

Machine Learning for Derivative Pricing

Gaussian processes vs. Gradient boosting

Sofie Reyners

Joint work with Jesse Davis, Laurens Devos, Jan De Spiegeleer, Dilip Madan and Wim Schoutens

Derivative pricing is time-consuming...

- ▶ Sophisticated models
- ▶ Exotic products
 - we need to rely on numerical methods, e.g. Monte Carlo simulations
- ▶ Portfolio valuation
 - Sensitivity analysis
 - VaR/ES calculations

... and markets are moving!

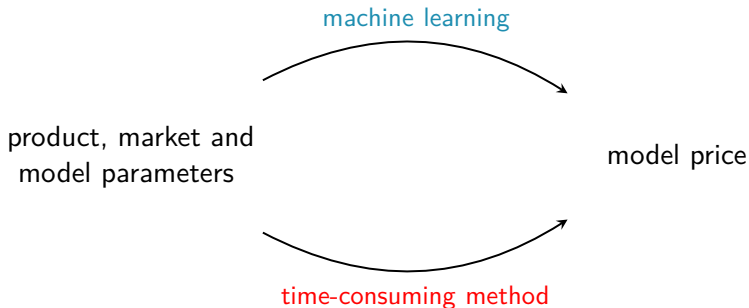
time-consuming algorithms



continuously moving markets

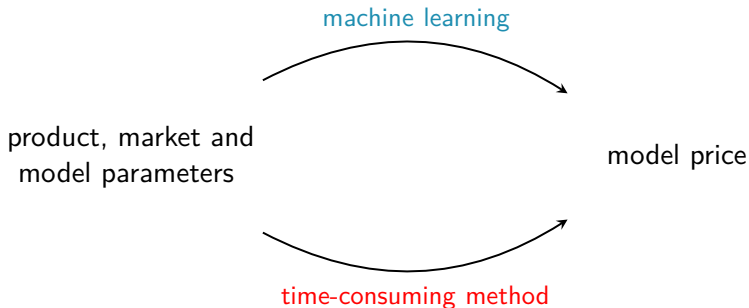
→ prices are outdated when available, risk calculations cannot be completed in time, ...

Let a machine learn the pricing function



Expensive pricing function is summarized with machine learning.

Let a machine learn the pricing function



When training is completed, prediction is extremely fast!

Literature

- ▶ Neural networks:
 - Ferguson & Green (2018), Deeply learning derivatives.
 - Liu, Oosterlee & Bohte (2019), Pricing options and computing implied volatilities using neural networks.
 - Horvath, Muguruza & Tomas (2019), Deep learning volatility.
 - ...

Literature

- ▶ Gaussian process regression:
 - De Spiegeleer, Madan, Reyners & Schoutens (2018), Machine learning for quantitative finance: fast derivative pricing, hedging and fitting.
 - Crépey & Dixon (2019), Gaussian process regression for derivative portfolio modeling and application to CVA computations.
 - Sousa, Esquivel & Gaspar (2012), Machine Learning Vasicek Model Calibration with Gaussian Processes.
- ▶ Gradient boosting machines:
 - Davis, Devos, Reyners & Schoutens, Gradient Boosting for Quantitative Finance. Working paper.

Gaussian process regression

Gaussian process regression (GPR)

Consider a training set $(X, \mathbf{y}) = \{(\mathbf{x}_i, y_i) \mid i = 1, \dots, n\}$ and assume that

$$y_i = f(\mathbf{x}_i) + \varepsilon_i$$

- ▶ $f(\mathbf{x})$ is a Gaussian process, characterized by two functions:
 - mean function: $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$
 - kernel function: $k(\mathbf{x}, \mathbf{x}') = \text{Cov}(f(\mathbf{x}), f(\mathbf{x}'))$
- ▶ $\varepsilon_i \sim \mathcal{N}(0, \sigma_n^2)$ are i.i.d. random variables representing the noise in the data.

Gaussian process

- ▶ $f(\mathbf{x}) \sim GP(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$
- ▶ If $(X, \mathbf{f}) = \{(\mathbf{x}_i, f_i) \mid i = 1, \dots, n\}$ is a sample from $f(\mathbf{x})$, then

$$\mathbf{f} \sim \mathcal{N}(M(X), K(X, X))$$

with

$$M(X) = \begin{bmatrix} m(\mathbf{x}_1) \\ \vdots \\ m(\mathbf{x}_n) \end{bmatrix}, \quad K(X, X) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}.$$

GPR: a Bayesian method

- ▶ Don't model the relation as one function, but as a distribution over functions.
- ▶ Procedure:
 - 1 Start from a prior GP (with zero mean)
 - prior knowledge: smooth function, periodic function, ...
 - prior distribution over functions
 - 2 Include observed data points
 - 3 Compute a posterior GP

Posterior distribution

Only consider functions that agree with the data.

- ▶ Take new inputs X_* , with corresponding (unknown) function values f_* .
- ▶ Joint distribution of training outputs and function values:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix} \right)$$

Posterior distribution

- ▶ Condition on the observations:

$$\mathbf{f}_* | X_*, X, \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

with

$$\boldsymbol{\mu} = K(X_*, X) [K(X, X) + \sigma_n^2 I]^{-1} \mathbf{y}$$

$$\boldsymbol{\Sigma} = K(X_*, X_*) - K(X_*, X) [K(X, X) + \sigma_n^2 I]^{-1} K(X, X_*)$$

- ▶ Point prediction: $\boldsymbol{\mu}$

Kernel function

Squared exponential kernel function (with ARD):

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp \left(-\frac{1}{2} \sum_{j=1}^d \frac{|x_j - x'_j|^2}{\ell_j^2} \right)$$

with hyperparameters σ_f and ℓ :

- ▶ σ_f^2 : signal variance
- ▶ ℓ_1, \dots, ℓ_d : characteristic length-scale parameters

→ **Hyperparameters** (including σ_n) are estimated on the training set, usually with MLE.

Machine learning for derivative pricing

- ▶ Construct a **training set**:

product, market and
model parameters

time-consuming method

model price

Machine learning for derivative pricing

- ▶ Construct a **training set**:

product, market and model parameters $\xrightarrow{\text{time-consuming method}}$ model price



sample n random combinations x_i



compute n corresponding prices y_i

Machine learning for derivative pricing

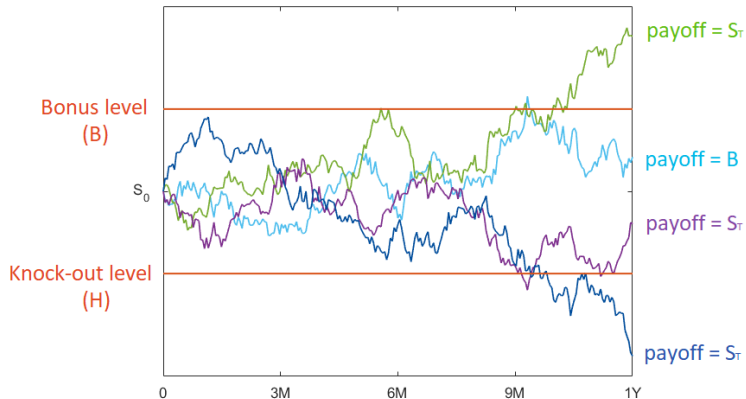
- ▶ Construct a **training set**:

product, market and time-consuming method model price
model parameters

- ▶ Fit a GPR model.
- ▶ Fast prediction of new model prices.

Case study:
pricing a structured product

Bonus certificate's payoffs



→ Price according to Heston model, using Monte Carlo simulation.

Build a training set

Parameter ranges:

Product/market	Heston model (*)
$B \in [105\%, 155\%]$	$\kappa \in [0.2, 1.6]$
$H \in [55\%, 95\%]$	$\rho \in [-0.95, -0.25]$
$T \in [11M, 1Y]$	$\theta \in [0.15, 0.65]$
$r \in [2\%, 3\%]$	$\eta \in [0.01, 0.25]$
$q \in [0\%, 5\%]$	$v_0 \in [0.01, 0.25]$

→ sample parameter combinations + calculate corresponding prices.

(*) κ = rate of mean reversion, ρ = correlation stock - vol, θ = vol-of-vol,
 η = long run variance, v_0 = initial variance.

Pricing a bonus certificate

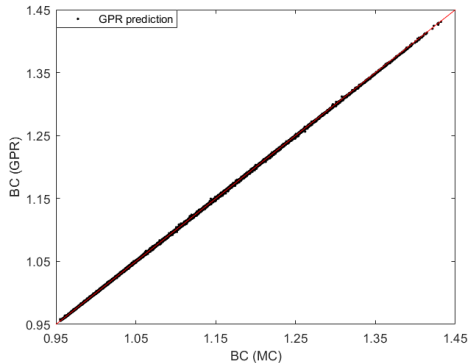
- ▶ Train a GPR model

$$B, H, T, r, q, \frac{\text{learn this relation}}{\kappa, \rho, \theta, \eta, v_0} \quad \text{BC price}$$

→ matrix inversion + hyperparameter optimization

- ▶ Construct a test set (100 000 instances):
 - similarly to training set construction
 - slightly smaller parameter intervals

GPR's predictions



Maximal error

absolute: 0.0058

relative: 0.51%

Prediction: 84 seconds

Speed-up × 2900

→ GPR model trained on 10 000 samples, tested on 100 000 samples.

Gradient Boosting Machines

Gradient boosting machine (GBM)

= ensemble of regression trees h

$$\hat{y}_i = F_K(\mathbf{x}_i) = \sum_{k=0}^K f_k(\mathbf{x}_i) = c + \sum_{k=1}^K \nu h_k(\mathbf{x}_i)$$

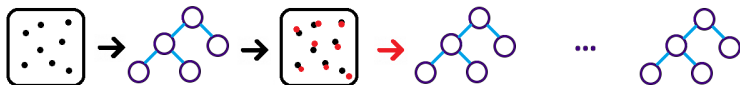
with $0 < \nu \leq 1$ a shrinkage parameter.

→ Many simple models f_k are combined into one strong model.



Gradient BOOSTING machine

- ▶ GBM is trained **stage-wise**:

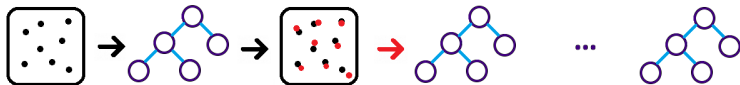


→ Trees are added sequentially:

f_k tries to correct for the **errors** made by model F_{k-1} .

→ Previously fitted trees are not readjusted.

GRADIENT boosting machine



How to measure the errors?

- ▶ Use pseudo-residuals

$$g_{i,k} = - \left[\frac{\partial \mathcal{L}(y_i, \hat{y})}{\partial \hat{y}} \right]_{\hat{y}=F_{k-1}(\mathbf{x}_i)} \quad i = 1, \dots, n$$

where \mathcal{L} is a loss function that measures how well the data is captured by the model.

Case study:
pricing a structured product

Boosting derivative pricing

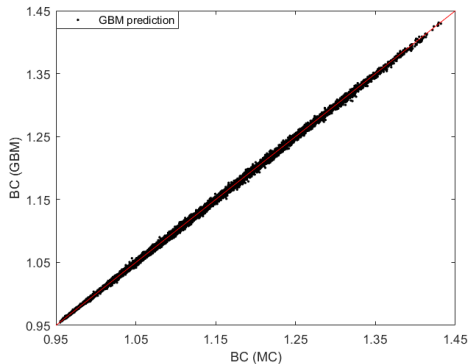
- ▶ Construct a training set.
- ▶ Train a GBM model:

$$\begin{array}{l} B, H, T, r, q, \\ \kappa, \rho, \theta, \eta, v_0 \end{array} \xrightarrow{\text{learn this relation}} \text{BC price}$$

→ LightGBM implementation.

- ▶ Fast prediction of new model prices.

GBM's predictions



Maximal error

absolute: 0.0146

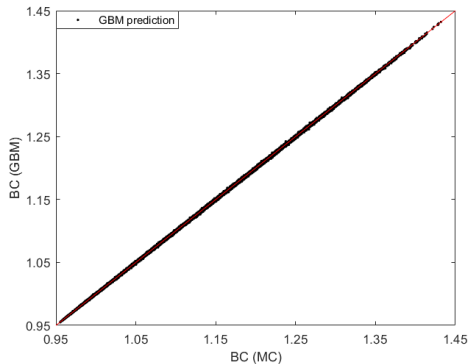
relative: 1.11%

Prediction: 42 seconds

Speed-up × 5800

→ GBM model trained on 10 000 samples, tested on 100 000 samples.

GBM's predictions



Maximal error

absolute: 0.0052

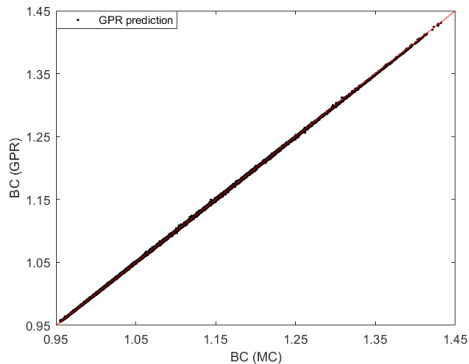
relative: 0.41%

Prediction: 42 seconds

Speed-up × 5800

→ GBM model trained on 100 000 samples, tested on 100 000 samples.

GPR's predictions



Maximal error

absolute: 0.0058

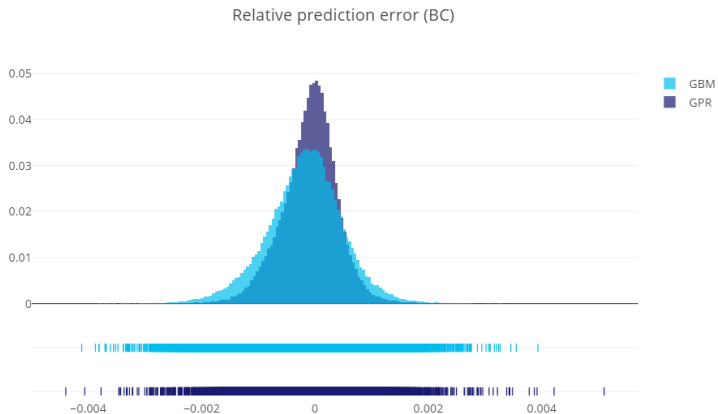
relative: 0.51%

Prediction: 84 seconds

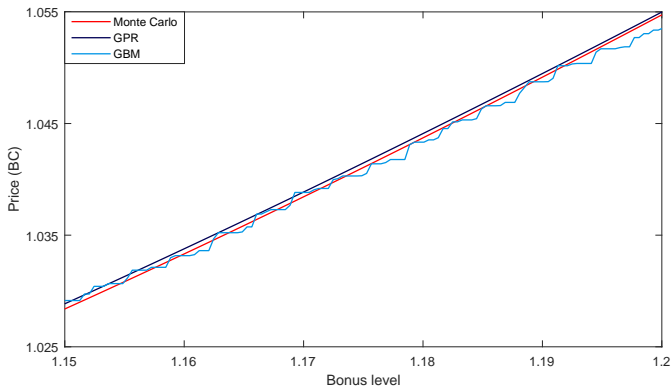
Speed-up × 2900

→ GPR model trained on 10 000 samples, tested on 100 000 samples.

GBM vs. GPR



GPR predictions are smoother



Fixed parameters: $H = 80\%$, $T = 1Y$, $r = 2.5\%$, $q = 3\%$, $\kappa = 1$, $\eta = 0.05$,
 $\theta = 0.5$, $\rho = -0.7$, $v_0 = 0.05$.

GPR has built-in analytic derivatives

- ▶ GPR's point prediction for \mathbf{x}_* :

$$f^*(\mathbf{x}_*) = K(\mathbf{x}_*, X) \underbrace{[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{y}}_{\boldsymbol{\alpha}}$$

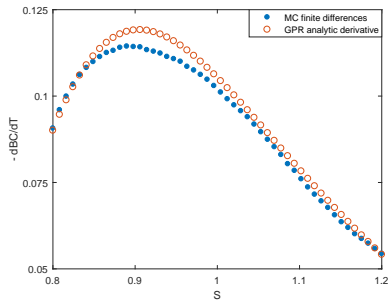
- ▶ Derivatives w.r.t. \mathbf{x}_* :

$$\frac{\partial f^*(\mathbf{x}_*)}{\partial \mathbf{x}_*} = \frac{\partial K(\mathbf{x}_*, X)}{\partial \mathbf{x}_*} \boldsymbol{\alpha} = -\Lambda^{-1} \tilde{X}_*^T [K(\mathbf{x}_*, X)]^T \circ \boldsymbol{\alpha}$$

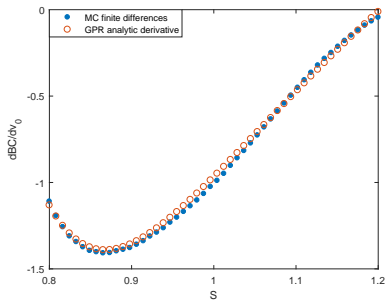
with

$$\tilde{X}_* = [\mathbf{x}_* - \mathbf{x}_1, \dots, \mathbf{x}_* - \mathbf{x}_n]^T, \quad \Lambda = \begin{bmatrix} \ell_1^2 & & 0 \\ & \ddots & \\ 0 & & \ell_d^2 \end{bmatrix}$$

Smooth Greek profiles for free!



(a) GPR's Theta



(b) GPR's Vega

→ Monte Carlo simulations based on 100 million price paths.

Conclusion

- ▶ Time-consuming pricing methods
- ▶ Machine learning methods
 - Gaussian process regression (GPR)
 - Gradient boosting machines (GBM)
- ▶ Pricing a structured product:
 - Speed-ups of several orders of magnitude.
 - GBM prediction is faster, training scales more easily.
 - GPR provides smoother price predictions and Greeks.

Thank you!

Contact: sofie.reyners@kuleuven.be

Hyperparameters

ntrees = 10 000

learning_rate = 0.2

num_leaves = 16

min_obs_in_leaf = 24

max_bin = 255

bagging_freq = 1

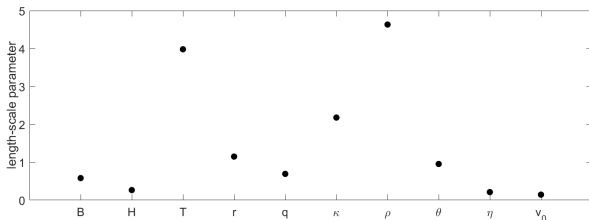
bagging_fraction = 0.5

L2_regularization = 0.8

booster = DART

drop_rate = 0.1

GPR



$$\sigma_f = 0.4580$$

$$\sigma_n = 0.0011$$