

**University of Warsaw**  
Faculty of Economic Sciences

**Krzysztof Wojdalski**

Nr albumu: 310284

# **Reinforcement Learning for Trading**

**Praca magisterska**  
**na kierunku ECONOMICS**

The thesis written under the supervision of  
**dr Pawel Sakowski**

2025-05-03

Data

Podpis opiekuna

### **Statement of the Supervisor on Submission of the Thesis**

I hereby certify that the thesis submitted has been prepared under my supervision and I declare that it satisfies the requirements of submission in the proceedings for the award of a degree.

Date

Signature of the Supervisor:

### **Oświadczenie autora (autorów) pracy**

Data

Podpis autora (autorów) pracy

### **Statement of the Author(s) on Submission of the Thesis**

Aware of legal liability I certify that the thesis submitted has been prepared by myself and does not include information gathered contrary to the law.

I also declare that the thesis submitted has not been the subject of proceedings resulting in the award of a university degree.

Furthermore I certify that the submitted version of the thesis is identical with its attached electronic version.

Date

Signature of the Author(s) of the thesis

# Table of contents

## 1. Introduction

1.1. Statistical Arbitrage . . . . .	
1.1.1. Statistical Arbitrage vs Other Types of Arbitrage . . . . .	
1.1.2. Critical Parameter Selection and Optimization . . . . .	
1.2. Scope and Objectives of the Research . . . . .	
1.2.1. Hypothesis . . . . .	
1.2.2. Objectives of the Research . . . . .	
1.2.3. Methodological Framework . . . . .	
1.3. Thesis Organization and Chapter Overview . . . . .	

## 2. Classical Finance Context and Literature Review

2.1. Modern Portfolio Theory and CAPM . . . . .	
2.2. The Modern Portfolio Theory . . . . .	
2.3. Efficient Market Hypothesis . . . . .	
2.4. Factor Models . . . . .	
2.5. Modern Approaches in Algorithmic Trading . . . . .	
2.6. Critical Analysis of Traditional Financial Models . . . . .	
2.6.1. Criticism of the Efficient Market Hypothesis and Modern Portfolio Theory	

## 3. Machine Learning

3.1. Why is machine learning important? . . . . .	
3.2. Classification of machine learning algorithms . . . . .	
3.2.1. Supervised Learning . . . . .	
3.2.2. Unsupervised Learning . . . . .	
3.2.3. Semi-Supervised Learning . . . . .	
3.2.4. Reinforcement Learning . . . . .	
3.2.5. Selected aspects of reinforcement learning . . . . .	
3.2.6. Rule-Based Trading Strategies . . . . .	
3.3. Trading based on Forecasts . . . . .	
3.4. Algorithmic Trading Systems . . . . .	

## 4. Reinforcement Learning

## 5. Design of the Trading Agent

5.1. Action Space . . . . .	
5.1.1. Book Pressure-Related Variables . . . . .	
5.1.2. Volume-Related Variables . . . . .	
5.1.3. Last Trade Price-Related Variables . . . . .	
5.1.4. Time-Related Variables . . . . .	
5.1.5. Technical Indicators . . . . .	
5.1.6. Autoregressive Variables . . . . .	
5.2. Reward Function . . . . .	
5.2.1. Differential Sharpe Ratio . . . . .	
5.2.2. Other Reward Functions . . . . .	
5.3. Value Function . . . . .	
5.4. Policy . . . . .	
5.5. Step Size . . . . .	

- 6. Implementation of the Trading Agent
  - 6.1. Data Preparation
    - 6.1.1. Data Collection
    - 6.1.2. Data Preprocessing
    - 6.1.3. Feature Engineering
    - 6.1.4. Data Splitting
  - 6.2. Code
    - 6.2.1. Code Structure
    - 6.2.2. Code Implementation
- 7. Empirical Evaluation and Performance Analysis
  - 7.1. Statistical Validation
  - 7.2. Robustness Assessment
  - 7.3. Performance Evaluation
- 8. Conclusions and Future Work
  - 8.1. Summary of Findings
  - 8.2. Limitations and Future Research
  - 8.3. Implications for Trading Systems
  - 8.4. Recommendations for Practitioners
  - 8.5. Conclusion

**Bibliography**

## **Abstract**

This thesis examines the application of reinforcement learning techniques to foreign exchange market trading strategies. The investigation commences with a comprehensive analysis of the FX market structure, including its organizational framework, key participants, and recent evolutionary developments. The author provides a critical assessment of prevailing trends and potential future trajectories in this domain.

The subsequent section establishes the theoretical foundation for the research, presenting a detailed exposition of relevant financial models and artificial intelligence methodologies that underpin the analytical framework.

The empirical component in the third chapter implements a Q-learning based reinforcement learning approach, yielding inconclusive results that warrant further investigation.

# 1. Introduction

The conventional representation of trading floors as chaotic environments characterized by vocal trader interactions was historically accurate approximately four decades ago. However, the financial industry has since undergone substantial structural transformation.

The financial sector has consistently functioned as an early adopter of computational innovations. This integration represents a strategic imperative—competitive advantage through technological implementation frequently correlates with financial performance. The industry has systematically pioneered the deployment of state-of-the-art technologies, ranging from sophisticated Bloomberg terminal infrastructure in the 1990s to contemporary blockchain applications and ultra-low latency systems, all oriented toward maximizing operational efficiency<sup>1</sup>.

Trading entities, functioning as economic agents, aim to optimize market-derived profits. Multiple methodological approaches exist to achieve this objective. For instance, Warren Buffett, a statistically significant outlier in investment performance metrics, continues to employ a passive buy-and-hold strategy.

In recent years, methodologies predicated on artificial intelligence and machine learning algorithms have demonstrated increasing significance. This phenomenon correlates with computational processing capacity, decreasing infrastructure expenditure requirements, and empirical recognition of cognitive biases in human decision-making processes (Arnold2017?). A growing consensus within the literature suggests that algorithmic systems should supplement or potentially supersede human involvement in decision-making and execution processes Turner (2015). The progressive automation of trading activities will likely continue its trajectory in subsequent temporal periods, contradicting the aforementioned stereotypical trading floor representation.

Although machine learning's theoretical foundations date to the 1950s, its explicit implementation in trading contexts remains relatively limited. For example, in institutional foreign exchange trading, only entities with substantial capital resources and sophisticated quantitative infrastructure have developed effective machine learning trading systems Masic (2017).

This research investigates potential applications of machine learning—specifically reinforcement learning methodologies—in developing trading systems capable of generating statistically significant positive outcomes.

Currently, most systems documented in academic literature aim to maximize either absolute trading profits or risk-adjusted performance metrics. Despite numerous methodological approaches to create consistently profitable systems utilizing variables derived from financial econometrics, fundamental analysis, or machine learning algorithms, many have demonstrated suboptimal performance due to several factors, including:

- Frequent large drawdowns resulting in excessive performance volatility
- Prohibitive transaction costs rendering strategies economically unfeasible
- Excessive computational complexity, particularly problematic for high-frequency trading implementations

Even when empirical research demonstrates exceptional results, upon publication, any com-

---

<sup>1</sup>In empirical literature, the banking sector is identified as a primary industry implementing blockchain technology.

petitive advantage tends to diminish through market efficiency mechanisms. Strategies with persistent alpha-generating capabilities must remain proprietary to maintain their effectiveness.

#### **1.0.0.1. Work Structure**

This thesis is structured to provide a comprehensive analytical framework for understanding reinforcement learning algorithms applied to trading contexts.

The first section provides a detailed examination of the intersection between artificial intelligence methodologies and financial markets, exploring the historical relationship between quantitative finance and computational science.

The second chapter presents a systematic review of selected literature from quantitative finance, examining both classical equilibrium models such as CAPM (the established paradigm in equity research) and contemporary approaches. This section evaluates the implicit advantages and limitations of various financial models, with particular emphasis on algorithmic trading methodologies employed in comparable research contexts.

The third section analyzes machine learning frameworks, providing a theoretical basis for why reinforcement learning may represent an optimal approach for certain trading applications. It presents a taxonomic comparison of major machine learning categories to elucidate their methodological distinctions, introduces key reinforcement learning concepts with illustrative examples, and addresses potential limitations and implementation challenges associated with these algorithms.

The fourth part details the experimental methodology, including research objectives, data characteristics, experimental design parameters, and empirical results. The primary objective was to develop trading agents capable of statistically outperforming established benchmarks on risk-adjusted performance measures in the foreign exchange market—agents characterized by statistical robustness, adaptive learning capability, and consistent performance metrics. This chapter presents the mathematical formulations and procedural implementations leading to the empirical results, examining each component of the trading system.

The implemented algorithms utilize a dynamic optimization approach. Beyond a value function based on Differential Sharpe Ratio, the system incorporates various technical indicators such as Relative Strength Index to inform algorithmic decision-making processes. The methodology incorporates transaction cost models to simulate realistic trading conditions.

The value function integrates multiple statistical measures, including Sharpe and Differential Sharpe Ratios, to capture both risk and return dimensions. The algorithm outputs agent actions in the discrete action space  $\{-1,0,1\}$ . The final section of this chapter evaluates the reinforcement learning-based trading system against two benchmark methodologies:

- A buy-and-hold strategy (maintaining consistent long positions in selected currency pairs)
- Random action generation—producing stochastic values in the domain of  $\{-1,0,1\}$  to determine positions in underlying pairs. This benchmark excludes transaction costs, as such a strategy would incur prohibitive cumulative costs with position changes occurring in approximately two-thirds of states.

The concluding section presents a comparative analysis with similar research and proposes directions for future investigation, addressing research questions such as:



- What additional implementations could enhance performance metrics?
- What methodological limitations were encountered and how might they be addressed in subsequent research?

## 1.1. Statistical Arbitrage

In financial markets the term arbitrage means the practice of exploiting riskless opportunities to make profit. In the purest form, if an asset is priced differently at two different places and a trader has access to both then he would have an arbitrage opportunity. In practice financial markets are designed in such a way that any deterministic arbitrage opportunities vanish, i.e., are being exploited by market participants. Still there are patterns in the market which seem to be persistent. Strategies which use statistical techniques to extract highly reliable patterns and then carry out trades to exploit them are called statistical arbitrage strategies.

Statistical arbitrage is a general term in which assets are put into baskets by market-based similarities. Once the relation between the co-moving financial instruments is found, any deviation from it can be exploited. If that relation is true then either one or both stocks are mispriced and a trader can take a long position on the undervalued one(s) and short position on the overvalued asset(s) to bet on their return to the prevalent relative price behaviour. Statistical arbitrage also fittingly resembles the fact that the strategy hedges risk from market movements. As the pair trade goes short and long on similar assets, the overall risk balance should cancel out and the strategy should be profitable without market risk exposure.

To determine if a security is over- or undervalued one can either use some theoretic approach to determine what the absolute price of the security should be or one could use the idea of relative pricing. In this approach it is only of interest how one security is priced in relation to another security. Now to form a simple mean-reverting pair trading strategy the only thing that remains is a definition of what constitutes a reliable price relationship between the asset prices. There are many ways to identify such security pairs. One approach would be to measure fundamentals (such as some accountancy variables) of the firms. The statistical approach, which also makes the whole concept easily applicable in a system is to use cointegration.

### 1.1.1. Statistical Arbitrage vs Other Types of Arbitrage

While statistical arbitrage relies on probabilistic relationships between securities, several other arbitrage types exist in financial markets, each with distinct characteristics:

#### 1.1.1.1. Classical Arbitrage

Classical arbitrage involves simultaneously trading two or more fungible instruments and converting between them to capture price differentials. At least one leg might involve derivatives:

- **ETF arbitrage:** Buying an ETF while selling its constituent basket, then redeeming the ETF
- **ADR arbitrage:** Selling an ADR while buying its underlying foreign stock and currency, then creating the ADR
- **Futures arbitrage:** Buying a future while selling its underlying, holding both until expiration

#### **1.1.1.2. Latency Arbitrage**

Latency arbitrage leverages speed advantages to be first to complete trades triggered by discrete events. Examples include racing to take stale bids/offers on one exchange after observing trades on another, or using private fill prices to trade ahead of other participants.

#### **1.1.1.3. Time Arbitrage**

Time arbitrage captures spreads between buyers and sellers demanding liquidity at different periods. US equity wholesalers practice this by guaranteeing retail buy orders, expecting offsetting sell orders to follow. Opportunities also arise when linked instruments trade on exchanges with different operating hours.

#### **1.1.1.4. Microstructure Arbitrage**

This captures price dislocations using exchange-specific idiosyncrasies like matching semantics, fee treatment, and specialized order types. Common opportunities arise from non-continuous matching mechanisms such as opening/ closing auctions and their associated order types.

#### **1.1.1.5. Statistical Arbitrage**

Statistical arbitrage trades spreads between instruments based on probabilistic estimations of convergence (mean reversion) or divergence (momentum). The “dispersion trade” bets that index option prices should converge to weighted baskets of options on the index’s components. Unlike classical arbitrage, statistical arbitrage involves probability-based relationships that may break down, introducing additional risk elements.

### [A Taxonomy of Arbitrage Trading](#)

#### **1.1.2. Critical Parameter Selection and Optimization**

### **1.2. Scope and Objectives of the Research**

#### **1.2.1. Hypothesis**

The research is based on the following hypotheses:

1. AI-driven algorithms can deliver superior value by exceeding benchmark performance metrics in terms of both risk management and return generation
2. Enhanced performance is particularly evident in high-frequency trading scenarios as well as across extended time horizons
3. These algorithms demonstrate the capability to implicitly learn market patterns that would be impossible for humans to identify and efficiently select securities. In consequence, they generate alpha

#### **1.2.2. Objectives of the Research**

The research has three main objectives:

##### **1. Design and Implementation**

- Design and implement a basic and an extended Reinforcement Learning trading agent for pair trading

- Develop robust testing frameworks for both approaches
- 2. Testing and Comparison**
- Test and compare both approaches on out-of-sample test data
  - Evaluate performance using selected performance measures
  - Analyze the robustness of the approaches under different market conditions
- 3. Interpretation and Implications**
- Analyze how the agents evolve their learning processes
  - Study how their developed policies integrate conventional pair trading approaches
  - Investigate transaction cost management strategies
  - Identify and suggest potential avenues for future research

### 1.2.3. Methodological Framework

## 1.3. Thesis Organization and Chapter Overview

# 2. Classical Finance Context and Literature Review

## 2.1. Modern Portfolio Theory and CAPM

The following chapter introduces articles that correspond with the subject of the current thesis and are considered as fundamentals of modern finance. Specifically, the beginning contains financial market models. The next subchapter includes basic investment effectiveness indicators that implicitly or explicitly result from the fundamental formulas from the first subchapter.

### 2.1.0.1. Capital Asset Pricing Model

Works considered as a fundament of quantitative finance and investments are Sharpe (1964), Lintner (1965), and Mossin (1966). All these authors, almost simultaneously, formulated Capital Asset Pricing Model (CAPM) that describes dependability between rate of return and its risk, risk of the market portfolio, and risk premium.

Assumptions in the model are as follows:

- Decisions in the model regard only one period,
- Market participants has risk aversion, i.e. their utility function is related with plus sign to rate of return, and negatively to variance of portfolio rate of return,
- Risk-free rate exists,
- Asymmetry of information non-existent,
- Lack of speculative transactions,
- Lack of transactional costs, taxes included,
- Market participants can buy a fraction of the asset,
- Both sides are price takers,
- Short selling exists,

Described by the following model formula is as follows:

$$E(R_P) = R_F + \frac{\sigma_P}{\sigma_M} \times [E(R_M) - R_F]$$

where:

- $E(R_P)$  – the expected portfolio rate of return,

- $E(R_M)$  – the expected market rate of return,
- $R_F$  – risk-free rate,
- $\sigma_P$  – the standard deviation of the rate of return on the portfolio,
- $\sigma_M$  – the standard deviation of the rate of return on the market portfolio.

$E(R_P)$  function is also known as Capital Market Line (CML). Any portfolio lies on that line is effective, i.e. its rate of return corresponds to embedded risk. The next formula includes all portfolios, single assets included. It is also known as Security Market Line (SML) and is given by the following equation:

$$E(R_i) = R_F + \beta_i \times [E(R_M) - R_F]$$

where:

- $E(R_i)$  – the expected  $i$ -th portfolio rate of return,
- $E(R_M)$  – the expected market rate of return,
- $R_F$  – risk-free rate,
- $\beta_i$  – Beta factor of the  $i$ -th portfolio.

## 2.2. The Modern Portfolio Theory

The following section discuss the Modern Portfolio Theory developed by Henry Markowitz Stulz (1995). The author introduced the model in which the goal (investment criteria) is not only to maximize the return but also to minimize the variance. He claimed that by combining assets in different composition it is possible to obtain the portfolios with the same return but different levels of risk. The risk reduction is possible by diversification, i.e. giving proper weights for each asset in the portfolio. Variance of portfolio value can be effectively reduced by analyzing mutual relations between returns on assets with use of methods in statistics (correlation and covariance matrices). It is important to say that any additional asset in portfolio reduces minimal variance for a given portfolio but it is the correlation what really impacts the magnitude. The Markowitz theory implies that for any assumed expected return there is the only one portfolio that minimizes risk. Alternatively, there is only one portfolio that maximizes return for the assumed risk level. The important term, which is brought in literature, is the effective portfolio, i.e. the one that meets conditions above. The combination of optimal portfolios on the bullet.

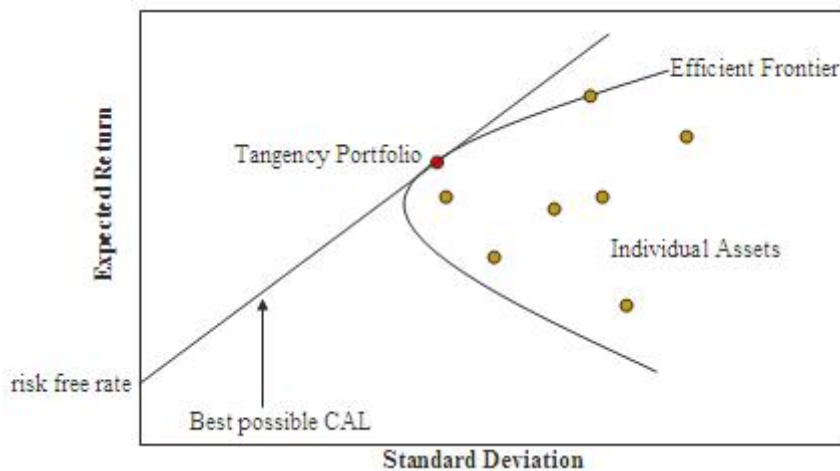


Figure 1: Efficient Frontier

The Markowitz concept is determined by the assumption that investors are risk-averse. This observation is described by the following formula:

$$E(U) < U(E(X))$$

where:

- $E(U)$  – the expected value of utility from payoff;
- $U(E(X))$  – utility of the expected value of payoff.

The expected value of payoff is given by the following formula:

$$E(U) = \sum_{i=1}^n \pi_i U(c_i)$$

where:

- $\pi_i$  – probability of the  $c_i$  payoff,
- $U(c_i)$  – utility from the  $c_i$  payoff.

One of the MPT biggest flaws is the fact that it is used for ex post analysis. Correlation between assets changes overtime so results must be recalculated. Real portfolio risk may be underestimated. Also, time window can influence the results.

## 2.3. Efficient Market Hypothesis

In 1965, Eugene Fama introduced the efficient market term. Fama claimed that an efficient market is one that instantaneously discounts new information arrival in the market price of a given asset. Because this definition applies to financial markets, it determined the further belief that it is not possible to beat the market because assets are correctly priced. Also, if this hypothesis were true, market participants could not be better or worse. Their portfolio return would be a function of new, unpredictable information. In that respect, the only role of an investor is to manage assets so that the risk is acceptable. Fama (1965)

It is highly unlikely that EMH exists in its strongest form due to successful quantitative hedge funds that consistently beat the markets. For instance, Renaissance Technologies hedge fund generated on average 40% per annum in the last 30 years Shen (2017).

Formally, Efficient Market Hypothesis states that a market is efficient with respect to information set  $F_t$  if it is impossible to make economic profits by trading on the basis of that information set. In other words, it is not possible to achieve any better than risk-adjusted average rate of return. In its essence, that claim is consistent with classical price theory Weber (2012). Over time, other versions (forms) of the EMH have been introduced - weak, semi-strong, and strong Fama, E. F.; Malkiel (1970).

**Definition 1.** Weak Form of the EMH  $F_t$  represents only the information contained in the past price history of the market as of time  $t$

According to this form, investors cannot achieve above-normal returns by analyzing historical price patterns and trading volumes to forecast future price movements. Nevertheless, fundamental analysis may still yield exceptional results since markets are not entirely efficient at identifying under or overvalued securities. Therefore, market participants can potentially discover profitable investment opportunities through thorough examination of companies' financial reports.

**Definition 2.** Semi-Strong Form of the EMH  $F_t$  represents all information publicly available at time  $t$

It states that neither technical nor fundamental analysis can be exploited for gaining superior returns, and only non-public material information might help in achieving above average results.

**Definition 3.** Strong Form of the EMH  $F_t$  represents all information (public and private) known to anyone at time  $t$ .

The strong form rejects the idea of any possibility of consistently beating the market. According to this idea, any kind of information, public or non-public, is completely embedded into current financial asset prices. In other words, there is no advantage for anyone in the market. Returns that deviate from expected values are attributed to pure randomness.

### 2.3.0.1. Critic of strong form of the EMH

There are at least a few documented anomalies that contradicts with efficient market hypothesis. For example, price/earnings (P/E) measure can help in systematically outperforming stocks Malkiel (2003). The neglected firm effect claims that “uninteresting” companies, often ignored by market analysts are sometimes incorrectly priced, and offer investors potentially fruitful opportunities. Another phenomenon that cannot be explained by the strong form of EMH is so called the January effect Haug and Hirschey (2006). According to the authors of “The January Effect” working paper, returns reached in January has predictive power for the upcoming 11 months. It persists for both small and large cap companies.

Although the strongest form in its essence is justified, logically correct, it is rather unlikely that it explains the reality, even due to the effects mentioned above.

## 2.4. Factor Models

## 2.5. Modern Approaches in Algorithmic Trading

## 2.6. Critical Analysis of Traditional Financial Models

### 2.6.1. Criticism of the Efficient Market Hypothesis and Modern Portfolio Theory

The Efficient Market Hypothesis (EMH) posits that investors cannot earn excess returns on a risk-adjusted basis. Complementarily, Markowitz’s Modern Portfolio Theory (MPT) provides a framework for constructing portfolios that optimize risk-return profiles based on available investment options.

If markets are truly efficient as EMH suggests, portfolios constructed using the Markowitz framework would not generate excess returns. However, if market inefficiencies exist, these portfolios could potentially earn returns exceeding market benchmarks. Notably, MPT allows investors to incorporate subjective views regarding expected asset returns, contributing to its widespread applicability and Markowitz’s Nobel Prize recognition.

Both theories gained popularity synergistically and share similar criticisms. Despite MPT becoming an industry standard, it faces several significant limitations:

1. **Problematic Risk Measurement:** MPT treats upside and downside price movements equally in risk calculations. It also assumes constant volatility, whereas real-world markets demonstrate volatility clustering across various time horizons.

2. **Simplified Cash Flow Assumptions:** The model ignores transaction costs and tax implications. It also fails to account for investor preferences regarding dividend income.
3. **Unrealistic Liquidity Assumptions:** MPT assumes investors can take arbitrarily large positions at current market prices without affecting those prices, contradicting market impact realities.
4. **Idealized Investor Behavior:** The theory presumes all market participants are rational, risk-averse, and share identical investment time horizons.

### 3. Machine Learning

In this chapter, the term **machine learning** and its subfields are explained. The discussion also encompasses potential applications for trading financial instruments.

As the field evolves, numerous definitions of machine learning emerge from various sources. In this subchapter, the author has selected definitions that accurately capture the essence of the discipline.

What is machine learning? The most widely accepted definitions are as follows:

- “Field of study that gives computers the ability to learn without being explicitly programmed.” - Arthur Samuel, a pioneer in machine learning and computer gaming Samuel (1959)
- “A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .” - Tom Mitchell, a computer scientist and E. Fredkin University Professor at Carnegie Mellon University (CMU) Mitchell (1997)

The latter is particularly regarded as an elegant and modern definition. Less formal, but equally relevant observations come from textbook authors in the discipline:

- “Pattern recognition has its origins in engineering, whereas machine learning grew out of computer science. However, these activities can be viewed as two facets of the same field...” - Christopher Bishop
- “One of the most interesting features of machine learning is that it lies on the boundary of several different academic disciplines, principally computer science, statistics, mathematics, and engineering. ...machine learning is usually studied as part of artificial intelligence, which puts it firmly into computer science ...understanding why these algorithms work requires a certain amount of statistical and mathematical sophistication that is often missing from computer science undergraduates.” - Stephen Marsland Marsland (2009)

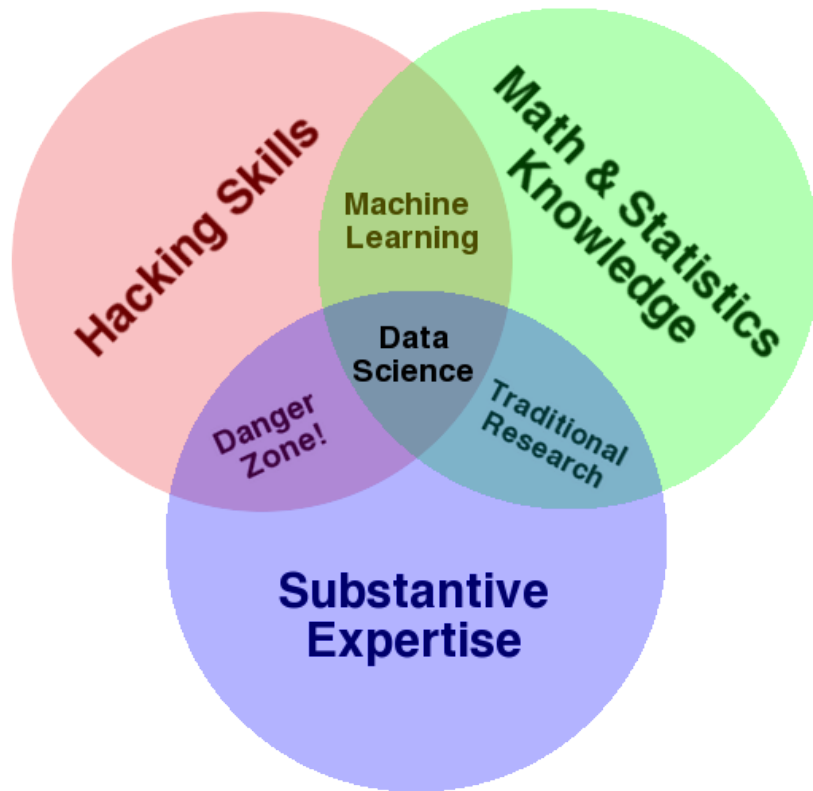


Figure 2: Data Science Graph

Despite numerous concepts and perspectives on what machine learning entails, the general objective remains consistent: Machine learning involves building models that sufficiently resemble reality, are optimal with respect to a value function, and can subsequently be utilized for predictions on new data.

### 3.1. Why is machine learning important?

Machine learning facilitates solving problems that are difficult or impossible to address deterministically Jason (2013). Variables may be missing or observed values may contain embedded errors. Traditional models are often susceptible to being under- or overdetermined. They may fail to generalize adequately or may be excessively general. An appropriate machine learning model should provide an approximate solution incorporating only relevant components.

### 3.2. Classification of machine learning algorithms

In machine learning (ML), tasks are categorized based on how learning/feedback ( $P$ ) is received and/or the type of problem they address. The following categories can be distinguished:

- Supervised Learning - the complete set  $(Y_t; X_{t,1}, \dots, X_{t,n})$  is available. The objective is to model the target variable  $Y_t$  using a subset of  $X_t$  variables, i.e., find a functional relationship  $Y_t = f(\mathbb{X}_t)$  between input and output variables that minimizes a predefined loss function  $g(f(\mathbb{X}_t); Y_t)$ . The structural form of this relationship is constrained by the



class of functions considered. For example, assuming a linear relationship between input and output variables with a square loss function, the problem becomes:

$$\min_{b_1 \dots b_n} \mathbb{E}[(Y_t - (b_1 X_{t,1} + \dots + b_n X_{t,n}))^2]$$

The estimation method above is known as the least squares method for linear regression. Despite its simplicity, it often yields sufficient results. Other popular supervised learning methods include:

- K-nearest neighbors, Neural Networks
  - SVM - Support Vector Machines
  - Random Forests
- Unsupervised learning - this category deals exclusively with the  $\mathbb{X}_t$  set. The goal is to identify patterns within the dataset and categorize observations. The most prevalent methods include:
  - Clustering - based on finding groups of instances that are as similar as possible to observations within the same group while being as different as possible from observations in other groups
  - Feature extraction - this subcategory comprises methods for extracting relevant variables from a set  $\mathbb{X}_t$ . Often, a subset of a dataset can contain a similar amount of information as the original while reducing dimensionality, thereby enhancing computational efficiency and improving the model in accordance with Occam's Razor.
  - Anomaly detection - this approach aids in identifying outlier observations that warrant careful investigation. Variables may require transformation or invalid observations may need removal.
- Reinforcement Learning - this is perhaps the most intuitive category of ML in terms of what is commonly associated with artificial intelligence. According to (**Silver2017?**), it incorporates influences from engineering, economics, mathematics, neuroscience, psychology, and computer science. Reinforcement learning algorithms maximize long-term cumulative rewards and **interact with the environment**, making them suitable for non-stationary problems.

The two distinctive features of reinforcement learning algorithms are trial-and-error and delayed rewards, meaning this type of ML evaluates actions rather than providing definitive instructions. This distinguishes reinforcement learning from supervised learning and constitutes one reason why it is considered a distinct subfield in ML.

Moreover, it does not rely on a training set of labeled examples. In supervised learning, each observation specifies precisely what an algorithm should do. For instance, if blue balls according to the model should be in a blue basket, they will invariably end up there.

The goal of supervised learning is to generalize effectively from training data so that the formula works for test data as well. While important and extensively researched, this approach is insufficient when interaction between an agent and an environment occurs. In such scenarios, an agent must learn from its own actions, sense states, and accumulate experience.

Reinforcement learning must also be distinguished from unsupervised learning. Unsupervised learning focuses on identifying structures not explicitly given in collections of unlabeled datasets. While this sounds similar, it differs fundamentally from RL, where the objective is to maximize the sum of reward signals. Finding data patterns may be useful (as noted in the section on unsupervised learning), but it does not address a reinforcement learning problem. Hence, the approach analyzed in this thesis should be considered a separate paradigm.

The only feedback an agent receives is a scalar reward. Its objective is to maximize a long-run value function comprising summed (discounted) rewards in subsequent states. The agent aims to learn through trial-and-error which actions maximize long-run rewards. The environment changes stochastically and sometimes interacts with the agent. The agent must select a policy that optimizes the rewards received. The design must account for this by adjusting the agent to avoid purely greedy behavior, i.e., it should explore new actions rather than exclusively exploiting existing optimal (possibly sub-optimal) solutions.

### 3.2.1. Supervised Learning

Supervised learning represents a fundamental paradigm in machine learning where algorithms learn from labeled training data to make predictions or decisions without explicit programming. At its core, supervised learning involves a dataset consisting of input-output pairs, where each example contains features (input variables) and their corresponding target values or labels (output variables). The primary objective is to learn a mapping function that can accurately predict the output value for new, previously unseen inputs.

Formally, given a dataset of pairs  $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ , where  $X_i$  represents the feature vector and  $Y_i$  is the corresponding target value, supervised learning aims to find a function  $f$  such that  $Y = f(X)$  which minimizes a predefined loss function  $L(f(X), Y)$ . The function  $f$  is constrained by the model class chosen for the task, which determines the complexity and expressiveness of the relationships that can be captured.

The mathematical foundation of supervised learning lies in statistical learning theory and optimization. For instance, in a linear regression model, we seek parameters  $\beta$  that minimize the sum of squared errors:  $\min_{\beta} \sum (Y_i - X_i\beta)^2$ . In more complex models like neural networks, we optimize weights and biases across multiple layers using gradient-based methods to minimize error functions across the entire training dataset.

Supervised learning tasks typically fall into two main categories: - Classification: When the output variable  $Y$  is categorical or discrete (e.g., spam/not spam, fraud/legitimate, image categories) - Regression: When the output variable  $Y$  is continuous (e.g., stock prices, temperature, house prices)

The choice of algorithm depends on various factors including the nature of the problem, dataset characteristics, computational resources, and the desired balance between model interpretability and predictive performance. Common supervised learning algorithms include:

- Linear models: Linear regression for regression tasks and logistic regression for classification
- Tree-based methods: Decision trees, random forests, and gradient boosting machines
- Support vector machines: Effective for both classification and regression with high-dimensional data

- **K-nearest neighbors:** A non-parametric method that makes predictions based on similarity measures
- **Neural networks:** Deep learning architectures capable of capturing complex non-linear relationships
- **Ensemble methods:** Combining multiple models to improve overall performance and robustness

A critical aspect of supervised learning is the bias-variance tradeoff. Simple models may underfit the data (high bias), failing to capture important patterns, while overly complex models may overfit (high variance), learning noise rather than underlying relationships. Techniques such as regularization, cross-validation, and ensemble methods help manage this tradeoff to create models that generalize well to unseen data.

In financial applications, supervised learning has become important for various tasks such as price prediction, risk assessment, credit scoring, fraud detection, and market sentiment analysis. For instance, in predicting stock prices, historical market data with known outcomes serves as the training set, where features might include technical indicators, fundamental data, and macroeconomic variables, while the target variable could be future price movements or returns.

The effectiveness of supervised learning models for financial markets is often challenged by the non-stationary nature of markets, where relationships between variables change over time and in the presence of noise. This necessitates continuous model updating and validation against recent data. Additionally, feature engineering—the process of creating relevant variables from raw data—plays a crucial role in financial applications, often requiring domain expertise to identify meaningful predictors.

### 3.2.2. Unsupervised Learning

Unsupervised learning represents a paradigm in machine learning where algorithms learn patterns and structures from unlabeled data without explicit guidance. Unlike supervised learning, there are no target outputs or labels to guide the learning process. Instead, the algorithm discovers the inherent structure within the data itself.

Formally, in unsupervised learning, the dataset consists of a collection of unlabeled examples  $X_1, X_2, \dots, X_N$ , where each  $X_i$  represents a feature vector. The primary objective is to create a model that processes these feature vectors to either transform them into another representation or extract meaningful patterns that can solve practical problems.

The mathematical foundation of unsupervised learning involves finding patterns, relationships, or structures within the data space. For instance, in clustering algorithms, we seek to minimize within-cluster distances while maximizing between-cluster distances, often expressed as optimization problems such as minimizing  $\sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$  where  $C_i$  represents clusters and  $\mu_i$  their centroids.

Unsupervised learning encompasses several key application areas:

- **Clustering:** Identifying natural groupings within data where instances within the same cluster exhibit high similarity while being dissimilar to instances in other clusters. Common algorithms include K-means, hierarchical clustering, and DBSCAN.
- **Dimensionality Reduction:** Transforming high-dimensional data into lower-dimensional representations while preserving essential information. Techniques such as

Principal Component Analysis (PCA), t-SNE, and autoencoders fall into this category.

- **Density Estimation:** Modeling the probability distribution that generates the observed data, enabling better understanding of data characteristics and generation of new samples. Methods include Gaussian Mixture Models and kernel density estimation.
- **Anomaly Detection:** Identifying instances that deviate significantly from the norm or expected patterns within the dataset. These outliers often represent rare events, errors, or fraudulent activities that warrant special attention.
- **Feature Learning:** Automatically discovering useful representations from raw data that can subsequently enhance performance in downstream tasks, including supervised learning problems.

In financial applications, unsupervised learning proves valuable for market segmentation, identifying trading patterns, detecting fraudulent transactions, and discovering latent factors driving market movements. The absence of labeled data makes unsupervised learning particularly useful in exploratory analysis and when dealing with novel or evolving financial phenomena where historical classifications may not exist or apply.

The effectiveness of unsupervised learning in finance is often measured by metrics such as silhouette scores for clustering quality, explained variance for dimensionality reduction, or business impact metrics like improved portfolio diversification or fraud detection rates. As markets evolve and data complexity increases, unsupervised learning provides tools for uncovering hidden structures and relationships in financial data.

### 3.2.3. Semi-Supervised Learning

Semi-supervised learning represents a hybrid approach in machine learning where the dataset contains both labeled and unlabeled examples. Formally, the explanatory variables  $X_i$  are available for all observations, but the labels  $Y_i$  are only available for a subset of the data. Typically, the quantity of unlabeled examples significantly exceeds the number of labeled examples.

The primary objective, similar to supervised learning, is to discover the relationship  $Y = f(X)$ . This is generally accomplished through a strategic combination of supervised and unsupervised learning techniques. The underlying principle involves labeled observations effectively “diffusing” their labels to unlabeled observations that exhibit high similarity according to specific criteria.

The appeal of semi-supervised learning lies in its ability to leverage large amounts of unlabeled data, which is often more abundant and less costly to obtain than labeled data. By incorporating these unlabeled examples, the learning algorithm can potentially develop a more robust and generalizable model than would be possible using only the limited labeled examples.

Several key techniques in semi-supervised learning include:

- **Self-training:** An iterative process where a model trained on labeled data makes predictions on unlabeled data, then adds high-confidence predictions to the training set.

- **Co-training:** Using multiple views or feature subsets to train separate models that teach each other by labeling unlabeled examples for one another.
- **Graph-based methods:** Constructing similarity graphs where nodes represent data points and edges represent similarities, allowing label propagation through the graph structure.
- **Generative models:** Using techniques like Gaussian Mixture Models to model the joint distribution of features and labels.
- **Semi-supervised SVMs (S3VMs):** Extensions of Support Vector Machines that incorporate unlabeled data by seeking decision boundaries that avoid dense regions of unlabeled points.

### 3.2.4. Reinforcement Learning

Reinforcement learning represents a distinct paradigm in machine learning where an agent interacts with an environment, learning optimal behaviors through trial and error. In this framework, the agent perceives the environment's state as a feature vector and can execute actions in each state. These actions yield varying rewards and potentially transition the agent to new states.

The fundamental objective in reinforcement learning is to develop a policy—a function that maps states to actions. Formally, this policy function  $f$  (analogous to supervised learning models) takes a state's feature vector as input and outputs the optimal action for that state. An action is considered optimal when it maximizes the expected cumulative reward over time.

Key characteristics of reinforcement learning include:

- **Sequential Decision Making:** Unlike other learning paradigms, reinforcement learning addresses problems where decisions occur in sequence and have long-term implications.
- **Delayed Rewards:** The agent must learn to take actions that may not immediately yield rewards but contribute to greater cumulative rewards in the future.
- **Exploration vs. Exploitation:** Agents must balance exploring new actions to discover potentially better strategies against exploiting known rewarding actions.

Reinforcement learning is particularly well-suited for non-stationary environments where relationships between variables evolve over time. In these contexts, the agent continuously adapts its policy to changing conditions, making it valuable for financial applications where market dynamics constantly shift.

Applications of reinforcement learning span diverse domains including game playing, robotics, resource management, logistics, and increasingly, financial trading strategies such as statistical arbitrage. The target variable in reinforcement learning can be viewed as a specialized output designed to solve specific optimization problems within dynamic environments.

This type of Machine Learning is discussed in detail in the following section.

### 3.2.5. Selected aspects of reinforcement learning

In the following section the author discussed relevant aspects and challenges of the paradigm.

### 3.2.5.1. Exploration/exploitation

One of the most important problems in RL is the trade-off between exploration and exploitation. To maximize cumulated rewards an agent should take actions that worked in the past and caused bigger payoffs (exploit). At the very beginning of learning process it never knows what works well, though. Hence, it needs to discover desirable actions for its state (explore). The dilemma is unresolved as of now, there are at least a few approaches to tackle the problem, though. In the next subsection the author presents that possible methods on the example of Bandit problem.

#### 3.2.5.1.1. $\epsilon$ -greedy policy

The simplest version is to behave greedily most of the time, i.e. an agent selects such action ( $A_t$ ) that maximizes the used value function (e.g.  $Q_t(a)$ ), but sometimes, with probability of  $\epsilon$  pick up a random action from those available, apart from the action value estimates. Such an algorithm guarantees that every action for every state will be explored and eventually  $Q_t(a) = q_*(a)$ . It implies that probability of choosing the most optimal action will converge to more than  $1 - \epsilon$ , to near certainty. The disadvantage of this simple method is that it says very little of its practical effectiveness. Asymptotic guarantee might take too long in a real environment. It can be shown that small  $\epsilon$  causes the agent to gain more reward at initial steps, but tends to underperform against larger  $\epsilon$  values when number of steps is getting larger.

#### 3.2.5.1.2. Optimistic initial values

One of the techniques to improve agent's choices is based on the idea of encouraging the agent to explore. Why is that? If the actual reward is smaller than initially set up action-value methods, an agent is more likely to pick up actions that potentially can stop getting rewards that constantly worsen value function  $q(a)$ . Eventually, the system does a lot more exploration even if greedy actions are selected all the time.

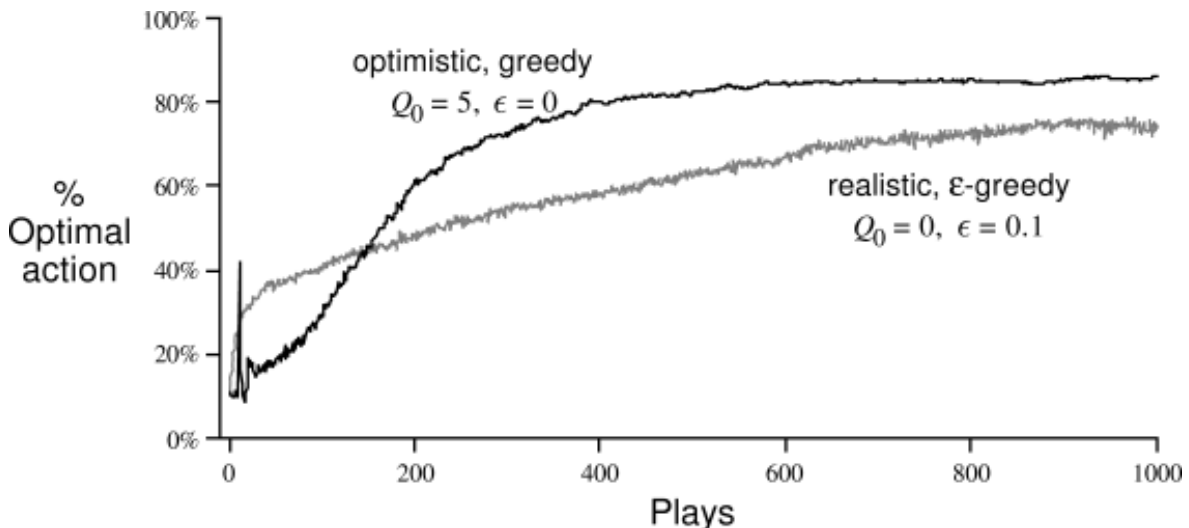


Figure 3: The effect of optimistic initial action-value estimates on the 10-armed testbed

#### 3.2.5.1.3. Upper-Confidence-Bound Action Selection

The other method for handling the exploration/exploitation problem is by using the special bounds that narrow with the number of steps taken. The formula is as follows:

$$A_t = \arg \max_a [Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}}$$

where:

- $\ln t$  is the natural logarithm of  $t$
- $N_t(a)$  - the number of times that action  $a$  has been selected prior to time  $t$
- $c$  - the exploration rate

The idea of this solution is that the square-root part is an uncertainty measure or variance in the  $a$  estimation. The use of the natural logarithm implies that overtime square-root term, so does the confidence interval, is getting smaller. All actions will be selected at some point, but the ones with non-optimal values for  $Q(a)$  are going to be selected much less frequently over time. UCB performs well, but it is harder to apply (generalize) for a broader amount of problems than  $\epsilon$ -greedy algorithm. Especially, when one is dealing with nonstationary problems. In such situations, algorithms more complex than those presented in this subsection should be selected.

Reinforcement learning algorithms can be classified into three general subcategories:

- Model Based - they are based on the idea that an model of the environment is known. Actions are chosen by searching and planning in this model. Markov Decision Process (MDP) is a typical example of such method since it requires knowledge of the Markov transition matrix and reward function.
- Model-free - it uses experience to learn in a direct way from state-action values or policies. They can achieve the same behaviour, but without any knowledge on the world model an agent acts in. In practical examples, reinforcement learning is primarily used for environments where a transition matrix is not known. Given a policy, a state has some value, which is defined as cumulated utility (reward) starting from the state. Model-free methods are generally less efficient than model-based ones because information about the environment is combined with possibly incorrect estimates about state values Dayan and Niv (2008).

### 3.2.5.2. Model-based Methods in Reinforcement Learning

- Value iteration - a model-based algorithm that computes the optimal state value function by improving the estimate of  $V(s)$ . It starts with initializing arbitrary values and then updates  $Q(s, a)$  and  $V(s)$  values until they converge. The pseudocode is as follows:
- Policy iteration - while in the previous bullet the algorithm is improving value function, policy iteration is based on the different approach. Concretely, there are two functions to be optimized  $V^\pi(s)$  and  $\pi'(s)$ . This method is based on the premise that a RL agent cares about finding out the right policy. Sometimes, it is more convenient to directly use policies as a function of states. Below is the pseudocode:

### 3.2.5.3. Model Free Learning

Model Free learning is a subcategory of reinforcement learning algorithms which are used when a model is not known. The agent improves its accuracy in choosing right actions by

interacting with the environment without explicit knowledge of the underlying transition matrix. It fits trading conditions - in financial markets it is impossible to know what the model is and what probabilities in the transition matrix are (they are not stationary). Hence, value or policy iteration algorithm can not be used directly.

Even though, Markov Decision Process and its element are not visible, the agent can gain experience from sampled states. It is assumed that eventually the distribution of sampled states will converge to the one in the transition matrix. So do  $Q(s, a)$  converge to  $Q^*$  and  $\pi^*$  to the optimal policy. The conditions required by the convergence is that all state-action pairs were visited infinite times and the agent is greedy once it finds the best action in every state.

#### 3.2.5.4. Components of an reinforcement learning system

Reinforcement learning systems are developed to solve sequential decision making problems, to select such actions that eventually maximize cumulative discounted future rewards. In the following section the author explained components of reinforcement learning on the example of game of chess and trading. The subsection was partially inspired and based on Sutton and Barto (2017).

- Environment ( $E$ ) - it defines what states and actions are possible. In the game of chess it is the whole set of rules and possible combination of figures on the chessboard. It must be stated that some states are not available and will be never reached. In trading such rules might constitute that for instance the only position an agent can take is 0 or 1, or that weights of assets in a portfolio must sum up to 1.
- State ( $s$ ) - can be seen as a snapshot of the environment. It contains a set of information in time  $t$  that a RL agent uses to pick the next action. States can be terminal, i.e. the agent will no longer be able to choose any action. In such scenario they end an episode (epoch), a sequence of state-action pairs from the start to the end of the game. For a trading application, a state in time  $t$  can be a vector of different financial measures, such as rate of return, implied/realized volatility, moving averages, economics measures, technical indicators, market sentiment measures, etc.
- Action ( $a$ ) - given a current state the agent chooses an action which directs him into a new state, either deterministically or stochastically. The action choice process itself may also be deterministic or based on probability distributions. In the game of chess analogy, an action is to move a figure in accordance to the game's rules. In trading it could be for instance going long, short, staying flat, outweighing.
- Policy ( $\pi$ ) - a policy is a mapping from state of the environment to action to be taken in that state. In psychology it is called a set of stimulus, i.e. response rules or associations. The policy might be a lookup table or a simple function (e.g. linear), but not necessarily. Especially in trading where variables are often continuous extensive computations to set up a satisfying outcome take place. The policy is the most essential part of a reinforcement learning agent because they determine how it behaves. It may be stochastic. Policies do not imply deterministic nature of the mapping. Even after countless number of episodes and states, there is a chance that an efficient RL algorithm will explore other states rather than by exploiting the then-optimal action
- Value Function - it is a prediction of future, usually discounted rewards. Value functions are used for determining how much a state should be desired by the agent. They depend



on initial states ( $S_0$ ), and a policy that is picked up by the agent. Every state should have an associated value, even if the path it is part of was never explored - in such cases they usually equal to zero. The general formula for value function is as follows:

$$V^\pi = \mathbb{E}_\pi \left[ \sum_{k=1}^{\infty} \gamma^k r_{t+k} | s_t = s \right]$$

where  $\gamma$  is a discount factor from the range  $[0; 1]$ . It measures how much more instant rewards are valued. The smaller it is the more immediate values are relatively more relevant and cause algorithm to be more greedy. Sometimes  $\gamma$  is equal to 1 if it is justified by the design of the whole agent.

Value estimation, as a area of research in RL is probably the most vital one in the last decade. The most important distinction between different RL algorithms lies as to how it is calculated, in what form, and what variables it incorporates.

- Reward ( $r$ ) - rewards are the essence of reinforcement learning predictions. Value function, as previously stated, is a sum of, often discounted, rewards. Without them, as components of value function, an agent would not be able to spot (or optimal) better policies actions are based on. Hence, it is assumed rewards are the central point, required element, of every RL algorithm. Rewards are always given as a scalar, single value that is retrieved from the environment, that is easy to observe and interpret. With value function it is much harder since it can be obtained only by calculating a sequence of observations a RL agent makes over its lifetime.
- Model ( $m$ ) - a model shows the dynamics of environment, how it will evolve from  $S_{t-1}$  to  $S_t$ . The model helps in predicting what the next state and next reward will be. They are used for planning, i.e. trial-and-error approach is not needed in order to achieving the optimal policy. Formally, it is a set of transition matrices:

$$\begin{aligned} \mathbb{P}_{ss'}^a &= \mathbb{P}[s' | s, a] \\ \mathbb{R}_s^a &= \mathbb{E}[r | s, a] \end{aligned}$$

where:

- $\mathbb{P}_{ss'}^a$  is a matrix of probability of transitions from state  $s$  to state  $s'$  when taking action  $a$ . Analogously,  $\mathbb{R}_s^a$  is an expected value of reward when an agent is in state  $s$  and taking action  $a$

### 3.2.5.5. Limitations

Reinforcement learning is not a panacea for all kinds of ML problems, they should be heavily associated with problems based on states, some policy to determine and defined value function. A state is just a signal that reaches the agent, a snapshot of the environment at a time. Most of pure reinforcement learning methods are oriented about states and their values. Even though it is useful for simpler environments, for more sophisticated ones it is not as easy. First of all, tabular data is not a good way to store information about expected values for states/states-actions. They would not fit into memory as variables are continuous. Hence, some additional approaches must be used in order to solving the problem efficiently.

#### **3.2.5.6. Data Handlers**

#### **3.2.5.7. Reference Data**

#### **3.2.5.8. Optimization**

#### **3.2.5.9. Execution**

#### **3.2.5.10. UI / UX**

### **3.2.6. Rule-Based Trading Strategies**

A rule-based trading strategy uses explicit, human-defined rules for buying and selling assets. These interpretable rules don't rely on forecasting models or machine-generated predictions.

These strategies employ conditional logic where trading decisions follow deterministic paths based on market conditions. Rules typically derive from established market heuristics or technical analysis principles with proven effectiveness.

The clear causal relationship between market conditions and trading actions provides transparency, allowing traders to understand exactly why positions change. Decisions respond to current market states rather than attempting to predict future movements.

Rule-based systems apply consistently across all market conditions, eliminating psychological biases and ensuring reproducible trading decisions. This deterministic nature enables straightforward performance evaluation

The framework adapts to changing markets while maintaining its structure. Parameter optimization enhances performance while preserving core logic, allowing systematic evolution through iterative refinement.

### **3.3. Trading based on Forecasts**

In such a system forecast is optimized to produce price forecasts from a set of inputs. Supervised learning techniques are used, basing on statistics such as mean squared error. Minimizing an error is only an intermediate step - the outcome of Forecasting System is then used for buying or selling decisions for analyzed asset(s) inside Trading Rules. Then the latter module is subject to evaluation module (Profits/Losses  $U(\theta, \theta')$ ) which consists of some financial measure. It must be noted that Forecasting module outputs only that predicted price while leaving the original inputs. It can be inefficient since Trading Rules sometimes might need more than just price. <sup>2</sup>

#### **3.3.0.1. Training a Trading System on labeled Data**

Such a trading system is based on the idea of direct integration between Trading System and Input Series. Trades (signals) are based on labelled trades (training set), and actual trades take place basing on input (Input Series).

---

<sup>2</sup>In empirical literature, the banking sector is identified as a primary industry implementing blockchain technology.

Its efficiency and effectiveness rely on how well Trading Module can utilize information from Input Series and Labeled Trades  $\theta'$ . Since it's Trading System, not utility function  $U(\theta, \theta')$  optimized, also in this case the system tends to be sub-optimal.

### **3.3.0.2. Direct Optimization of Performance**

In this modern approach, there is no intermediate step and labeled data is not given. The environment is observed,  $X_t$ , the system carries out an action, and receives a scalar reward for its past activities, representing the trading performance in some form (e.g. rate of return). Based on this reward, the system alters the way it behaves in subsequent episodes and steps.

## **3.4. Algorithmic Trading Systems**

In its essence, the investor's or trader's main goal is to optimize some measure of trading system performance, such as risk-adjusted return, economic utility, or simple profit. In the work the author presented direct optimization methods, based on reinforcement learning. It is flexible and can work out for trading a single asset (+ risk-free instrument), but also a portfolio consisting of n-instruments. There are also other options which often direct to non-optimal solutions. In the following section the author brought up different types of algorithmic systems and outlines advantages, disadvantages and differences between them.

## **4. Reinforcement Learning**

## **5. Design of the Trading Agent**

### **5.1. Action Space**

#### **5.1.1. Book Pressure-Related Variables**

#### **5.1.2. Volume-Related Variables**

#### **5.1.3. Last Trade Price-Related Variables**

#### **5.1.4. Time-Related Variables**

#### **5.1.5. Technical Indicators**

#### **5.1.6. Autoregressive Variables**

### **5.2. Reward Function**

#### **5.2.1. Differential Sharpe Ratio**

#### **5.2.2. Other Reward Functions**

### **5.3. Value Function**

### **5.4. Policy**

### **5.5. Step Size**

## **6. Implementation of the Trading Agent**

The following assumptions are established for the trading mechanisms and algorithmic simulation framework: (b) Zero latency exists between market data reception and trade execution

when a favorable signal is present. (a) The absence of market competition enables execution at counterparty prices (selling at bid price and purchasing at ask price). (c) Position constraints limit holdings to  $x$  contract(s) at any given moment, with transactions restricted to whole contract units. (d) A transaction fee of  $x\%$  of the execution price is applied to all trades.

## **6.1. Data Preparation**

### **6.1.1. Data Collection**

### **6.1.2. Data Preprocessing**

### **6.1.3. Feature Engineering**

### **6.1.4. Data Splitting**

## **6.2. Code**

### **6.2.1. Code Structure**

### **6.2.2. Code Implementation**

## **7. Empirical Evaluation and Performance Analysis**

### **7.1. Statistical Validation**

### **7.2. Robustness Assessment**

### **7.3. Performance Evaluation**

## **8. Conclusions and Future Work**

### **8.1. Summary of Findings**

### **8.2. Limitations and Future Research**

### **8.3. Implications for Trading Systems**

### **8.4. Recommendations for Practitioners**

### **8.5. Conclusion**

**List of Figures**

- 1. Efficient Frontier . . . . .
- 2. Data Science Graph . . . . .
- 3. The effect of optimistic initial action-value estimates on the 10-armed testbed

## Bibliography

- Dayan, Peter, and Yael Niv. 2008. "Reinforcement learning: The Good, The Bad and The Ugly." <https://doi.org/10.1016/j.conb.2008.08.003>.
- Fama, E. F.;Malkiel, B. G. 1970. "Efficient capital markets: A review of theory and empirical work." *Journal of Finance* 25 (2): 383–417.
- Fama, Eugene F. 1965. "The Behavior of Stock-Market Prices." *The Journal of Business* 38 (1): 34. <https://doi.org/10.1086/294743>.
- Haug, Mark, and Mark Hirschey. 2006. "The january effect." <https://doi.org/10.2469/faj.v62.n5.4284>.
- Jason, Brownlee. 2013. "What is Machine Learning?" <https://machinelearningmastery.com/what-is-machine-learning/>.
- Lintner, John. 1965. "Security prices, risk, and maximal gains from diversification." *Jf* 20 (4): 587–615. <https://doi.org/10.1111/j.1540-6261.1965.tb02930.x>.
- Malkiel, Burton G. 2003. "The Efficient Market Hypothesis and Its Critics." *Journal of Economic Perspectives* 17 (1): 59–82. <https://doi.org/10.1257/089533003321164958>.
- Marsland, Stephen. 2009. *Machine Learning: An Algorithmic Perspective*. [https://doi.org/10.1111/j.1751-5823.2010.00118\\_11.x](https://doi.org/10.1111/j.1751-5823.2010.00118_11.x).
- Mitchell, Tom M. 1997. *Machine Learning*. 1. <https://doi.org/10.1145/242224.242229>.
- Mosic, Ranko. 2017. "Deep Reinforcement Learning Based Trading Application at JP Morgan Chase." *Medium*. <https://medium.com/@ranko.mosic/reinforcement-learning-based-trading-application-at-jp-morgan-chase-f829b8ec54f2>.
- Mossin, Jan. 1966. "Equilibrium in a capital asset market." *Econometrica* 34 (4): 768–83. <https://doi.org/10.2307/1910098>.
- Samuel, Arthur L. 1959. "Some Studies in Machine Learning Using the Game of Checkers Some Studies in Machine Learning Using the Game of Checkers." *IBM Journal* 3. <https://www.cs.virginia.edu/%7B~%7DDevans/greatworks/samuel1959.pdf>.
- Sharpe, William F. 1964. "Capital Asset Prices: A Theory of Market Equilibrium under Conditions of Risk." *The Journal of Finance* 19 (3): 425–42. <https://doi.org/10.2307/2329297>.
- Shen, Lucinda. 2017. "Here's How Much the Top Hedge-Fund Manager Made Last Year." <http://fortune.com/2017/05/16/hedge-fund-james-simons-renaissance-technologies/>.
- Stulz, Rene M. 1995. "American Finance Association, Report of the Managing Editor of the Journal of Finance for the Year 1994." *The Journal of Finance* 50 (3): 1013. <https://doi.org/10.2307/2329297>.
- Sutton, R S, and A G Barto. 2017. "Reinforcement Learning: An Introduction." *Neural Networks IEEE Transactions on* 9 (2): 1054. <https://doi.org/10.1109/TNN.1998.712192>.
- Turner, Matt. 2015. "The robot revolution is coming for Wall Street traders." <http://www.businessinsider.com/robots-to-replace-wall-street-traders-2015-8?IR=T>.
- Weber, Thomas A. 2012. "Price Theory in Economics." In *The Oxford Handbook of Pricing Management*. <https://doi.org/10.1093/oxfordhb/9780199543175.013.0017>.