# Solving The Dynamic Volatility Fitting Problem:

A Deep Reinforcement Learning Approach in Continuous State Action Spaces

Flow Equity Derivatives.

Emmanuel Gnabeyeu<sup>1,2</sup> and Omar Karkar <sup>2</sup>

October 14, 2025







<sup>&</sup>lt;sup>1</sup>Ecole Polytechnique, Paris, France.

<sup>&</sup>lt;sup>2</sup>EMEA Head of EQA, Citigroup Global Markets Limited, London, UK.

#### Table of Contents

- General Overview of the Volatility Fitting Problem
- 2 A Reinforcement Learning Framework and Algorithm
- 3 Numerical Results
- 4 Conclusion and Prospects:

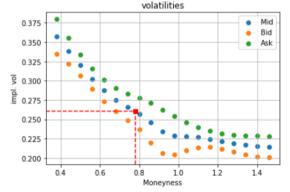
- General Overview of the Volatility Fitting Problem
  - Volatility Fitting: Problem Statement
  - Traditional approach for the Volatility Fitting Problem
- A Reinforcement Learning Framework and Algorithm
- 3 Numerical Results
- Conclusion and Prospects

# Volatility Fitting Problem and why?

Mark-to-market quoted price P:  $\left\{P_{\textit{Call/Put}}^{\textit{Ask}}, P_{\textit{Call/Put}}^{\textit{Bid}}\right\} \xrightarrow{\phi^{\mathsf{Model}}} \left\{\sigma_{\textit{Call/Put}}^{\textit{Ask}}, \sigma_{\textit{Call/Put}}^{\textit{Bid}}\right\} \xrightarrow{\phi} \left\{\sigma^{\textit{Ask}}, \sigma^{\textit{Bid}}\right\}$  for different moneyness  $(\kappa_1, ..., \kappa_n)$  where  $\phi$  is a non unique mapping and  $\phi^{\mathsf{Model}} = \mathsf{Black}\text{-Scholes}$  for Ex.

# Volatility Fitting Problem and why?

Mark-to-market quoted price P:  $\left\{P_{Call/Put}^{Ask}, P_{Call/Put}^{Bid}\right\} \xrightarrow{\phi^{\mathsf{Model}}} \left\{\sigma_{Call/Put}^{Ask}, \sigma_{Call/Put}^{Bid}\right\} \xrightarrow{\phi} \left\{\sigma_{Call/Put}^{Ask}, \sigma_{Call/Put}^{Bid}\right\} \xrightarrow{\phi} \left\{\sigma_{Call/Put}^{Ask}, \sigma_{Call/Put}^{Bid}\right\}$  for different moneyness  $(\kappa_1, ..., \kappa_n)$  where  $\phi$  is a non unique mapping and  $\phi^{\mathsf{Model}} = \mathsf{Black-Scholes}$  for Ex.



The volatility fitting problems aims to match the market implied volatility. It is used:

- To price structured products with volatility exposure.
- For risk management
- etc.

#### Encoding the market data:

Many methods in the industry to encode the implied volatility information <sup>3</sup>

- Non-Parametric grid of impl. Vol then interpolation via simple spline.
- Modelling the implied density of the asset directly ( Carr and Wu approach)
- Use a diffusion style model ( (L)SV models)

$$\Psi_{\vec{\theta}}^{SVI}(k)^2 = a + b(\rho(k-m) + \sqrt{(k-m)^2 + \sigma^2})$$

<sup>&</sup>lt;sup>3</sup>flexible and compatible with arbitrage constraints (call spread, butterfly and calendar spread)

<sup>&</sup>lt;sup>4</sup>Ex: Stochastic volatility inspired (SVI) parameterization by Gatheral and Antoine Jacquier:  $\vec{\theta} = (a, b, \rho, m, \sigma)$ 

#### Encoding the market data:

Many methods in the industry to encode the implied volatility information <sup>3</sup>

- Non-Parametric grid of impl. Vol then interpolation via simple spline.
- Modelling the implied density of the asset directly ( Carr and Wu approach)
- Use a diffusion style model ( (L)SV models)
- Surface-Level Parametrizations <sup>4</sup>
- $\rightarrow$  In this setting, we model the vol with a parametric functional  $\Psi^{vol}_{\vec{\theta}}(\kappa)$  for a given maturity.

$$\Psi_{\vec{\theta}}^{SVI}(k)^2 = a + b(\rho(k-m) + \sqrt{(k-m)^2 + \sigma^2})$$

<sup>&</sup>lt;sup>3</sup>flexible and compatible with arbitrage constraints (call spread, butterfly and calendar spread)

<sup>&</sup>lt;sup>4</sup>Ex: Stochastic volatility inspired (SVI) parameterization by Gatheral and Antoine Jacquier:  $\vec{\theta} = (a, b, \rho, m, \sigma)$ 

- General Overview of the Volatility Fitting Problem
  - Volatility Fitting: Problem Statement
  - Traditional approach for the Volatility Fitting Problem
- A Reinforcement Learning Framework and Algorithm
- 3 Numerical Results
- Conclusion and Prospects

# Traditional approach to volatility fitting with Parametrisation:

⇒ Choice of the target objective function:<sup>5</sup> : Among Others, the squared-error loss

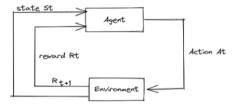
$$\xi(\vec{\theta}_{t_i}) := \sqrt{\frac{\sum_{j=1}^{n} w_j \ell(t_i, \kappa_j)^2}{\sum_{j=1}^{n} w_j}} \text{ with } \ell(t_i, \kappa_j) := \frac{\sigma^{\textit{Mid}}(t_i, \kappa_j) - \Psi^{\textit{vol}}_{\vec{\theta}_{t_i}}(\kappa_j)}{\sigma^{\textit{spread}}(t_i, \kappa_j)}, \ w_j = \textit{Black-Scholes Vega}$$
 (1)

 $\Rightarrow$  Given vol. surf. param.  $\Psi_{\vec{\theta}_{t_i}}$  and  $\xi(\vec{\theta}_{t_i})$ , adjust  $\vec{\theta}_{t_i}$  to match the mid market observable impl. vol.  $\sigma^{Mid}$  Drawbacks:

- $\bullet \ \ \, \text{Large set of conditions (liquidity, macro-events, term-structure, stability)} \rightarrow \text{optimizer too heavy to run}$
- Imply an optimisation at each move of the market (speed and reliability).
- Cannot re-use past accumulated experiences (every output is flashed out).
- Do not handle very well a case with market data limitations and irregularities.
- ⇒ We rely on a fully Al-based method to learn market regimes and discover optimal behaviours.

- General Overview of the Volatility Fitting Problem
- 2 A Reinforcement Learning Framework and Algorithm
  - Background on Reinforcement Learning
  - Problem Statement in Reinforcement Learning
  - Market environment
  - Deep Deterministic Policy Gradient Algorithm
- 3 Numerical Results
- 4 Conclusion and Prospects:

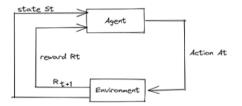
## Towards a Novel Framework: Background on Reinforcement Learning.



An agent reinforced with awards based on its interaction with the environment would learn to operate optimally in that environment to maximize the rewards.

• The agent observes a current state  $s_t \in \mathcal{S} \subset \mathbb{R}^d$ , takes an action  $a_t \in \mathcal{A} \subset \mathbb{R}^n$  drawn from the policy  $\Pi \supset \pi : \mathcal{S} \to \mathcal{A}$  (resp.  $\mathcal{P}(\mathcal{A})$ ) and receives a reward  $r_t = r(a_t, s_t)$ .

## Towards a Novel Framework: Background on Reinforcement Learning.



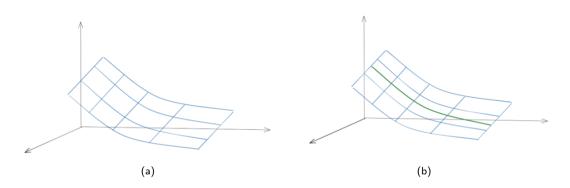
An agent reinforced with awards based on its interaction with the environment would learn to operate optimally in that environment to maximize the rewards.

- The agent observes a current state  $s_t \in \mathcal{S} \subset \mathbb{R}^d$ , takes an action  $a_t \in \mathcal{A} \subset \mathbb{R}^n$  drawn from the policy  $\Pi \supset \pi : \mathcal{S} \to \mathcal{A}$  (resp.  $\mathcal{P}(\mathcal{A})$ ) and receives a reward  $r_t = r(a_t, s_t)$ .
- Set  $\gamma \in [0,1]$ , define the Return from a state from following the policy  $\pi$ ,  $R_t = \sum_{i=t}^T \gamma^{(i-t)} r(s_i, a_i)$
- state-action value function  $Q^{\pi}(s, a) = \mathbb{E}^{\pi}[R_t|s_t = s, a_t = a]$ .
- Mathematical formulation of the problem:  $\pi^* = \arg\max_{\pi \in \Pi} Q^{\pi}$



- General Overview of the Volatility Fitting Problem
- 2 A Reinforcement Learning Framework and Algorithm
  - Background on Reinforcement Learning
  - Problem Statement in Reinforcement Learning
  - Market environment
  - Deep Deterministic Policy Gradient Algorithm
- 3 Numerical Results
- Conclusion and Prospects

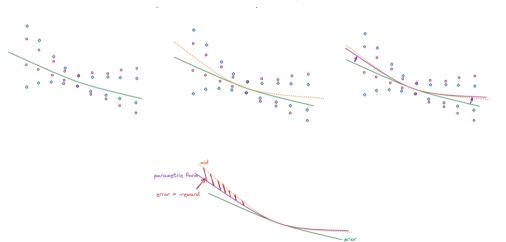
# Problem Statement in Reinforcement Learning:



- Fixed expiry T, and observe the market at time  $t_i$ , i.e. a collection  $\left\{\sigma^{Ask}_{Call/Put}(t_i,\kappa_j),\sigma^{Bid}_{Call/Put}(t_i,\kappa_j)\right\}_{j\in[1,...,n]}$  of markets quotes for different moneyness  $(\kappa_1,...,\kappa_n)$ .
- we consider a parametrisation  $\Psi^{vol}_{\vec{\theta}}$  of the volatility slice with a vector of K parameters.

# RL Framework for Fitting:

• Starting from a prior parametrized volatility surface  $\Psi^{vol}_{\vec{\theta}_r}$  ( i.e. the latest fit is the prior input).



# Goodness of Fit in RL-Based Volatility Fitting

• Mean squared error (MSE) <sup>6</sup> defined by :

$$\xi(\vec{\theta}_{t_i}) := \sum_{i=1}^n \left( \sigma^{Mid}(t_i, \kappa_j) - \Psi^{vol}_{\vec{\theta}_{t_i}}(k_i) \right)^2 \tag{3}$$

• The problem is to find the shifted vector  $\Delta \vec{\theta}_{t_i} = \left(\Delta \theta_{t_i}^1, ..., \Delta \theta_{t_i}^K\right)$  to bump the old parameters  $\vec{\theta}_{t_i}$  in order to maximize our selected reward function.  $\Delta \vec{\theta}_{t_i}^* = \arg\max_{\Delta \vec{\theta}_{t_i} \in \mathbb{R}^K} Q^{\pi}(s_{t_i}, \Delta \vec{\theta}_{t_i}).$ 

with 
$$Q^{\pi}(s_{t_i}, \Delta \vec{\theta}_{t_i}) = \sum_{t=t_i}^{T} \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} \left[ \gamma^{(t-t_i)} r(s_t, \Delta \vec{\theta}_t) \right]$$
 and  $\underline{r(s_t, a_t = \Delta \vec{\theta}_t) = -\xi(\theta_t)}$ 

<sup>6</sup>Alternative. Setting  $\sigma^{\text{spread}}(t_i, \kappa_j) := \sigma_{t_i, \kappa_i}^{Ask} - \sigma_{t_i, \kappa_i}^{Bid}$  the volatility spread, we define:

$$\xi(\vec{\theta}_{t_i}) := \sum_{i=1}^{n} \mathsf{vega}_{BS}(\kappa_j) \left( \frac{\sigma^{Mid}(t_i, \kappa_j) - \Psi^{vol}_{\vec{\theta}_{t_i}}(\kappa_j)}{\sigma^{\mathsf{spread}}(t_i, \kappa_j)} \right)^2 \tag{2}$$

- General Overview of the Volatility Fitting Problem
- 2 A Reinforcement Learning Framework and Algorithm
  - Background on Reinforcement Learning
  - Problem Statement in Reinforcement Learning
  - Market environment
  - Deep Deterministic Policy Gradient Algorithm
- 3 Numerical Results
- 4 Conclusion and Prospects:

#### Market environment and Non Linear functional approximation by FCNNs:

States are the collection of observable market quotes:

$$oldsymbol{s}_{t_{i+1}} = \left(\sigma^{\textit{Bid}}(t_{i+1}, \kappa_j), \sigma^{\textit{Ask}}(t_{i+1}, \kappa_j), heta_{t_i}^1, ..., heta_{t_i}^K
ight)_j$$

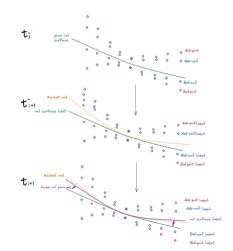
**Actions** are the bumps to apply in order to shift the prior volatility parametrization:  $a_{ti,t} = \Delta \vec{\theta}_{ti,t} = \vec{\theta}_{ti,t} - \vec{\theta}_{t}$ 

**1** Rewards are the opposite of well designed errors. 
$$r(s_{ti}, a_{ti}) = \Delta \vec{\theta}_{ti} = -\xi(\theta_{ti})$$

$$Q^{\pi}(s_{t_{i+1}}, a_{t_{i+1}}) = -\sum_{k=t_{i+1}}^{T} \mathbb{E}_{(s_k, a_k) \sim 
ho_{\pi}}[\gamma^{(k-t_{i+1})} \xi(\vec{ heta}_k)].$$

$$a_{t_{i+1}}^*(s_{t_{i+1}}) = \operatorname{arg\,max}_{a_{t_{i+1}} \in \mathcal{A}} \, Q^{\pi}(s_{t_{i+1}}, a_{t_{i+1}})$$

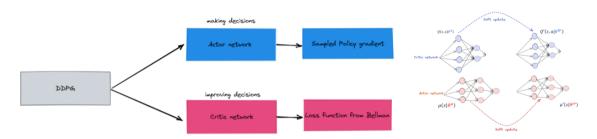
**Remark:**  $\Delta \vec{\theta}$  (the action vectors) can take any values in  $\mathbb{R}^K$  space:  $\Rightarrow$  continuous actions spaces.



- General Overview of the Volatility Fitting Problem
- 2 A Reinforcement Learning Framework and Algorithm
  - Background on Reinforcement Learning
  - Problem Statement in Reinforcement Learning
  - Market environment
  - Deep Deterministic Policy Gradient Algorithm
- 3 Numerical Results
- 4 Conclusion and Prospects:

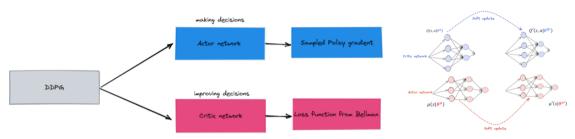
# Actor Critic: Deep Deterministic Policy Gradient.(Lillicrap et al., 2015)

• off-policy actor-critic method that learns a deterministic policy in continuous domain.



# Actor Critic: Deep Deterministic Policy Gradient.(Lillicrap et al., 2015)

• off-policy actor-critic method that learns a deterministic policy in continuous domain.

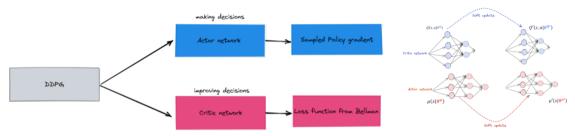


Action network: DPG= Deep Policy gradient

$$a_t \sim \pi^D(s_t, \theta^{\pi}) + \epsilon_t$$
 with  $\epsilon_t \sim \mathcal{N}(0, \sigma_n^2 I_K)$  and  $\sigma_n = \max(\sigma_0(1 - \frac{n}{N})^4, \sigma_{\min})$ 

# Actor Critic: Deep Deterministic Policy Gradient. (Lillicrap et al., 2015)

off-policy actor-critic method that learns a deterministic policy in continuous domain.



Action network: DPG= Deep Policy gradient

$$a_t \sim \pi^D(s_t, \theta^\pi) + \epsilon_t$$
 with  $\epsilon_t \sim \mathcal{N}(0, \sigma_n^2 I_K)$  and  $\sigma_n = \max(\sigma_0(1 - \frac{n}{N})^4, \sigma_{\min})$ 

• Critic network: Q-Learning and Bellman equation.

$$\begin{cases} R_t = \sum_{i=t}^T \gamma^{(i-t)} r(s_i, a_i) \\ Q^{\pi}(s_t, a_t) = \mathbb{E}[R_t | s_t, a_t] \end{cases} \Rightarrow L(\theta^Q) = \mathbb{E}[(Q^{\pi}(s_t, a_t; \theta^Q) - [r(s_t, a_t) + \gamma Q^{\pi}(s_{t+1}, a_{t+1}; \theta^Q)])^2]$$

- General Overview of the Volatility Fitting Problem
- 2 A Reinforcement Learning Framework and Algorithm
- 3 Numerical Results
  - Static Market, MSE Reward
  - Quasi Dynamic Scenario, MSE Reward
- 4 Conclusion and Prospects

# Numerical results: Convergence of DDPG Algorithm.

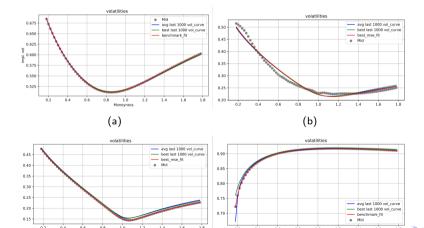
- To assess the performance of the RL algorithms for the volatility fitting, we consider toy problems of increasing complexity and perform a Monte Carlo on 5 differents random seeds:
  - Static scenario:
    - we explore the scenarios where the market conditions remain constant.

▶ Play Video



# Static Market, DDPG algorithm, MSE Reward:

Below, best response of the agent amongst the last 1000 episodes (*blue*) and the mean of the last 1000 volatilities slices ( *red*) (a) High Smile (b) Skew (c) High Skew (d) Inverse Smile.



# DDPG: Convergence and Consistency.

- Performance of the RL algorithms: Monte Carlo on 5 differents random seeds:
  - Static scenario:
    - RL agents perform as-good as traditional approaches in fitting volatility curve.

Table: DDPG MSE Rewards

type	Skew	High Smile	Inv. Smile
Bench	-0.011913	-0.0000220	-0.0000436
seeds' avg	-0.012145	-0.000163	-0.000315

- Final rewards achieved is close to the one coming from the optimizer.
- Performance is stable across the different market configuration considerered



- General Overview of the Volatility Fitting Problem
- 2 A Reinforcement Learning Framework and Algorithm
- 3 Numerical Results
  - Static Market, MSE Reward
  - Quasi Dynamic Scenario, MSE Reward
- Conclusion and Prospects

# Quasi Dynamic Scenario: The market evolve freely.

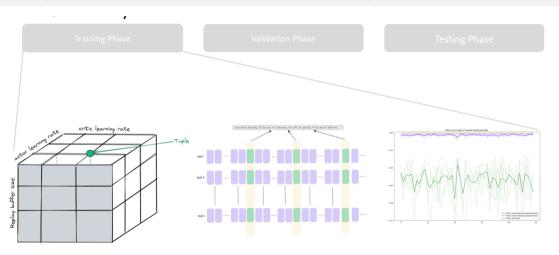


Figure: Training and Consistency over random seeds.



# Quasi Dynamic Scenario, MSE Reward

#### **2.Validation Phase:** We compare the performance of different successful agents under different seeds.

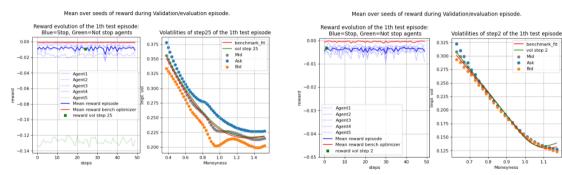


Figure: Evolution of mean episodic rewards(several for a wide spread stock with DDPG.

Figure: Evolution of mean episodic rewards(several random seeds) and Implied volatility(of one selected step)<sub>random seeds)</sub> and Implied volatility(of one selected step) for a tight spread stock with DDPG.

# Quasi Dynamic Scenario, MSE Reward, Testing Phase

**3.Testing Phase:** we assess the performance of the best agent in a test episode.

▶ Play Video

Vol fitting animation

## Conclusion and Prospects

- Sum up,
  - A wide range of problems can be tackled with Deep RL Actor-Critic methods while neural networks help to approximate the solution function.<sup>7</sup>
  - We showed that DRL algorithms are natively tailored for the volatility fitting problem:
    - native exploration
    - effective catalog of past experiences(re-use past accumulated experiences)
    - handle complex objective function (desk PnL).

## Conclusion and Prospects

- Sum up,
  - A wide range of problems can be tackled with Deep RL Actor-Critic methods while neural networks help to approximate the solution function.<sup>7</sup>
  - We showed that DRL algorithms are natively tailored for the volatility fitting problem:
    - native exploration
    - effective catalog of past experiences(re-use past accumulated experiences)
    - handle complex objective function (desk PnL).

- Euture work could be:
  - Reflects stylistic effects in volatility surface term-structure.
  - Careful rewards shaping and environment design.

<sup>&</sup>lt;sup>7</sup>optimal execution, portfolio optimization, option pricing and hedging, market making, smart<sub>f</sub>order routing

# Thanks for your attention!



#### References.

- B. Hambly, R. Xu, and H. Yang, Recent Advances in Reinforcement Learning in Finance.
- J. Gatheral, A. Jacquier. Arbitrage-Free SVI Volatility Surfaces
- V. Zetocha, Sculpting implied volatility surfaces of illiquid assets
- Sutton, R. S. and Barto, A. G, Reinforcement learning: An introduction.
- T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, Continuous control with deep reinforcement learning
- R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, Policy gradient methods for reinforcement learning with function approximation