# A Novel Algorithm for Online Bilevel Optimization for Price Forecasting around Market Moving News

Jason Bohne[1,2], David Rosenberg[2], Gary Kazantsev [2], Paweł Polak[1]

[1]Stony Brook University [2] Bloomberg

**QR code for the full paper**

**Stony Brook University**

**Bloomberg**

## Abstract

**Bilevel optimization** methods are increasingly relevant in machine learning, especially for tasks such as hyperparameter optimization and meta-learning. Compared to the offline setting, online bilevel optimization offers a more dynamic framework by accommodating time-varying functions and sequentially arriving data in the problem formulation. We introduce a novel **online Bregman bilevel optimizer (OBBO)** with an improved theoretical guarantee for regret minimization and an efficient computational implementation via PyTorch. Empirically, we apply OBBO against established online/offline bilevel benchmarks in an online hyperparameter optimization for financial time series and display the **superior performance of OBBO** in terms of forecasting loss on an independent test set.

## Bilevel Optimization

**Canonical hyperparameter optimization problem**:
Optimization of a hyperparameter $\lambda$ on a validation set for optimal parameters $\widehat{\beta}(\lambda)$ on a training set.

$$\arg\min_{\lambda \in \mathbb{R}^+} \frac{1}{2}\left\|y^{val} - X^{val}\widehat{\beta}(\lambda)\right\|_2^2, \qquad \text{(outer level)}$$

$$\widehat{\beta}(\lambda) \in \arg\min_{\beta \in \mathbb{R}^m} \frac{1}{2}\left\|y^{train} - X^{train}\beta\right\|_2^2 + \lambda\|\beta\|_2^2 \qquad \text{(inner level)}$$

The above problem includes ridge regression, smoothing spline regression, and kernel ridge regression. Hyperparameter optimization is a special case of a **bilevel optimization**, where there is an outer level optimization parameterized by the optimal solution of an inner level optimization:

$$\arg\min_{\lambda \in X \subseteq \mathbb{R}^{d_1}}\left\{F(\lambda) \triangleq f\left(\lambda, \widehat{\beta}(\lambda)\right)\right\}, \qquad \text{(outer level)}$$

$$\widehat{\beta}(\lambda) \in \arg\min_{\beta \in \mathbb{R}^{d_2}} g(\lambda, \beta) \qquad \text{(inner level)}$$

Other special cases include meta-learning, neural architecture search, dataset distillation, and RLHF.

## The Hypergradient

One can differentiate through bilevel optimization problems under a few technical assumptions – leading to **bilevel specific gradient descent algorithms**. With the chain rule, the gradient of the outer level objective can be decomposed into a direct and indirect gradient term.

$$\nabla F(\lambda) = \underbrace{\nabla_\lambda f(\lambda, \widehat{\beta}(\lambda))}_{(a)\text{ Hyperparameter Direct Gradient}} + \overbrace{\underbrace{\nabla\widehat{\beta}(\lambda)}_{(b)\text{ Best-Response Jacobian}} \underbrace{\nabla_\beta f(\lambda, \widehat{\beta}(\lambda))}_{(c)\text{ Parameter Direct Gradient}}}^{\text{Hyperparameter Indirect Gradient}}$$

Direct gradient terms (**a**,**c**) are computable if the objective is differentiable w.r.t. hyperparameters and parameters, e.g., neural networks. For any hyperparameter values $\lambda$, the best-response Jacobian (**b**) is the typically **unknown** gradient of the corresponding optimal parameters. However with **implicit function theorem**, one can derive the best-response Jacobian in terms of computable gradients as

$$\nabla\widehat{\beta}(\lambda) = -\underbrace{\nabla^2_{\lambda,\beta}g(\lambda,\widehat{\beta}(\lambda))}_{\text{Training Partials}}\overbrace{\left(\nabla^2_{\beta,\beta}g(\lambda,\widehat{\beta}(\lambda))\right)}^{\text{Training Hessian}}{}^{-1}$$

## An Improved Bilevel Optimizer

Our online Bregman bilevel optimizer (**OBBO**) generalizes the gradient step in known online bilevel optimizers through the application of a Bregman Divergence $\mathcal{D}_\phi(\cdot, \cdot)$. This offers a generalization from the squared Euclidean distance, as in Lin et al. 2024; Tarzanagh et al. 2024. Given a continuously differentiable $\rho$-strongly convex function $\phi(\lambda)$, a Bregman Divergence $\mathcal{D}_\phi(\cdot, \cdot)$ is defined for all $\lambda_1, \lambda_2 \in \mathcal{X}$ as:

$$\mathcal{D}_\phi(\lambda_2, \lambda_1) := \phi(\lambda_2) - \phi(\lambda_1) - \langle\nabla\phi(\lambda_1), \lambda_2 - \lambda_1\rangle.$$

For a Bregman divergence $\mathcal{D}_\phi(\cdot, \cdot)$, our generalized gradient step then has the form of

$$\lambda^+ = \arg\min_{\lambda \in \mathcal{X}}\left\{\langle q, \lambda\rangle + \frac{1}{\alpha}\mathcal{D}_\phi(\lambda, u)\right\},$$

where $\alpha > 0$ is a step size, and $q, u \in \mathbb{R}^{d_1}$. Our analysis shows **OBBO** achieves an improved sublinear rate of bilevel local regret– a measure of stationarity for online bilevel algorithms from Tarzanagh et al. 2024. For a window smoothing parameter $w \geq 1$, the bilevel local regret is defined for a smooth $F_t(\lambda)$ as

$$BLR_w(T) := \sum_{t=1}^{T}\left\|\nabla F_{t,w}(\lambda_t)\right\|^2, \quad F_{t,w}(\lambda_t) := \frac{1}{w}\sum_{i=0}^{w-1}F_{t-i}(\lambda_{t-i}),$$

Specifically, the condition number of the inner objective $g(\lambda, \beta)$ is $\kappa_g > 1$. **OBBO** achieves a $\kappa_g^2$ dependency, whereas benchmarks of SOBOW (Lin et al. 2024) and OAGD (Tarzanagh et al. 2024) only achieve $\kappa_g^3$ and $\kappa_g^4$ respectively.

- One special case of the generalized gradient step given adaptive matrix $\mathbf{H}_t$ is **Adagrad** – which can better capture the underlying geometry via use of adaptive learning rates.

- Another special case of the generalized gradient step is the reduction to gradient descent, when $\phi(\lambda) = \frac{1}{2}\|\lambda\|^2$ and $\mathcal{X} = \mathbb{R}^{d_1}$.

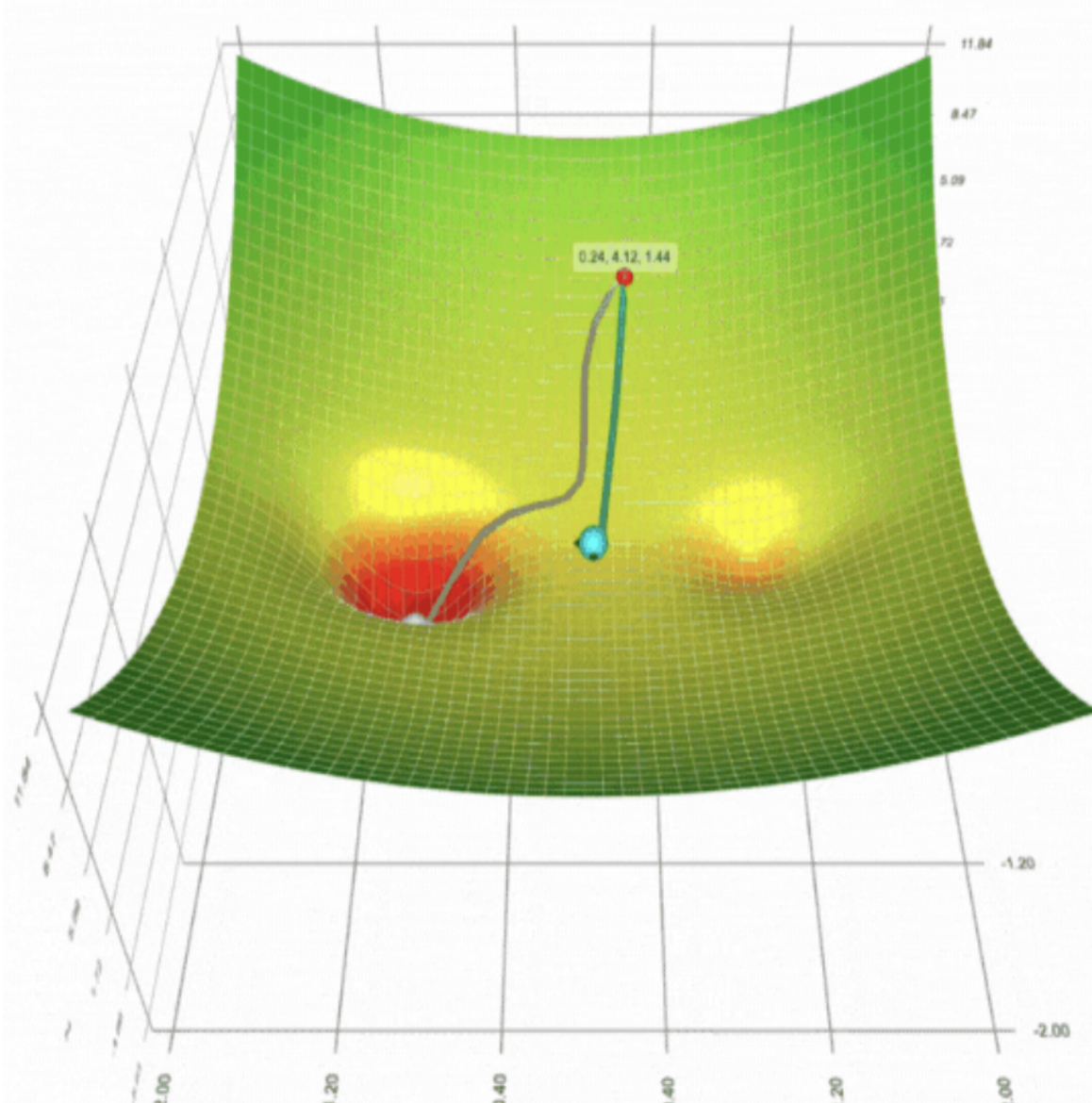- Empirically, we update hyperparameters with an Adagrad step (2nd order information) vs. a gradient descent step (only 1st order).

Figure 1: For local minimum in above, see the improvement of Adagrad (gray) relative to Gradient Descent (teal).



| Example | Bregman Function |
|---|---|
| **Gradient Descent** | $\frac{1}{2}\|\lambda\|_2^2$ |
| Adagrad | $\frac{1}{2}\lambda^T H_t \lambda$ |

Table 2: Example Bregman Functions

Full paper can be found on https://arxiv.org/abs/2409.10470

## Application to Hyperparameter Optimization

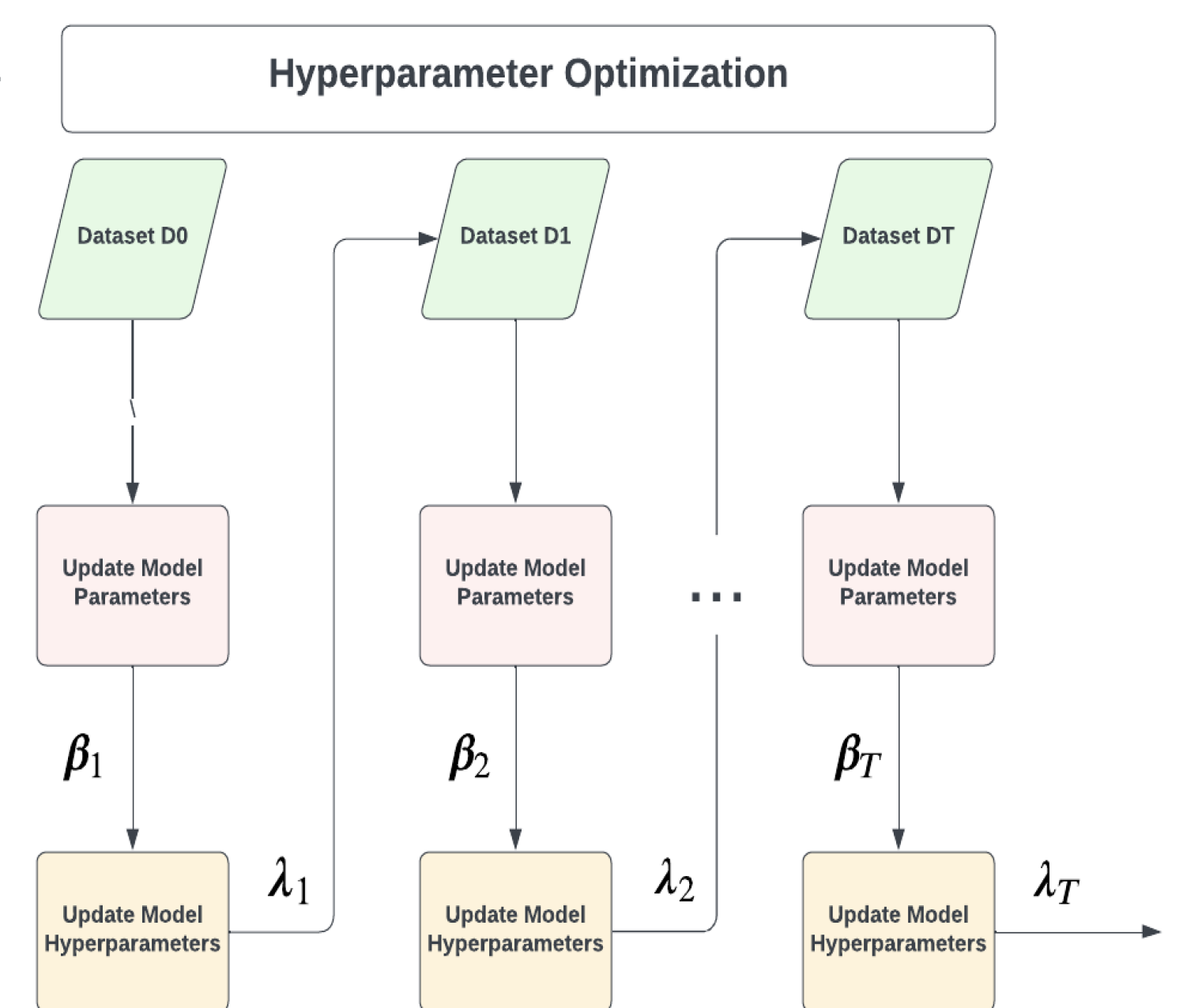Consider **hyperparameter optimization** where:

- New training-validation datasets sequentially arrive such that optimal values for (hyper)-parameters can vary over datasets.

- The goal is to update hyperparameters $\lambda_t$ for optimal parameters $\widehat{\beta}_t(\lambda_t)$ on dataset $D_t$.

- Ex: Linear smoothing spline model with B-spline coefficients as parameters $\beta_t$ and positive regularization hyperparameter $\lambda_t$.

For datasets $D_t := (X_t^{train}, y_t^{train}, X_t^{val}, y_t^{val})$, formulate our hyperparameter optimization $\forall t$ as

$$\arg\min_{\lambda \in \mathbb{R}^+}\frac{1}{2}\left\|y_t^{val} - X_t^{val}\widehat{\beta}_t(\lambda)\right\|_2^2,$$

$$\widehat{\beta}_t(\lambda) \in \arg\min_{\beta \in \mathbb{R}^m}\frac{1}{2}\left\|y_t^{train} - X_t^{train}\beta\right\|_2^2 + \lambda\|\beta\|_2^2$$

Figure 2: Diagram of Hyperparameter Optimization



## Price Forecasting around Market Moving News

**Market Moving News Dataset**:

- Significant price events and news stories annotated via segmentation algorithm.

- Subset of the RAY index (U.S. equities) between January 1, 2021 and June 1, 2022.

**Experiment Setup**:

- 440 samples of equity time series partitioned into a rolling window of training-validation data.

- Separate test set post annotated price event to evaluate forecasting mean-squared error.

- Model choice of linear smoothing spline with B-spline coefficients as parameters $\beta_t$ and regularization hyperparameter $\lambda_t$ results in forecasts satisfying assumptions of annotation pipeline.
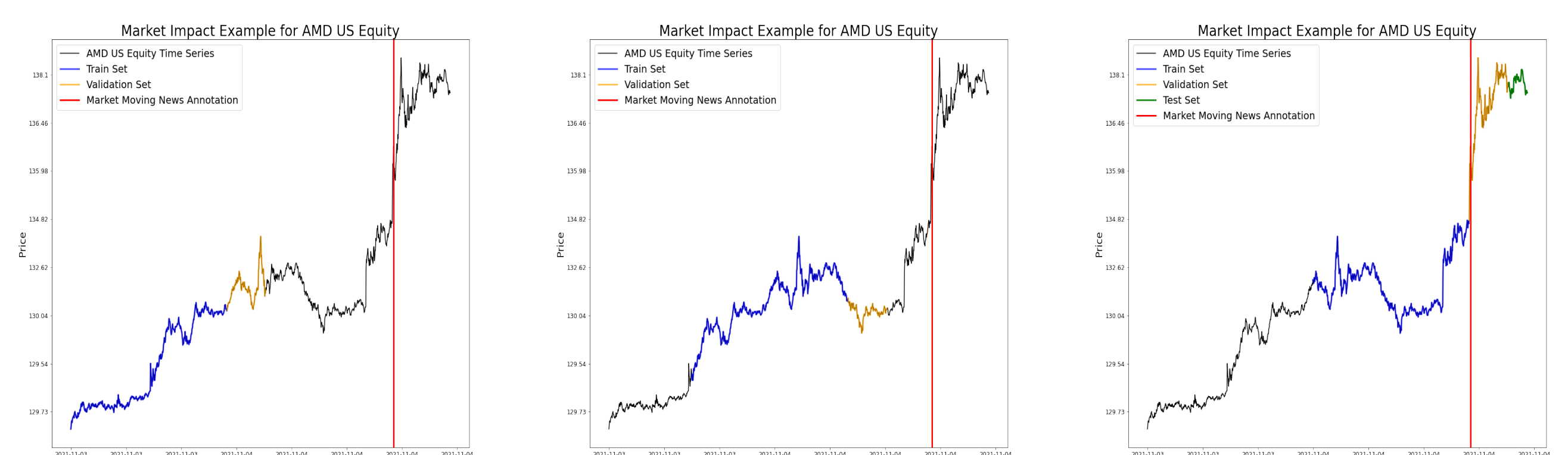


Figure 3: Sample training-validation subsets for AMD U.S. Equity and market event on 11-08-2021.

Benchmark algorithms include gradient bilevel algorithms OAGD (Tarzanagh et al. 2024) and SOBOW (Lin et al. 2024), limited to a gradient descent step. Further benchmarked by general purpose optimizers; ADAM (Kingma 2014), SGDM.
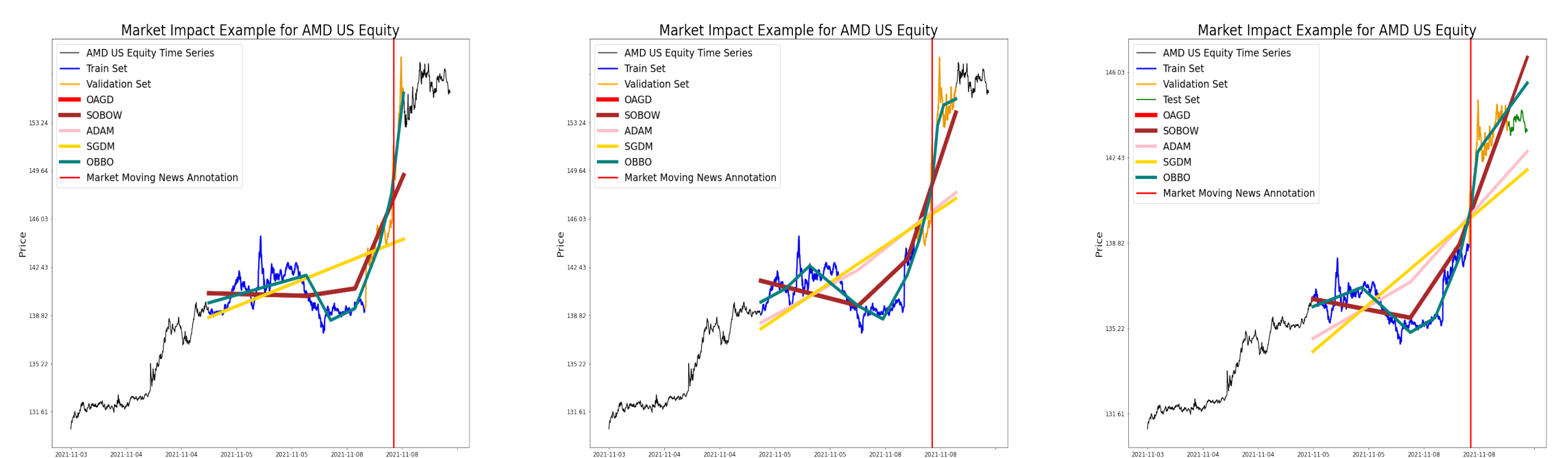


Figure 4: Example forecasts generated with **OBBO** vs. benchmark algorithms.

Better convergence of **OBBO** results in quicker adaptability of our underlying model to annotated event and improved forecasting MSE on a test set post market event— statistics provided below.

| Algorithm | Forecasting Loss across U.S. Markets | |
|---|---|---|
| | Mean Loss ± Standard Error | Median Loss ± Median Absolute Deviation |
| **OBBO** | **0.661** ± 0.055 | **0.205** ± 0.150 |
| **OAGD** | 0.707 ± 0.053 | 0.265 ± 0.209 |
| **SOBOW** | 0.689 ± 0.053 | 0.273 ± 0.215 |
| **Adam** | 1.265 ± 0.176 | 0.267 ±0.230 |
| **SGDM** | 0.872 ± 0.078 | 0.401 ± 0.286 |

Table 3: Statistics of forecasting mean-squared error across U.S. markets.

## Conclusions and Extensions

- Hyperparameter optimization (HO) is actually a special case of bilevel optimization.

- Provide an **improved algorithm for general bilevel optimization** problems.

- Empirically show benefit of our improved algorithm in HO for financial time series forecasting.

- Determine if our algorithm offers an empirical improvement in **other special cases**, e.g., RLHF.

**References:**

Kingma, Diederik P (2014). "Adam: A Method for Stochastic Optimization". In: *arXiv preprint arXiv:1412.6980*.

Lin, Sen et al. (2024). "Non-Convex Bilevel Optimization with Time-Varying Objective Functions". In: *Advances in Neural Information Processing Systems* 36.

Tarzanagh, Davoud Ataee et al. (2024). "Online Bilevel Optimization: Regret Analysis of Online Alternating Gradient Methods". In: *International Conference on Artificial Intelligence and Statistics*.