

# Reinforcement Learning for Finance:

## Discrete-time Black-Scholes model

Igor Halperin

NYU Tandon School of Engineering

Brooklyn NY 2017

1

---

<sup>1</sup>Opinions presented in these notes are author's only, and not necessarily of his employer. The standard disclaimer applies.

There is nothing more  
practical than a good theory.

---

V. Vapnik (?)

# In these notes:

- **Can a RL agent learn to trade an option from *directly from trading data*?**
- **Plan:**
  - Take a discrete-time version of the Black-Scholes model for option pricing
  - Convert it to a Markov Decision Process (MPD) model
  - Apply Reinforcement Learning methods to this model, using simulated data

# Discrete-time BSM model

We start with a discrete-time version of the BSM model. The problem of option hedging and pricing in this formulation amounts to a sequential risk minimization. To define risk in an option, we follow a local risk minimization approach pioneered in the work of Föllmer and Schweizer, Schweizer (1994). A similar method was developed by physicists Potters and Bouchaud (2001), see also the work by Kapoor et. al. (2010). We use a version of this approach suggested in a Ph.D. thesis by Grau (2007). In this approach, we take the view of a seller of a European option (e.g. a put option) with maturity  $T$  and the terminal payoff of  $H_T(S_T)$  at maturity, that depends on the final stock price  $S_T$  at that time. To hedge the option, the seller use the proceeds of the sale to set up a replicating (hedge) portfolio  $\Pi_t$  made of the stock  $S_t$  and a risk-free bank deposit  $B_t$ . The value of hedge portfolio at any time  $t \leq T$  is

$$\Pi_t = u_t S_t + B_t \quad (1)$$

where  $u_t$  is a stock position at time  $t$ , taken to hedge risk in the option.

# Hedge portfolio evaluation

The replicating portfolio tries to exactly match the option price in all possible future states of the world. Start at  $t = T$ : the option position is closed, the hedge  $u_t$  is closed ( $u_t = 0$ ), and therefore

$$\Pi_T = B_T = H_T(S_T) \quad (2)$$

This sets a terminal condition for  $B_T$  for all future states of the world at  $T$ . To find  $B_t$  for previous times  $t < T$ , we impose the *self-financing constraint* which requires that all future changes in the hedge portfolio should be funded from an initially set bank account, without any cash infusions or withdrawals over the lifetime of the option.

$$u_t S_{t+1} + e^{r\Delta t} B_t = u_{t+1} S_{t+1} + B_{t+1} \quad (3)$$

This can be expressed as a recursive relation for  $B_t$  at any time  $t < T$  using its value at the next time instance:

$$B_t = e^{-r\Delta t} [B_{t+1} + (u_{t+1} - u_t) S_{t+1}] , \quad t = T - 1, \dots, 0 \quad (4)$$

# Hedge portfolio evaluation

Plugging this into Eq.(1) produces a recursive relation for  $\Pi_t$  in terms of its values at later times, which can therefore be solved backward in time, starting from  $t = T$  with the terminal condition (2), and continued all the way to the current time  $t = 0$ :

$$\Pi_t = e^{-r\Delta t} [\Pi_{t+1} - u_t \Delta S_t], \quad \Delta S_t = S_{t+1} - e^{r\Delta t} S_t, \quad t = T-1, \dots, 0 \quad (5)$$

Eqs.(4) and (5) imply that both  $B_t$  and  $\Pi_t$  are not measurable at any  $t < T$ , as they depend on the future. Respectively, their values today  $B_0$  and  $\Pi_0$  will be random quantities with some distributions. We can compute them using Monte Carlo!

# Hedge portfolio evaluation with Monte Carlo

For any given hedging strategy  $\{u_t\}_{t=0}^T$ , these distributions can be estimated using Monte Carlo simulation:

- **Forward pass** : Simulate  $N$  paths of the underlying  $S_1 \rightarrow S_2 \rightarrow \dots \rightarrow S_N$ ,
- **Backward pass**: Evaluate  $\Pi_t$  going backward on each path.

As the choice of a hedge strategy does not affect the evolution of the underlying, such simulation of forward paths should only be performed once. Alternatively, we could use real historical data for stock prices, together with a pre-determined hedging strategy  $\{u_t\}_{t=0}^T$  and a terminal condition (2).

But first, we need a hedge strategy  $u_t$  to implement this Monte Carlo!

# Optimal hedging in the discrete-time BSM model

The optimal hedge  $u^*(S_t)$  in this model is obtained from the requirement that the variance of  $\Pi_t$  across all simulated MC paths at time  $t$  is minimized when conditioned on the currently available *cross-sectional* information  $\mathcal{F}_t$ , i.e.

$$\begin{aligned} u_t^*(S_t) &= \arg \min_u \text{Var} [\Pi_t | \mathcal{F}_t] \\ &= \arg \min_u \text{Var} [\Pi_{t+1} - u_t \Delta S_t | \mathcal{F}_t] , \quad t = T-1, \dots, 0 \quad (6) \end{aligned}$$

The first expression in Eq.(7) implies that all uncertainty in  $\Pi_t$  is due to uncertainty regarding the amount  $B_t$  needed to be held at the bank account at time  $t$  in order to be able to cover the future obligations at the option maturity  $T$ . This means that an optimal hedge should minimize the cost of hedge capital at each time step  $t$ .



# Optimal hedging in the discrete-time BSM model

The optimal hedge  $u^*(S_t)$  :

$$\begin{aligned} u_t^*(S_t) &= \arg \min_u \text{Var} [\Pi_t | \mathcal{F}_t] \\ &= \arg \min_u \text{Var} [\Pi_{t+1} - u_t \Delta S_t | \mathcal{F}_t] , \quad t = T-1, \dots, 0 \end{aligned} \quad (7)$$

The optimal hedge can be found analytically by setting the derivative of (7) to zero. This gives

$$u_t^*(S_t) = \frac{\text{Cov} (\Pi_{t+1}, \Delta S_t | \mathcal{F}_t)}{\text{Var} (\Delta S_t | \mathcal{F}_t)} , \quad t = T-1, \dots, 0 \quad (8)$$

How one compute one-step expectations depends on whether we deal with a continuous or a discrete state space.

We will use a general notation as in Eq.(8) to denote similar conditional expectations where  $\mathcal{F}_t$  stands for cross-sectional information set at time  $t$ .

# Option pricing in discrete-time BSM

A *mean* option price  $\hat{C}_t$  is defined as a time- $t$  expected value of the hedge portfolio  $\Pi_t$ :

$$\hat{C}_t = \mathbb{E}_t [\Pi_t | \mathcal{F}_t] \quad (9)$$

Using the equation for  $\Pi_t$  and the tower law of conditional expectations, we obtain

$$\begin{aligned} \hat{C}_t &= \mathbb{E}_t \left[ e^{-r\Delta t} \Pi_{t+1} \middle| \mathcal{F}_t \right] - u_t(S_t) \mathbb{E}_t [\Delta S_t | \mathcal{F}_t] \\ &= \mathbb{E}_t \left[ e^{-r\Delta t} \mathbb{E}_{t+1} [\Pi_{t+1} | \mathcal{F}_{t+1}] \middle| \mathcal{F}_t \right] - u_t(S_t) \mathbb{E}_t [\Delta S_t | \mathcal{F}_t] \quad (10) \\ &= \mathbb{E}_t \left[ e^{-r\Delta t} \hat{C}_{t+1} \middle| \mathcal{F}_t \right] - u_t(S_t) \mathbb{E}_t [\Delta S_t | \mathcal{F}_t], \quad t = T-1, \dots, 0 \end{aligned}$$

# Fair Option pricing in discrete-time BSM

The dealer cannot just ask the mean price  $\hat{C}_0$  when selling the option, as she has to compensate for risk of exhausting the bank account  $B_t$  some time in the future, which would require cash infusions into the hedge portfolio, after any fixed amount  $\hat{B}_0 = \mathbb{E}_0[B_0]$  is put in the bank account at time  $t = 0$  right after selling the option.

One possible specification of a risk premium is to add the cumulative expected discounted variance of the hedge portfolio along all time steps  $t = 0, \dots, N$ , with a risk-aversion parameter  $\lambda$ :

$$C_0^{(ask)}(S, u) = \mathbb{E}_0 \left[ \Pi_0 + \lambda \sum_{t=0}^T e^{-rt} \text{Var} [\Pi_t | \mathcal{F}_t] \middle| S_0 = S, u_0 = u \right] \quad (11)$$

The option seller should *minimize* this option price to be competitive.

# Maximization problem for option pricing

The problem of *minimization* of a fair (to the dealer) option price (11) can be equivalently expressed as the problem of *maximization* of its negative  $V_t = -C_t^{(ask)}$ , where

$$V_t(S_t) = \mathbb{E}_t \left[ -\Pi_t - \lambda \sum_{t'=t}^T e^{-r(t'-t)} \text{Var} [\Pi_{t'} | \mathcal{F}_{t'}] \middle| \mathcal{F}_t \right] \quad (12)$$

# Hedging and pricing in the BS limit

The framework presented above provides a smooth transition to the strict BS limit  $\Delta t \rightarrow 0$ . In this limit, the BSM model dynamics under the physical measure  $\mathbb{P}$  is described by a continuous-time Geometric Brownian motion with a drift  $\mu$  and volatility  $\sigma$ :

$$\frac{dS_t}{S_t} = \mu dt + \sigma dW_t \quad (13)$$

where  $W_t$  is a standard Brownian motion.

# Hedging in the BS limit

Consider first the optimal hedge strategy in the BS limit  $\Delta t \rightarrow 0$ . Using the first-order Taylor expansion

$$\hat{C}_{t+1} = C_t + \frac{\partial C_t}{\partial S_t} \Delta S_t + O(\Delta t) \quad (14)$$

in (??), we obtain

$$u_t^{BS}(S_t) = \lim_{\Delta t \rightarrow 0} u_t^*(S_t) = \frac{\partial C_t}{\partial S_t} \quad (15)$$

which is the correct optimal hedge in the continuous-time BSM model.

# Pricing in the BS limit

First compute the limit of the second term in the equation for  $\hat{C}_t$ :

$$\lim_{\Delta t \rightarrow 0} u_t(S_t) \mathbb{E}_t [\Delta S_t | \mathcal{F}_t] = \lim_{dt \rightarrow 0} u_t^{BS} S_t (\mu - r) dt = \lim_{dt \rightarrow 0} (\mu - r) S_t \frac{\partial C_t}{\partial S_t} dt$$

To evaluate the first term in the equation for  $\hat{C}_t$ , we use the second-order Taylor expansion:

$$\begin{aligned} \hat{C}_{t+1} &= C_t + \frac{\partial C_t}{\partial t} dt + \frac{\partial C_t}{\partial S_t} dS_t + \frac{1}{2} \frac{\partial^2 C_t}{\partial S_t^2} (dS_t)^2 + \dots \\ &= C_t + \frac{\partial C_t}{\partial t} dt + \frac{\partial C_t}{\partial S_t} S_t (\mu dt + \sigma dW_t) \\ &\quad + \frac{1}{2} \frac{\partial^2 C_t}{\partial S_t^2} S_t^2 (\sigma^2 dW_t^2 + 2\mu\sigma dW_t dt) + O(dt^2) \end{aligned}$$

# Pricing in the BS limit

We had the recursive formula for the option price

$$\hat{C}_t = \mathbb{E}_t \left[ e^{-r\Delta t} \hat{C}_{t+1} \middle| \mathcal{F}_t \right] - u_t(S_t) \mathbb{E}_t [\Delta S_t | \mathcal{F}_t], \quad t = T-1, \dots, 0$$

Use Eqs.(15) and (16) for two terms in this expression, use  $\mathbb{E}[dW_t] = 0$  and  $\mathbb{E}[dW_t^2] = dt$ , and simplify. We see that the stock drift  $\mu$  drops out from the problem.

This becomes the celebrated **Black-Scholes equation** in the limit  $dt \rightarrow 0$ :

$$\frac{\partial C_t}{\partial t} + rS_t \frac{\partial C_t}{\partial S_t} + \frac{1}{2} \sigma^2 S_t^2 \frac{\partial^2 C_t}{\partial S_t^2} - rC_t = 0 \quad (16)$$

Therefore, if the world is lognormal, both our hedging and pricing formulae become the original formulae of the Black-Scholes-Merton model in the strict limit  $\Delta t \rightarrow 0$ .



# Markov Decision Process model for discrete-time BS model

Now we will re-formulate the discrete-time BSM model as a Markov Decision Processes (MDP) model.

Our discrete-time MDP model works directly in the physical measure  $\mathbb{P}$ , and addresses the problem of option pricing and hedging as a problem of stochastic optimal control in discrete time, where the system being controlled is a hedge portfolio, and control is a stock position in this hedge portfolio. The problem is then solved by a sequential maximization of "rewards" (negatives of hedge portfolio one-step variances times the risk-aversion  $\lambda$ , plus a drift term).

Two cases are possible:

- The model is *known*
- The model is *unknown*

Respectively, we will have *two ways* to solve a Bellman equation for these two cases.

# State variables

We first define a new variable  $X_t$  by the following relation:

$$X_t = - \left( \mu - \frac{\sigma^2}{2} \right) t + \log S_t \quad (17)$$

This implies that

$$dX_t = - \left( \mu - \frac{\sigma^2}{2} \right) dt + d \log S_t = \sigma dW_t \quad (18)$$

Therefore,  $X_t$  is a standard Brownian motion, scaled by volatility  $\sigma$ . For a given  $X_t$  in a MC scenario, the corresponding value of  $S_t$  is

$$S_t = e^{X_t + \left( \mu - \frac{\sigma^2}{2} \right) t} \quad (19)$$

As  $X_t$  is a martingale, i.e.  $\mathbb{E}[dX_t] = 0$ , on average it should not run too far away from  $X_0$  during the lifetime of an option. The state variable  $X_t$  is time-uniform, unlike the stock price  $S_t$  that has a drift. But Eq.(19) can always be used in order to map the dynamics of  $S_t$  into the dynamics of  $X_t$ .

# Finite-state approximation

We can discretize the set of admissible values of state variables  $X_t$  defined in Eq.(17), while keeping the relation (19) that expresses the stock price  $S_t$  in terms of the (now discretized) state variable  $X_t$ .

Can use a discrete-time, discrete-state Markov Chain model as an approximation to a continuous-state BSM dynamics (Duan and Simonato, 2001).

A state-space discretization is *not* necessary within our framework. Below we assume a general continuous-state case in a Monte Carlo setting.

Simplifications due to discretization:

- Can use simple algorithms for discrete-state MDPs (such as Q-Learning)
- Discrete state formulation simplifies calculation of one-step conditional expectations

# Value function

Now we reformulate our risk minimization procedure in a language of MDP problems.

Express the dynamics in terms of variables  $X_t$  using Eq.(19). Actions  $u_t = u_t(S_t)$  in terms of stock prices are then obtained by the substitution

$$u_t(S_t) = a_t(X_t(S_t)) = a_t\left(\log S_t - \left(\mu - \frac{\sigma^2}{2}\right)t\right) \quad (20)$$

Actual hedging decisions  $a_t(x_t)$  are determined by a time-dependent *policy*  $\pi(t, X_t)$ . We consider *deterministic policies*, i.e.

$$\pi : \{0, \dots, T-1\} \times \mathcal{X} \rightarrow \mathcal{A} \quad (21)$$

This deterministic policy maps the time  $t$  and the current state  $X_t = x_t$  into the action  $a_t \in \mathcal{A}$ :

$$a_t = \pi(t, x_t) \quad (22)$$

# Bellman equation

First re-write the value maximization problem of Eq.(12) in terms of a new state variable  $X_t$ , and with an upper index to denote its dependence on the policy  $\pi$ :

$$\begin{aligned} V_t^\pi(X_t) &= \mathbb{E}_t \left[ -\Pi_t(X_t) - \lambda \sum_{t'=t}^T e^{-r(t'-t)} \text{Var} [\Pi_{t'}(X_{t'}) | \mathcal{F}_{t'}] \middle| \mathcal{F}_t \right] \\ &= \mathbb{E}_t \left[ -\Pi_t(X_t) - \lambda \text{Var} [\Pi_t] - \lambda \sum_{t'=t+1}^T e^{-r(t'-t)} \text{Var} [\Pi_{t'}(X_{t'}) | \mathcal{F}_{t'}] \middle| \mathcal{F}_t \right] \end{aligned} \quad (23)$$

The last term in this expression can be expressed in terms of  $V_{t+1}$  using the definition of the value function with a shifted time argument:

$$- \lambda \mathbb{E}_{t+1} \left[ \sum_{t'=t+1}^T e^{-r(t'-t)} \text{Var} [\Pi_{t'} | \mathcal{F}_{t'}] \right] = \gamma (V_{t+1} + \mathbb{E}_{t+1} [\Pi_{t+1}]) \quad (24)$$

Here  $\gamma \equiv e^{-r\Delta t}$  is a discrete-time discount factor

# Bellman equation

Substitute this into (23), re-arrange terms, and use the portfolio process Eq.(5).

This produces the Bellman equation for our model:

$$V_t^\pi(X_t) = \mathbb{E}_t^\pi [R(X_t, a_t, X_{t+1}) + \gamma V_{t+1}^\pi(X_{t+1})] \quad (25)$$

Here  $R(X_t, a_t, X_{t+1})$  is a one-step time-dependent random reward

$$\begin{aligned} R_t(X_t, a_t, X_{t+1}) &= \gamma a_t \Delta S_t(X_t, X_{t+1}) - \lambda \text{Var}[\Pi_t | \mathcal{F}_t] \quad (26) \\ &= \gamma a_t \Delta S_t(X_t, X_{t+1}) - \lambda \gamma^2 \mathbb{E}_t \left[ \hat{\Pi}_{t+1}^2 - 2a_t \Delta \hat{S}_t \hat{\Pi}_{t+1} + a_t^2 (\Delta \hat{S}_t)^2 \right] \end{aligned}$$

where we used Eq.(5) in the second line, and  $\hat{\Pi}_{t+1} \equiv \Pi_{t+1} - \bar{\Pi}_{t+1}$ , where  $\bar{\Pi}_{t+1}$  is the sample mean of all values of  $\Pi_{t+1}$ , and similarly for  $\Delta \hat{S}_t$ . For  $t = T$ , we have  $R_T = -\lambda \text{Var}[\Pi_T]$  where  $\Pi_T$  is determined by the terminal condition (2).

## Expected rewards

Note that Eq.(26) implies that the expected reward  $R_t$  at time step  $t$  is *quadratic* in the action variable  $a_t$ :

$$\begin{aligned}\mathbb{E}_t [R_t (X_t, a_t, X_{t+1})] &= \gamma a_t \mathbb{E}_t [\Delta S_t] \\ &\quad - \lambda \gamma^2 \mathbb{E}_t [\hat{\Pi}_{t+1}^2 - 2a_t \Delta \hat{S}_t \hat{\Pi}_{t+1} + a_t^2 (\Delta \hat{S}_t)^2]\end{aligned}\tag{27}$$

As we will see, this simple quadratic dependence on  $a_t$  is very useful for a solution of the MDP dynamics in this model.

Note: when  $\lambda \rightarrow 0$ , the expected reward is *linear* in  $a_t$ , so it does *not* have a maximum when  $\lambda \rightarrow 0$  (i.e. there is no risk aversion).

In our framework, quadratic risk is incorporated in a standard (risk-neutral) MDP formulation.

# Bellman optimality equation

The *optimal policy*  $\pi_t^*(\cdot|X_t)$  is determined as a policy that maximizes the value function  $V_t^\pi(X_t)$ :

$$\pi_t^*(X_t) = \arg \max_{\pi} V_t^\pi(X_t) \quad (28)$$

The optimal value function satisfies the Bellman optimality equation

$$V_t^*(X_t) = \mathbb{E}_t^{\pi^*} [R_t(X_t, u_t = \pi_t^*(X_t), X_{t+1}) + \gamma V_{t+1}^*(X_{t+1})] \quad (29)$$

If the system dynamics are *known*, the Bellman optimality equation can be solved using methods of Dynamic Programming such as Value Iteration. If dynamics are *unknown*, the optimal policy should be computed using *samples*. This is a setting of Reinforcement Learning. A formalism based on an action-value function provides a better framework for a RL setting.