

Regime-Switching Reinforcement Learning for Portfolio Allocation in Pairs Trading

Author: Tsvetelina Ilieva

Responsible Professor: Frans Oliehoek

Supervisor: Fenghui Yu

1. Introduction

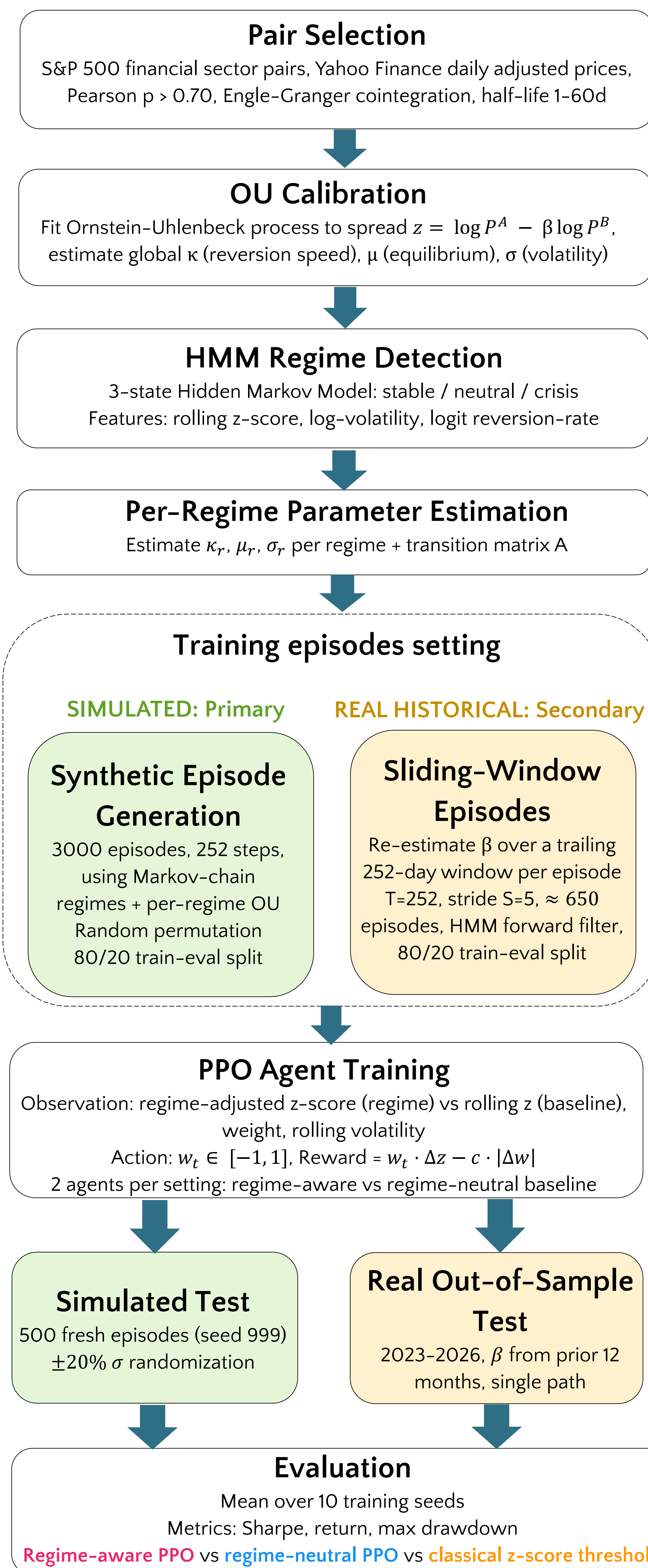
- **Pairs trading** is a statistical arbitrage strategy that exploits the **spread** between two historically **cointegrated** stocks, betting that temporary divergences from their long-run equilibrium will revert.
- That equilibrium **isn't stable**. Under different market conditions the spread exhibits different levels of volatility and mean-reverting behavior.
- Classical pairs-trading and RL approaches largely ignore that the spread's dynamics shift across **market regimes**. Existing regime-aware RL targets broad asset-class allocation, not the **cointegration structure specific to pairs trading**.
- We propose a **regime-aware RL agent** that infers the latent market regime and centers its trading signal on the **per-regime** equilibrium rather than a stale rolling mean, making regime information directly actionable **without increasing model complexity**.

2. Research Questions

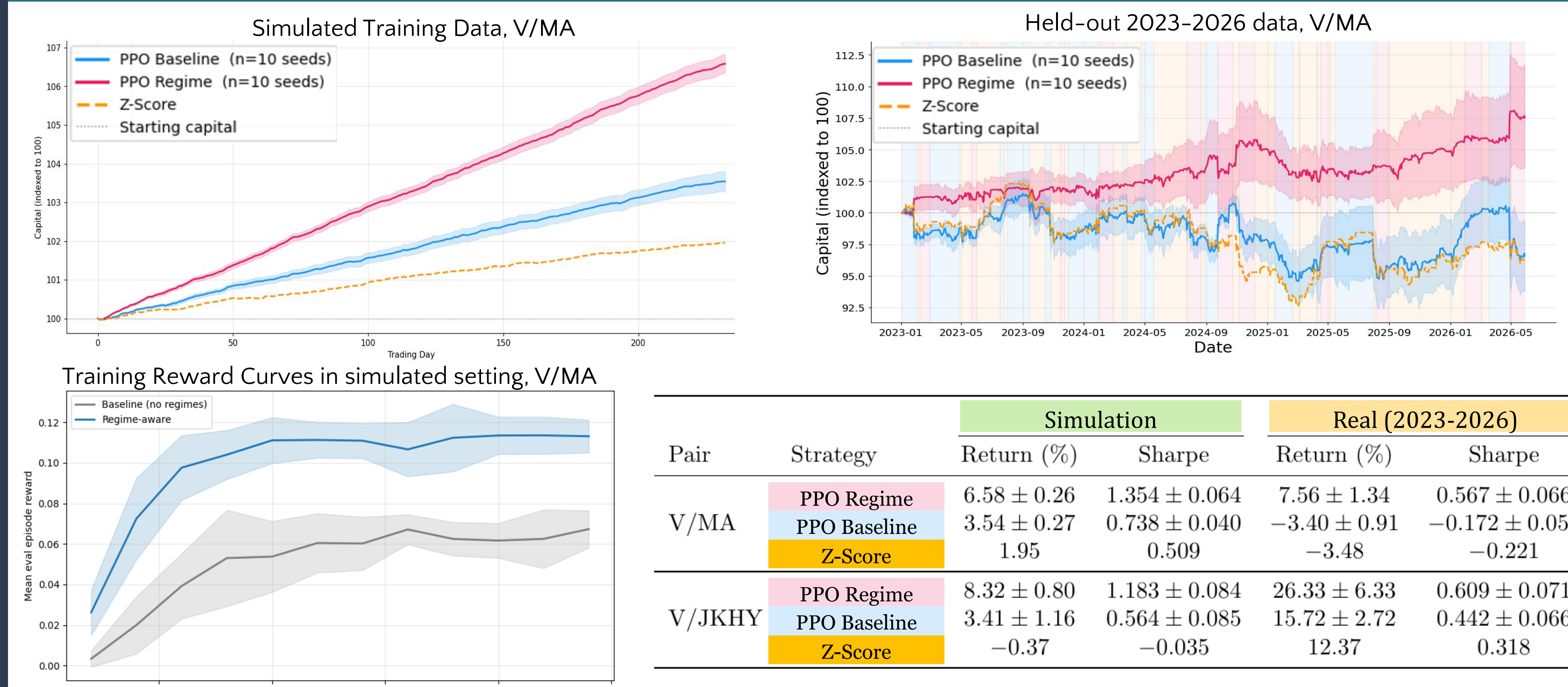
Does a regime-aware reinforcement learning agent outperform regime-neutral alternatives in portfolio allocation for pairs trading?

- Does a regime-aware Proximal Policy Optimization (PPO) agent outperform a regime-neutral PPO baseline trained on the same data?
- Does a regime-aware PPO agent outperform a classical z-score threshold strategy?
- Does the performance advantage of the regime-aware agent observed in simulation persist when both agents are trained and tested on real historical data?

3. Methodology



4. Results



- The regime-aware agent **outperforms the baseline and the z-score strategy in both settings**, though with **more variable results** in the real-data setting. The regime-aware agent reaches **higher mean episode reward during training**.
- In **simulation** the regime-aware agent nearly **doubles the baseline's Sharpe ratio** on V/MA (+0.616) at a lower drawdown, just as large as on V/JKHY (+0.619).
- The advantage **transfers to real 2023-2026 data**: on V/MA the regime agent earns a **positive Sharpe (0.567)** while both the regime-neutral baseline and the z-score **lose money** (-0.172 and -0.221). On V/JKHY all three strategies profit, yet the regime agent **still leads (0.609 vs 0.442)**.

5. Discussion

- The **per-regime equilibrium** signal, rather than a stale rolling mean, lets the agent spot high-confidence entries it would otherwise miss, without added model complexity.
- The holds only under a **stable cointegrating relationship**. Both pairs sit in payments/fintech, largely insulated from the 2022-2026 shocks, whereas sectors hit by **structural breaks** would likely degrade the framework regardless of the agent.
- The **sim-to-real gap** can be narrowed by enriching the simulator with **structural breaks** and **volatility clustering**. Replace Gaussian HMM emissions with **Student-t** to make regime detection robust to fat tails.
- Further extension include adding a **cash allocation** to exit non-cointegrated periods, a short-selling **margin constraint** for realistic risk, and **implicit regime learning** that removes the explicit HMM.