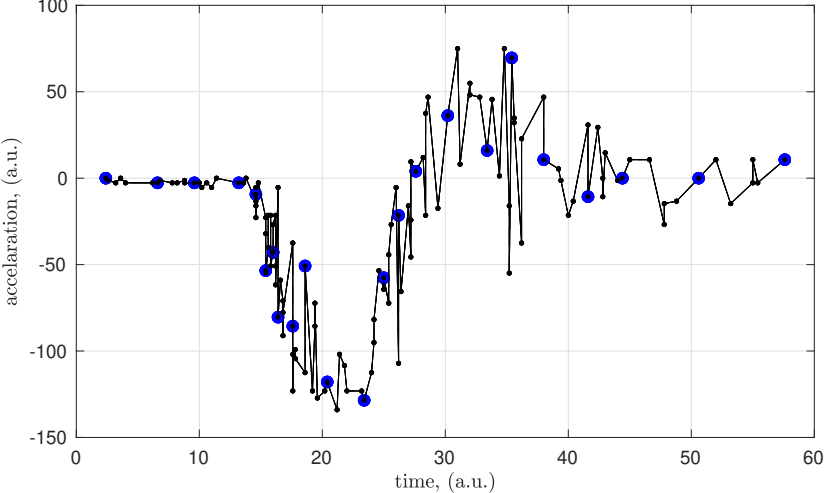
- ▶ interpolation if $p > 2r - 1$;
- ▶ smoothing splines of order r is $p = 2r - 1$;
- ▶ penalized polynomial regression if $p < 2r - 1$ and even ;
- ▶ polynomial regression of order $\lfloor p/2 + 1 \rfloor$ if $p < 2r - 1$ and odd.

Same kind of results in the multivariate case ! and stuffs like :

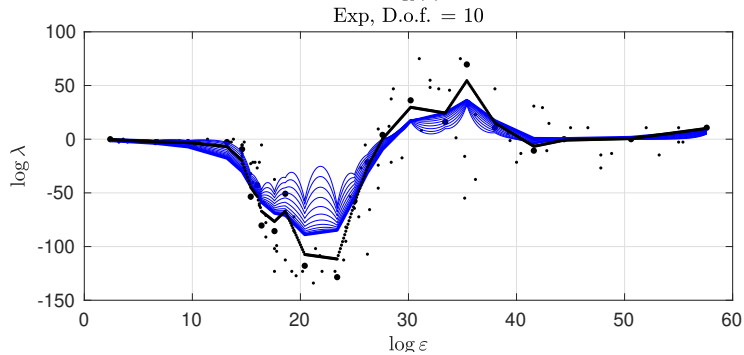
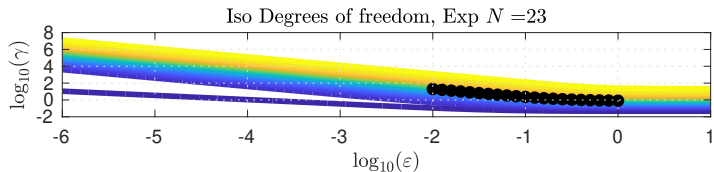
as $\varepsilon \rightarrow 0$, GPR for $\varepsilon^{-p} \exp(\varepsilon\|\mathbf{x} - \mathbf{y}\|^2)$ behaves as GPR with $(\mathbf{x}^\top \mathbf{y})^{p/2-1}$ (poly. reg., penalized or not)

Illustrations

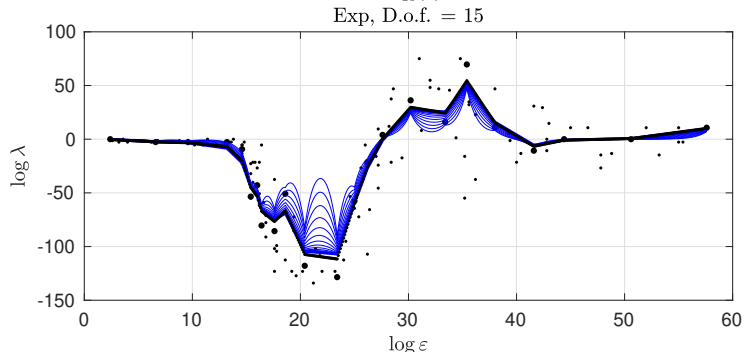
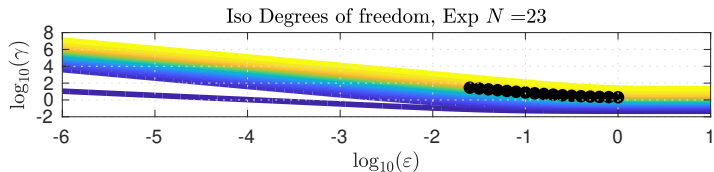
acceleration data in a motorcycle crash, from Wang 2011



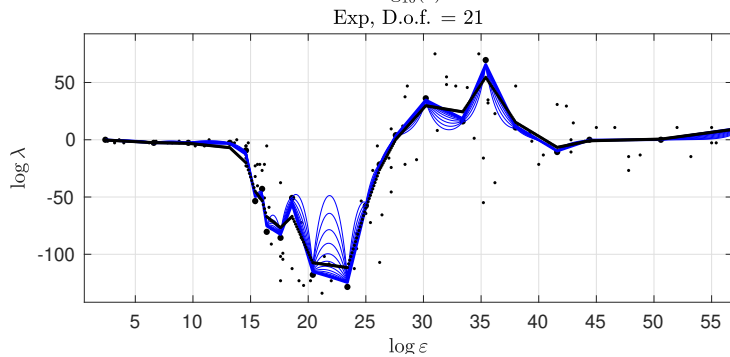
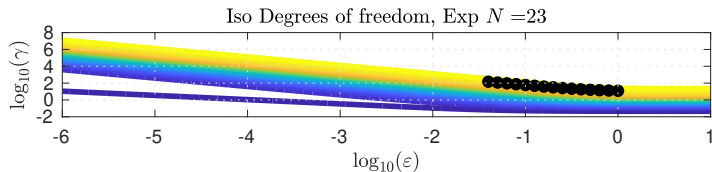
Exponential kernel, $r = 1$



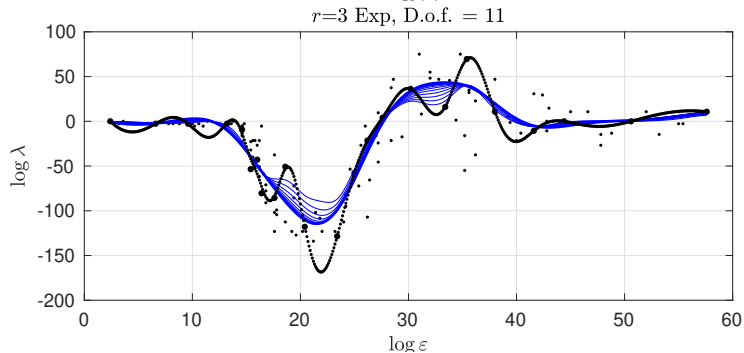
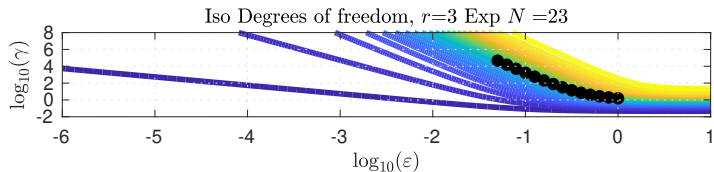
Exponential kernel, $r = 1$



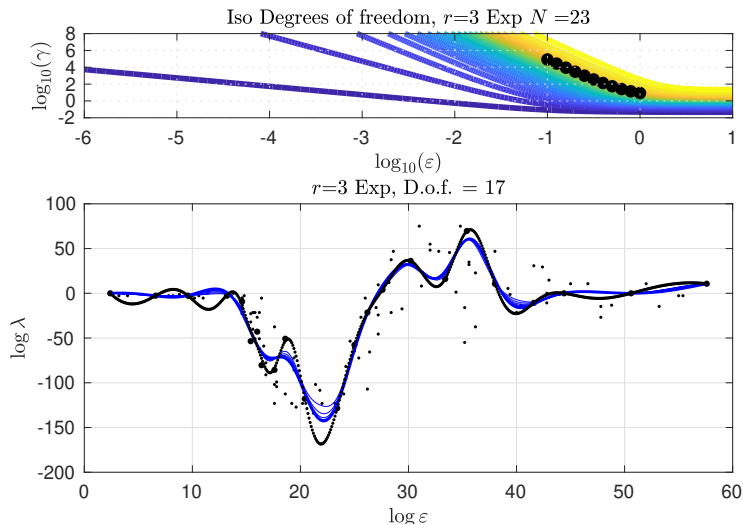
Exponential kernel, $r = 1$



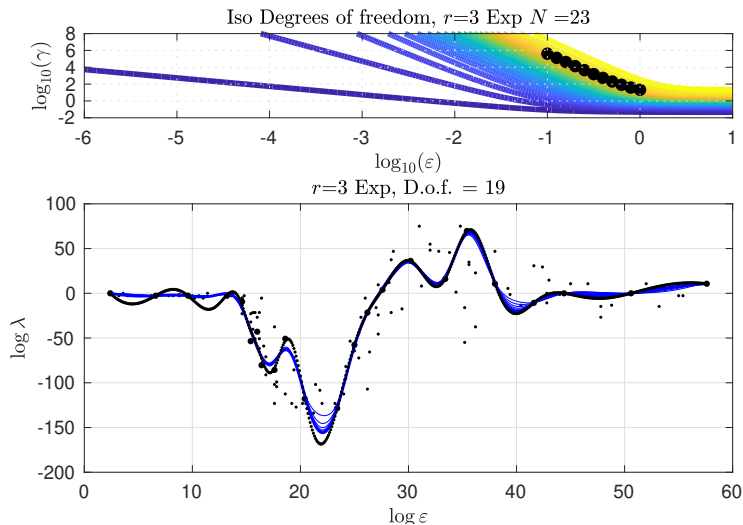
Exponential kernel, $r = 3$



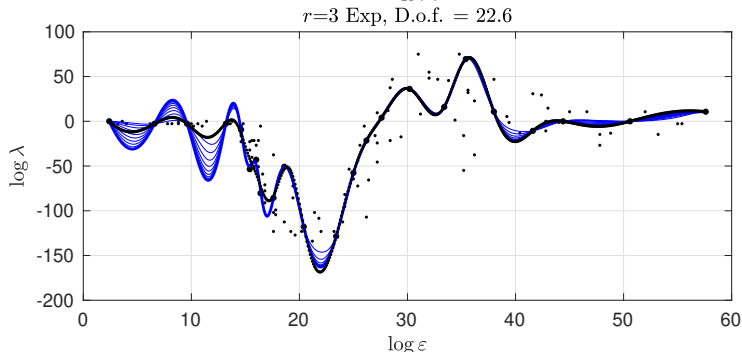
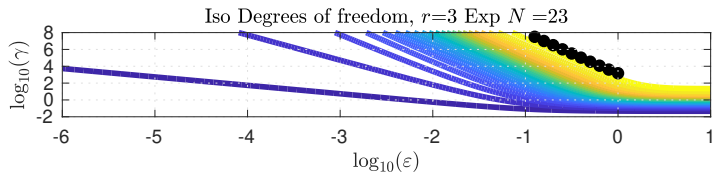
Exponential kernel, $r = 3$



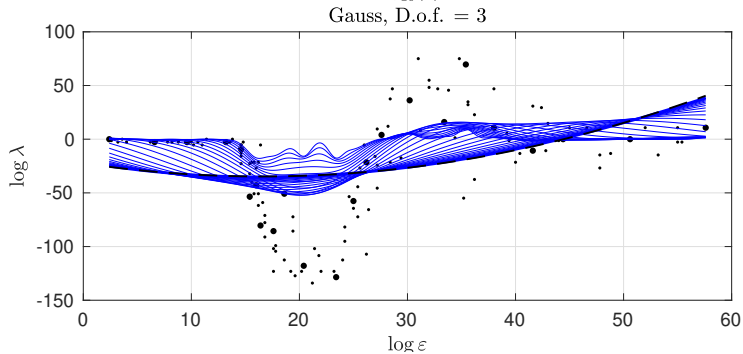
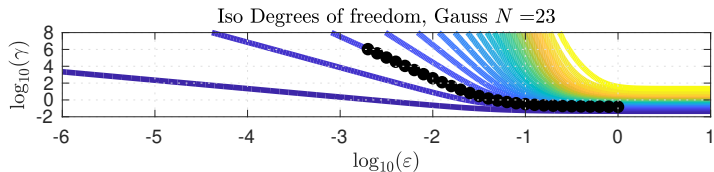
Exponential kernel, $r = 3$



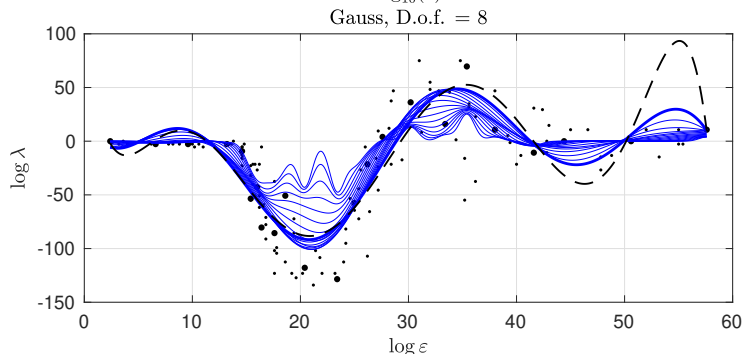
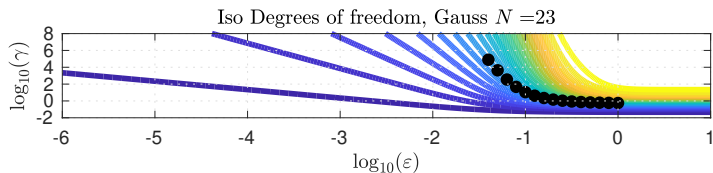
Exponential kernel, $r = 3$



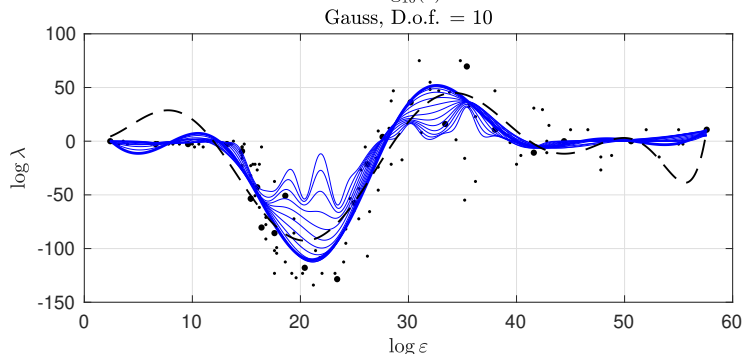
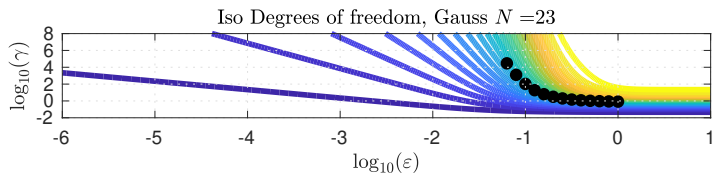
Gaussian kernel



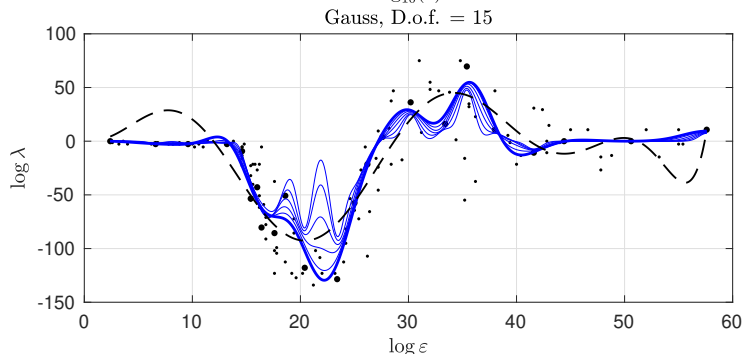
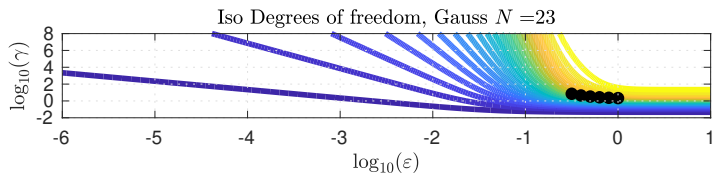
Gaussian kernel



Gaussian kernel



Gaussian kernel



To conclude

- ↪ Deep links between GP & splines/orthogonal polynomials in the flat limit
- ↪ Flat limit may be obtain quite rapidly
- ↪ Flat limits attractive (*e.g.* polyharmonic splines) since less hyperparameter to control
- ↪ We've studied the flat limit in DPP's!

More in :

Bathelmé&Usevitch, SIMAX 2020 and our ArXiv papers currently submitted (or to be)