

International Journal of Pattern Recognition and Artificial Intelligence  
 © World Scientific Publishing Company

## Hedging Portfolio with Multi-Agent Reinforcement Learning

Quoc Luu

*Quantitative and Computational Finance, John von Neumann Institute,  
 Ho Chi Minh, Vietnam  
 quoc.luu2015@qcf.jvn.edu.vn*

Uyen Pham\*

*Economic Mathematics, University of Economics and Law  
 Ho Chi Minh, Vietnam  
 uyenph@uel.edu.vn*

Hien Tran

*School of Engineering, Tan Tao University,  
 Long An, Vietnam  
 hien.tran@ttu.edu.vn*

Developing a hedge strategy to reduce risk of losses for a given set of stocks in portfolio is a difficult task due to cost of the hedge. In Vietnam stock market, cross-hedge is involved to hedge a long position of a stock because there is no put option for the stock. There is only VN30 stock index futures contracts are traded on Hanoi Stock Exchange (HNX). Inspired by recently achievement of deep reinforcement learning (DRL), we explore feasibility to construct a hedging strategy automatically by leveraging cooperative multi-agent in reinforcement learning techniques without advanced domain knowledge. In this work, we use 10 popular stocks on Ho Chi Minh Stock Exchange (HSX), and VN30F1M (VN30 Index Futures contract within one month settlement) to develop a stock market simulator (including transaction fee, tax, and settlement date of transactions) for RL agent training. We use daily return as input data for training process. Results suggest that the agent can learn trading and hedging policy to make profit and reduce losses. Furthermore, we also find that our agent can protect portfolio and make positive profit in case market collapses systematically. In practical, this work can help Vietnam's stock market investors to improve performance and reduce losses in trading.

*Keywords:* deep reinforcement learning, hedging, trading, portfolio

### 1. Introduction

Hedging a position in stock is an attractive topic for both academics and practitioners. The objective of hedging is to minimize market risk due to price fluctuation, maximize profit by speculation on the basis, and construct a portfolio with reduced risk [1]. Portfolio managers have used stock index futures as a means to adjust

\*corresponding author: uyenph@uel.edu.vn

## 2 Authors' Names

desired return of a portfolio and potential loss since the 1980s. Main advantage of index futures as a major hedging tool is liquidity, lower transaction costs [2]. However, hedge strategies are not always effective as expected because relationship of cash price and price of a futures contract is usually not perfect, or hedged position in stock is different from the underlying portfolio for the index contract [3], [1]. It is possible to increase risk of potential loss that leads to negative return. Hence, hedgers have to determine the optimal hedge ratio to control the risk of the portfolio.

Differ from supervised learning and unsupervised learning, reinforcement learning mainly relies on experience of repeated interaction to learn optimal policy in order to make sequential decisions to maximum rewards in a given environment [4]. In a complex and dynamic environment, it may require huge amounts of computational power over a long period of time to train. With revolution of deep learning techniques and computer hardware, reinforcement learning has become more feasible by using deep neural network as a function estimator. From long time horizons with high dimensional observation and action spaces in real-time strategy controls game, to self-driving vehicles, data center cooling systems, deep reinforcement learning has been more and more applied to solve many complex challenges in real-world [5], [6], [7]. In finance, deep reinforcement learning is also widely adopted. The algorithm is used to design trading strategies for continuous futures contracts [8]. They use technical indicators such as Moving Average Convergence Divergence (MACD), Relative Strength Index (RSI) as a part of input features. The agent shows that it can deliver profit even under heavy transaction costs. Multi-agent dealer market was developed for market marking with different competitive scenarios, market price conditions [9]. The research suggests that trained agent can learn to manage inventory and its competitor's policy for pricing.

However, there are not many researches carried out to investigate feasibility of stock and futures trading to hedge portfolio at the same time using deep reinforcement learning. In this study, we selected 10 popular stocks on HSX, and one stock index futures contract on HNX to build a simulation of stock market environment with real market data to study learning performance of our agent. Our objective in this work is to investigate if cooperation of multi-agent to determine optimal hedging strategy to protect stock portfolio is achievable.

## 2. Related Work

### 2.1. Hedging Effectiveness of Stock Index Futures

Motivated by risk reduction, hedging a stock portfolio with index futures has been an active research topic since it was introduced [3], [2], [10]. Hedger supposes that return of a hedged position (e.g. stock portfolio) can be close to risk-free interest rate. Optimal hedge ratio  $hr$  was estimated using one-to-one hedge, the beta hedge, and minimum variance hedge ratio (MVHR) [11], [12]. With cash prices and futures moving closely together assumption, one-to-one hedge strategy suggests  $hr = -1$ .



Fig. 1. Daily Movement of VN30 Index from 2020-01-17 to 2020-04-08

Beta hedge strategy uses negative of the beta cash portfolio as  $hr$ . The hedger expects the overall beta of the portfolio is zero. However, in practice, change of price of spot and futures does imperfectly correlated. Especially in case of cross-hedge, one-to-one and beta hedge may not reduce risk. In contrast, futures hedging can lead to unexpected loss. The MVHR was introduced to workaround for the problem by taking the imperfect relationship of prices into account and determine the optimal ratio  $hr$ :

$$hr = -\frac{Cov(R_s, R_f)}{Var(R_f)} \quad (1)$$

Furthermore, by using ordinary least squares (OLS) regression to estimate minimum risk hedge, it was found that hedging effectiveness of a large capitalization portfolio can yield “fairly good” for a one week holding period [3] (p. 663). However, with diversified portfolio of small stocks, the effectiveness is reduced significantly. Basis risk is also not negligible even if the spot is hedged with index futures itself. When basis risk arises, it can generate profit or loss. It is suggested that one day hedge strategy can potentially increase basis risk and reduce risk effectiveness than one week hedge.

Stating that traditional methods to estimate optimal hedge ratio are misspecified, error correction model (ECM) was proposed to estimate optimal hedge ratio and forecast out-of-sample for evaluation as in [2]. Firstly, it carries out cointegration test. Secondly, it use OLS regression to estimate error correction model. The model incorporates relationship of the long-run equilibrium as well as the short-run dynamics. The result shows that optimal hedge ratio is significantly improved with adjusted  $R^2$  from ECM is higher than traditional methods. Also by comparing root

mean squared error (RMSE), out-of-sample forecasts from the ECM are found to be better than other methods.

Beyond variance and standard deviation, value at risk (VaR) and conditional conditional value at risk (CVaR) are extensively applied to measure market risk for hedging strategies of portfolio [13], [14]. VaR was introduced by J.P. Morgan in the 1990s and widely adopted to summarize risk of an entire portfolio at the end of each day [15]. However, VaR is not a coherent risk measure. To be coherent, it must be monotonicity, positive homogeneity, translation invariance, and subadditivity [16], [17]. CVaR was constructed with these properties as a new valid practical alternative to VaR [18]. Result showed that CVaR is applicable to a wide range of derivatives portfolio including American options and exotic options [19]. In addition, it is found that CVaR risk metric is suitable for asymmetric return distributions and expected loss of portfolio can be minimized in many circumstances [20].

## **2.2. Deep Reinforcement Learning in Trading**

Reinforcement learning was proposed to train trading systems to make profit and adjust risk [21], [22]. Recurrent learning and Q-Learning with neural networks were used to optimize financial performance functions including risk-adjusted return, immediate utility for online learning [21]. Furthermore, portfolio with continuous quantities of multiple assets was consider. The result shows that reinforcement learning can avoid large losses when market crashed.

Basis risk hedging strategy was developed using reinforcement learning as in [23]. Without assets modeling requirement, State-action-reward-state-action (SARSA) based algorithm was applied to find an optimal trading policy to hedge a non-traded asset. Q-Learning is proposed to extend Black-Scholes-Merton (BSM) model for option pricing and hedging in [24].

In an attempt to escape Greeks and complete market assumptions in risk management, by leveraging deep reinforcement learning, a greek-free approach is proposed to focus on on realistic market dynamics and out-sample testing performance for optimizing hedging of a portfolio of derivatives [25]. Deep reinforcement learning is further investigated for hedging a portfolio of over-the-counter derivatives under generic market frictions as in [26]. Trading costs and liquidity constraints are considered in the approach.

## **3. Background**

### **3.1. Deep Reinforcement Learning**

#### *3.1.1. Single Agent Reinforcement Learning*

For a given stochastic environment  $\varepsilon$ , an agent interacts with the environment by choosing to take a legal action  $a_t$  from many actions at time step  $t$ ,  $a_t \in A \equiv \{1, \dots, L\}$ . Action space can be discrete or continuous. When the selected action is passed to the environment  $\varepsilon$ , internal state  $s_t$  is switched to another state in many

states  $S$ . In other words, the process of sequential interactions between the agent and the environment is result of mapping from perceived states  $s_t$  to actions  $a_t$  by policy  $\pi$ . For instance, in Dota 2 game, internal state can be all the available information for human player including positions, health, map [5]. In this research, internal state is asset return in percentage, position of each asset. In return, the agent receives reward  $r_t$  of the passed action as feedback, and new internal state  $s_{t+1}$  for each time step until reaching terminate state. The ultimate goal of deep reinforcement learning is to find an policy  $\pi$  that can select optimal action to maximize reward signal for each state  $s_t$ .

Value of a state measures total expected return by predicting future reward with discount-rate  $\gamma \in [0, 1]$ . The total accumulated discounted return  $G_t$  from time  $t$  with  $k$  time steps in the future is defined:

$$G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \quad (2)$$

The state value  $V_\pi(s)$  is defined as in [4]:

$$\begin{aligned} V_\pi(s) &\doteq \mathbb{E}_\pi[G_t | S_t = s] \\ &= \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right] \end{aligned} \quad (3)$$

Similarly, action value  $Q_\pi(s, a)$  is the expected return for state  $s$  from selecting action  $a$  following policy  $\pi$ .

$$\begin{aligned} Q_\pi(s, a) &= \mathbb{E}_\pi[G_t | S_t = s, A_t = a] \\ &= \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right] \end{aligned} \quad (4)$$

In value-based reinforcement learning, off-policy Q-learning was introduced to estimate the action value function  $Q_\pi(s, a)$ , defined as [27].

$$\begin{aligned} Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \\ &\alpha \left[ R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right] \end{aligned} \quad (5)$$

The algorithm directly approximate the optimal action value function  $Q_*(s, a) = \max_a Q(s, a)$ . By extending neural network as a function approximator, the function can be estimated as  $Q_*(s, a) \approx Q(s, a, \theta)$ . The approach is referred as Q-network with weights  $\theta$  [28].

In contrast to value-based methods, policy-based can select actions directly by parameterizing the policy  $\pi(a|s, \theta)$  and using gradient ascent to optimize  $\mathbb{E}[R_t]$  to find the best  $\theta$  that can produce the highest reward. In term of probability, we can

express the policy as  $\pi(a|s, \theta) = \Pr\{A_t = a | S_t = s, \theta_t = \theta\}$  for the probability of a given environment  $\varepsilon$  in state  $s$  at time  $t$  with parameter  $\theta$  to take action  $a$ . Actor-critic algorithms use both value and policy functions to learn approximations [29]. To improve performance, the critic learns a value function (e.g. state value) and is used to update policy parameters of actor.

### 3.1.2. Multi-Agent Reinforcement Learning

Extent from single agent, multi-agent learning is considered  $n$  agents interacting with the environment  $\varepsilon$ . At state  $s_t$  of time step  $t$ , each agent selects action  $a_t^i$  to react to the state and receive reward  $r_t^i$ , where  $i \in \{1, \dots, n\}$ . Hence, for any given joint policy  $\pi(a|s) \doteq \prod_i^n \pi^i(a^i|s)$  with state  $s \in S$ , value state function can be defined as in [30]:

$$V_{\pi^i, \pi^{-i}}^i(s) \doteq \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}^i \middle| a_t^i \sim \pi^i(\cdot | s_t), s_0 = s \right] \quad (6)$$

where  $-i$  indicates that all agents except agent  $i$ .

### 3.2. High Throughput Architecture with Importance Weighted Actor-Learner Architecture

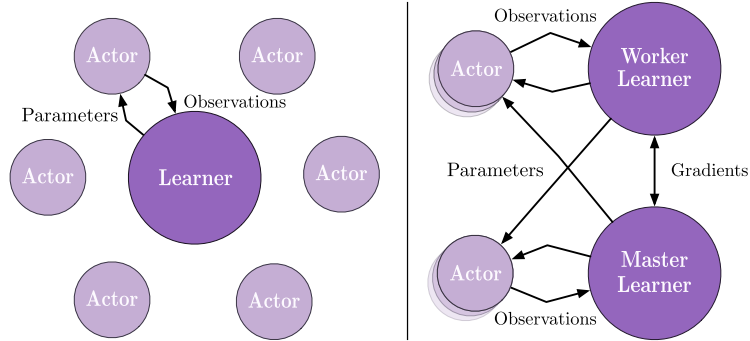


Fig. 2. Left: Single Learner. Right: Multiple Synchronous Learners. Adopted from [31]

Importance weighted actor-learner architecture (IMPALA) is a decoupled actor-critic style learner with introduction of V-trace off-policy to learn a policy  $\pi$  and a baseline function  $V^\pi$  that achieves stability, high data throughput and efficiency for agent training [31]. Moreover, deep neural networks can be train efficiently with IMPALA as suggested in Fig. 2.

Suppose at time  $t$ , a given local actor policy  $\mu$  generates trajectory  $(s_t, a_t, r_t)_{t=k}^{t=k+n}$ , The  $n$ -step V-trace target for value approximation  $V(s_k)$  at state  $s_k$ , is defined as:

$$\begin{aligned}
v_k &\doteq V(s_k) + \sum_{t=k}^{k+n-1} \gamma^{t-k} \left( \prod_{i=k}^{t-1} c_i \right) \delta_t V \\
\delta_t V &\doteq \rho_t(r_t + \gamma V(s_{t+1}) - V(s_t)) \\
\rho_t &\doteq \min(\bar{\rho}, \frac{\pi(a_t|s_t)}{\mu(a_t|s_t)}) \\
c_i &\doteq \min(\bar{c}, \frac{\pi(a_i|s_i)}{\mu(a_i|s_i)})
\end{aligned} \tag{7}$$

where  $\delta_t V$  is temporal difference for  $V$ ,  $\rho_t$  and  $c_i$  are truncated importance sampling. It is worth to note that the truncation levels are assumed  $\bar{c} \leq \bar{\rho}$ .

Furthermore, value function  $V_\theta$  and policy  $\pi_\omega$  with  $\theta$  and  $\omega$  parameters respectively, can be updated in the direction of:

$$\begin{aligned}
\Delta\theta &= (v_k - V_\theta(s_k)) \nabla_\theta V_\theta(s_k) \\
\Delta\omega &= \rho_k \nabla_\omega \log \pi_\omega(a_k|s_k) (r_k + \gamma v_{k+1} - V_\theta(s_k)) - H(\omega) \\
H(\omega) &= \nabla_\omega \sum_a \pi_\omega(a|s_k) \log \pi_\omega(a|s_k).
\end{aligned} \tag{8}$$

Entropy  $H(\omega)$  is added to avoid immature convergence and encourage exploration in agent training process.

IMPALA algorithm can be used to concurrently train for multiple tasks with one set of weights due to efficiency of the architecture.

## 4. Empirical Result

### 4.1. Data Preprocessing and Network Architecture

Training data from raw inputs rather than handcrafted features is often recommended for feeding data to deep neural networks to achieve higher performance [32]. Likewise, researches show that reinforcement learning can exceed human capabilities without human expert data or domain knowledge [28], [33]. As a result, in this study, instead of applying advanced quantitative finance theories to develop trading and hedging strategy, daily return data of each asset collected from HSX and HNX exchanges is used as main components of environment observation. Specifically, we use data from 25 September 2017 to 16 January 2020 for training, and from 17 January 2020 to 08 April 2020 for evaluating. We also provide position and unrealized profit of each asset (P&L) in portfolio to our neural network.

Actor-critic based algorithms use policy and value networks. We can use policy network and value network separately, or combined these two networks. We choose the combined architecture due to improvement of computational efficiency. Furthermore, we use LSTM [34] beside dense layer in the shared network.

As suggested in Fig. 3, we use a shallow network architecture for this study. In detail, in term of shared network, a trajectory is feed to the first fully-connected

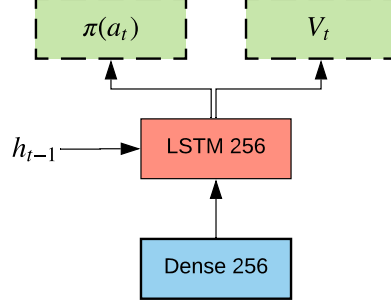


Fig. 3. Model architecture for policy and value networks

hidden layer with 256 units and apply *tanh* activation function. The next hidden layer is stateful LSTM with 256 unit. The tanh function is also applied to the LSTM layer. Finally, policy and value heads are fully-connected linear layer for single output of each action and state value respectively.

#### 4.2. Experiments

We use same network architecture and hyper-parameters all trading agent without tweaking. The agent learn to decide to long buy or short sell assets on its own. For instance, agent can cut loss or hold positions overnight without any constraint.

**Input** : Selected valid action  $\pi(a)$  of learned agent  $\pi$ , Transaction fee  $F$ ,  
Price of asset  $m$ , Asset  $p$  in portfolio  $P$ , Total time step  $T$

**Output:** Portfolio value  $PV$

**for**  $t=1$  to  $t=T$  **do**

**for**  $p$  in  $P$  **do**

**if**  $\pi(a_t) \leftarrow 0$  **then**

            Hold position;

**else**

**if**  $\pi(a_t) \leftarrow 1$  **then**

                Trade return  $r_p \leftarrow$  Long buy at  $m_t$ ;

**else**

                Trade return  $r_p \leftarrow$  Short sell at  $m_t$ ;

**end**

$r_p -= F$ ;

$PV += r_p$ ;

**end**

**end**

**end**

**Algorithm 1:** Simulation of stock market environment



Our experiment use discrete actions. For equity, trading agent can hold, buy, or sell stocks without considering amount of volume. Stocks are only sold after T+2 settlement. Likewise, derivatives trading agent can hold, long, or short futures contracts. However, the agent can trade continuously as it is T+0 settlement market. We include transaction fees for every trade return (see Algorithm. 1). For every agent, reward can be 1 if overall profit of an episode is positive and -1 in case of negative profit.

Finally, we use RLLib [35] with 32 workers to train the agent. RMSProp algorithm is used as optimizer for training.

### 4.3. Trading Results

Learned deep RL agent was deployed to trade out-of-sample market data from 17 January 2020 to 08 April 2020. The result shows that the learned agent can protect portfolio by short selling. In addition, our agent can cut loss, make positive profit in equity market even market plunged as traders had panic-sold out of COVID-19 pandemic fear (see Table. 1).

Table 1. Performance of each trading agent

	BID	BVH	CTG	FPT	GAS	MSN
<i>Number orders</i>	6	8	4	10	10	6
<i>Trading return</i>	-31.1%	-26.5%	-8.8%	6.6%	-12.2%	<b>2.5%</b>
<i>Market return</i>	-33.6%	-28.1%	-22.7%	-19.7%	-25.8%	9.2%

Table 1. Performance of each trading agent (Continue)

	MWG	SSI	VCB	VNM	VN30F1M
<i>Number orders</i>	10	11	7	7	6
<i>Trading return</i>	-14.7%	-1.6%	-7.9%	1.8%	209.1%
<i>Market return</i>	-42.4%	-17.9%	-20.1%	-15.1%	-25.8%

Specifically, VN30 Index was lost nearly 200 points (22.2%) during the period (see Fig. 1). In term of equity trading, nearly every stock in portfolio had negative market return. With many orders for buying and selling, after transaction and commission fees, almost every trading agent has achieves better performance than market performance. In futures market, our agent executed orders for opening and closing position to hedge equity. Especially, the portfolio profit did not decrease when the market rebounded due to dynamic hedge strategy. In case of MSN stock,

## 10 REFERENCES

our trading agent can create positive profit. However, the agent can not achieve higher performance than return of market.

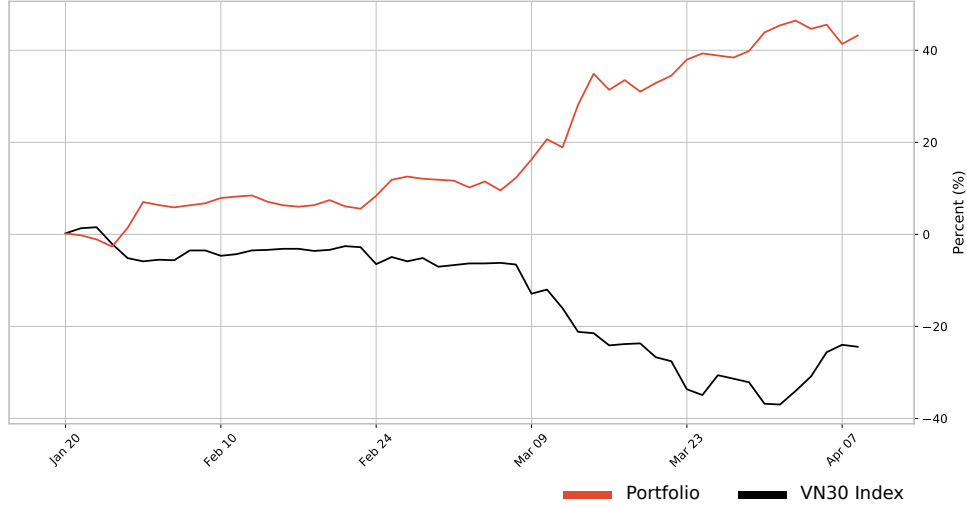


Fig. 4. Uncompounded cumulative return of portfolio and VN30 Index in percentage from 2020-01-17 to 2020-04-08

As a result, we show that our method can reduce losses and achieve positive profit in trading. Furthermore, the trading data also suggests that dynamic hedging strategy for equity in portfolio is feasible in cross-hedging case. The futures trading agent generated far profit than losses. Overall, our deep RL agent earned over 40% profit of portfolio value in case market collapsed systematically.

## 5. Conclusion

This study proposed a feasible approach for cross-hedging in trading without domain knowledge by applying deep reinforcement learning. Our result also suggests that the approach can cut loss efficiently when market is in selling panic. Overall, the proposed method can generate positive profit with dynamic hedge strategy. In future work, we should further study stability and safety in reinforcement learning for trading.

## References

- [1] Christos Floros\* and Dimitrios V Vougas. “Hedge ratios in Greek stock index futures market”. In: *Applied Financial Economics* 14.15 (2004), pp. 1125–1136.

- [2] Asim Ghosh. “Hedging with stock index futures: Estimation and forecasting with error correction model”. In: *Journal of futures markets* 13.7 (1993), pp. 743–752.
- [3] Stephen Figlewski. “Hedging performance and basis risk in stock index futures”. In: *The Journal of Finance* 39.3 (1984), pp. 657–669.
- [4] RS Sutton and AG Barto. *Reinforcement learning: An introduction*. 2nd. MIT press, 2018.
- [5] Christopher Berner et al. “Dota 2 with Large Scale Deep Reinforcement Learning”. In: *arXiv preprint arXiv:1912.06680* (2019).
- [6] Matthew O’Kelly et al. “Scalable end-to-end autonomous vehicle testing via rare-event simulation”. In: *Advances in Neural Information Processing Systems*. 2018, pp. 9827–9838.
- [7] Richard Evans and Jim Gao. “Deepmind AI reduces Google data centre cooling bill by 40%”. In: *DeepMind blog* 20 (2016), p. 158.
- [8] Zihao Zhang, Stefan Zohren, and Roberts Stephen. “Deep Reinforcement Learning for Trading”. In: *The Journal of Financial Data Science* (2020).
- [9] Sumitra Ganesh et al. “Reinforcement Learning for Market Making in a Multi-agent Dealer Market”. In: *arXiv preprint arXiv:1911.05892* (2019).
- [10] Darren Butterworth and Phil Holmes. “The hedging effectiveness of stock index futures: evidence for the FTSE-100 and FTSE-Mid250 indexes traded in the UK”. In: *Applied Financial Economics* 11.1 (2001), pp. 57–68.
- [11] Louis H Ederington. “The hedging performance of the new futures markets”. In: *The journal of finance* 34.1 (1979), pp. 157–170.
- [12] Stewart L Brown. “A reformulation of the portfolio model of hedging”. In: *American Journal of Agricultural Economics* 67.3 (1985), pp. 508–512.
- [13] Zhiguang Cao, Richard DF Harris, and Jian Shen. “Hedging and value at risk: A semi-parametric approach”. In: *Journal of Futures Markets: Futures, Options, and Other Derivative Products* 30.8 (2010), pp. 780–794.
- [14] Markus Huggenberger, Peter Albrecht, and Alexandr Pekelis. “Tail risk hedging and regime switching”. In: *Asian Finance Association (AsianFA) 2015 Conference Paper*. 2016.
- [15] Michael B Miller. *Quantitative Financial Risk Management*. John Wiley & Sons, 2018.
- [16] Ph Artzner et al. “Thinking coherently. RISK, 10”. In: *November*, 68 71 (1997).
- [17] Philippe Artzner et al. *Coherent measures of risk, mathematical Finance 9 (3): 203–228*. 1999.
- [18] Carlo Acerbi and Dirk Tasche. “Expected shortfall: a natural coherent alternative to value at risk”. In: *Economic notes* 31.2 (2002), pp. 379–388.
- [19] Siddharth Alexander, Thomas F Coleman, and Yuying Li. “Derivative portfolio hedging based on CVaR”. In: *New Risk Measures in Investment and Regulation: Wiley* (2003).

## 12 REFERENCES

- [20] Nikolas Topaloglou, Hercules Vladimirov, and Stavros A Zenios. “CVaR models with selective hedging for international asset allocation”. In: *Journal of Banking & Finance* 26.7 (2002), pp. 1535–1561.
- [21] John Moody et al. “Reinforcement learning for trading systems and portfolios: Immediate vs future rewards”. In: *Decision Technologies for Computational Finance*. Springer, 1998, pp. 129–140.
- [22] John Moody et al. “Performance functions and reinforcement learning for trading systems and portfolios”. In: *Journal of Forecasting* 17.5-6 (1998), pp. 441–470.
- [23] Samuel Watts. *Hedging basis risk using reinforcement learning*. Tech. rep. Working Paper, University of Oxford, 2015.
- [24] Igor Halperin. “Qlbs: Q-learner in the black-scholes (-merton) worlds”. In: *Available at SSRN 3087076* (2017).
- [25] Hans Buehler et al. “Deep hedging”. In: *Quantitative Finance* 19.8 (2019), pp. 1271–1291.
- [26] Hans Buehler et al. “Deep Hedging: Hedging Derivatives Under Generic Market Frictions Using Reinforcement Learning-Machine Learning Version”. In: *Available at SSRN* (2019).
- [27] Christopher John Cornish Hellaby Watkins. “Learning from delayed rewards”. In: (1989).
- [28] Volodymyr Mnih et al. “Playing atari with deep reinforcement learning”. In: *arXiv preprint arXiv:1312.5602* (2013).
- [29] Vijay R Konda and John N Tsitsiklis. “Actor-critic algorithms”. In: *Advances in neural information processing systems*. 2000, pp. 1008–1014.
- [30] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. “Multi-Agent Reinforcement Learning: A Selective Overview of Theories and Algorithms”. In: *arXiv preprint arXiv:1911.10635* (2019).
- [31] Lasse Espeholt et al. “Importance weighted actor-learner architectures: Scalable distributed deep-rl with importance weighted actor-learner architectures”. In: *arXiv preprint arXiv:1802.01561* (2018).
- [32] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. “Imagenet classification with deep convolutional neural networks”. In: *Advances in neural information processing systems*. 2012, pp. 1097–1105.
- [33] David Silver et al. “Mastering the game of go without human knowledge”. In: *nature* 550.7676 (2017), pp. 354–359.
- [34] Sepp Hochreiter and Jürgen Schmidhuber. “Long short-term memory”. In: *Neural computation* 9.8 (1997), pp. 1735–1780.
- [35] Eric Liang et al. “RLlib: Abstractions for distributed reinforcement learning”. In: *arXiv preprint arXiv:1712.09381* (2017).