
000
001
002
003
004
005
006
007 **BSMRL: A RISK-SENSITIVE AND TIME-AWARE**
008 **FRAMEWORK FOR REINFORCEMENT LEARNING**
009 **(DRAFT VER.)**

010
011 **Anonymous authors**
012 Paper under double-blind review

013 **ABSTRACT**
014

015 In this paper, we propose BSMRL, a novel reinforcement learning frame-
016 work that leverages the Black-Scholes-Merton (BSM) model for reward
017 shaping. By integrating financial mathematics into the RL paradigm,
018 BSMRL aims to enhance learning efficiency and stability in environments
019 with sparse and delayed rewards. We further extend our approach by incor-
020 porating the Merton Jump-Diffusion Model to account for sudden changes
021 in the environment, providing a risk-sensitive and time-aware framework.
022 Experimental results demonstrate that BSMRL outperforms traditional RL
023 methods in various benchmark tasks, showcasing its potential for broader
024 applications in software testing, quantitative finance, embodied AI, and
025 beyond.

026 **CONTENTS**
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

054 1 INTRODUCTION
055

056 In the original Temporal-Difference (TD) learning framework, the agent learns to estimate
057 the value function based on one-step transitions. However, in real-world scenarios, rewards
058 can be sparse and delayed, making it challenging for agents to learn effective policies.
059

060 Here, we propose BSMRL, a novel reinforcement learning framework that leverages the
061 Black-Scholes-Merton (BSM) model for reward shaping. By integrating financial mathe-
062 matics into the RL paradigm, BSMRL aims to enhance learning efficiency and stability in
063 environments with sparse and delayed rewards. The BSM model provides a closed-form
064 solution for European-style options, which can be adapted to shape rewards based on the
065 agent's current state and time-to-go. This approach provides more informative feedback to
066 the agent, guiding its learning process. We treat the option price as a advanced exploration
067 bonus for the agent, which encourages it to explore states randomly and discover potential
068 rewards, even when immediate rewards are sparse.

069 However, we note that the BSM model assumes that the noise in the environment is Gaus-
070 sian, which may not hold when the environment exhibits sudden changes or jumps. To ad-
071 dress this limitation, we extend our approach by incorporating the Merton Jump-Diffusion
072 Model, which accounts for jumps in the underlying asset price dynamics. This method how-
073 ever, lead to the uncertainty in the convergence of the value estimates, thus we introduce an
074 adapted decaying mechanism for the jump intensity parameter in the Merton model, which
075 helps stabilize the learning process, guaranteeing the convergence of the value estimates.
076

077 2 RELATED WORKS
078

079 3 BSMRL
080

081 3.1 BLACK-SCHOLES-MERTON MODEL
082

083 In the field of financial mathematics, the Black-Scholes-Merton (BSM) model?? is a founda-
084 tional framework for (European) option pricing. It assumes that the price of the underlying
085 asset follows a geometric Brownian motion with constant volatility and drift. The BSM
086 model provides a closed-form solution for European-style options.

087 They derived a partial differential equation (PDE) that the option price must satisfy, known
088 as the Black-Scholes equation.

089 **Theorem 3.1** (Black-Scholes Equation). *The price of a European call option $C(S, t)$ on a
090 non-dividend-paying stock satisfies the following PDE:*

091
$$\frac{\partial C}{\partial t} + \frac{1}{2}\sigma^2 S^2 \frac{\partial^2 C}{\partial S^2} + rS \frac{\partial C}{\partial S} - rC = 0 \quad (1)$$

092

093 where $C(S, t)$ is the price of the option at time t when the underlying asset price is S , σ is
094 the volatility of the underlying asset, and r is the risk-free interest rate.
095

096 By applying Ito's Lemma and constructing a riskless portfolio, they eliminated the stochastic
097 component and derived the pricing formula for European call options.
098

099 **Theorem 3.2** (Black-Scholes-Merton Formula). *The price of a European call option $C(S_t, t)$
100 on a non-dividend-paying stock is given by:*

101
$$C(S_t, t) = S_t \Phi(d_1) - K e^{-r(T-t)} \Phi(d_2) \quad (2)$$

102

103 where:
104

$$d_1 = \frac{\ln(S_t/K) + (r + \frac{1}{2}\sigma^2)(T - t)}{\sigma\sqrt{T - t}}, \quad d_2 = d_1 - \sigma\sqrt{T - t} \quad (3)$$

105 Here, $C(S_t, t)$ is the price of a European call option at time t , S_t is the current price of
106 the underlying asset, K is the strike price, r is the risk-free interest rate, σ is the volatility
107 of the underlying asset, T is the time to maturity, and $\Phi(\cdot)$ is the cumulative distribution
108 function of the standard normal distribution.
109

108 **Rationale** In Reinforcement Learning, the agent interacts with an environment to maximize cumulative rewards. However, in real-world scenarios, rewards can be sparse and delayed, making it challenging for agents to learn effective policies, leading to high variance in value estimates and slow convergence. To address this, we propose using the BSM model to shape rewards based on the agent’s current state and time-to-go, providing more informative feedback and guiding the agent’s learning process.

114 In a RL task, we have to estimate the value function $V(s)$ or action-value function $Q(s, a)$,
115 usually within a certain period of time. For instance, in Monte-Carlo methods, we estimate
116 the expected return over an episode, while in Temporal-Difference (TD) learning, we estimate
117 the value function based on one-step transitions. This is analogous to option pricing,
118 where the option’s value depends on the underlying asset price and time to maturity. These
119 method often suffer from slow convergence rate due to sparse rewards, e.g., in website bug
120 mining, the agent only receives a reward when a bug is found, which may take a long time.

121 By mimicing the idea of option pricing, we use the BSM formula and its derivatives to shape
122 the expected accumulated reward, providing power for the agent to discover although the
123 immediate reward is sparse. The method will provide a time-sensitive potential shaping
124 reward, which can be viewed as an additional reward signal that guides the agent towards
125 more promising states and actions.

126 However, such model assumes that the noise in the environment is Gaussian, which may not
127 hold when the environment exhibits sudden changes or jumps. To address this limitation,
128 we extend our approach by incorporating the Merton Jump-Diffusion Model, which accounts
129 for jumps in the underlying asset price dynamics.

131 3.2 MERTON JUMP MODEL

133 In financial markets, asset prices often exhibit sudden and significant changes, known as
134 jumps, which cannot be captured by the standard Black-Scholes model. To address this
135 limitation, Robert C. Merton extended the Black-Scholes framework by incorporating jump
136 processes into the asset price dynamics, leading to the Merton Jump-Diffusion Model. The
137 Merton Jump-Diffusion Model assumes that the underlying asset price follows a stochastic
138 process that combines both continuous diffusion and discrete jumps. The asset price dynamics
139 under the Merton model can be described by the following stochastic differential equation (SDE):

$$140 \quad dS_t = \mu S_t dt + \sigma S_t dW_t + J_t S_t dN_t \quad (4)$$

141 where:

- 143 • S_t is the asset price at time t .
- 144 • μ is the drift rate of the asset price.
- 145 • σ is the volatility of the continuous component.
- 146 • W_t is a standard Brownian motion.
- 147 • N_t is a Poisson process with intensity λ , representing the number of jumps up to time t .
- 148 • J_t is the jump size, typically modeled as a log-normal random variable.

151 **Remark 3.1.** Such model can be also viewed as a Levy process, which generalizes Brownian motion by allowing for jumps.

154 The Merton Jump-Diffusion Model leads to a modified option pricing formula that accounts
155 for the possibility of jumps in the underlying asset price.

156 **Theorem 3.3** (Merton Jump-Diffusion Option Pricing Formula). *The price of a European call option $C(S_t, t)$ under the Merton Jump-Diffusion Model is given by:*

$$159 \quad C(S_t, t) = \sum_{n=0}^{\infty} \frac{e^{-\lambda(T-t)} (\lambda(T-t))^n}{n!} C_{BS}(S_t, t; \sigma_n) \quad (5)$$

161 where:

-
- 162 • $C_{BS}(S_t, t; \sigma_n)$ is the Black-Scholes price of the option with adjusted volatility $\sigma_n =$
 163 $\sqrt{\sigma^2 + \frac{n\delta^2}{T-t}}$, where δ is the standard deviation of the jump size.
 164
 165 • λ is the jump intensity.
 166
 167 • T is the time to maturity.
 168
 169 • n is the number of jumps.
 170
 171

172
 173 **Rationale** While the Black-Scholes model assumes continuous price movements, real-world environments often exhibit sudden changes or jumps, leading to fat-tailed reward distributions.
 174
 175

176 To better capture these dynamics, we propose using the Merton Jump-Diffusion Model for
 177 reward shaping in Reinforcement Learning. By incorporating jump processes, the Merton
 178 model provides a more accurate representation of environments with abrupt changes.
 179

180
 181 3.3 ADAPTED DECAYING FOR POISSON JUMP INTENSITY
 182

183 We have assumed that the jump intensity λ is constant in the Merton model. However,
 184 such assumption may lead to the overestimation of jump effects, resulting in non-convergent
 185 value estimates. To address this, we introduce an adapted decaying mechanism for the jump
 186 intensity λ .
 187

188 We use the random counting process $N(t)$ to model the number of jumps occurring up to
 189 time t . When a significant jump is detected in the reward signal, we increase the jump
 190 intensity λ to reflect the higher likelihood of future jumps. As the $N(t) \sim P(\lambda)$, then
 191 $t \sim \Gamma(n, \lambda)$, where n is the number of the observed jumps. After time T , where the option
 192 is bound to expire, we may count how many jumps with intensity λ . If there are
 193 too few jumps, then it means that the jump intensity is overestimated, thus we decay the λ
 194 accordingly. Specifically, we update λ as follows:
 195

$$\lambda \leftarrow \lambda \cdot \exp(-\beta \cdot (N(T) - \lambda T)) \quad (6)$$

196 And we sample the observed time period T from $\Gamma(1, \lambda)$.
 197
 198

199 **Rationale** If in real environment, within times T , we observe less jumps than expected
 200 λT , it indicates that the jump intensity λ is overestimated. To address this, we introduce
 201 an adapted decaying mechanism for λ . By adjusting λ based on observed jump counts, we
 202 ensure that the model remains responsive to actual environmental dynamics, thus stabilizing
 203 the learning process and improving convergence of value estimates. Then, we sample T from
 204 $\Gamma(1, \lambda)$ to reflect the updated jump intensity.
 205
 206

207 4 ALGORITHM FRAMEWORK
 208
 209

210 The overall algorithm framework for BSMRL using Merton Potential Shaping is outlined in
 211 Algorithm 1. We introduce a method from ?, which allows us to decompose the reward into
 212 a predictable component and a chaotic component, thus estimate the volatility rate of the
 213 chaotic component, which will be used to compute the Merton potential shaping reward.
 214
 215

Algorithm 1 BSMRL

```

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

1: Hyperparameters: Learning rate  $\alpha$ , Discount factor  $\gamma$ , Moving average rate  $\eta$  (for statistics), Maturity  $T$ , Risk-free rate  $r = -\ln \gamma$ , Jump mean  $\mu_J$ , Jump std  $\delta_J$ , Expansion terms  $K$ .
2: Initialize:
3:  $Q(s, a)$  arbitrarily.
4:  $\bar{R}(s, a) \leftarrow 0$  ▷ Predictable Reward Component
5:  $\Sigma_{chaos}^2(s, a) \leftarrow 0$  ▷ Chaotic Variance (Doob Volatility) [cite: 139]
6:  $\lambda_{jump}(s, a) \leftarrow \lambda_{init}$  ▷ Estimated Jump Intensity for Levy Process
7: Function BLACKSCHOLES( $S, K_{strike}, T, r, \sigma$ ):
8:  $d_1 \leftarrow \frac{\ln(S/K_{strike}) + (r + 0.5\sigma^2)T}{\sigma\sqrt{T}}$ 
9:  $d_2 \leftarrow d_1 - \sigma\sqrt{T}$ 
10: return  $S \cdot \Phi(d_1) - K_{strike}e^{-rT} \cdot \Phi(d_2)$ 
11: End Function
12: Function MERTONPRICE( $S_0, T, r, \sigma_{chaos}, \lambda, \mu_J, \delta_J$ ):
13:  $Price \leftarrow 0$ 
14:  $k \leftarrow e^{\mu_J + 0.5\delta_J^2} - 1$  ▷ Expected jump size
15:  $\lambda' \leftarrow \lambda(1 + k)$ 
16: for  $n = 0$  to  $K$  do
17:    $w_n \leftarrow \frac{e^{-\lambda'T}(\lambda'T)^n}{n!}$  ▷ Poisson weight for n jumps
18:    $\sigma_n \leftarrow \sqrt{\sigma_{chaos}^2 + \frac{n\delta_J^2}{T}}$  ▷ Effective volatility blending Doob noise & Jumps
19:    $r_n \leftarrow r - \lambda k + \frac{n \ln(1+k)}{T}$ 
20:    $Val_n \leftarrow \text{BLACKSCHOLES}(S_0, S_0, T, r_n, \sigma_n)$  ▷ ATM Call Option as Enhanced Value
21:    $Price \leftarrow Price + w_n \cdot Val_n$ 
22: end for
23: return  $Price$ 
24: End Function
25: loop
26: Initialize state  $s$ 
27: while  $s$  is not terminal do
28:   Choose action  $a$  (e.g.,  $\epsilon$ -greedy based on  $Q$ )
29:   Take action  $a$ , observe reward  $R_{obs}$  and next state  $s'$ 
30:    $\delta_{pred} \leftarrow R_{obs} - \bar{R}(s, a)$  ▷ Chaotic innovation
31:    $\bar{R}(s, a) \leftarrow \bar{R}(s, a) + \eta \cdot \delta_{pred}$  ▷ Update Predictable Component
32:    $\Sigma_{chaos}^2(s, a) \leftarrow (1 - \eta)\Sigma_{chaos}^2(s, a) + \eta \cdot (\delta_{pred})^2$  ▷ Update Chaotic Variance
33:    $\lambda_{jump}(s) \leftarrow \lambda_{jump}(s) \exp(-\beta(\frac{(V_t - V_{t-x})}{\lambda_{jump}(s)} - \lambda_{jump}(s)T))$ 
34:    $T \leftarrow \Gamma(1, \lambda_{jump}(s))$  ▷ Adapted Decaying for Jump Intensity
35:    $a'_{best} \leftarrow \arg \max_{a'} Q(s', a')$ 
36:    $S_{underlying} \leftarrow Q(s', a'_{best})$  ▷ Underlying asset is the naive value
37:    $\sigma_{input} \leftarrow \sqrt{\Sigma_{chaos}^2(s', a'_{best})}$  ▷ Use Doob Volatility as input
38:    $V_{enhanced}(s') \leftarrow \text{MERTONPRICE}(S_{underlying}, T, r, \sigma_{input}, \lambda_{jump}(s'), \mu_J, \delta_J)$ 
39:    $Y_{target} \leftarrow R_{obs} + \gamma \cdot V_{enhanced}(s')$  ▷ Maximize Option-Adjusted Return
40:    $Q(s, a) \leftarrow Q(s, a) + \alpha(Y_{target} - Q(s, a))$ 
41:    $s \leftarrow s'$ 
42: end while
43: end loop

```

270 5 THEORETICAL ANALYSIS
271
272 5.1 MDP FORMULATION FOR BSMRL
273
274 5.2 CONVERGENCE ANALYSIS
275
276 5.3 VARIANCE ANALYSIS
277
278 6 EXPERIMENTS
279
280 7 DISCUSSION AND FUTURE WORK
281
282 8 CONCLUSION
283
284 A PROOFS
285
286 B EXPERIMENT DETAILS
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323