

1 Introduction

The purpose of this document is to detail the architecture of the universal pricing machine in the case of the local volatility model.

2 Local volatility model in a nutshell

Assuming that it pays no dividend and that the interest rates are zero, we find the following dynamics under the risk neutral probability measure \mathbb{P}^* :

$$\frac{dS_t}{S_t} = \sigma(S_t, t) dW_t \quad (1)$$

where W_t is the standard Brownian motion.

What matters for us is that, in that general diffusive model, one can relate the model parameter $\sigma(S, t)$ to the prices $C(K, T)$ of the call options of maturity T and strike price K :

$$\Lambda(K, T) = \frac{K}{2} \sigma^2(K, T) = \frac{\frac{\partial C(T, K)}{\partial T}}{K \frac{\partial^2 C(T, K)}{\partial K^2}} \quad (2)$$

That is going to be the corner stone of our machine learning approach.

3 Challenge of the pricing machine

The goal of the machine is to determine the probability density function of the stock $p_S(s, t)$ under the measure \mathbb{P}^* i.e. the probability of ending up in the vicinity of s at time t given that the stock was at S_0 at time 0.

That density is related to the option price by:

$$p_S(K, T) = \frac{\partial^2 C(K, T)}{\partial K^2} \quad (3)$$

For that, we must use as input the model parameter $\nu(K, T)$ and determine from it a solution meeting the Dupire's formula Eq.2

To complete the picture, we are going to consider that the solution can be represented by some Gaussian mixture. That is not really restrictive since any function can be decomposed as a series of Gaussian radial basis functions.

4 Gaussian radial basis functions

4.1 Expression

Gaussian radial basis functions are natural interpolating functions and, in the world of information theory, they correspond to a mixture of Gaussian variables. As the stock price S is log-normal, it is natural to introduce the log price X_t as:

$$X_t \triangleq \ln(S_t/S_0) \quad (4)$$

And we have:

$$p_X(x, T) = p_S(S_0 e^x, T) S_0 e^x \quad (5)$$

$$p_S(S, T) = p_X(\ln(S/S_0), T)/S \quad (6)$$

The Gaussian mixture model reads:

$$p_X(x, t) = \sum_{i=0}^I a_i(t) g(x - c_i(t), b_i(t)) \quad (7)$$

where

$$g(x, b) = \frac{1}{\sqrt{2\pi b}} e^{-\frac{x^2}{2b}} \quad (8)$$

The parameters $a_i(t) \geq 0$, $c_i(t), b_i(t) \geq 0$ are in turn functions of the time and $T_i(t)$ must be an increasing function. Those functions are typically going to be splines or coefficients of polynomials. We will see later on how to can formulated more explicitly in terms of neural network architecture.

4.2 Rationales of the choice

The reasons for that choice are manifold:

- If the local volatility surface ν is interpolated from a Gaussian radial basis function then it is clear from the Fokker Planck equation that once the solution is in that class it is going to stay within it (if we accept the idea that the number of centroids may grow over time). Clearly, the initial condition being a Delta distribution, only one single centroid is enough for small maturities.
- From the Chapman Kolmogorov identity, we can consider that if the probability measure at time T is approximated by some discretized measure (which all amounts to saying that the solutions can be quantized and remarking that discrete measures are degenerated cases of the Gaussian mixture), the probability density at the next step is of the form 7.
- It is conditionally Gaussian just as the local volatility model.

The interpretation of 7 in terms of financial market model is clear. The probabilities a_i are the discrete big picture scenarios under which the market behaves like a Black Scholes model of predefined volatility. Or, instead of saying that there is a single constant volatility coefficient as we do when plugging a Black Scholes model in, it all amounts to saying that there are I economics experts giving their views on the market and the model is the combination of those alternative market views, weighted by the credibility of each expert.

4.3 Black Scholes connection

4.4 General case

The Gaussian mixture being a mixture of Black Scholes models, we get for free the European option prices:

$$C(T, K) = \int_{\mathbb{R}} (S_0 e^{\xi} - K)_+ p_X(\xi, T) d\xi \quad (9)$$

$$= \sum_{i=0}^I a_i(T) \int_{\mathbb{R}} (S_0 e^{\xi} - K)_+ g(\xi - c_i(T), b_i(T)) d\xi \quad (10)$$

$$= \sum_{i=0}^I a_i(T) \int_{\mathbb{R}} (S_0 e^{c_i(T)} e^x - K)_+ g(\xi, b_i(T)) dx \quad (11)$$

$$= S_0 \sum_{i=0}^I a_i(T) e^{c_i(T)} \psi(b_i(T), e^{-c_i(T)} K / S_0) \quad (12)$$

where

$$\psi(b, \kappa) = e^{b/2} N(d^+) - \kappa N(d^-) \quad (13)$$

$$d^+ = d^- + \sqrt{b} \quad (14)$$

$$d^- = -\frac{\ln(\kappa)}{\sqrt{b}} \quad (15)$$

4.5 Martingale

We are going to focus on the martingale framework where:

$$c_i(t) = \underbrace{c_i(0)}_{\gamma_i} - \frac{b_i(t)}{2} \quad (16)$$

In that case, equations Eq.9 becomes:

$$C(T, K) = \sum_{i=0}^I a_i(T) \Omega_i(K, T) \quad (17)$$

$$\Omega_i(K, T) = S_0 e^{\gamma_i} N(d_i^+(T, K)) - K N(d_i^-(T, K)) \quad (18)$$

where

$$d_i^+(T, K) = d_i^-(T, K) + \sqrt{b_i(T)} \quad (19)$$

$$d_i^-(T, K) = \frac{\delta_i(K, T)}{\sqrt{b_i(T)}} \quad (20)$$

$$\delta_i(K, T) = \gamma_i + \ln(S_0/K) - b_i(T)/2 \quad (21)$$

We then have:

$$\frac{\partial C(K, T)}{\partial T} = \sum_{i=0}^I \frac{da_i(T)}{dT} \Omega_i(K, T) + a_i(T) \frac{\partial \Omega_i(K, T)}{\partial b_i} \frac{db_i(T)}{dT} \quad (22)$$

with:

$$\frac{\partial \Omega_i(K, T)}{\partial b_i} = \frac{K}{2} \underbrace{g(\delta_i(K, T), b_i(T))}_{g_i(K, T)} \quad (23)$$

4.6 Network architecture

Crucially, $a_i(t), b_i(t), c_i(t)$ must be the outputs of the neural network.

Contrary to the weights of the mixture and the width of the Gaussian kernels, the initial locations of the centroids $c_i(0)$ is not expected to be a causal process so we make it depend on the whole local volatility surface.

Hence, we have:

$$a_i(t) = A_i(t, \mathcal{W}_a(\mathbf{\Lambda}_{0,t})) \quad (24)$$

$$b'_i(t) = B_i(t, \mathcal{W}_b(\mathbf{\Lambda}_{0,t})) \quad (25)$$

$$c_i(0) = C_i(t, \mathcal{W}_c(\mathbf{\Lambda}_{0,T})) \quad (26)$$

where $\mathbf{\Lambda}_{t_1, t_2}$ is the collection of all the local volatilities $\Lambda(K, \tau)$ with τ between t_1 and t_2 .

It means that the probability of the stock to be at some value at some date t must not depend on the volatility of the stock after t .

The general idea would be then to find the parameters of the networks $\mathcal{W}_a, \mathcal{W}_b, \mathcal{W}_c$ minimizing the gap to Dupire identity:

$$E_{\mathcal{W}_{a,b,c}}^2 = \int_0^T \int_{K_{\min}}^{K_{\max}} \varepsilon(k, t)^2 dk dt \quad (27)$$

where

$$\varepsilon(k, t) = \frac{\partial C(k, t)}{\partial t} - \Lambda(k, t) k p_S(k, t) \quad (28)$$

and

$$k p_S(k, t) = \sum_{i=0}^I a_i(t) g(\ln(k/S_0) - c_i(t), b_i(t)) = \sum_{i=0}^I a_i(t) g_i(k, t) \quad (29)$$

Using Eq17, we get:

$$\varepsilon(k, t) = \sum_{i=0}^I \varepsilon_i(k, t) \quad (30)$$

where

$$\varepsilon_i(k, t) = \frac{da_i(t)}{dt} \Omega_i(k, t) + \frac{k}{2} a_i(t) g_i(k, t) \left(\frac{db_i(t)}{dt} - \sigma^2(k, t) \right) \quad (31)$$

5 Time Block approach

For clarification sake, we are going to slice the time domain into blocks and demand that the time dependence of each parameter be piecewise linear.

6 First Time block

To make it more concrete, let us now detail the first time block of the neural network architecture.

That is going to correspond to times between 0 and T_1 . T_1 could typically be 1 year.

The initial condition is a Dirac distribution in $X_0 = 1$.

$$b_i(0) = 0 \quad (32)$$

And

$$b_i(t) = tb_i(1) \quad (33)$$

The weights are going to be extrapolated flat for the first block:

$$a_i(0) = a_i(1) \quad (34)$$

We are then left with the following 2 vectors to predict:

$$\mathbf{a}(1) = [a_0(1), \dots, a_I(1)] \quad (35)$$

$$\mathbf{b}(1) = [b_0(1), \dots, b_I(1)] \quad (36)$$

So the only thing we have now to determine is:

$$\mathbf{a}(1) = \mathbf{A}(1, \mathcal{W}_a(\mathbf{\Lambda}_{0,T_1})) \triangleq \mathcal{W}_a^1(\mathbf{\Lambda}_{0,T_1}) \quad (37)$$

$$\mathbf{b}(1) = \mathbf{B}(1, \mathcal{W}_b(\mathbf{\Lambda}_{0,T_1})) \triangleq \mathcal{W}_b^1(\mathbf{\Lambda}_{0,T_1}) \quad (38)$$

with

$$\mathcal{W}_a^1(\mathbf{\Lambda}_{0,T_1}) \geq 0 \quad (39)$$

$$\sum_{i=0}^I \mathcal{W}_a^1(\mathbf{\Lambda}_{0,T_1})_i = 1 \quad (40)$$

Let us write:

$$\mathbf{\Lambda}_{0,T_1} = [\Lambda(k_m, t_n)]_{m \in \{0, \dots, M\}, n \in \{0, \dots, N\}} \quad (41)$$

The matching error given by Eq.27 now reads:

$$E_{\mathcal{W}_{a,b,c}^1}^2 = \sum_{m=1}^{M-1} \sum_{n=1}^N \varepsilon(k_m, t_n)^2 \Delta k_m \Delta t_n \quad (42)$$

$$\Delta k_m \triangleq k_m - k_{m-1} \quad (43)$$

$$\Delta t_n \triangleq t_n - t_{n-1} \quad (44)$$

In that case, we also have a simple expression for $\varepsilon(k_m, t_n)$:

$$\varepsilon(k, t) \triangleq \sum_{i=0}^I \varepsilon_i(k, t) \quad (45)$$

$$\varepsilon_i(k, t) = ka_i(1)g_i(k, t)(b_i(1) - \nu(k, t)) \quad (46)$$

7 Next time blocks

For $t \in [T_Y, T_{Y+1}]$, let us posit $\Delta_Y t = t - t_Y$. We generalize inductively as follows:

$$\mathbf{a}(t) = \mathbf{a}(T_Y) + \Delta_Y t (\mathbf{a}(T_{Y+1}) - \mathbf{a}(T_Y)) \quad (47)$$

$$\mathbf{b}(t) = \mathbf{b}(T_Y) + \Delta_Y t \mathbf{b}'(T_{Y+1}) \quad (48)$$

We must determine

$$\mathbf{a}(T_{Y+1}) = \mathcal{W}_a^{Y+1}(\mathbf{\Lambda}_{T_Y, T_{Y+1}}, \boldsymbol{\pi}_{T_Y}) \quad (49)$$

$$\mathbf{b}'(T_{Y+1}) = \mathcal{W}_b^{Y+1}(\mathbf{\Lambda}_{T_Y, T_{Y+1}}, \boldsymbol{\pi}_{T_Y}) \quad (50)$$

where $\boldsymbol{\pi}_{T_Y}$ is the representation of the Gaussian mixture at time T_Y :

$$\boldsymbol{\pi}_{T_Y} = [\mathbf{a}(T_Y), \mathbf{b}(T_Y)] \quad (51)$$

We must have:

$$\mathcal{W}_a^{Y+1}(\mathbf{\Lambda}_{T_Y, T_{Y+1}}, \boldsymbol{\pi}_{T_Y}) \geq 0 \quad (52)$$

$$\sum_{i=0}^I \mathcal{W}_a^{Y+1}(\mathbf{\Lambda}_{T_Y, T_{Y+1}}, \boldsymbol{\pi}_{T_Y})_i = 1 \quad (53)$$

$$\mathcal{W}_b^{Y+1}(\mathbf{\Lambda}_{T_Y, T_{Y+1}}, \boldsymbol{\pi}_{T_Y}) \geq 0 \quad (54)$$

As for the first block, we need to determine the parameters of each neural network to minimize Eq. 2 encoding the matching gap to Dupire's formula.

8 Initial centroids

The training of the initial centroids γ_i must be made at a higher level. The loss of each depends on the centroids so need to optimize that too. The initialization of the centroid must be rich enough. Taking them all equal to zero is probably not a good idea.

9 Remarks

- The time step within each block can be taken equal to one month.
- We might have to do something to prevent centroids c_i from cycling while moving from one time block to the next one. Penalizing large variations could then be considered.
- Topology matters for both that problem the concept of distance in the local volatility surface must be kept. The parameter b_i is both a time and the typical radius of interaction of a centroid. Said differently, I would expect the update of the parameters to be essentially driven by the local volatility in the region located within a distance b_i of the centroid c_i . I would even expect more precisely that b_i must depend on the average level of the local volatility while the change in a_i and c_i must rather be sensitive to the gradient of the local volatility.