

# Beyond Surrogate Modeling: Learning the Local Volatility Via Shape Constraints\*

Marc Chataigner<sup>†</sup>, Areski Cousin<sup>‡</sup>, Stéphane Crépey<sup>§</sup>, Matthew Dixon<sup>¶</sup>, and Djibril Gueye<sup>‡</sup>

**Abstract.** We explore the abilities of two machine learning approaches for no-arbitrage interpolation of European vanilla option prices, which jointly yield the corresponding local volatility surface: a finite dimensional Gaussian process (GP) regression approach under no-arbitrage constraints based on prices, and a neural net (NN) approach with penalization of arbitrages based on implied volatilities. We demonstrate the performance of these approaches relative to the SSVI industry standard. The GP approach is proven arbitrage-free, whereas arbitrages are only penalized under the SSVI and NN approaches. The GP approach obtains the best out-of-sample calibration error and provides uncertainty quantification. The NN approach yields a smoother local volatility and a better backtesting performance, as its training criterion incorporates a local volatility regularization term.

**Key words.** Gaussian Processes; Local Volatility; Option pricing; Neural Networks; No-arbitrage.

**1. Introduction.** There have been recent surges of literature about the learning of derivative pricing functions by machine learning surrogate models, i.e. neural nets and Gaussian processes that are respectively surveyed in [11] and [4, Section 1]. There has, however, been relatively little coverage of no-arbitrage constraints when interpolating prices, and of the ensuing question of extracting the corresponding local volatility surface.

Tegnér & Roberts [12, see their Eq. (10)] first attempt the use of GPs for local volatility modeling by placing a Gaussian prior directly on the local volatility surface. Such an approach leads to a nonlinear least squares training loss function, which is not obviously amenable to gradient descent (stochastic or not), so the authors resort to a MCMC optimization. Zheng et al. [13] introduce shape constraint penalization via a multi-model gated neural network, which uses an auxiliary network to fit the parameters. The gated network is interpretable and lightweight, but the training is

---

*Acknowledgements:* The authors are thankful to Antoine Jacquier and Tahar Ferhati for useful hints regarding the SSVI method, and to an anonymous referee for stimulating comments.

\*Single-file demos Master.html (with the results of the paper) and Master.ipynb (for dynamic execution of all scripts) are available on [https://github.com/mChataign/Beyond-Surrogate-Modeling-Learning-the-Local-Volatility-Via-Shape-Constraints]. Note that, due to github size limitations, the file Master.html file must be downloaded locally (and then opened with a browser) to be displayed.

<sup>†</sup>LaMME, Université d'Evry, CNRS, Université Paris-Saclay; marc.chataigner@univ-evry.fr. The PhD thesis of Marc Chataigner is co-funded by the Research Initiative "Modélisation des marchés actions, obligations et dérivés", financed by HSBC France under the aegis of the Europlace Institute of Finance, and by the public grant ANR-11-LABX-0056-LLH LabEx LMH.

<sup>‡</sup>Institut de Recherche en Mathématique Avancée, Université de Strasbourg, 7 rue René Descartes, 67084 Strasbourg, cedex; a.cousin@unistra.fr

<sup>§</sup>LPSM, Université de Paris; Stephane.Crepey@lpsm.paris. The research of S. Crépey benefited from the support of the Chair Stress Test, RISK Management and Financial Steering, led by the French Ecole polytechnique and its Foundation and sponsored by BNP Paribas.

<sup>¶</sup>Department of Applied Mathematics, Illinois Institute of Technology, Chicago; matthew.dixon@iit.edu.

30 expensive and there is no guarantee of no-arbitrage. They do not consider the local  
 31 volatility and the associated regularization terms, nor do they assess the extent to  
 32 which no-arbitrage is violated in a test set.

33 Maatouk & Bay [9] introduce finite dimensional approximation of Gaussian pro-  
 34 cesses (GP) for which shape constraints are straightforward to impose and verify.  
 35 Cousin et al. [3] apply this technique to ensure arbitrage-free and error-controlled  
 36 yield-curve and CDS curve interpolation.

37 In this paper, we propose an arbitrage-free GP option price interpolation, which  
 38 jointly yields the corresponding local volatility surface, with uncertainty quantifica-  
 39 tion. Another contribution of the paper is to introduce a neural network approx-  
 40 imation of the implied volatility surface, penalizing arbitrages on the basis of the  
 41 Dupire formula, which is also used for extracting the corresponding local volatility  
 42 surface. This is all evidenced on an SPX option dataset.

43 Throughout the paper we consider European puts on a stock (or index)  $S$  with  
 44 dividend yield  $q$ , in an economy with interest rate term  $r$ , with  $q$  and  $r$  constant in  
 45 the mathematical description and deterministic in the numerics.

46 Given any rectangular domain of interest in time and space, we tacitly rescale  
 47 the inputs so that the domain becomes  $\Omega = [0; 1]^2$ . This rescaling avoids any  
 48 one independent variable dominating over another during any fitting of the market  
 49 prices.

50 **2. Gaussian process regression for learning arbitrage-free price surfaces.** We  
 51 denote by  $P_*(T; K)$  the time-0 market price of the put with maturity  $T$  and strike  
 52  $K$  on  $S$ , observed for a finite number of pairs  $(T; K)$ . Our first goal is to construct,  
 53 by Gaussian process regression, an arbitrage-free and continuous put price surface  
 54  $P : \mathbb{R}_+ \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$ , interpolating  $P_*$  up to some error term, and to retrieve the  
 55 corresponding local volatility surface  $(\cdot; \cdot)$  by the Dupire formula.

56 In terms of the reduced prices  $\rho(T; k) = e^{qT} P(T; K)$ ; where  $k = K e^{-(r-q)T}$ ; the  
 57 Dupire formula [5] reads (assuming  $\rho$  of class  $C^{1,2}$  on  $fT > 0g$ ):

$$58 \quad (2.1) \quad \frac{\partial^2(T; K)}{2} = \frac{\partial_T \rho(T; k)}{k^2 \partial_k^2 \rho(T; k)} =: \text{dup}(T; k):$$

59 Obviously, for this formula to be meaningful, its output must be nonnegative, which  
 60 holds if the interpolating map  $\rho$  exhibits nonnegative derivatives w.r.t.  $T$  and second  
 61 derivative w.r.t.  $k$ , i.e.

$$62 \quad (2.2) \quad \partial_T \rho(T; k) \geq 0; \quad \partial_k^2 \rho(T; k) \geq 0;$$

63 In this section, we consider a zero-mean Gaussian process prior on the mapping  
 64  $\rho = \rho(x)_{x \in \Omega}$  with correlation function  $c$  given, for any  $x = (T; k); x' = (T'; k') \in \Omega$ ,  
 65 by

$$66 \quad (2.3) \quad c(x; x') = \tau(T - T') \kappa(k - k');$$

67 Here  $(\tau; \kappa) = (\cdot; \cdot)^2$  correspond to length scale and variance hyper-parameters  
 68 of the kernel function  $c$ , whereas the functions  $\tau$  and  $\kappa$  are kernel correlation  
 69 functions.

70 Without consideration of the conditions (2.2), (unconstrained) prediction and  
71 uncertainty quantification are made using the conditional distribution  $p_j | \rho(\mathbf{x}) + \boldsymbol{\epsilon} =$   
72  $\mathbf{y}$ , where  $\mathbf{y} = [y_1; \dots; y_n]^\top$  are  $n$  noisy observations of the function  $\rho$  at input points  
73  $\mathbf{x} = [x_1; \dots; x_n]^\top$ , corresponding to observed maturities and strikes  $x_i = (T_i; k_i)$ ;  
74 the additive noise term  $\boldsymbol{\epsilon} = [\epsilon_1; \dots; \epsilon_n]^\top$  is assumed to be a zero-mean Gaussian  
75 vector, independent from  $\rho(\mathbf{x})$ , and with an homoscedastic covariance matrix given  
76 as  $\delta^2 I_n$ , where  $I_n$  is the identity matrix of dimension  $n$ . Note that bid and ask prices  
77 are considered here as (noisy) replications at the same input location.

78 **2.1. Imposing the no-arbitrage conditions.** To deal with the constraints (2.2),  
79 we adopt the solution of Cousin et al. [3] that consists in constructing a finite di-  
80 mensional approximation  $p^h$  of the Gaussian prior  $\rho$  for which these constraints can  
81 be imposed in the entire domain  $\Omega$  with a finite number of checks. One then recov-  
82 ers the (non Gaussian) constrained posterior distribution by sampling a truncated  
83 Gaussian process.

84 **Remark 1.** *Switching to a finite dimensional approximation can also be viewed*  
85 *as a form of regularization, which is also required to deal with the ill-posedness of*  
86 *the (numerical differentiation) Dupire formula.*

We first consider a discretized version of the (rescaled) input space  $\Omega = [0; 1]^2$   
as a regular grid  $(\{h\})_\ell$ , where  $\ell = (i; j)$ , for a suitable mesh size  $h$  and indices  $i; j$   
ranging from 0 to  $1/h$  (taken in  $\mathbb{N}^2$ ). For each knot  $\ell = (i; j)$ , we introduce the hat  
basis functions  $\varphi_\ell$  with support  $[(i-1)h; (i+1)h] \times [(j-1)h; (j+1)h]$  given, for  
 $x = (T; k)$ , by

$$\varphi_\ell(x) = \max(1 - \frac{jT - ihj}{h}; 0) \max(1 - \frac{jk - jhj}{h}; 0);$$

87 We take  $V = H^1(\Omega) = \{u \in L_2(\Omega) : D u \in L_2(\Omega)\}$ , where  $D u$  is a  
88 weak derivative of order  $j$ , as the space of (the realizations of)  $\rho$ . Let  $V^h \subset V$   
89 denote the finite dimensional linear subspace spanned by the  $M$  linearly independent  
90 basis functions  $\varphi_\ell$ . The (random) surface  $\rho$  in  $V$  is projected onto  $V^h$  as

91 (2.4) 
$$p^h(x) = \sum_{\ell} \rho(\{h\})_\ell \varphi_\ell(x); \quad \forall x \in \Omega;$$

92 If we denote  $\%_\ell = \rho(\{h\})_\ell$ , then  $\% = (\%_\ell)_\ell$  is a zero-mean Gaussian column vector  
93 (indexed by  $\ell$ ) with  $M \times M$  covariance matrix  $\Gamma^h$  such that  $\Gamma_{\ell; \ell'}^h = c(\{h\}; \ell, \ell')$ , for any  
94 two grid nodes  $\ell$  and  $\ell'$ . Let  $\varphi(x)$  denote the vector of size  $M$  given by  $\varphi(x) =$   
95  $(\varphi_\ell(x))_\ell$ . The equality (2.4) can be rewritten as  $p^h(x) = \varphi(x) \%$ . Denoting by  
96  $p^h(\mathbf{x}) = [p^h(x_1); \dots; p^h(x_n)]^\top$  and by  $\Phi(\mathbf{x})$  the  $n \times M$  matrix of basis functions  
97 where each row  $i$  corresponds to the vector  $\varphi(x_i)$ , one has  $p^h(\mathbf{x}) = \Phi(\mathbf{x}) \%$ . By  
98 application of the results of [9]:

- 99 **Proposition 2.** (i) *The finite dimensional process  $p^h$  converges uniformly to  $\rho$  on*  
100  *$\Omega$  as  $h \rightarrow 0$ , almost surely,*  
101 (ii)  *$p^h(T; k)$  is a nondecreasing function of  $T$  if and only if  $\%_{i+1; j} \geq \%_{i; j}; \forall (i; j)$ ,*  
102 (iii)  *$p^h(T; k)$  is a convex function of  $k$  if and only if  $\%_{i; j+2} \geq \%_{i; j+1} + \%_{i; j+1}$*   
103  *$\%_{i; j}; \forall (i; j)$ . ■*

104 In view of (i), denoting by  $\mathcal{I}$  the set of 2d continuous positive functions which are  
 105 nondecreasing in  $T$  and convex in  $k$ , we choose as constrained GP metamodel for  
 106 the put price surface the law of  $\rho^h$  conditional on

$$107 \quad \begin{cases} \rho^h(\mathbf{x}) + \boldsymbol{\mu} = \mathbf{y} \\ \rho^h \in \mathcal{I} \end{cases}$$

108 In view of (ii)-(iii),  $\rho^h \in \mathcal{I}$  if and only if  $\% \in \mathcal{I}^h$ ; where  $\mathcal{I}^h$  corresponds to the set of  
 109 ( $\ell$  indexed) vectors  $\% = (\rho_{i;j})_{i,j}$  such that  $\rho_{i+1;j} \geq \rho_{i;j}$  and  $\rho_{i;j+2} \geq \rho_{i;j+1} \geq \rho_{i;j}$   
 110  $\mathcal{S}(i;j)$ . Hence, our GP metamodel for the put price surface can be reformulated as  
 111 the law of  $\%$  conditional on

$$112 \quad (2.5) \quad \begin{cases} \Phi(\mathbf{x}) \% + \boldsymbol{\mu} = \mathbf{y} \\ \% \in \mathcal{I}^h \end{cases}$$

113 **2.2. Hyper-parameter learning.** Hyper-parameters consist in the length scales  
 114 and the variance parameter  $\sigma^2$  in (2.3), as well as the noise variance  $\delta$ . Up to a  
 115 constant, the so called marginal log likelihood of  $\%$  at  $\boldsymbol{\mu} = [\boldsymbol{\mu}; \delta]^\top$  can be expressed  
 116 as (see e.g. [10, Section 15.2.4, p. 523]):

$$117 \quad L(\boldsymbol{\mu}) = \frac{1}{2} \mathbf{y}^\top (\Phi(\mathbf{x}) \Gamma^h \Phi(\mathbf{x})^\top + \delta^2 I_n)^{-1} \mathbf{y} - \frac{1}{2} \log \left( \det (\Phi(\mathbf{x}) \Gamma^h \Phi(\mathbf{x})^\top + \delta^2 I_n) \right);$$

118 We maximize  $L$  for learning the hyper-parameters (MLE estimation).

119 **Remark 3.** *The above expression does not take into account the inequality con-*  
 120 *straints in the estimation. However, Bachoc et al. [1, see e.g. their Eq. (2)] argue*  
 121 *(and we observed empirically) that, unless the sample size is very small, condition-*  
 122 *ing by the constraints significantly increases the computational burden with negligible*  
 123 *impact on the MLE.*

**2.3. The most probable response surface and measurement noises.** We com-  
 124 pute the joint MAP  $(\hat{\rho}; \hat{\boldsymbol{\mu}})$  of the truncated Gaussian vector  $\%$  and of the Gaussian  
 125 noise vector  $\boldsymbol{\mu}$ ,

$$(\hat{\rho}; \hat{\boldsymbol{\mu}}) = \arg \max_{(\rho; \boldsymbol{\mu})} \text{Prob} \left( \% \in [\boldsymbol{\mu}; \boldsymbol{\mu} + d]; \boldsymbol{\mu} \in [\mathbf{e}; \mathbf{e} + d\mathbf{e}] \mid \Phi(\mathbf{x}) \% + \boldsymbol{\mu} = \mathbf{y}; \% \in \mathcal{I}^h \right)$$

124 (for the probability measure Prob underlying the GP model). As  $(\%; \boldsymbol{\mu})$  is Gaussian  
 125 centered with block-diagonal covariance matrix with blocks  $\Gamma^h$  and  $\delta^2 I_n$ ; this implies  
 126 that the MAP  $(\hat{\rho}; \hat{\boldsymbol{\mu}})$  is a solution to the following quadratic problem :

$$127 \quad (2.6) \quad \arg \min_{\Phi(\mathbf{x}) \cdot \rho + \boldsymbol{\mu} = \mathbf{y}; \rho \in \mathcal{I}^h} \left( \boldsymbol{\mu}^\top (\Gamma^h)^{-1} \boldsymbol{\mu} + \mathbf{e}^\top (\delta^2 I_n)^{-1} \mathbf{e} \right);$$

128 We define the most probable measurement noise to be  $\hat{\boldsymbol{\mu}}$  and the most probable  
 129 response surface  $\hat{\rho}^h(\mathbf{x}) = \Phi(\mathbf{x}) \hat{\rho}$ . Distance to the data can be an effect of arbitrage  
 130 opportunities within the data and/or misspecification / lack of expressiveness of the  
 131 kernel.

132 **2.4. Sampling finite dimensional Gaussian processes under shape constraints.**

133 The conditional distribution of  $\mathbf{y} | \Phi(\mathbf{x}) = \mathbf{y}$  is multivariate Gaussian with  
 134 mean  $\mathbf{y}(\mathbf{x})$  and covariance matrix  $\mathbf{C}_y(\mathbf{x})$  such that

135 (2.7) 
$$\mathbf{y}(\mathbf{x}) = \Gamma^h \Phi(\mathbf{x})^\top (\Phi(\mathbf{x}) \Gamma^h \Phi(\mathbf{x})^\top + \ell^2 I_n)^{-1} \mathbf{y}$$

136 (2.8) 
$$\mathbf{C}_y(\mathbf{x}) = \Gamma^h \Phi(\mathbf{x})^\top (\Phi(\mathbf{x}) \Gamma^h \Phi(\mathbf{x})^\top + \ell^2 I_n)^{-1} \Phi(\mathbf{x}) \Gamma^h$$

137 In view of (2.5), we thus face the problem of sampling from this truncated mul-  
 138 tivariate Gaussian distribution, which we do by Hamiltonian Monte Carlo, using  
 139 the MAP  $\hat{\mathbf{y}}$  of  $\mathbf{y}$  as the initial vector (which must verify the constraints) in the  
 140 algorithm.

141 **2.5. Local volatility.** Due to the shape constraints and to the ensuing finite-  
 142 dimensional approximation with basis functions of class  $C^0$  (for the sake of Proposi-  
 143 tion 2),  $\rho^h$  is not differentiable. Hence, exploiting GP derivatives analytics, as done  
 144 for the mean in [4, cf. Eq. (10)] and also for the covariance in [8], is not possible for  
 145 deriving the corresponding local volatility surface here. Computation of derivatives  
 146 involved in the Dupire formula is implemented by finite differences with respect to  
 147 a coarser grid (than the grid of basis functions). Another related solution would  
 148 be to formulate a weak form of the Dupire equation and construct a local volatility  
 149 surface approximation using a finite element method.

150 See Algorithm 2.1 for the main steps of the GP approach.

---

**Algorithm 2.1** The GP algorithm for local volatility surface approximation.

---

**Data:** Put price training set  $\rho^h$

**Result:**  $M$  realizations of the local volatility surface  $\hat{f}_{\text{dup}}^h; g_{i=1}^M$

- 1  $\hat{\mathbf{y}}$  Maximize the marginal log-likelihood of the put price surface  $\rho^h$  w.r.t.  
 // Hyperparameter fitting
  - 2  $(\hat{\mathbf{y}}; \hat{\Sigma})$  Minimize quadratic problem (2.6) based on  $\hat{\mathbf{y}}$  // Joint MAP estimate
  - 3  $\hat{\mathbf{y}}$  Initialize a Hamiltonian MC sampler
  - 4  $\rho_1^h; \dots; \rho_M^h$  Hamiltonian MC Sampler // Sampling price surfaces
  - 5  $\text{dup}_i^h$  Finite difference approximation using each  $\rho_i^h; i := 1 \dots M$
- 

151 **3. Neural networks implied volatility metamodeling.** Our second goal is to use  
 152 neural nets (NN) to construct an implied volatility (IV) put surface  $\Sigma : \mathbb{R}_+ \times \mathbb{R} \rightarrow$   
 153  $\mathbb{R}_+$ , interpolating implied volatility market quotes  $\Sigma_*$  up to some error term, both  
 154 being stated in terms of a put option maturity  $T$  and log-(forward) moneyness  
 155  $\kappa = \log\left(\frac{K}{S_0}\right) = \log\left(\frac{K}{S_0}\right) - (r - q)T$ . The advantage of using implied volatilities  
 156 rather than prices (as previously done in [2]), both being in bijection via the Black-  
 157 Scholes put pricing formula as well known, is their lower variability, hence better  
 158 performance as we will see.

159 The corresponding local volatility surface is given by the following local volatil-  
 160 ity implied variance formula, i.e. the Dupire formula stated in terms of the implied

161 total variance<sup>1</sup>  $\Theta(T; \cdot) = \Sigma^2(T; \cdot)T$  (assuming  $\Theta$  of class  $C^{1,2}$  on  $fT > 0g$ ):  
 (3.1)

$$162 \quad \Sigma^2(T; K) = \frac{\partial_T \Theta}{1 - \frac{\kappa}{\Theta} \partial_\kappa \Theta + \frac{1}{4} \left( \frac{1}{4} - \frac{1}{\Theta} + \frac{\kappa^2}{\Theta^2} \right) (\partial_\kappa \Theta)^2 + \frac{1}{2} \partial_{\kappa^2} \Theta} (T; \cdot) =: \frac{\text{cal}_T(\Theta)}{\text{butt}_k(\Theta)} (T; \cdot);$$

We use a feedforward NN with weights  $\mathbf{W}$ , biases  $\mathbf{b}$  and smooth activation functions for parameterizing the implied volatility and total variance, which we denote by

$$\Sigma = \Sigma_{\mathbf{W}, \mathbf{b}}; \Theta = \Theta_{\mathbf{W}, \mathbf{b}};$$

163 The terms  $\text{cal}_T(\Theta_{\mathbf{W}, \mathbf{b}})$  and  $\text{butt}_k(\Theta_{\mathbf{W}, \mathbf{b}})$  are available analytically, by automatic dif-  
 164 ferentiation, which we exploit below to penalize calendar spread arbitrages, i.e. neg-  
 165 ativity of  $\text{cal}_T(\Theta)$ , and butterfly arbitrage, i.e. negativity of  $\text{butt}_k(\Theta)$ .

166 The training of NNs is a non-convex optimization problem and hence does not  
 167 guarantee convergence to a global optimum. We must therefore guide the NN opti-  
 168 mizer towards a local optima that has desirable properties in terms of interpolation  
 169 error and arbitrage constraints. This motivates the introduction of an arbitrage pen-  
 170 alty function into the loss function to select the most appropriate local minima. An  
 171 additional challenge is that maturity-log moneyness pairs with quoted option prices  
 172 are unevenly distributed and the NN may favor fitting to a cluster of quotes to the  
 173 detriment of fitting isolated points. To remedy this non-uniform data fitting prob-  
 174 lem, we re-weight the observations by the Euclidean distance between neighboring  
 175 points. More precisely, given  $n$  observations  $\mathbf{x}_i = (T_i; \cdot)_i$  of maturity-log moneyness  
 176 pairs and of the corresponding market implied volatilities  $\Sigma_*(\cdot)_i$ , we construct the  
 177  $n \times n$  distance matrix with general term  $d(\mathbf{x}_i; \mathbf{x}_j) = \sqrt{(T_j - T_i)^2 + (\cdot_j - \cdot_i)^2}$ : We  
 178 then define the loss weighting  $w_i$  for each point  $\mathbf{x}_i$  as the distance  $w_i = \min_{j: j \neq i} d(\mathbf{x}_i; \mathbf{x}_j)$ :  
 179 with the closest point. These modifications aim at reducing error for any isolated  
 180 points. In addition, in order to avoid linear saturation of the neural network, we  
 181 apply a further log-maturity change of variables (adapting the partial derivatives  
 182 accordingly).

183 Learning the weights  $\mathbf{W}$  and biases  $\mathbf{b}$  to the data subject to no arbitrage soft  
 184 constraints (i.e. with penalization of arbitrages) then takes the form of the following  
 185 (nonconvex) loss minimization problem:

$$186 \quad (3.2) \quad \arg \min_{\mathbf{W}, \mathbf{b}} \sqrt{\frac{1}{n} \sum_i \left( w_i \frac{\Sigma_{\mathbf{W}, \mathbf{b}}(\cdot)_i - \Sigma_*(\cdot)_i}{\Sigma_*(\cdot)_i} \right)^2} + \frac{w}{m} \sum_{\epsilon \in \Omega_h} \mathbf{1}^\top R(\Theta_{\mathbf{W}, \mathbf{b}})(\epsilon);$$

187 where  $\mathbf{1} = [1; 2; 3]^\top \in \mathbb{R}_+^3$  and

$$188 \quad R(\Theta) = [\text{cal}_T^-(\Theta); \text{butt}_k^-(\Theta); \left( \frac{\text{cal}_T}{\text{butt}_k}(\Theta) - \bar{a} \right)^+ + \left( \frac{\text{cal}_T}{\text{butt}_k}(\Theta) - \underline{a} \right)^-]^\top$$

189 is a regularization penalty vector evaluated over a penalty grid  $\Omega_h$  with  $m$  nodes  
 190 as detailed below. The error criterion is calculated as a root mean square error on  
 191 relative difference, so that it does not discriminate high or low implied volatilities.

---

<sup>1</sup>This follows from the Dupire formula by simple transforms detailed in [6, p.13].

192 The first two elements in the penalty vector  $\mathcal{R}(\Theta)$  favor the no-arbitrage conditions  
193 (2.2) and the third element favors desired lower and upper bounds  $0 < \underline{a} < \bar{a}$   
194 (constants or functions of  $T$ ) on the estimated local variance  $\sigma^2(T; \mathcal{K})$ . In order to  
195 adjust the weight of penalization, we multiply our penalties by the weighting mean  
196  $w := \frac{1}{m} \sum_i w_i$ . Suitable values of the ‘‘Lagrange multipliers’’  $w_i$ ; ensuring the right  
197 balance between fit to the market implied volatilities and the constraints, is then  
198 obtained by grid search. Of course a soft constraint (penalization) approach does  
199 not fully prevent arbitrages. However, for large  $m$ , arbitrages are extremely unlikely  
200 to occur, except perhaps very far from  $\Omega$ . With this in mind, we use a penalty grid  
201  $\Omega_h$  that extends well beyond the domain of the IV interpolation. This is intended  
202 so that the penalty term penalizes arbitrages outside of the domain used for IV  
203 Interpolation.  
204 See Algorithm 3.1 for the pseudo-code of the NN approach.

---

**Algorithm 3.1** The NN-IV algorithm for local volatility surface approximation.

---

**Data:** Market implied volatility surface  $\Sigma_*$

**Result:** The local volatility surface  $\sqrt{\frac{\text{cal}_T}{\text{butt}_k}}(\Theta_{\hat{\mathbf{W}}; \hat{\mathbf{b}}})$

1  $(\hat{\mathbf{W}}; \hat{\mathbf{b}})$  Minimize the penalized training loss (3.2) w.r.t.  $(\mathbf{W}; \mathbf{b})$ ;

2  $\sqrt{\frac{\text{cal}_T}{\text{butt}_k}}(\Theta_{\hat{\mathbf{W}}; \hat{\mathbf{b}}})$  AAD differentiation of the trained NN implied vol. surface

---

## 205 4. Numerical results.

206 **4.1. Experimental design.** Our training set is prepared using SPX European  
207 puts with different available strikes and maturities ranging from 0.005 to 2.5 years,  
208 listed on 18th May 2019, with  $S_0 = \$2859:53$ . Each contract is listed with a bid/ask  
209 price and an implied volatility corresponding to the mid-price. The associated in-  
210 terest rate is constructed from US treasury yield curve and dividend yield curve  
211 rates are then obtained from call/put parity applied to the option market prices  
212 and forward prices. We preprocess the data by removing the shortest maturity  
213 options, with  $T < 0.055$ , and the numerically inconsistent observations for which  
214 the gap between the listed implied volatility and the implied volatility calibrated  
215 from mid-price with our interest/dividend curves exceeds 5% of the listed implied  
216 volatility. But we do not remove arbitrable observations. The preprocessed training  
217 set is composed of 1720 market put prices. The testing set consists of a disjoint set  
218 of 1725 put prices.

219 All results for the GP method are based on using Matern  $\nu = 5/2$  kernels over a  
220  $[0; 1]^2$  domain with fitted kernel standard-deviation hyper-parameter  $\hat{\kappa} = 185:7611$ ,  
221 length-scale hyper-parameters  $\hat{\tau} = 0.3282$  and  $\hat{\tau} = 0.2211$ , and homoscedastic  
222 noise standard deviation,  $\hat{\ell} = 0.6876$ .<sup>2</sup> The grid of basis functions for constructing  
223 the finite-dimensional process  $p^h$  has 100 nodes in the modified strike direction and  
224 25 nodes in the maturity direction. The Matlab interior point convex algorithm  
225 quadprog is used to solve the MAP quadratic program (2.6).

---

<sup>2</sup>When re-scaled back to the original input domain, the fitted length scale parameters of the 2D Matern  $\nu = 5/2$  are  $\hat{\theta}_\kappa = 973.1901$  and  $\hat{\theta}_\tau = 0.5594$ .

226 Regarding the NN approach, we use a three layer architecture similar to the one  
 227 based on prices (instead of implied volatilities in Section 3) in [2], to which we refer  
 228 the reader for implementation details. We use a penalty grid  $\Omega_h$  with  $m = 50 \dots 100$   
 229 nodes. In the moneyness and maturity coordinates, the domain of the penalty grid  
 230 is  $[0;0.005;10] \times [0;5;2]$ .

231 **4.2. Arbitrage-free SVI.** We benchmark the machine learning results with the  
 232 industry standard provided by the arbitrage free stochastic volatility inspired (SVI)  
 233 model of [7]. Under the “natural parameterization”  $\text{SVI} = (\Delta; \rho; \nu; \theta; \sigma)$ , the implied  
 234 total variance is given, for any fixed  $T$ , by

$$235 \quad (4.1) \quad \Theta_{\text{SVI}}(\kappa) = \Delta + \frac{\nu}{2} \left( 1 + \kappa \rho + \sqrt{(\kappa \rho + \nu)^2 + (1 - \rho^2)} \right) :$$

236 Our SSVI parameterization of a surface corresponds to  $\text{SVI}_T = (0; 0; \nu; \Theta_T; (\Theta_T))$   
 237 for each  $T$ , where  $\Theta_T$  is the at-the-money total implied variance and we use for  $\rho$  a  
 238 power law function  $\rho(\kappa) = \frac{\nu}{\nu + \kappa}$ . [7, Remark 4.4] provides sufficient conditions  
 239 on SSVI parameters  $(\nu, \rho, \sigma)$  with  $\rho = 0.5$  that rule out butterfly arbitrage,  
 240 whereas SSVI is free of calendar arbitrage when  $\Theta_T$  is nondecreasing.

241 We calibrate the model as in [7]:<sup>3</sup> First, we fit the SSVI model; Second, for  
 242 each maturity in the training grid, the five SVI parameters are calibrated, (starting  
 243 in each case from the SSVI calibrated values. The implied volatility is obtained  
 244 for new maturities by a weighted average of the parameters associated with the  
 245 two closest maturities in the training grid,  $T$  and  $U$ , say, with weights determined  
 246 by  $\Theta_T$  and  $\Theta_U$ . The corresponding local volatility is extracted by finite difference  
 247 approximation of (3.1).

248 As, in practice, no arbitrage constraints are implemented for SSVI by penaliza-  
 249 tion (see [7, Section 5.2]), in the end the SSVI approach is in fact only practically  
 250 arbitrage-free, much like our NN approach, whereas it is only the GP approach that  
 251 is proven arbitrage-free.

252 **4.3. Calibration results.** Training times for SSVI, GP, and NNs are reported  
 253 in the last row of Table 1 which, for completeness, also includes numerical results  
 254 obtained by NN interpolation of the prices as per [2]. Because price based NN results  
 255 are outperformed by IV based NN results we only focus on the IV based NN in the  
 256 figures that follow, referring to [2] for every detail on the price based NN approach.  
 257 We recall that, in contrast to the SSVI and NNs which fit to mid-quotes, GPs fit to  
 258 the bid-ask prices.

259 The GP implementation is in Matlab whereas the SSVI and NN approaches  
 260 are implemented in Python. On our (large) dataset, the constrained GP has the  
 261 longest training time. Training is longer for constrained SSVI than for unconstrained  
 262 SSVI because of the ensuing amendments to the optimization routine. There are  
 263 no arbitrage violations observed for any of the constrained methods in neither the  
 264 training or the testing grid. Unconstrained methods yield 18 violations with NN and  
 265 177 with SSVI on the testing set, out of a total of 1725 testing points, i.e. violations  
 266 in 1.04% and 10.26% of the test nodes. The unconstrained GP approach yields  
 267 constraint violations on 12.5% of the basis function nodes  $\{h\}$ . The NN penalizations

---

<sup>3</sup>Building on <https://www.mathworks.com/matlabcentral/profile/authors/4439546>.



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
CN 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 1: The IV and price RMSEs of the SSVI, GP and NN approaches. Last row: computation times (in seconds).

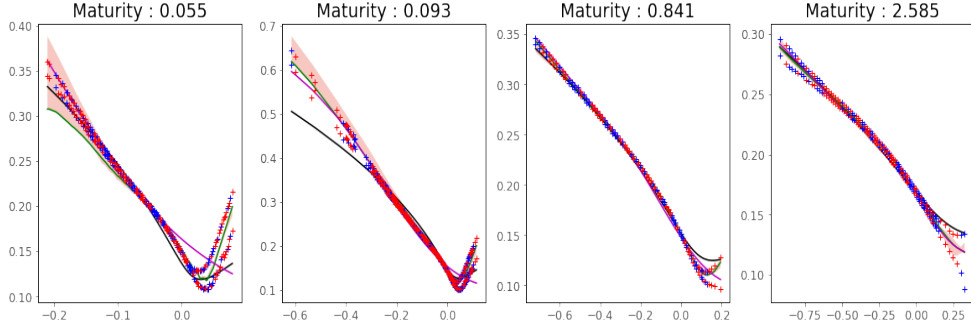
268  $(\text{cal}_T)^-$  and  $(\text{butt}_k)^-$  vanish identically on the penalty grid  $\Omega_h$  in the constrained  
269 case, whereas in the unconstrained case their averages across grid nodes in  $\Omega_h$  are  
270  $(\text{cal}_T)^- = 3.91 \cdot 10^{-6}$  and  $(\text{butt}_k)^- = 1.60 \cdot 10^{-2}$  with the IV based NN.

271 Fig. 1(a-b) respectively compare the fitted IV surfaces and their errors with  
272 respect to the market mid-implied volatilities, among the constrained methods. The  
273 surface is sliced at various maturities (more slices are available in the github) and  
274 the IVs corresponding to the bid-ask price quotes are also shown – the blue and red  
275 points respectively denote training and test observations.

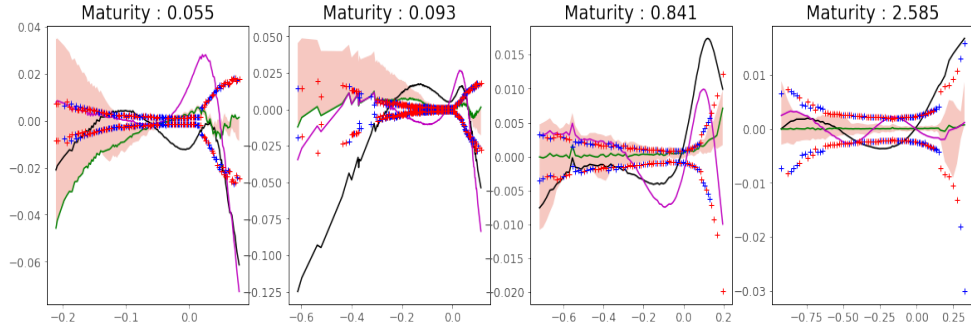
276 We generally observe good correspondence between the models and that each  
277 curve typically falls within the bid-ask spread, except for the shortest maturity con-  
278 tracts where there is some departure from the bid-ask spreads for observations with  
279 the lowest log-moneyness values. We see on Fig. 1(b) that the GP IV errors are  
280 small and mostly less than 5 volatility points, whereas NN and SSVI exhibit IV  
281 error that may exceed 15 volatility points. The green line and the red shaded en-  
282 velopes respectively denote the GP MAP estimates and the posterior uncertainty  
283 bands under 100 samples per observation. The support of the posterior GP process  
284 assessed on the basis of 100 simulated paths of the GP captures the majority of  
285 bid-ask quotes. The GP MAP estimate occasionally corresponds to the boundary  
286 of the support of the posterior simulation. This indicates that the posterior trun-  
287 cated Gaussian distribution is heavily skewed for some points, and that the MAP  
288 estimate consequently saturates the arbitrage constraints. This indicates a tension  
289 between these constraints and the calibration requirement, which cannot be fully  
290 reconciled, most likely because some of the (short maturity) data are arbitrable  
291 (they are at least illiquid and hence noisy). See notebook for location of arbitrages  
292 in the unconstrained approach.

293 Fig. 1(a-b) suggest that the data may exhibit arbitrage at the lowest maturities  
294 where the methods depart from the bid-ask spreads. This is further supported  
295 in Fig. 2(a-b) which shows the corresponding methods without the no-arbitrage  
296 constraints. In Fig. 2(a-b) we observe that the estimated IVs now fall within close  
297 proximity of the bid-ask spreads—all methods exhibit an error typically less than 5  
298 volatility points. Note that the y-axis has been scaled for each plot in Fig. 2(b) to

299 accommodate the wide uncertainty band of the posterior for the unconstrained GP.  
 300 Whereas the uncertainty band of the constrained GP spanned at most 10 volatility  
 301 points, the uncertainty band of the unconstrained GP is an order of magnitude  
 302 larger, sometimes spanning more than 100 volatility points.



(a) Implied volatilities.

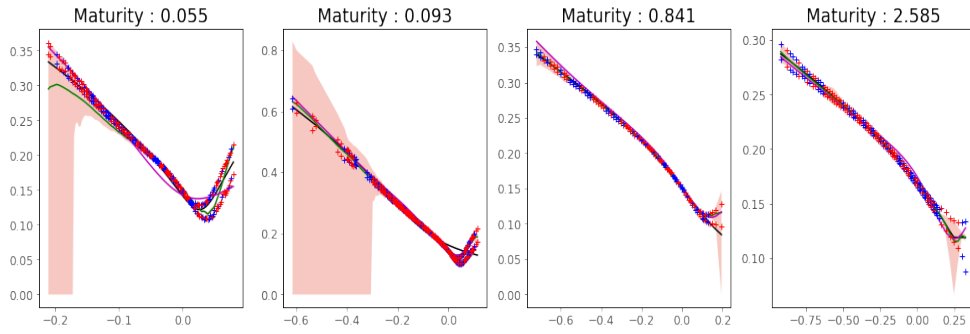


(b) Fitted IV errors with respect to mid-price IVs.

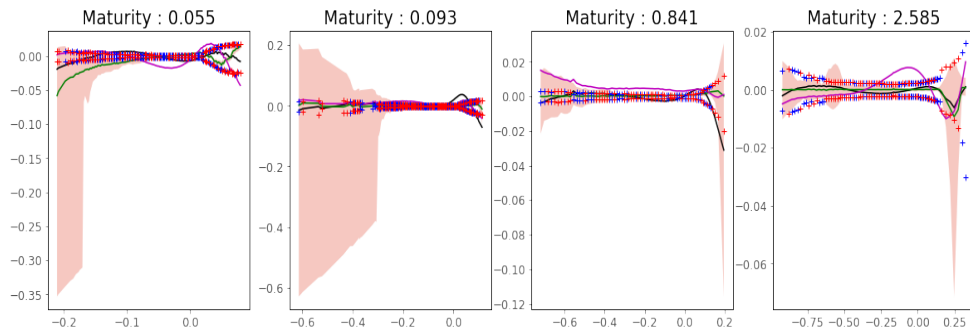
Figure 1: Slices of constrained GP (green), NN (purple), and SSVI (black) models of SPX puts with training bid-asks IVs (+) and testing bid-asks IVs as a function of log forward moneyness (+) (the bid-ask IVs are reconstructed numerically from the corresponding bid-ask market prices). The shaded envelopes show 100 paths of the constrained GP's posterior.

303 Fig. 3 shows the local volatility surfaces that stem from the three constrained  
 304 approaches. Fig. 3(a) shows the spiky local volatility surface generated by SSVI,  
 305 capped at the 200% level for scaling convenience. Fig. 3(b) shows the capped local  
 306 volatility surface constructed from the GP MAP price estimate. Fig. 3(c) shows the  
 307 (complete) NN local volatility surface.

308 **4.4. In-sample and out-of-sample calibration errors.** The error between the  
 309 prices of the calibrated models and the market data are evaluated on both the  
 310 training and the out-of-sample data set. The first two rows of Table 1 compare the  
 311 in-sample and out-of-sample RMSEs of the prices and implied volatilities across the  
 312 different approaches. The differences between the training and testing RMSEs are  
 313 small, suggesting that all approaches are not over-fitting the training set. The GP

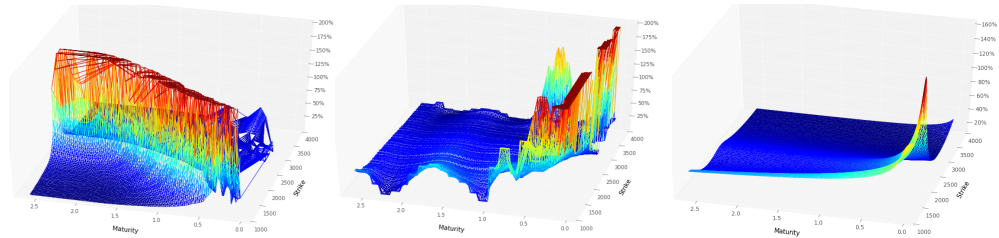


(a) Implied volatilities.



(b) Fitted IV errors with respect to mid-price IVs.

Figure 2: Same as Figure 1 but for unconstrained GP, NN and SSVI.



(a) The local volatility surface generated by SSVI with finite differences, capped at the 200% level. (b) The MAP estimate of the GP local volatility surface, capped at the 200% level. (c) The implied volatility based NN local volatility surface (with the local volatility penalization).

Figure 3: The GP, SSVI, and NN local volatility estimate.

314 exhibits the lowest price RMSEs.

315 **4.5. Backtesting results.** The first repricing backtest estimates the prices of  
 316 the European options corresponding to the testing set, by Monte Carlo sampling in  
 317 each calibrated local volatility model (same methodology as in [2, Section 7.2]). The  
 318 second approach uses finite differences to price the options with the calibrated local  
 319 volatility surfaces. The pricing PDEs with local volatility are discretized using a

320 Crank-Nicolson (CN) scheme implemented on a  $100 \times 100$  backtesting grid. The last  
321 two rows in Table 1 compare the resulting price backtest RMSEs across the different  
322 approaches. The NN fitted to implied volatilities exhibit significantly lower errors in  
323 the backtests, followed by NN based on prices, SSVI and GP. To quantify discretiza-  
324 tion error in these backtesting results (as opposed to the part of the error stemming  
325 from a wrong local volatility), we ran the same backtests in a Black-Scholes model  
326 with 20% volatility and the associated prices. The corresponding Monte Carlo and  
327 Crank-Nicholson backtesting IV(price) RMSEs are 2.90%(1.56) and 0.846%(4.10),  
328 confirming the significance of the above results.

329

330 **Conclusion.** We approach the option quote fitting problem from two perspectives:  
331 (i) the GP approach assumes noisy data and hence the existence of a latent function.  
332 The mid-prices are not considered, rather the GP calibrates to bid-ask quotes; and  
333 (ii) the NN and SSVI approaches fit to the mid-prices under a noise-free assumption.  
334 While these two approaches are important to distinguish on theoretical grounds, in  
335 practice there are other factors which are more important for, in particular, local  
336 volatility modeling. In line with classical inverse problems theory, we find that  
337 regularization of the local volatility is critical for backtesting performance.

338

### References.

- 339 [1] François Bachoc, Agnes Lagnoux, Andrés F López-Lopera, et al. Maximum  
340 likelihood estimation for Gaussian processes under inequality constraints. *Elec-*  
341 *tronic Journal of Statistics*, 13(2):2921–2969, 2019.
- 342 [2] Marc Chataigner, Stéphane Crépey, and Matthew Dixon. Deep local volatility.  
343 *Risks*, 8(3):82, 2020.
- 344 [3] Areski Cousin, Hassan Maatouk, and Didier Rullière. Kriging of financial term-  
345 structures. *European J. Oper. Res.*, 255(2):631–648, 2016.
- 346 [4] Stéphane Crépey and Matthew Dixon. Gaussian process regression for deriv-  
347 ative portfolio modeling and application to CVA computations. *Journal of*  
348 *Computational Finance*, 24(1):47–81, 2020.
- 349 [5] Bruno Dupire. Pricing with a smile. *Risk*, 7:18–20, 1994.
- 350 [6] Jim Gatheral. *The volatility surface: a practitioner’s guide*. Wiley, 2011.
- 351 [7] Jim Gatheral and Antoine Jacquier. Arbitrage-free SVI volatility surfaces.  
352 *Quantitative Finance*, 14(1):59–71, 2014.
- 353 [8] Mike Ludkovski and Yuri Saporito. Kri hedge: Gaussian process surrogates for  
354 delta hedging, 2020. arXiv:2010.08407.
- 355 [9] Hassan Maatouk and Xavier Bay. Gaussian process emulators for computer  
356 experiments with inequality constraints. *Math. Geosci.*, 49(5):557–582, 2017.
- 357 [10] K. Murphy. *Machine Learning: A Probabilistic Perspective*. MIT Press, 2012.
- 358 [11] Johannes Ruf and Weiguan Wang. Neural networks for option pricing and  
359 hedging: a literature review. *Journal of Computational Finance*, 24(1), 2020.
- 360 [12] Martin Tegnér and Stephen Roberts. A probabilistic approach to nonparametric  
361 local volatility. *arXiv preprint arXiv:1901.06021*, 2019.
- 362 [13] Yu Zheng, Yongxin Yang, and Bowei Chen. Gated neural networks for implied  
363 volatility surfaces, 2020. arXiv:1904.12834.