



# Reinforcement Learning for Portfolio Management

Thesis Presentation

---

Angelos Filos

Supervisor: Professor Danilo Mandic

Electrical & Electronic Engineering, Imperial College London

June 26, 2018

# Table of Contents

## 1. Background

- 1.1. Financial Signal Processing
- 1.2. Portfolio Optimization
- 1.3. Reinforcement Learning

## 2. Innovation

- 2.1. Financial Markets: Stochastic Dynamical Systems
- 2.2. Model-based Recurrent Neural Agent
- 2.3. Model-free Mixture of Score Machines
- 2.4. Pre-Training to Quadratic Programming

## 3. Experiments

- 3.1. Deterministic Signals
- 3.2. Market Simulation: Surrogate Data
- 3.3. S&P 500 and EURO STOXX 50

## **Background**

---

# Financial Terms & Concepts

## Asset

An **asset** is an item of economic value, such as: cash (i.e., in bank or in hand), stocks and loans. Useful abbreviations (Kennedy, 2016) include AAPL for Apple and GE for General Electric.

## Portfolio

A **portfolio** is a collection of multiple financial assets (i.e., master asset), characterized by its:

- **Constituents:**  $M$  assets of which it consists;
- **Portfolio vector,**  $\mathbf{w}_t$ : its  $i$ -th component represents the ratio of the total budget invested to the  $i$ -th asset, given by:

$$\mathbf{w}_t = [w_{1,t}, \dots, w_{M,t}]^T \in \mathbb{R}^M$$

such that:  $\sum_{i=1}^M w_{i,t} = 1$  (1)

## Problem Definition

The aim of this report is to investigate the effectiveness of context-independent Reinforcement Learning agents on sequential asset allocation\*. A finite universe of financial instruments, assets, such as stocks, is selected and the role of an agent is to construct an internal representation (model) of the market, allowing it to determine how to *optimally allocate funds* of a finite budget to those assets.

The agent is trained on both synthetic and real market data. Then, its performance is compared with standard portfolio management algorithms on an out-of-sample dataset; data that the agent has not been trained on (i.e., test set).

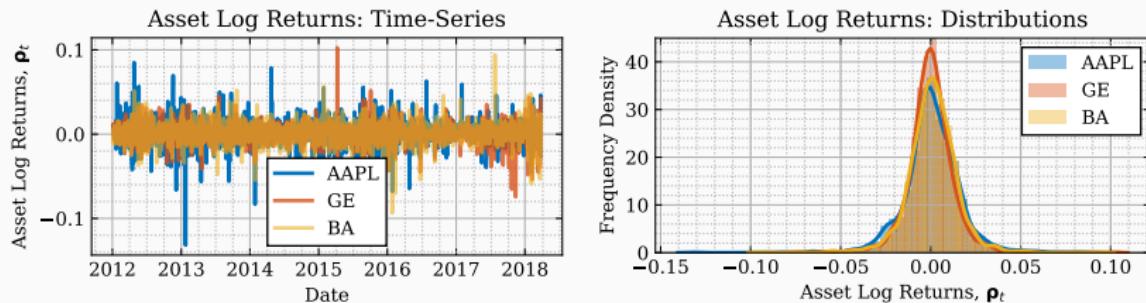
---

\* The terms *Asset Allocation* and *Portfolio Management* are used interchangeably throughout the presentation (Zhang and Wang, 2017).

# Financial Time-Series

Let  $\rho_{i,t}$  the **log return**, or the difference of the log prices,  $\log(p_{i,t})$ , of asset  $i$  from time  $(t - 1)$  to time  $t$ , given by:

$$\rho_{i,t} \triangleq \log\left(\frac{p_{i,t}}{p_{i,t-1}}\right) \in \mathbb{R} \quad (2)$$



**Figure 1:** Log returns time-series (left) and distributions (right).

The log returns are recorded per asset at each time step, hence they form a multivariate time-series,  $\vec{P} \in \mathbb{R}^{M \times T}$ .

The **expected value** of log returns:

$$\mathbb{E}[\rho] = \mu_\rho \approx \frac{1}{T} \sum_{t=1}^T \rho_t^{(\text{portfolio})} = \mathbf{w}_t^T \boldsymbol{\mu}_\rho \in \mathbb{R} \quad (3)$$

## Greedy Criterion

For the same level of risk, choose the portfolio that maximizes the expected returns (Wilmott, 2007).

The **variance**, or **squared volatility** of the log returns:

$$\text{Var}[\rho] = \sigma_\rho^2 \approx \frac{1}{T-1} \sum_{t=1}^T (\rho_t - \mu_\rho)^2 = \mathbf{w}_t^T \boldsymbol{\Sigma}_\rho \mathbf{w}_t \in \mathbb{R} \quad (4)$$

## Risk-Aversion Criterion

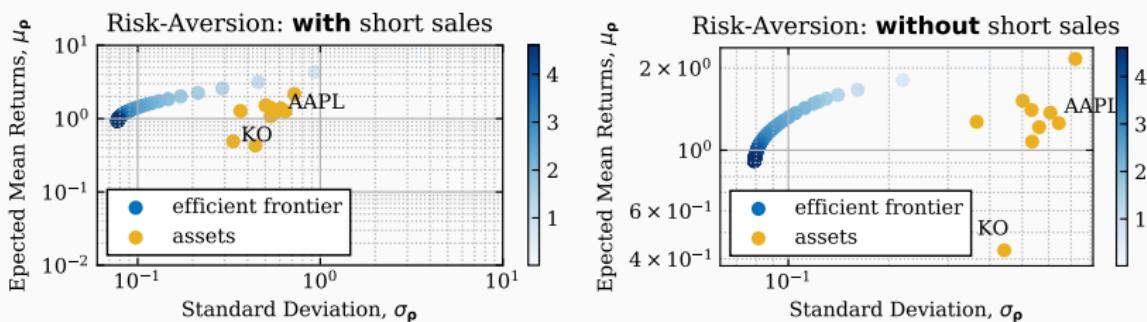
For the same expected returns, choose the portfolio that minimizes the volatility (Wilmott, 2007).

# Sequential Portfolio Theory

## Risk-Aversion with Transaction Costs

Construct a portfolio out of single assets so that to: balance profitability-risk trade-off and eliminate transaction cost effects:

$$\begin{aligned} & \underset{\mathbf{w}_t}{\text{maximize}} \quad \underbrace{\mathbf{w}_t^T \mu_p}_{\mu_p} - \underbrace{\mathbf{1}^T \beta \|\mathbf{w}_{t-1} - \mathbf{w}_t\|_1}_{\text{transaction consts}} - \alpha \underbrace{\mathbf{w}_t^T \Sigma_p \mathbf{w}_t}_{\sigma_p^2} \\ & \text{and } \mathbf{1}^T \mathbf{w}_t = 1 \\ & \text{and } \mathbf{w}_t \succeq 0 \end{aligned} \quad (5)$$



**Figure 2:** Risk-aversion with transaction costs efficient frontier for various  $\alpha$  values and fixed coefficient  $\beta = 0.0025$ .

## Innovation

---

# Financial Market as Discrete-Time Stochastic Dynamical System

- **Observation**,  $o_t \in \mathbb{O}$ , asset prices at  $t$ ;
- **Action**,  $a_t \in \mathbb{A}$ , portfolio vector at  $(t+1)$ ;
- **Reward signal**,  $r_t \in \mathbb{R}$ , reward signal at  $t$ ;
- **Environment state**,  $s_t \in \mathbb{S} \Rightarrow$  unobservable;
- **Reward generating function**,  $\mathcal{R} \Rightarrow$  framework hyperparameter;
- **Environment dynamics**  $\mathcal{P}_{ss'}^a \Rightarrow$  intractable.

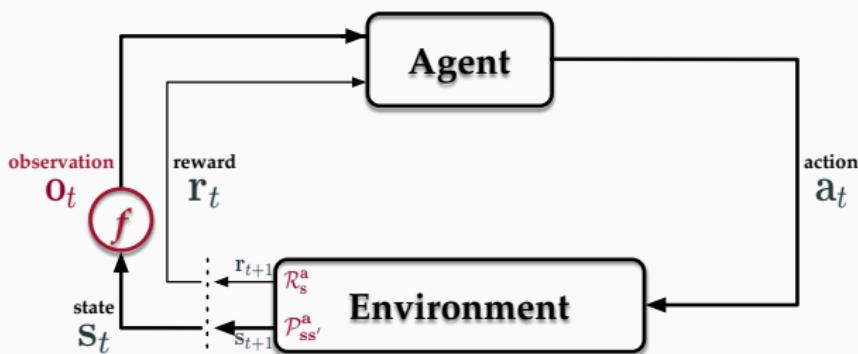


Figure 3: Partially Observable Markov Decision Process.

# Trading Agent

---

**Algorithm 1:** General setup for trading agents.

---

**inputs** : trading universe of  $M$ -assets; initial portfolio vector  $\mathbf{w}_1 = \mathbf{a}_0$ ; initial asset prices  $\mathbf{p}_0 = \mathbf{o}_0$ ; objective function  $\mathcal{J}$ .

**output**: optimal agent parameters  $\theta_*, \varphi_*$

```
1 repeat
2   | for  $t = 1, 2, \dots, T$  do
3   |   | observe 2-tuple  $\langle \mathbf{o}_t, r_t \rangle$ 
4   |   | calculate gradients  $\nabla_{\theta} \mathcal{J}(r_t)$  and  $\nabla_{\varphi} \mathcal{J}(r_t)$            // BPTT
5   |   | update agent parameters  $\theta, \varphi$ 
6   |   |   | using adaptive gradient optimizers                         // ADAM
7   |   | get estimate of agent state:  $s_t \approx f(\cdot, o_t)$ 
8   |   | sample and take action:  $\mathbf{a}_t \sim \pi(\cdot | s_t; \varphi)$     // portfolio rebalance
9   | end
10 until convergence
11 set  $\theta_*, \varphi_* \leftarrow \theta, \varphi$ 
```

---

# Model-based Reinforcement Learning

- System Identification such that:

$$\underbrace{\hat{P}_{ss'}^a}_{\text{model}} \approx \underbrace{P_{ss'}^a}_{\text{environment}} \quad (6)$$

- Predictive model & planning:

$$s_t \xrightarrow{\mathcal{P}_{ss'}^a} s_{t+1} \xrightarrow{\mathcal{P}_{ss'}^a} \cdots \xrightarrow{\mathcal{P}_{ss'}^a} s_{t+L}, \quad a_{t+1} \equiv \max_{a \in \mathbb{A}} \mathcal{J}(a | a_t, s_t, \dots, s_{t+L}) \quad (7)$$

- Implementations: vector autoregressive processes (VAR), recurrent neural network (RNN) and Gaussian processes (GP).

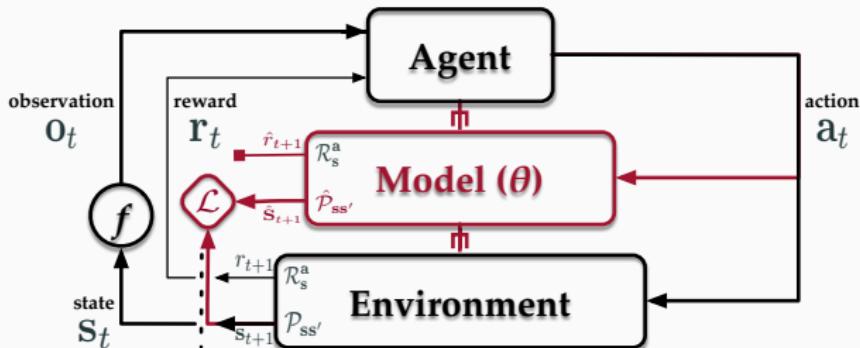


Figure 4: General setup for model-based reinforcement learning.

# Model-free Reinforcement Learning

- Multi-stage decision making:

$$a_{\geq t} = \underset{a \in \mathbb{A}}{\operatorname{argmax}} \sum_{\tau=1}^{\infty} \gamma^{\tau} r_{t+\tau}, \quad \gamma \in [0, 1] \quad (8)$$

- Direct parametrization of policy  $\pi$  or/and action value function  $q$ :

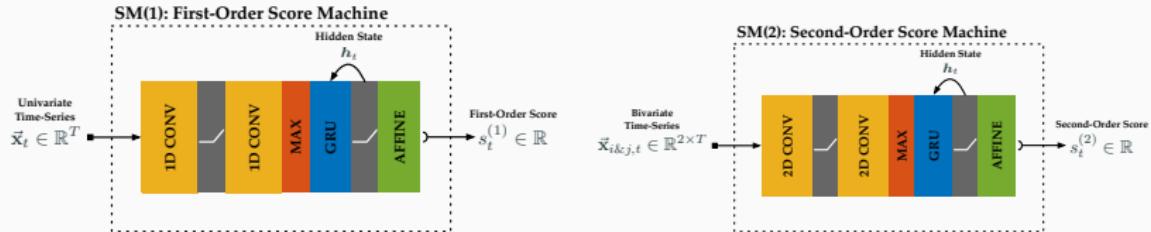
$$\pi(\cdot | s; \theta) : \mathbb{A} \times \mathbb{S} \rightarrow [0, 1] \quad (9)$$

$$q(\cdot | s; \theta) : \mathbb{A} \times \mathbb{S} \rightarrow \mathbb{R} \quad (10)$$

- Given objective  $\mathcal{J}$ , follow gradient:

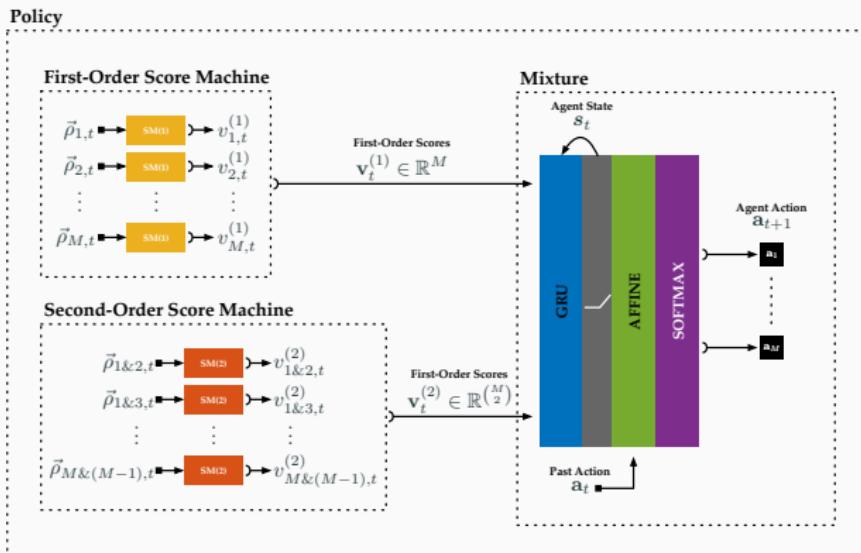
$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \mathcal{J}(\theta) \quad (11)$$

# Score Machine



**Figure 5:** Score Machine (SM) neural network architecture. Convolutional layers followed by non-linearities (e.g., ReLU) and Max-Pooling (Giusti *et al.*, 2013) construct a feature map, which is selectively stored and filtered by the Gate Recurrent Unit (GRU) layer. Finally, a linear layer combines the GRU output components to a single scalar value, the score. (*Left*) First-Order Score Machine SM(1); given a univariate time-series  $\vec{x}_t$  of  $T$  samples, it produces a scalar score value  $s_t^{(1)}$ . (*Right*) Second-Order Score Machine SM(2); given a bivariate time-series  $\vec{x}_{i\&j,t}$ , with components  $\vec{x}_{i,t}$  and  $\vec{x}_{j,t}$ , of  $T$  samples each, it generates a scalar value  $s_t^{(2)}$ .

# Mixture of Score Machines



**Figure 6:** Mixture of Score Machines (MSM) architecture. The historic log returns  $\vec{p}_{t-T \rightarrow t} \in \mathbb{R}^{M \times T}$  processes by the score machines SM(1) and SM(2), which assign scores to each asset ( $\mathbf{v}_t^{(1)}$ ) and pair of assets ( $\mathbf{v}_t^{(2)}$ ), respectively. The scores concatenated and passed to the mixture network, which combines them with the past action (i.e., current portfolio vector) and the generated agent state (i.e., state manager hidden state) to determine the next optimal action.

## Weaknesses

Trading Agents Comparison Matrix: 12-assets of S&P 500					
	Model-Based		Model-Free		
	VAR	RNN	DSRQN	REINFORCE	MSM
Non-Stationarity	✗	✓	✓	✓	✓
Long-Term Memory	✗	✓	✓	✓	✓
Non-Linear Model	✗	✓	✓	✓	✓
End-to-End	✗	✗	✗	✓	✓
Linear Scaling	✗	✗	✗	✗	✓
Universality	✗	✗	✗	✗	✓
Low Variance	✗	✗	✗	✗	✗
Short Sales	✓	✓	✗	✗	✗

Table 1: Comprehensive comparison of trading agent weaknesses.

## Pre-Training i

---

**Algorithm 2:** Pre-training supervised dataset generation.

---

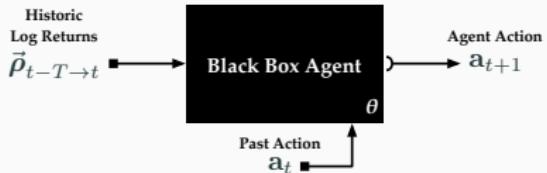
**inputs :** number of pairs to generate  $N$   
number of assets in portfolio  $M$   
look back window size  $T$   
transaction costs coefficient  $\beta$

**output:** dataset  $\{(\mathbf{X}_i, \mathbf{y}_i)\}_{i=1}^N$

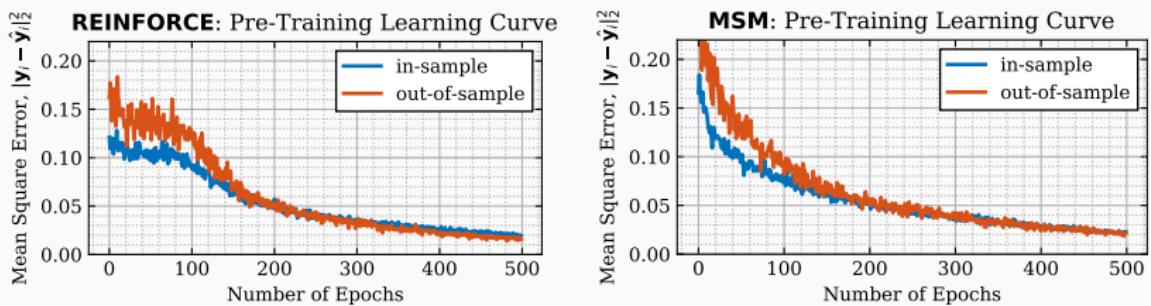
```
1 for  $i = 1, 2, \dots, N$  do
2   sample valid random initial portfolio vector  $\mathbf{w}_{t_i}$ 
3   sample random lower triangular matrix  $\mathbf{L} \in \mathbb{R}^{M \times M}$            // Cholesky
   decomposition
4   sample randomly distributed log returns:  $\vec{\rho}_{t_i - T \rightarrow t_i} \sim \mathcal{N}(\mathbf{1}, \mathbf{L}\mathbf{L}^T)$ 
5   calculate empirical mean vector of log returns:  $\mu = \mathbb{E}[\vec{\rho}_{t_i - T \rightarrow t_i}]$ 
6   calculate empirical covariance matrix of log returns:  $\Sigma = \text{Cov}[\vec{\rho}_{t_i - T \rightarrow t_i}]$ 
7   determine  $\mathbf{a}_{t_i+1}$  by solving quadratic program (??)
8   set  $\mathbf{X}_i = [\vec{\rho}_{t_i - T \rightarrow t_i}, \mathbf{a}_{t_i}]$  and  $\mathbf{y}_i = \mathbf{a}_{t_i+1}$ 
9 end
```

---

## Pre-Training ii



**Figure 7:** Interfacing with policy gradient agents as black boxes, with inputs (1) historic log returns  $\vec{\rho}_{t-T \rightarrow t}$  and (2) past action (i.e., current portfolio vector)  $a_t$  and output the next agent actions  $a_{t+1}$ .

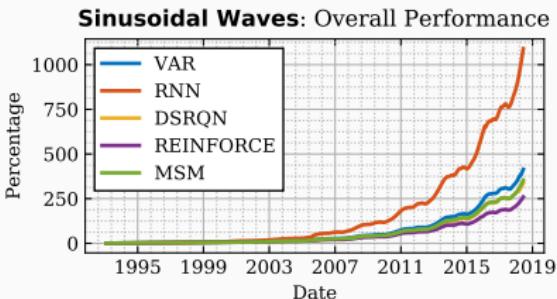
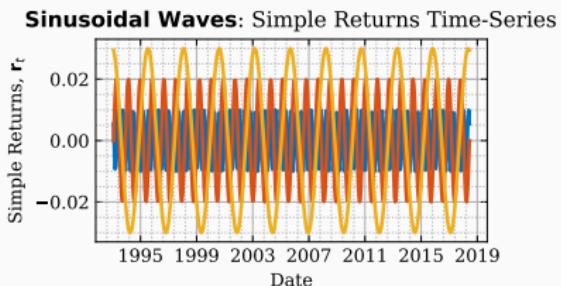


**Figure 8:** Pre-training mean square error (MSE) for policy gradient agents.

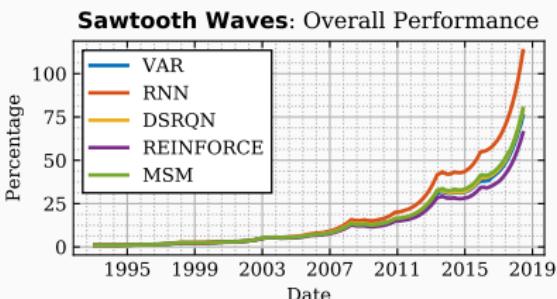
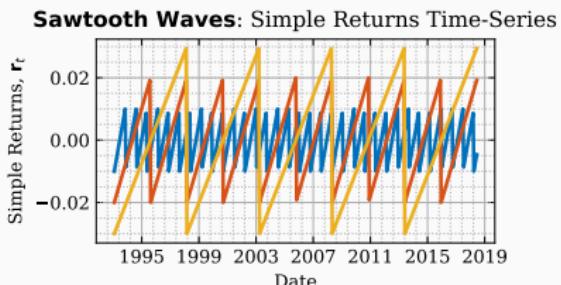
## Experiments

---

# Deterministic Synthetic Data



**Figure 9:** Synthetic universe of deterministic sinusoidal waves. (*Left*) Example series from universe. (*Right*) Cumulative returns of reinforcement learning trading agents.



**Figure 10:** Synthetic universe of deterministic sawtooth waves. (*Left*) Example series from universe. (*Right*) Cumulative returns of reinforcement learning trading agents.

# Surrogate Data

---

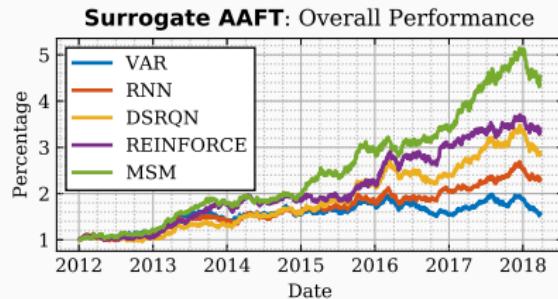
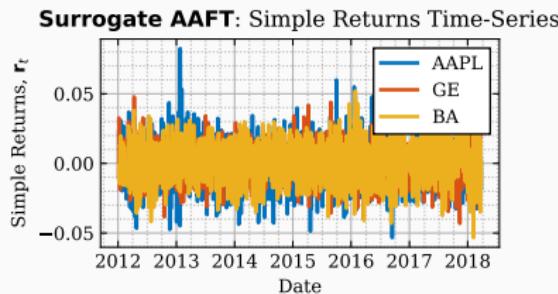
**Algorithm 3:** Amplitude Adjusted Fourier Transform (AAFT).

---

**inputs :**  $M$ -variate original time-series  $\vec{X}$   
**output:**  $M$ -variate synthetic time-series  $\hat{\vec{X}}$

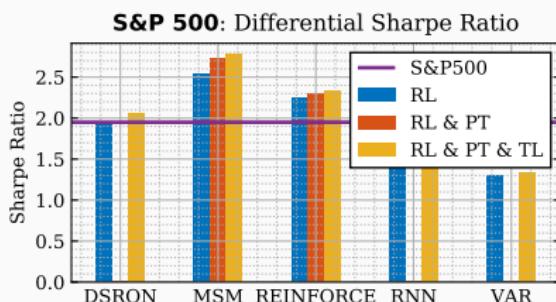
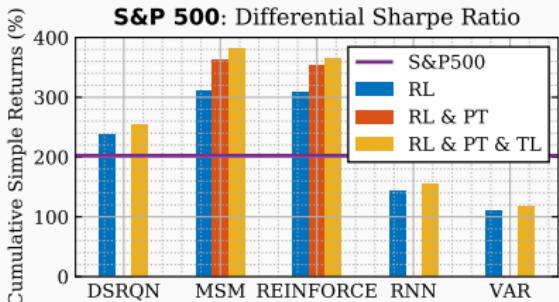
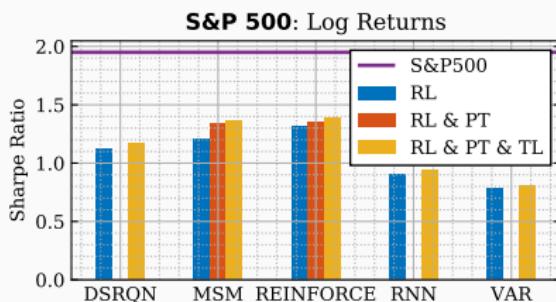
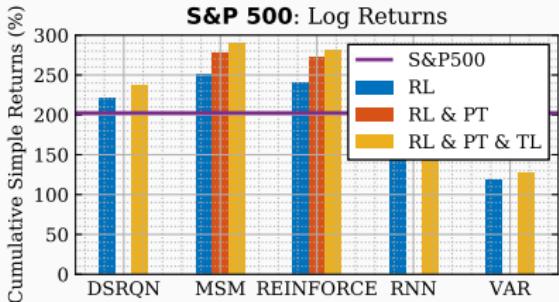
```
1 for  $i = 1, 2, \dots, M$  do
2     calculate Fourier Transform of univariate series  $\mathfrak{F}[\vec{X}_{:,i}]$ 
3     randomize phase component
4     calculate Inverse Fourier Transform of unchanged amplitude and
        randomized phase  $\hat{\vec{X}}_{:,i}$ 
5 end
```

---



**Figure 11:** Prices time-series (left) and distributions (right).

# Standard & Poor's 500



**Figure 12:** Comparison of reinforcement learning trading agents on S&P 500.

Trading Agents Comparison Matrix: S&P 500				
Reward Generating Function	Differential Sharpe Ratio		Log Returns	
	Cumulative Returns (%)	Sharpe Ratio	Cumulative Returns (%)	Sharpe Ratio
SPY	202.4	1.95	202.4	1.95
RNN	142.3	1.49	146.2	0.91
DSRQN	237.1	1.96	221.5	1.12
MSM	310.8	2.53	251.9	1.21
REINFORCE & PT	353.6	2.29	272.7	1.33
MSM & PT & TL	381.7	2.77	291.0	1.36

**Table 2:** Comprehensive comparison of evaluation metrics of reinforcement learning trading algorithms and their variants, namely pre-training (PL) and transfer learning (TL). Improvement 9.2% in annualized cumulative returns and 13.4% in annualized Sharpe Ratio, compared to the most recent models in (Jiang et al., 2017).

## Conclusions

1. Model-based Reinforcement Learning is appropriate when predictive models are accurate, while model-free Reinforcement Learning when environment dynamics are chaotic, but true objective is tractable;
2. Suggested a Unified and versatile framework for training agents and investment strategies, based on a mathematical formulation of financial markets as discrete-time stochastic dynamical systems;
3. A comprehensive account of reinforcement agents has been developed, including traditional, baseline agents, a universal agent by virtue of principles of parameter sharing (Bengio *et al.*, 2003) and transfer learning (Pan and Yang, 2010);
4. Overcame the barrier of training deep architectures due to lack of large datasets, via data augmentation (a generative approach) and a choice of pre-training strategies, both of which are validated against current state-of-the-art models.

## Future Work

1. **Interpretability:** Penalize model flexibility for understanding strategies (i.e., "gray-box" (Rico-Martinez *et al.*, 1994));
2. **Feature Engineering:** Extraction of financial signals and indicators to improve convergence properties;
3. **Reward Generating Functions:** Consider alternative risk metrics (e.g., VaR, CVaR);
4. **Probabilistic networks:** Replace point estimators with Bayesian Neural Networks (Nasrabadi, 2007);
5. **Exact Policy Gradient:** Fourier Policy Gradient (Fellows *et al.*, 2018).

## **Questions?**

**Email:** angelos.filos14@imperial.ac.uk;  
af814@ic.ac.uk;  
filos.angel@gmail.com

**Github:** [github.com/filangel/qtrader](https://github.com/filangel/qtrader);  
[FYP.Presentation.pdf](#);  
[FYP.Final-Report.pdf](#)

## References

---

-  Rico-Martinez, R, JS Anderson, and IG Kevrekidis (1994). "Continuous-time nonlinear signal processing: a neural network based approach for gray box identification". *Neural Networks for Signal Processing [1994] IV. Proceedings of the 1994 IEEE Workshop*. IEEE, pp. 596–605.
-  Bengio, Yoshua et al. (2003). "A neural probabilistic language model". *Journal of Machine Learning Research* 3.Feb, pp. 1137–1155.
-  Tsay, Ruey S (2005). *Analysis of financial time series*. Vol. 543. John Wiley & Sons.
-  Nasrabadi, Nasser M (2007). "Pattern recognition and machine learning". *Journal of Electronic Imaging* 16.4, p. 049901.
-  Wilmott, Paul (2007). *Paul Wilmott introduces quantitative finance*. John Wiley & Sons.

## Bibliography ii

-  Pan, Sinno Jialin and Qiang Yang (2010). "A survey on transfer learning". *IEEE Transactions on Knowledge and Data Engineering* 22.10, pp. 1345–1359.
-  Giusti, Alessandro et al. (2013). "Fast image scanning with deep max-pooling convolutional neural networks". *Image Processing (ICIP), 2013 20th IEEE International Conference on*. IEEE, pp. 4034–4038.
-  Goodfellow, Ian et al. (2016). *Deep learning*. Vol. 1. MIT press Cambridge.
-  Kennedy, Douglas (2016). *Stochastic financial models*. Chapman and Hall/CRC.
-  Jiang, Zhengyao, Dixin Xu, and Jinjun Liang (2017). "A deep reinforcement learning framework for the financial portfolio management problem". *arXiv preprint arXiv:1706.10059*.
-  Zhang, Xiao-Ping Steven and Fang Wang (2017). "Signal processing for finance, economics, and marketing: Concepts, framework, and big data applications". *IEEE Signal Processing Magazine* 34.3, pp. 14–35.

## Bibliography iii

-  Fellows, Matthew, Kamil Ciosek, and Shimon Whiteson (2018). "Fourier Policy Gradients". *CoRR*.

## **Disclaimer**

*This presentation is for informational purposes only and does not constitute an offer to sell, a solicitation to buy, or a recommendation for any security; nor does it constitute an offer to provide investment advisory or other services. Nothing contained herein constitutes investment advice or offers any opinion with respect to the suitability of any security, and any views expressed herein should not be taken as advice to buy, sell, or hold any security or as an endorsement of any security or company. In preparing the information contained herein, we have not taken into account the investment needs, objectives, and financial circumstances of any particular investor. Any views expressed and data illustrated herein were prepared based upon information, believed to be reliable, available to us at the time of publication. We make no guarantees as to their accuracy or completeness. All information is subject to change and may quickly become unreliable for various reasons, including changes in market conditions or economic circumstances.*