

Beyond Surrogate Modeling : Learning the Local Volatility via Shape Constraints

Areski Cousin

IRMA, Université de Strasbourg

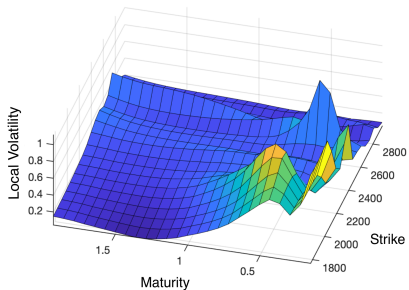
Joint work with Marc Chataigner, Stéphane Crépey, Matthew Dixon
and Djibril Gueye

Séminaire équipe Calisto - INRIA, July 2, 2021

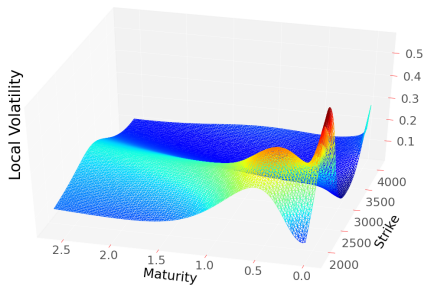


Motivation

- We explore 2 alternatives machine learning approaches for the construction of local volatility surfaces
- A no-arbitrage Gaussian process (GP) approach based on price and a neural net (NN) approach with penalization of arbitrages based on implied volatility
- Construction at a given cotation date, from observations of bid-ask prices for a set of liquidly traded European options



GP estimate



NN estimate

In this presentation, we focus on the GP approach

- Construct a GP surface estimate $P : \mathbb{R}^d \rightarrow \mathbb{R}$, interpolating noisy option prices $\{\mathbf{x}_i, y_i\}_{i=1}^n$ so that $y_i = P(\mathbf{x}_i) + \epsilon_i$ for some additive noise $\epsilon_i, \forall i$.
- For a given option, bid and ask prices considered as two (noisy) replicates of P
- Subject to shape constraints, i.e., non-decreasing in maturity and convex in strike (to avoid calendar spread and butterfly arbitrages).
- How to effectively enforce shape constraints in the GP regression ?
- Uncertainty quantification under shape constraints : how to derive no-arbitrage confidence bands ?

Connections between option price, IV and local volatility

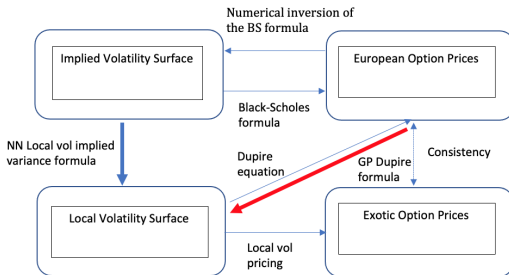


Figure: Mathematical connections between option prices, implied, and local volatility. The **bold red arrow** shows the route under the no-arbitrage GP approach based on price. The **bold blue arrow** shows the route under the NN approach based on implied volatility.

Absence of static arbitrage

The put price function $(T, K) \rightarrow P(T, K)$ is considered to be free of static arbitrage if there exists a measure \mathbb{Q} such that the discounted asset $e^{-(r-q)t} S_t$ is a \mathbb{Q} -martingale and $P(T, K) = \mathbb{E}_{\mathbb{Q}} [e^{-rT} (K - S_T)^+]$.

Absence of arbitrage [constant interest rate r and dividend yield q]

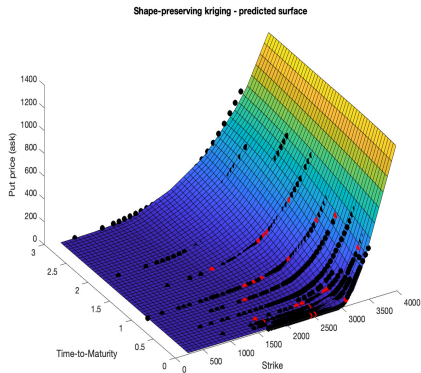
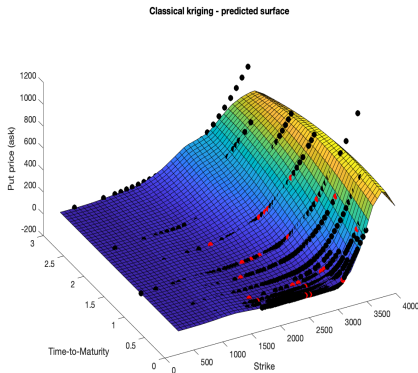
Let $p(T, k) := e^{qT} P(T, K)$ where $k = e^{-(r-q)T} K$ be the *reduced* put price. The put price surface $(T, K) \rightarrow P(T, K)$ is **free of static arbitrage** if and only if the reduced price function p is such that

- $p(\cdot, k)$ is a **non-decreasing** function and $p(0, k) = (k - S_0)^+$, for any $k \geq 0$
- $p(T, \cdot)$ is a **convex** function, $p(T, 0) = 0$, $\frac{\partial p}{\partial k}(T, 0) = 0$ and $\lim_{k \rightarrow \infty} p(T, k) = k - S_0$, for any $T \geq 0$

GP regression with shape constraints

Illustrative example

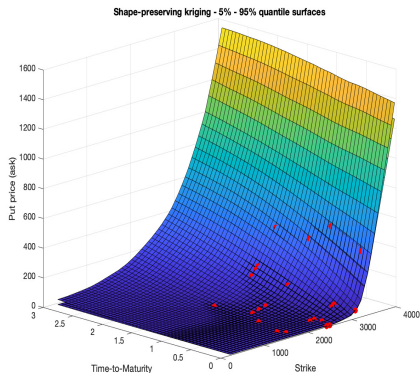
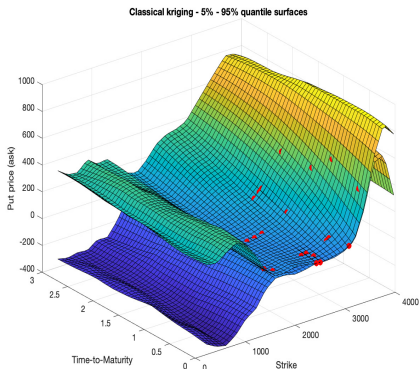
- Data : Euro Stoxx 50 put prices, January 10, 2019
- 5% of the data used (red points), zero-mean Gaussian prior with a Gaussian kernel
- Classical GP (left) vs GP with no-arbitrage constraints (right)



GP regression with shape constraints

Kriging of option price surface - quantification of uncertainty

- Pointwise 5% and 95% estimated quantiles of the fitted GP
- Classical kriging (left) vs kriging with no-arbitrage constraints (right)



GP regression with shape constraints

Assume that the a priori belief on the (reduced) price surface p is given as a GP prior process Y . The GP approach with shape constraints consists in estimating the conditional distribution of Y given

$$\begin{cases} \mathbf{y} = Y(X) + \varepsilon \\ Y \in \mathcal{M} \end{cases}$$

where $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$, $\mathbf{y} = (y_1, \dots, y_n)^\top$, ε is a zero-mean Gaussian noise in \mathbb{R}^n (independent of Y) and \mathcal{M} is a convex set of functions satisfying some shape properties.

For instance, \mathcal{M} can be :

- $\mathcal{M}_0^d := \{f \in \mathcal{C}([0, 1]^d, \mathbb{R}) \mid y_{\min} \leq f(x) \leq y_{\max}, \forall x \in D\}$
- $\mathcal{M}_1^1 := \{f \in \mathcal{C}([0, 1], \mathbb{R}) \mid f \text{ is non-decreasing}\}$
- $\mathcal{M}_2^1 := \{f \in \mathcal{C}([0, 1], \mathbb{R}) \mid f \text{ is convex}\}$
- $\mathcal{M}_{12}^2 := \{f \in \mathcal{C}([0, 1]^2, \mathbb{R}) \mid f \text{ is non-decreasing in } t \text{ and convex in } x\}$

Main issues :

- The posterior process is not Gaussian
- The shape condition is usually **infinite-dimensional**

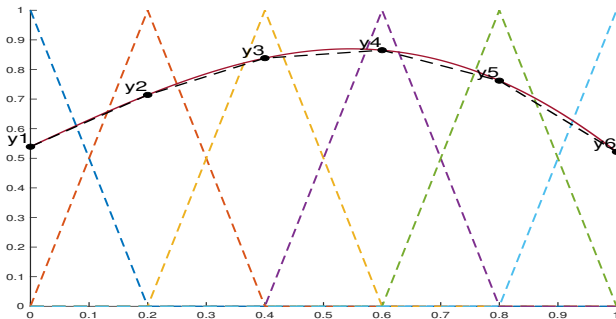
Proposed solutions :

- We construct a **finite-dimensional approximation** of Y for which the shape condition is easy to check.
- We consider the **mode of the posterior distribution** (as opposed to the posterior mean) as the response surface predictor
- Hyper-parameters are estimated using (unconstrained) MLE
- Sampling of the posterior (no-arbitrage) distribution by Hamiltonian Monte Carlo starting from the mode

Finite-dimensional approximation of GP (1d case)

As in [Maatouk and Bay \(2014\)](#), [Cousin et al. \(2016\)](#), [López et al. \(2018\)](#), we rely on basis function approximation.

- Input domain D is normalized on $[0, 1]$ and discretized on a regular subdivision $u_0 < \dots < u_N$ with a constant mesh δ .
- For each u_i , we consider hat functions $\phi_i(x) := \max\left(1 - \frac{|x-u_i|}{\delta}, 0\right)$
- Y is approximated on D by $Y^N(x) = \sum_{i=0}^N Y(u_i)\phi_i(x)$



Finite-dimensional approximation of GP (2d case)

- $D = [0, 1]^2$ is discretized on a $(N_t + 1) \times (N_x + 1)$ regular grid with knots (u_i, v_j) , $i = 0, \dots, N_t$, $j = 0, \dots, N_x$.
- For each knot (u_i, v_j) , we consider tensor product basis functions

$$\phi_{i,j}(t, x) := \max\left(1 - \frac{|t - u_i|}{\delta_t}, 0\right) \max\left(1 - \frac{|x - v_j|}{\delta_x}, 0\right)$$

- Y is approximated on D by

$$Y^N(t, x) = \sum_{i=0}^{N_t} \sum_{j=0}^{N_x} Y(u_i, v_j) \phi_{i,j}(t, x)$$

- $N = (N_t + 1)(N_x + 1)$ is the number of knots

Proposition

Let Y be a **zero-mean** GP with covariance function K and with almost surely continuous paths.

- The finite-dimensional process Y^N uniformly converges to Y on D as $N_t \rightarrow \infty$ and $N_x \rightarrow \infty$, almost surely.
- $Y^N(t, x) = \Phi(t, x)\xi$ where $\xi := (Y(u_0, v_0), Y(u_0, v_1), \dots, Y(u_{N_t}, v_{N_x}))^\top$ is a zero-mean Gaussian vector with $N \times N$ covariance matrix Γ^N such that $\Gamma^N = K((u_{i_1}, v_{j_1}), (u_{i_2}, v_{j_2}))$.

Shape-preserving conditions :

- Y^N is bounded on $[y_{\min}, y_{\max}]$ if and only if $y_{\min} \leq \xi_{i,j} \leq y_{\max}$
- $Y^N(t, x)$ is a non-decreasing function of t if and only if $\xi_{i+1,j} \geq \xi_{i,j}$
- $Y^N(t, x)$ is a convex function of x if and only if $\xi_{i,j+2} - \xi_{i,j+1} \geq \xi_{i,j+1} - \xi_{i,j}$
- ...

New formulation of the problem : Consider a zero-mean GP prior Y with covariance function K . The N -dimensional approximation Y^N of Y on D is such that, for any $x \in D$, $Y^N(x) = \Phi(x)\xi$, where ξ is a zero-mean Gaussian vector with covariance matrix Γ^N .

GP regression with shape constraints consists in finding the conditional distribution of $Y^N = \Phi(\cdot)\xi$ given that

$$\begin{cases} \mathbf{y} = \Phi(X) \cdot \xi + \varepsilon \\ \xi \in \mathcal{C}_{ineq} \end{cases}$$

where $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$, $\mathbf{y} = (y_1, \dots, y_n)^\top$, ε is a zero-mean Gaussian noise in \mathbb{R}^n (independent of ξ) and \mathcal{C}_{ineq} is a set of linear inequality constraints.

This is equivalent to finding the distribution of a **truncated Gaussian vector** $Z := [\xi \mid \mathbf{y} = \Phi(X) \cdot \xi + \varepsilon]$ given $\xi \in \mathcal{C}_{ineq}$

- We consider 2-dimensional **anisotropic stationary kernels** :

$$K(\mathbf{x}, \mathbf{x}') = \sigma^2 K_t(T - T'; \theta_t) K_x(k - k'; \theta_x)$$

where K_t, K_x are stationary kernels. ex : Gaussian, **Matérn 5/2**, Matérn 3/2, Exponential.

- Homoscedastic noise : $\varepsilon \sim \mathcal{N}(0, \Sigma_{noise})$ where $\Sigma_{noise} = \sigma_{noise}^2 \mathbb{I}_n$
- **Hyper-parameters** : $\lambda = (\sigma, \theta_1, \dots, \theta_d, \sigma_{noise})$

Following [López-Lopera et al \(2017\)](#), two MLE approaches can be considered

- **Unconditional likelihood** : Find λ that maximizes the Gaussian likelihood $\mathbb{P}(\Phi(X) \cdot \xi + \varepsilon = \mathbf{y} \mid \lambda)$ or log-likelihood

$$\mathcal{L}_N(\lambda) := -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |C| - \frac{1}{2} \mathbf{y}^\top C^{-1} \mathbf{y}$$

where $C := \Phi(X)\Gamma^N(\lambda)\Phi(X) + \Sigma_{noise}(\lambda)$

- **Conditional likelihood** : Find λ that maximizes the conditional probability $\mathbb{P}(\Phi(X) \cdot \xi + \varepsilon = \mathbf{y} \mid \xi \in \mathcal{C}_{ineq}, \lambda)$ or the log-likelihood

$$\mathcal{L}_{N,cond}(\lambda) := \mathcal{L}_N(\lambda) + \log \mathbb{P}(\xi \in \mathcal{C}_{ineq} \mid \Phi(X) \cdot \xi + \varepsilon = \mathbf{y}) - \log \mathbb{P}(\xi \in \mathcal{C}_{ineq})$$

We define the (a posteriori) **most probable response surface** and **measurement noises** as

$$\begin{cases} M_K^N(\mathbf{x}) := \Phi(\mathbf{x}) \cdot (\mathbf{c}_1^*, \dots, \mathbf{c}_N^*)^\top, \mathbf{x} \in D \\ \mathbf{e}^* := (\mathbf{e}_1^*, \dots, \mathbf{e}_n^*)^\top \end{cases}$$

where $(\mathbf{c}^*, \mathbf{e}^*)$ is the mode of the truncated Gaussian vector $(\boldsymbol{\xi}, \boldsymbol{\varepsilon})$ given the constraints, defined as solution of

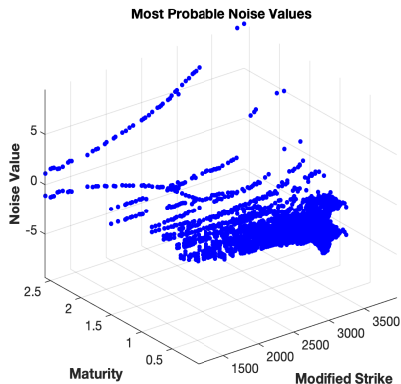
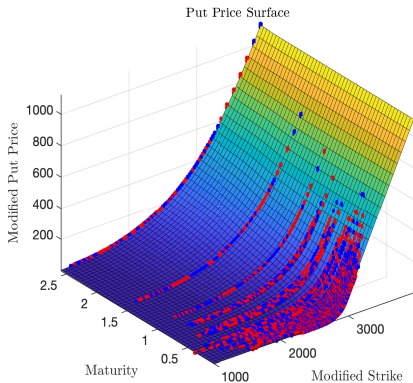
$$\max_{\mathbf{c}, \mathbf{e}} \mathbb{P}(\boldsymbol{\xi} \in [\mathbf{c}, \mathbf{c} + d\mathbf{c}], \boldsymbol{\varepsilon} \in [\mathbf{e}, \mathbf{e} + d\mathbf{e}] \mid \Phi(X) \cdot \boldsymbol{\xi} + \boldsymbol{\varepsilon} = \mathbf{y}, \boldsymbol{\xi} \in \mathcal{C}_{ineq}).$$

The mode $(\mathbf{c}^*, \mathbf{e}^*)$ is solution of a quadratic problem

$$\min_{\substack{\mathbf{c}^\top (\Gamma^N)^{-1} \mathbf{c} + \mathbf{e}^\top \Sigma_{noise}^{-1} \mathbf{e} \\ \Phi(X) \cdot \mathbf{c} + \mathbf{e} = \mathbf{y}, \mathbf{c} \in \mathcal{C}_{ineq}}} \left(\mathbf{c}^\top (\Gamma^N)^{-1} \mathbf{c} + \mathbf{e}^\top \Sigma_{noise}^{-1} \mathbf{e} \right)$$

Mode estimator

- Data : S&P 500 bid ask put prices as of of May 18, 2019
- Fitted Matérn 5/2 kernel using uncond. MLE, 3340 training data in blue, 1755 testing data in red, $N_t = 25$, $N_x = 100$.
- Most probable surface (left) vs most probable noise values (right)



First remark that the distribution of ξ given $\Phi(X) \cdot \xi + \varepsilon = y$ is multinormal $\mathcal{N}(\mu_{cond}, \Sigma_{cond})$ where

$$\begin{cases} \mu_{cond} = \Gamma^N \Phi(X)^\top (\Phi(X) \Gamma^N \Phi(X)^\top + \Sigma_{noise})^{-1} b \\ \Sigma_{cond} = \Gamma^N - \Gamma^N \Phi(X)^\top (\Phi(X) \Gamma^N \Phi(X)^\top + \Sigma_{noise})^{-1} \Phi(X) \Gamma^N \end{cases}$$

Following [López-Lopera et al \(2017\)](#), we consider the Hamiltonian Monte Carlo method introduced by [Pakman and Paninski \(2013\)](#) for sampling truncated multivariate Gaussian :

$$\mathcal{TN}(\mu_{cond}, \Sigma_{cond}, \mathcal{C}_{ineq})$$

MCMC initialized using the mode estimator since it satisfies the inequality constraints.

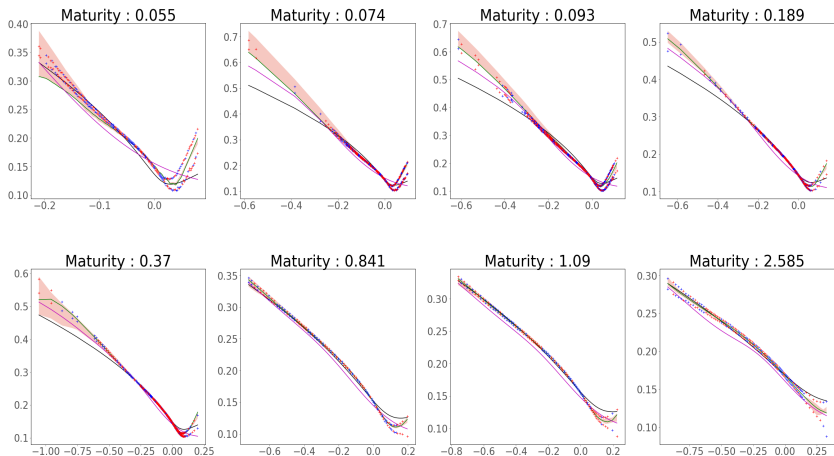
Comparison of RMSE and backtest results for the different approaches

SSVI, GP, NN based on IV, NN based on price, and unconstrained versions

IV RMSE (Price RMSE)	SSVI	GP	IV based NN	Price based NN	SSVI Unconstr.	GP Unconstr.	IV based NN Unconstr.	Price based NN Unconstr.
Calibr. fit on the training set	1.37% (2.574)	0.58% (0.338)	1.23% (2.897)	13.70% (9.851)	1.04% (2.691)	0.60% (0.321)	0.84% (2.163)	5.65 % (2.456)
Calibr. fit on the testing set	1.52% (2.892)	0.57% (0.355)	1.29% (2.966)	14.27% (10.347)	1.09% (2.791)	0.57% (0.477)	0.86% (2.045)	6.14% (2.888)
MC backtest	8.69% (22.826)	19.76% (74.017)	2.95% (4.989)	6.37% (11.764)	N/A	N/A	N/A	N/A
FD backtest	6.88% (33.545)	7.86% (35.270)	3.43% (11.976)	5.56% (26.785)	N/A	N/A	N/A	N/A
Comput. time (seconds)	33	856	191	185	1	16	76	229

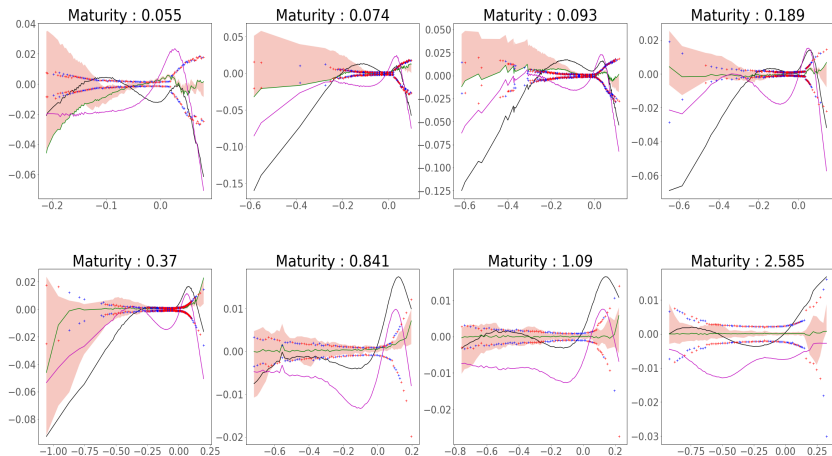
Table: The IV and price RMSEs of the SSVI, GP and NN approaches. Last line : computation times.

Slices of the IV Surface



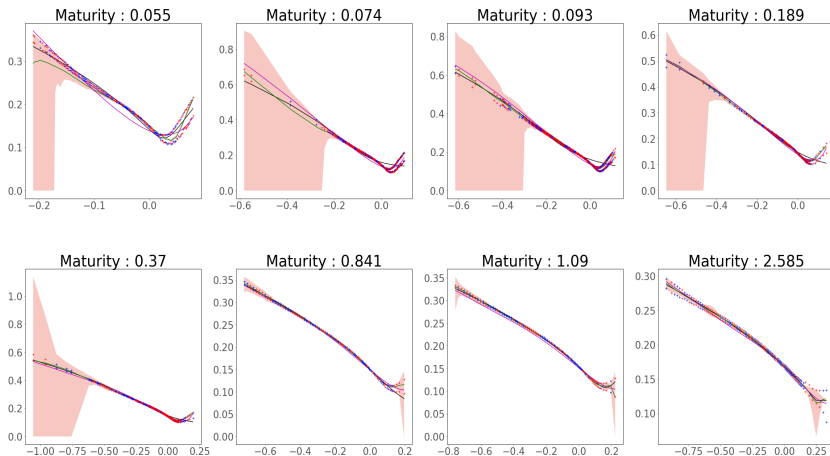
Slices of constrained GP (green), NN (purple), and SSVI (black) models of SPX puts with training bid-asks IVs (+) and testing bid-asks IVs (+) (the bid-ask IVs are reconstructed numerically from the corresponding bid-ask market prices). The shaded envelopes show 100 paths of the constrained GP's posterior.

Slices of fitted IV errors with respect to mid-price IVs



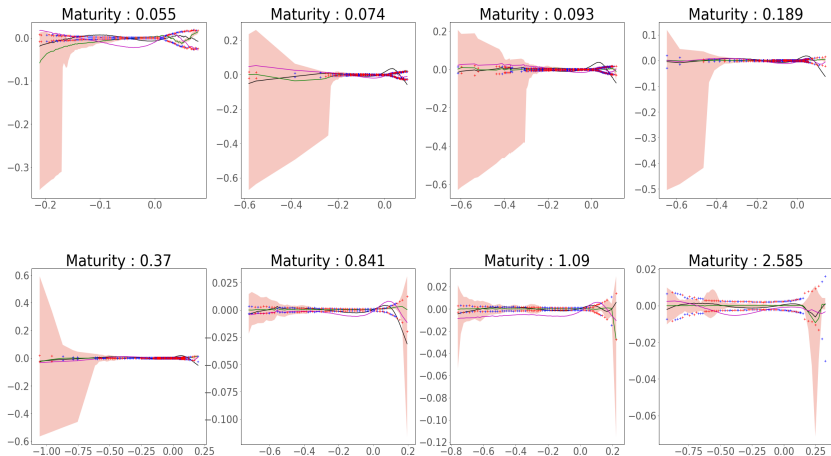
Slices of constrained GP (green), NN (purple), and SSVI (black) models of SPX puts with training bid-asks IVs (+) and testing bid-asks IVs (+) (the bid-ask IVs are reconstructed numerically from the corresponding bid-ask market prices). The shaded envelopes show 100 paths of the constrained GP's posterior.

Slices of the IV Surface - unconstrained case



Slices of constrained GP (green), NN (purple), and SSVI (black) models of SPX puts with training bid-asks IVs (+) and testing bid-asks IVs (+) (the bid-ask IVs are reconstructed numerically from the corresponding bid-ask market prices). The shaded envelopes show 100 paths of the constrained GP's posterior.

Slices of fitted IV errors with respect to mid-price IVs - unconstrained case



Slices of constrained GP (green), NN (purple), and SSVI (black) models of SPX puts with training bid-asks IVs (+) and testing bid-asks IVs (+) (the bid-ask IVs are reconstructed numerically from the corresponding bid-ask market prices). The shaded envelopes show 100 paths of the constrained GP's posterior.

Thanks for your attention.



Asgharian, H., Hess, W., and Liu, L. (2013).
A spatial analysis of international stock market linkages.
Journal of Banking & Finance, 37(12) :4738–4754.



Bay, X., Grammont, L., and Maatouk, H. (2016).
Generalization of the Kimeldorf-Wahba correspondence for constrained interpolation.



Cousin, A., Maatouk, H., and Rullière, D. (2016).
Kriging of financial term-structures.
European J. Oper. Res., 255(2) :631–648.



Cressie, N. (1990).
The origins of kriging.
Math. Geol., 22(3) :239–252.



da Barrosa, M. R., Salles, A. V., and Ribeiro, C. d. O. (2016).
Portfolio optimization through kriging methods.
Applied Economics, 48(50) :4894–4905.



De Spiegeleer, J., Madan, D. B., Reyners, S., and Schoutens, W. (2018).
Machine learning for quantitative finance : fast derivative pricing, hedging and fitting.
Quantitative Finance, 18(10) :1635–1643.



Dixon, M. F. and Crépey, S. (2018).
Multivariate gaussian process regression for derivative portfolio modeling : Application to
cva.



Liu, M. and Staum, J. (2010).

Stochastic kriging for efficient nested simulation of expected shortfall.
Journal of Risk, 12(3) :3.



López-Lopera, A. F., Bachoc, F., Durrande, N., and Roustant, O. (2018).

Finite-dimensional gaussian approximation with linear inequality constraints.
SIAM/ASA Journal on Uncertainty Quantification, 6(3) :1224–1255.



Ludkovski, M. (2018).

Kriging metamodels and experimental design for bermudan option pricing.
Journal of Computational Finance, 22(1) :37–77.



Ludkovski, M. and Risk, J. (2018).

Sequential design and spatial modeling for portfolio tail risk measurement.
SIAM Journal of Financial Mathematics.



Ludkovski, M., Risk, J., and Zail, H. (2018).

Gaussian process models for mortality rates and improvement factors.
ASTIN Bulletin : The Journal of the IAA, 48(3) :1307–1347.



Maatouk, H. and Bay, X. (2014).

Gaussian Process Emulators for Computer Experiments with Inequality Constraints.
in revision SIAM/ASA J. Uncertainty Quantification.



Pakman, A. and Paninski, L. (2014).

Exact hamiltonian monte carlo for truncated multivariate gaussians.
Journal of Computational and Graphical Statistics, 23(2) :518–542.



Roberts, S., Osborne, M., Ebden, M., Reece, S., Gibson, N., and Aigrain, S. (2013).
Gaussian processes for time-series modelling.
Philosophical Transactions of the Royal Society A : Mathematical, Physical and Engineering Sciences, 371(1984) :20110550.



Wahba, G. (1990).
Spline models for observational data, volume 59.
Siam.

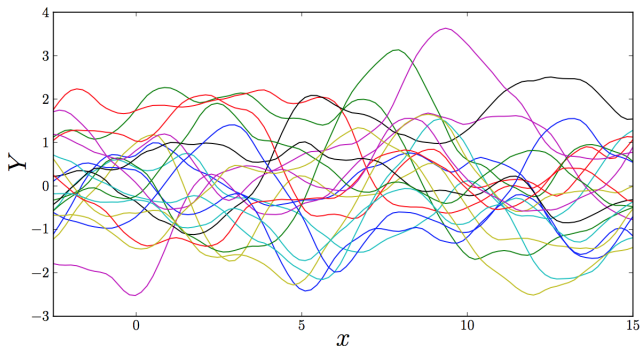


Williams, C. K. and Rasmussen, C. E. (2006).
Gaussian processes for machine learning.
the MIT Press, 2(3) :4.

Classical kriging

Estimation of the unknown function f using Bayesian statistics

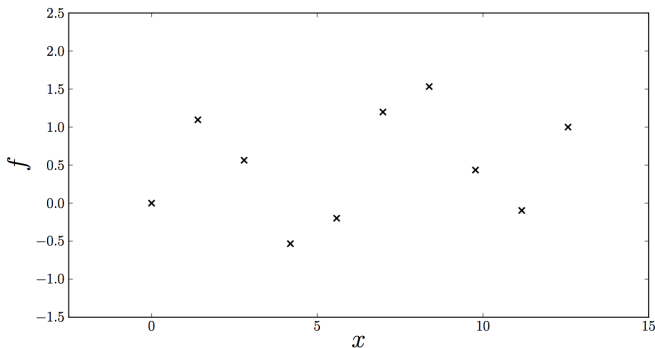
Our first belief in f is given as a Gaussian process prior Y



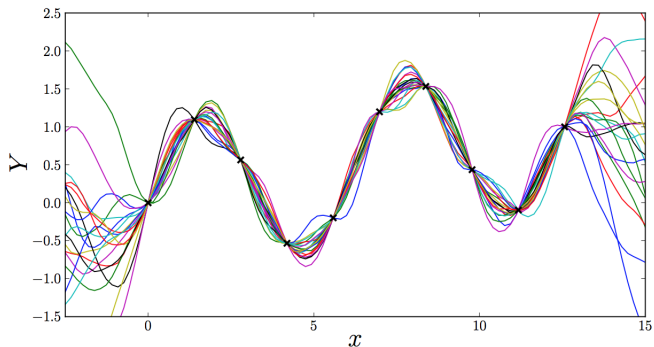
Classical kriging

The function f is known at some input points x^1, \dots, x^n :

$$f(x^1) = y^1, \dots, f(x^n) = y^n.$$



This belief is updated given that $Y(x_1) = y_1, \dots, Y(x_n) = y_n$



Source : presentation of N. Durrande

Definition : Gaussian process (GP) or Gaussian random field

A Gaussian process is a collection of random variables, any finite number of which have (consistent) joint Gaussian distributions.

A Gaussian process ($Y(x), x \in \mathbb{R}^d$) is characterized by its **mean function**

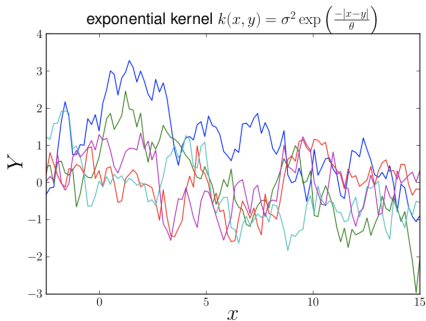
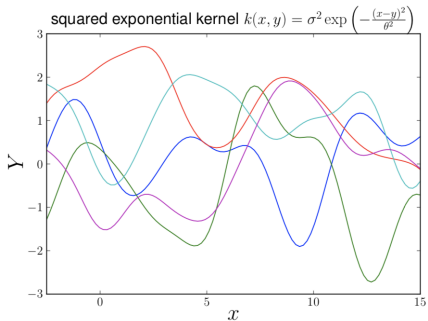
$$\mu : x \in \mathbb{R}^d \longrightarrow \mathbb{E}(Y(x)) \in \mathbb{R}.$$

and its **covariance function**

$$K : (x, x') \in \mathbb{R}^d \times \mathbb{R}^d \longrightarrow \text{Cov}(Y(x), Y(x')) \in \mathbb{R}.$$

1D kriging kernel	$K(x, x')$	Class
Gaussian	$\sigma^2 \exp\left(-\frac{(x-x')^2}{2\theta^2}\right)$	\mathcal{C}^∞
Matérn 5/2	$\sigma^2 \left(1 + \frac{\sqrt{5} x-x' }{\theta} + \frac{5(x-x')^2}{3\theta^2}\right) \exp\left(-\frac{\sqrt{5} x-x' }{\theta}\right)$	\mathcal{C}^2
Matérn 3/2	$\sigma^2 \left(1 + \frac{\sqrt{3} x-x' }{\theta}\right) \exp\left(-\frac{\sqrt{3} x-x' }{\theta}\right)$	\mathcal{C}^1
Exponential	$\sigma^2 \exp\left(-\frac{ x-x' }{\theta}\right)$	\mathcal{C}^0

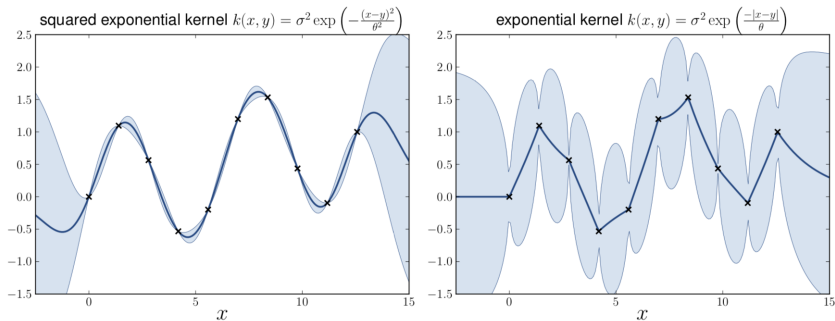
Changing the kernel K means changing the initial belief on f (i.e., the prior).



Source : presentation of N. Durrande

Given the observations, the model is entirely defined by the kernel.

Changing the kernel K has a huge impact on the model



Source : presentation of N. Durrande

- $\mathbf{X} = (x_1, \dots, x_n)^\top \in \mathbb{R}^{n \times d}$: some design points
- $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$: observed values of f at these points
- $Y(\mathbf{X}) = (Y(x_1), \dots, Y(x_n))^\top$: vector composed of Y at point \mathbf{X}

The conditional process is still a Gaussian Process

Let Y be a GP with mean μ and covariance function K . The conditional process $Y \mid Y(\mathbf{X}) = \mathbf{y}$ is a GP with mean function

$$\eta(x) = \mu(x) + \mathbf{k}(x)^\top \mathbb{K}^{-1}(\mathbf{y} - \boldsymbol{\mu}), \quad x \in \mathbb{R}^d$$

and covariance function \tilde{K} given by

$$\tilde{K}(x, x') = K(x, x') - \mathbf{k}(x)^\top \mathbb{K}^{-1} \mathbf{k}(x'), \quad x, x' \in \mathbb{R}^d$$

where $\boldsymbol{\mu} = \mu(\mathbf{X}) = (\mu(x_1), \dots, \mu(x_n))^\top$, \mathbb{K} is the covariance matrix of $Y(\mathbf{X})$ and $\mathbf{k}(x) = (K(x, x_1), \dots, K(x, x_n))^\top$

Note that **computational complexity** of \mathbb{K}^{-1} is $O(n^3)$.

In some applications, the unknown function f is not directly observed at X . But, if $f(X)$ is known up to solving a **linear equality system**, kriging can still be applied without loss of efficiency :

$$A \cdot f(X) = \mathbf{b}, \quad (1)$$

where

- A is a given matrix of dimension $n \times m$
- $\mathbf{b} = (b_1, \dots, b_n)^\top \in \mathbb{R}^n$
- $X = (x_1, \dots, x_m)^\top \in \mathbb{R}^{m \times d}$: some design points
- $f(X) = (f(x_1), \dots, f(x_m))^\top \in \mathbb{R}^m$

Classical kriging - indirect observations

The GP prior Y is updated given $AY(X) = \mathbf{b}$ where

- $Y(X) = (Y(x_1), \dots, Y(x_m))$: vector composed of Y at point X

The conditional process is still a Gaussian Process

Let Y be a GP with mean μ and covariance function K . The conditional process $Y \mid AY(X) = \mathbf{b}$ is a GP with mean function

$$\eta(x) = \mu(x) + (\mathbf{A}\mathbf{k}(x))^\top (\mathbf{A}\mathbb{K}\mathbf{A}^\top)^{-1} (\mathbf{b} - \mathbf{A}\boldsymbol{\mu}), \quad x \in \mathbb{R}^d$$

and covariance function \tilde{K} given by

$$\tilde{K}(x, x') = K(x, x') - (\mathbf{A}\mathbf{k}(x))^\top (\mathbf{A}\mathbb{K}\mathbf{A}^\top)^{-1} \mathbf{A}\mathbf{k}(x'), \quad x, x' \in \mathbb{R}^d$$

where $\boldsymbol{\mu} = \mu(X) = (\mu(x_1), \dots, \mu(x_m))^\top$, \mathbb{K} is the covariance matrix of $Y(X)$, $\mathbf{k}(x) = (K(x, x_1), \dots, K(x, x_m))^\top$

Assume that f is known up to solving a linear equality system with measurement errors :

$$A \cdot f(X) + \varepsilon = \mathbf{b}. \quad (2)$$

where

- A is a given matrix of dimension $n \times m$
- $\mathbf{b} = (b_1, \dots, b_n)^\top \in \mathbb{R}^n$
- $X = (x_1, \dots, x_m)^\top \in \mathbb{R}^{m \times d}$: some design points
- $f(X) = (f(x_1), \dots, f(x_m))^\top \in \mathbb{R}^m$
- ε is zero-mean Gaussian noise in \mathbb{R}^n with covariance matrix Σ_{noise}

Note that A is not necessarily a full-rank matrix in the presence of noise

Classical kriging - indirect observations with noise

The GP prior Y is updated given $AY(X) + \epsilon = \mathbf{b}$ where

- $Y(X) = (Y(x_1), \dots, Y(x_m))$: vector composed of Y at point X
- ϵ is assumed to be independent of Y

The conditional process is still a Gaussian Process

Let Y be a GP with mean μ and covariance function K . The conditional process $Y \mid AY(X) + \epsilon = \mathbf{b}$ is a GP with mean function

$$\eta(x) = \mu(x) + (A\mathbf{k}(x))^T \left(A\mathbb{K}A^T + \Sigma_{\text{noise}} \right)^{-1} (\mathbf{b} - A\boldsymbol{\mu}), \quad x \in \mathbb{R}^d$$

and covariance function \tilde{K} given by

$$\tilde{K}(x, x') = K(x, x') - (A\mathbf{k}(x))^T \left(A\mathbb{K}A^T + \Sigma_{\text{noise}} \right)^{-1} A\mathbf{k}(x'), \quad x, x' \in \mathbb{R}^d$$

where $\boldsymbol{\mu} = \mu(X) = (\mu(x_1), \dots, \mu(x_m))^T$, \mathbb{K} is the covariance matrix of $Y(X)$, $\mathbf{k}(x) = (K(x, x_1), \dots, K(x, x_m))^T$

Training a Gaussian process

- In general, we do not have enough information to define - a priori - the mean function and the kernel function
- Training the GP consists in selecting the kernel function (and the corresponding hyper-parameters) that best represents the data structure
- Hyper-parameters : $\mathbf{p} = (\mu(\cdot), \sigma, \theta, \Sigma)$
- Estimation of \mathbf{p} by **maximizing the likelihood function**
- As $Y(X)$ is a Gaussian vector with mean $\boldsymbol{\mu}$ and covariance matrix \mathbb{K} , the **log likelihood function** is given by :

$$\log \mathbb{P}(Y(X) = \mathbf{y} \mid \mathbf{p}) = -\frac{1}{2} \log \det(\mathbb{K}) - \frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \mathbb{K}^{-1} (\mathbf{y} - \boldsymbol{\mu}) - \frac{n}{2} \log(2\pi)$$