
Differentiable Multi-Period Portfolio Optimization (DMPO): Bridging Deep Sequence Modeling and Convex Control with Hard Constraints

Author One¹ Author Two¹

Abstract

Dynamic portfolio optimization in real-world financial markets requires agents to balance conflicting objectives—maximizing risk-adjusted returns while strictly adhering to physical constraints such as transaction costs, turnover limits, and leverage restrictions. Traditional two-stage frameworks, which first predict asset returns and then solve a mean-variance optimization problem, suffer from the “optimizer’s curse,” where estimation errors are magnified by the optimization step. Conversely, model-free Deep Reinforcement Learning (DRL) approaches often fail to guarantee feasibility for hard constraints, relying instead on unstable reward shaping or myopic projection layers.

To bridge this gap, we introduce **Differentiable Multi-Period Portfolio Optimization (DMPO)**, a hybrid end-to-end learning framework that integrates deep sequence modeling with differentiable convex optimization. We formulate the portfolio management problem as a Receding Horizon Control (RHC) task, where a Long Short-Term Memory (LSTM) network learns to output the parameters (e.g., expected returns and risk aversion) of a quadratic programming (QP) layer. By leveraging the implicit function theorem, DMPO analytically computes the gradients of the optimal solution with respect to the input parameters through the Karush-Kuhn-Tucker (KKT) conditions, enabling efficient backpropagation. This architecture allows the model to learn a “decision-aware” market view that directly optimizes the Sharpe ratio while rigorously satisfying multi-period turnover and risk constraints. Theoretical analysis confirms that our method avoids the pitfalls of sub-optimal projection and achieves

convergence in non-convex policy landscapes.

1. Introduction

The fundamental challenge of quantitative finance lies in constructing dynamic portfolios that can adapt to non-stationary market regimes while managing trading frictions and risk constraints. Since the seminal work of Markowitz (Markowitz, 1952), Mean-Variance Optimization (MVO) has served as the bedrock of modern portfolio theory (MPT). However, classical MVO faces significant practical hurdles: it is highly sensitive to estimation errors in expected returns and covariance matrices, a phenomenon often referred to as the “optimizer’s curse” (Michaud, 1989).

To mitigate these issues, the industry has largely adopted a “predict-then-optimize” (two-stage) paradigm. In the first stage, machine learning models (e.g., ARIMA, XGBoost, or Transformers) are trained to minimize a prediction error metric, such as Mean Squared Error (MSE). In the second stage, these predictions are fed into a numerical solver to generate portfolio weights. While intuitive, this approach is fundamentally flawed because the loss function used for training (prediction error) is misaligned with the ultimate objective (portfolio utility). A small error in return prediction can lead to a drastically different and sub-optimal portfolio allocation due to the ill-conditioned nature of the covariance matrix (Bengio et al., 2020).

In recent years, **Deep Reinforcement Learning (DRL)** has emerged as a promising alternative, aiming to learn a direct mapping from market states to portfolio weights (end-to-end learning) (Jiang et al., 2017; Liang et al., 2018). DRL agents, trained via Policy Gradient (PG) or Q-learning, can theoretically capture complex patterns and optimize non-differentiable rewards like the Sharpe ratio. However, standard DRL algorithms (e.g., PPO, DDPG) struggle significantly with **hard constraints**. Financial portfolios must adhere to strict constraints—such as non-negativity (no short selling), budget constraints ($\sum w_i = 1$), and, crucially, turnover constraints to limit transaction costs. Existing DRL solutions typically handle these constraints via:

1. *Softmax/Normalization*: Ensures budget constraints

¹School of Mathematical Sciences, Peking University, Beijing, China. Correspondence to: Author One <email@pku.edu.cn>.

but cannot handle complex inequality constraints like turnover limits.

2. *Reward Shaping*: Adds penalties to the reward function for violations. This soft constraint approach provides no feasibility guarantees and often destabilizes training (Achiam et al., 2017).
3. *Projection*: Projects the network output onto the feasible set. While valid for simple constraints, orthogonal projection is myopic and ignores the gradient information of the constraint boundary, leading to sub-optimal policies.

Our Contribution. In this work, we propose **Differentiable Multi-Period Portfolio Optimization (DMPO)**, a novel framework that embeds a convex optimization solver directly into the deep learning architecture. Instead of predicting asset returns (Two-Stage) or outputting weights directly (Standard DRL), our neural network learns to generate *optimization parameters*—conceptually representing a “subjective view” of the market. These parameters define a Quadratic Program (QP) whose solution constitutes the portfolio weights.

By differentiating through the QP solver using the Implicit Function Theorem (Amos & Kolter, 2017; Agrawal et al., 2019), DMPO achieves the best of both worlds: the representational power of deep learning to extract features from raw data, and the rigor of convex optimization to strictly enforce hard constraints. Furthermore, we extend this to a multi-period setting using Model Predictive Control (MPC) logic, allowing the agent to plan a sequence of trades that balances immediate returns against future transaction costs.

2. Related Work

2.1. Machine Learning in Portfolio Management

Early applications of ML in finance focused on forecasting price trends using Support Vector Machines (SVM) or Recurrent Neural Networks (RNN) (Fischer & Krauss, 2018). These forecasts were treated as deterministic inputs to MVO solvers. More recently, “decision-focused learning” (Donti et al., 2017) has argued that models should be trained to minimize the downstream optimization cost rather than prediction error. Our work aligns with this philosophy by backpropagating the Sharpe ratio loss through the optimization layer.

2.2. Deep Reinforcement Learning (DRL)

DRL has been extensively applied to algorithmic trading. (Jiang et al., 2017) introduced the EIIE (Ensemble of Identical Independent Evaluators) topology for portfolio management. (Wang et al., 2019) proposed the AlphaStock model

using attention mechanisms. Despite their success in unconstrained settings, handling constraints remains an open challenge. Constrained Policy Optimization (CPO) (Achiam et al., 2017) attempts to solve this for general MDPs, but it is computationally expensive and difficult to scale to high-dimensional portfolio spaces.

2.3. Differentiable Optimization Layers

The integration of optimization layers into neural networks was pioneered by OptNet (Amos & Kolter, 2017), which showed how to differentiate through quadratic programs. (Agrawal et al., 2019) generalized this to convex cone programs (CvxpyLayers). In the context of finance, (Butler & Kwon, 2023) applied differentiable convex layers for hedging. However, few works have explicitly addressed the *multi-period turnover constraint* problem in a unified end-to-end framework. DMPO fills this gap by formulating the problem as a differentiable Model Predictive Control (MPC) task.

3. Methodology

We formulate the portfolio optimization problem as a Constrained Markov Decision Process (CMDP). The goal is to learn a policy π_θ that maximizes the risk-adjusted return over a finite horizon H , subject to transaction costs and hard constraints.

3.1. Problem Formulation

Consider a market with N risky assets. At time step t , the state s_t consists of the historical window of asset returns $X_t \in \mathbb{R}^{T_{\text{window}} \times N}$ and the current portfolio weights $w_{t-1} \in \mathbb{R}^N$. The action $a_t = w_t$ is the target portfolio weight vector for the next period.

The system dynamics are given by the price evolution of the assets. The portfolio return at time $t+1$ (ignoring costs) is $r_{p,t+1} = w_t^\top y_{t+1}$, where y_{t+1} is the realized asset return vector.

Transaction Costs. Real-world trading incurs costs (commissions, slippage). We model the transaction cost c_t as proportional to the turnover:

$$c_t(w_t, w_{t-1}) = \kappa \sum_{i=1}^N |w_{t,i} - w_{t-1,i}| = \kappa \|w_t - w_{t-1}\|_1 \quad (1)$$

where κ is the transaction cost rate (e.g., 10 basis points).

Objective. We aim to maximize the annualized Sharpe Ratio. However, for the optimization layer, we utilize a mean-variance utility proxy, while the outer loop (RL training) optimizes the Sharpe Ratio directly.

3.2. The DMPO Architecture

The DMPO policy $\pi_\theta(s_t)$ is a composite function consisting of a **Prediction Network (Encoder)** and a **Differentiable Optimization Layer (Decoder)**.

3.2.1. 1. ENCODER: LATENT MARKET VIEWS

We use an LSTM network to process the time-series features X_t . The network does not directly output weights. Instead, it outputs a vector $\hat{\mu}_t \in \mathbb{R}^N$ and a structured matrix $\hat{L}_t \in \mathbb{R}^{N \times N}$ (where $\hat{\Sigma}_t = \hat{L}_t \hat{L}_t^\top$).

$$h_t = \text{LSTM}(X_t; \theta_{enc}) \quad (2)$$

$$\hat{\mu}_t = \text{Linear}_\mu(h_t) \quad (3)$$

$$\hat{L}_t = \text{Softplus}(\text{Linear}_\Sigma(h_t)) \quad (4)$$

Here, $\hat{\mu}_t$ can be interpreted as the model's "subjective view" of future returns, and $\hat{\Sigma}_t$ as the "subjective risk". Crucially, these are *not* trained to match realized returns (MSE loss) but are latent parameters trained to maximize final portfolio utility.

3.2.2. 2. DECODER: DIFFERENTIABLE QP LAYER

The core of DMPO is the convex optimization layer. Given the previous weights w_{t-1} and the learned parameters $\hat{\mu}_t, \hat{\Sigma}_t$, the network solves the following Quadratic Program (QP) to determine the optimal action w_t^* :

$$w_t^* = \arg \min_{w \in \mathbb{R}^N} \frac{1}{2} w^\top \hat{\Sigma}_t w - \lambda \hat{\mu}_t^\top w \quad (5a)$$

$$\text{s.t. } \mathbf{1}^\top w = 1, \quad (\text{Budget}) \quad (5b)$$

$$w \succeq 0, \quad (\text{No Shorting}) \quad (5c)$$

$$\|w - w_{t-1}\|_1 \leq \delta_t. \quad (\text{Turnover Constraint}) \quad (5d)$$

Here, δ_t is a dynamic turnover limit (which can also be learned or fixed). The term $\|w - w_{t-1}\|_1 \leq \delta_t$ is a hard constraint that prevents the model from excessive trading.

Since Equation (5) is a convex QP, it has a unique global minimum. We utilize the Alternating Direction Method of Multipliers (ADMM) or interior-point methods (via CVXPY) to solve this in the forward pass.

3.3. Implicit Differentiation & Backward Pass

To train the encoder parameters θ_{enc} end-to-end, we need to compute the gradient of the loss \mathcal{L} (e.g., negative Sharpe Ratio) with respect to $\hat{\mu}_t$ and $\hat{\Sigma}_t$. The chain rule gives:

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial w_t^*} \cdot \underbrace{\frac{\partial w_t^*}{\partial \hat{\mu}_t}}_{\text{Jacobian}} \cdot \frac{\partial \hat{\mu}_t}{\partial \theta} \quad (6)$$

The challenge is computing the Jacobian of the solution w_t^* with respect to the problem parameters. Since w_t^* is defined implicitly as the solution to an optimization problem, we invoke the **Implicit Function Theorem**.

The KKT optimality conditions for the QP (5) can be written as a system of equations $G(z, \psi) = 0$, where $z = (w^*, \nu^*, \xi^*)$ contains the primal and dual optimal variables, and $\psi = (\mu_t, \hat{\Sigma}_t, w_{t-1})$ contains the parameters. Specifically, differentiating the KKT conditions with respect to ψ yields:

$$D_z G(z^*, \psi) dz + D_\psi G(z^*, \psi) d\psi = 0 \quad (7)$$

$$\frac{\partial z^*}{\partial \psi} = -(D_z G(z^*, \psi))^{-1} D_\psi G(z^*, \psi) \quad (8)$$

The term $D_z G$ is essentially the KKT matrix. By solving a linear system involving the KKT matrix, we can efficiently compute the gradients without unrolling the iterative solver. This allows us to backpropagate the Sharpe Ratio gradient through the hard turnover constraint, effectively telling the neural network: "*Adjust your market view $\hat{\mu}_t$ such that the resulting constrained portfolio has a higher Sharpe ratio.*"

3.4. Multi-Period Extension (MPC)

While the QP above is a single-step lookahead, DMPO can be extended to a Receding Horizon Control (RHC) framework. The network predicts a sequence of views $\{\hat{\mu}_{t+k}\}_{k=0}^H$. The QP layer then solves for a trajectory of weights $\{w_t, \dots, w_{t+H}\}$, but only the first action w_t is executed. The cost function in the QP becomes:

$$\sum_{\tau=t}^{t+H} \left(\frac{1}{2} w_\tau^\top \hat{\Sigma}_\tau w_\tau - \hat{\mu}_\tau^\top w_\tau + \gamma \|w_\tau - w_{\tau-1}\|_1 \right) \quad (9)$$

This allows the agent to strategically plan trades to minimize impact costs over time.

4. Algorithm

The training procedure for DMPO is summarized in Algorithm 1.

5. Theoretical Analysis

5.1. Convexity and Uniqueness

The optimization layer in DMPO is constructed to be strictly convex. The regularization term $\frac{1}{2} w^\top \hat{\Sigma} w$ ensures positive definiteness provided $\hat{\Sigma}$ is positive definite (which we enforce via the Cholesky parametrization $\hat{L} \hat{L}^\top + \epsilon I$). This guarantees that for any learned view $\hat{\mu}$, there exists a unique optimal portfolio w^* , ensuring the mapping from view to action is a well-defined function.

Algorithm 1 Differentiable Multi-Period Portfolio Optimization (DMPO)

Input: Historical data \mathcal{D} , Horizon H , Batch size B , Learning rate α

Initialize: Encoder parameters θ

while not converged **do**

- Sample a batch of time windows $\{(X_t^{(i)}, w_{t-1}^{(i)})\}_{i=1}^B$ from \mathcal{D}
- for** each sample i in batch **do**
- 1. Forward Pass (View Generation):**
 $\hat{\mu}_t^{(i)}, \hat{\Sigma}_t^{(i)} \leftarrow \text{Encoder}(X_t^{(i)}; \theta)$
- 2. Forward Pass (QP Layer):**
 Solve QP (5) to get primal optimal $w_t^{*(i)}$
 Store KKT matrix factorization for backward pass
- 3. Loss Calculation:**
 Compute realized return $r_{t+1}^{(i)} = (w_t^{*(i)})^\top y_{t+1}$
 Compute transaction cost $c_t^{(i)} = \kappa \|w_t^{*(i)} - w_{t-1}^{(i)}\|_1$
 $r_{net}^{(i)} = r_{t+1}^{(i)} - c_t^{(i)}$
- end for**
- Compute Batch Sharpe Ratio: $\mathcal{L}(\theta) = -\frac{\text{Mean}(r_{net})}{\text{Std}(r_{net})}$
- 4. Backward Pass (Implicit Diff):**
 Compute $\nabla_\theta \mathcal{L}$ by backpropagating through KKT system
- Update $\theta \leftarrow \theta - \alpha \nabla_\theta \mathcal{L}$

end while

5.2. Gradient Existence

A key theoretical concern is whether the gradients exist at the boundaries of the feasible set. (Agrawal et al., 2019) showed that the solution map of a convex program is differentiable almost everywhere with respect to its parameters. Specifically, points where the active set of constraints changes (e.g., a weight hitting exactly 0 or the turnover constraint binding exactly) are measure zero. In practice, this means gradient-based training is stable, unlike in discrete action spaces.

5.3. Feasibility Guarantee

Unlike soft-constrained RL methods (e.g., PPO with penalty), DMPO guarantees $w_t \in \mathcal{F}$ (the feasible set) at every step of training and inference. This is because the output is the solution to a constrained optimization problem. This property is crucial for financial applications where violating leverage or turnover limits is legally or operationally impermissible.

6. Conclusion

We presented DMPO, an end-to-end framework for dynamic portfolio optimization. By differentiating through a convex optimization layer, DMPO aligns the feature learning capa-

bility of deep neural networks with the rigorous constraint handling of mathematical programming. This approach effectively solves the "optimizer's curse" by learning market views that are specifically tuned to perform well under the subsequent optimization step.

References

- Achiam, J., Held, D., Tamar, A., and Abbeel, P. Constrained policy optimization. In *International conference on machine learning*, pp. 22–31. PMLR, 2017.
- Agrawal, A., Amos, B., Barratt, S., Boyd, S., Diamond, S., and Kolter, J. Z. Differentiable convex optimization layers. In *Advances in neural information processing systems*, volume 32, 2019.
- Amos, B. and Kolter, J. Z. Optnet: Differentiable optimization as a layer in neural networks. In *International Conference on Machine Learning*, pp. 136–145. PMLR, 2017.
- Bengio, Y., Lodi, A., and Prouvost, A. Decision-focused learning: The case of portfolio optimization. *arXiv preprint arXiv:2011.08697*, 2020.
- Butler, A. and Kwon, R. Introduction to convex optimization for quantitative finance. *Journal of Financial Econometrics*, 2023.
- Donti, P., Amos, B., and Kolter, J. Z. Task-based end-to-end model learning in stochastic optimization. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Fischer, T. and Krauss, C. Deep learning with long short-term memory networks for financial market predictions. *European Journal of Operational Research*, 270(2):654–669, 2018.
- Jiang, Z., Xu, D., and Liang, J. A deep reinforcement learning framework for the financial portfolio management problem. *arXiv preprint arXiv:1706.10059*, 2017.
- Liang, Z., Chen, H., Zhu, J., Jiang, K., and Li, Y. Adversarial deep reinforcement learning in portfolio management. *arXiv preprint arXiv:1808.09940*, 2018.
- Markowitz, H. Portfolio selection. *The journal of finance*, 7(1):77–91, 1952.
- Michaud, R. O. The markowitz optimization enigma: Is 'optimized' optimal? *Financial analysts journal*, 45(1):31–42, 1989.
- Wang, J., Zhang, Y., Tang, K., Wu, J., and Xiong, Z. Alpha-stock: A buying-winners-and-selling-losers investment strategy using attentive lstm. *Proceedings of the 25th*

*ACM SIGKDD International Conference on Knowledge
Discovery & Data Mining, 2019.*