

## Inference Instructions and Notes

**Pre-trained models can be executed on unseen test data provided in the file `new_data.parquet`.** This file may be replaced in-place with any new test dataset, and predictions will automatically be generated using the `main.py` script.

For demonstration purposes, the accompanying notebook generates a 10-day test sample randomly selected from the existing `training_data.parquet` dataset. This test sample is saved as `new_data.parquet` with its dates modified to appear consecutively after the maximum date in the training data.

All model development and analysis are documented in the Jupyter notebook `230P-HW1-sol.ipynb`. Final models have been pre-trained using the methodology outlined in the notebook and serialized as `.pkl` files.

The inference pipeline expects input in the same structure as `training_data.parquet`, excluding the target column `ret`. If the `ret` column is present in `new_data.parquet`, it will be automatically dropped during inference to avoid data leakage.

## Feature Engineering Strategy

The dataset consists of five numerical features (`macro2`, `price`, `firm1`, `firm2`, `firm3`), one temporal feature (`date`), one categorical feature (`macro1`), and the target variable to be predicted (`ret`).

The primary objective of the feature engineering process was to construct a rich and informative set of predictors by systematically expanding the input space. This involved generating all meaningful pairwise interactions among numerical variables, as well as constructing various price-based ratios to capture relative relationships between firms and market-level variables.

Additional transformations, such as lagged features, rolling statistics (mean, std, min, max), and polynomial expansions (squared, cubed), were layered on top of these interactions to further enrich the feature space. The goal was to extract as much predictive signal as possible, even if it meant introducing a high-dimensional feature matrix.

The use of regularized models such as LASSO (L1) was motivated by their ability to perform automatic feature selection, effectively eliminating redundant or uninformative variables. This allowed for an aggressive feature engineering approach without the risk of overfitting due to high dimensionality.

## Pipeline Overview

The complete feature engineering pipeline is composed of five sequential substeps, each of which is modularized and applied through method chaining using the `.pipe()` interface. These steps are described below and elaborated in the subsequent subsections.

```

1 def full_feature_engineering_pipeline(df):
2     return (df.pipe(feature_interactions)
3             .pipe(feature_one_hot_encoding)
4             .pipe(feature_transformations)
5             .pipe(rolling_features)
6             .pipe(firm_id_one_hot_encoding))

```

Listing 1: Feature Engineering Pipeline

1. Feature Interaction Engineering
2. Categorical Feature One-Hot Encoding (`macro1`)
3. Feature Transformations (lagged, polynomial, log returns)
4. Rolling Window Statistics
5. Firm Identifier One-Hot Encoding (`firm_id`)

## 1. Feature Interaction Engineering

This step generates pairwise interaction terms across the core numerical features to uncover nonlinear relationships and capture cross-effects between firm-specific, macroeconomic, and price-related variables. In addition to multiplicative interactions (e.g., `firm1 * macro2`), the pipeline also includes price ratio features (e.g., `price / macro2`) which may capture relative scaling effects or valuation-like patterns.

```

1 def feature_interactions(df):
2     """
3     Creates interaction terms between firm variables, macro variables, and price.
4     """
5     print('Creating feature interactions')
6
7     # Firm-to-firm interactions
8     df['firm1*firm2'] = df['firm1'] * df['firm2']
9     df['firm2*firm3'] = df['firm2'] * df['firm3']
10    df['firm1*firm3'] = df['firm1'] * df['firm3']
11
12    # Firm-to-macro interactions
13    df['firm1*macro2'] = df['firm1'] * df['macro2']
14    df['firm2*macro2'] = df['firm2'] * df['macro2']
15    df['firm3*macro2'] = df['firm3'] * df['macro2']
16
17    # Firm-to-price interactions
18    df['firm1*price'] = df['firm1'] * df['price']
19    df['firm2*price'] = df['firm2'] * df['price']
20    df['firm3*price'] = df['firm3'] * df['price']
21
22    # Macro-to-price interaction
23    df['macro2*price'] = df['macro2'] * df['price']
24
25    # Price ratio features
26    df['price/macro2'] = df['price'] / df['macro2']
27    df['price/firm1'] = df['price'] / df['firm1']
28    df['price/firm2'] = df['price'] / df['firm2']
29    df['price/firm3'] = df['price'] / df['firm3']

```

```
30
31 return df
```

Listing 2: Feature Interaction Engineering

## 2. Categorical Feature One Hot Encoding (feature macro1)

To accommodate the categorical nature of the feature `macro1`, one-hot encoding was applied. This transformation enables the feature to be ingested by models such as Elastic Net Regression and XGBoost. While recent versions of XGBoost support native categorical handling, explicit one-hot encoding was performed in accordance with project instructions to ensure model-agnostic compatibility and transparency of the feature space.

```
1
2 def feature_one_hot_encoding(df):
3     '''
4     One-hot encodes the 'macro1' column.
5     '''
6     print('One-hot encoding categorical feature: macro1')
7
8     # Create dummy columns
9     dummies = pd.get_dummies(df['macro1'], prefix='macro1', drop_first=False)
10
11    # Drop original and concatenate dummies
12    df = pd.concat([df.drop(columns=['macro1']), dummies], axis=1)
13    return df
```

Listing 3: One Hot Encoding macro1

## 3. Feature Transformations

The next step involved applying non-linear transformations to the engineered features in order to capture potential higher-order relationships. These included:

- **Lag features:** Created per firm to capture short-term temporal dependencies.
- **Log returns:** Computed for the `price` feature to express relative changes.
- **Polynomial expansions:** Squared and cubed versions of key features were added to introduce curvature and interaction flexibility into the feature space.

Firm-wise computation was critical to preserve logical temporal structure and avoid data leakage across entities.

```

1 def feature_transformations(df):
2     """
3     Adds lag, squared, cubed, and log-return transformations per firm.
4     """
5     # Exclude macro1 from transformations
6     macro1_features = [col for col in df.columns if col.startswith('macro1')]
7
8     # Select columns to transform
9     columns = [col for col in df.columns if col not in ['firm_id', 'date', 'ret'] +
10                macro1_features]
11
12     for firm_id in tqdm(df['firm_id'].unique(), desc='Feature transformations per
13                          firm_id'):
14         firm_mask = df['firm_id'] == firm_id
15
16         for col in columns:
17             # Lag-1 feature
18             df.loc[firm_mask, f'{col}_lag_1'] = df.loc[firm_mask, col].shift(1)
19
20             # Polynomial features
21             df[f'{col}_squared'] = df[col] ** 2
22             df[f'{col}_cubed'] = df[col] ** 3
23
24             # Log return for price (requires lag column to be generated first)
25             df.loc[firm_mask, f'price_log_ret'] = (
26                 np.log(df.loc[firm_mask, 'price'] / df.loc[firm_mask, 'price_lag_1'])
27             )
28
29     return df

```

Listing 4: Feature Transformations

## 4. Rolling Features

To incorporate temporal context, rolling window statistics were computed over a 20-period horizon for all numeric features. The rolling features included:

- **Mean, Min, Max, Std:** Basic summary statistics within the window.
- **High/Low Ratio:** Max divided by Min — useful as a proxy for volatility skew.
- **Range:** Difference between Max and Min — an amplitude signal.
- **Window-based Return:** Percentage return over the lag window.

These features help embed autocorrelation and volatility structure that may aid the model in predicting the target variable more effectively.

```

1 def rolling_features(df, window_list=[20]):
2     """
3     Adds rolling mean, std, min, max, range, and returns for each numeric feature
4     per firm.
5     """

```

```

5  # Identify macro features to exclude from rolling calculations
6  macro1_features = [col for col in df.columns if col.startswith('macro1')]
7
8  # Select all numeric columns excluding identifiers and macro1
9  columns = [col for col in df.columns if col not in ['firm_id', 'date', 'ret'] +
10             macro1_features]
11
12 # Iterate over each unique firm
13 for firm_id in tqdm(df['firm_id'].unique(), desc='Rolling features per firm_id'):
14     for col in columns:
15         firm_mask = df['firm_id'] == firm_id
16         for window in window_list:
17             # Rolling average
18             df.loc[firm_mask, f'{col}_rolling_avg_{window}'] = df.loc[firm_mask,
19                                col].rolling(window).mean()
20
21             # Rolling standard deviation
22             df.loc[firm_mask, f'{col}_rolling_std_{window}'] = df.loc[firm_mask,
23                                col].rolling(window).std()
24
25             # Rolling max
26             df.loc[firm_mask, f'{col}_rolling_max_{window}'] = df.loc[firm_mask,
27                                col].rolling(window).max()
28
29             # Rolling min
30             df.loc[firm_mask, f'{col}_rolling_min_{window}'] = df.loc[firm_mask,
31                                col].rolling(window).min()
32
33             # High/low ratio (max / min)
34             df.loc[firm_mask, f'{col}_rolling_high/low_{window}'] = (
35                 df.loc[firm_mask, f'{col}_rolling_max_{window}'] /
36                 df.loc[firm_mask, f'{col}_rolling_min_{window}']
37             )
38
39             # Range = max - min
40             df.loc[firm_mask, f'{col}_rolling_range_{window}'] = (
41                 df.loc[firm_mask, f'{col}_rolling_max_{window}'] -
42                 df.loc[firm_mask, f'{col}_rolling_min_{window}']
43             )
44
45             # Window-based return (current / lag - 1)
46             df.loc[firm_mask, f'{col}_rolling_ret_{window}'] = (
47                 df.loc[firm_mask, col] / df.loc[firm_mask, col].shift(window -
48                 1) - 1
49             )
50     return df

```

Listing 5: Rolling Features

## 5. Firm ID One Hot Encoding

To complete the pipeline, the `firm_id` column—representing a categorical identifier of the security—was also one-hot encoded. Since the numeric values of identifiers do not imply ordinal relationships, encoding was necessary to prevent the model from misinterpreting them. This process mirrors the transformation applied to `macro1`, ensuring consistency across all categorical inputs.

```

1 def firm_id_one_hot_encoding(df):
2     '''
3     One-hot encodes the 'firm_id' column.
4     '''
5     print('One-hot encoding firm_id')
6
7     # Create dummy columns
8     dummies = pd.get_dummies(df['firm_id'], prefix='firm_id', drop_first=False)
9
10    # Drop original and concatenate dummies
11    df = pd.concat([df.drop(columns=['firm_id']), dummies], axis=1)
12    return df

```

Listing 6: Firm ID One Hot Encoding

## ElasticNet

### 1. Baseline Elastic Net Model

- Fixed hyper-parameters (selected via the initial validation-set grid search):

$$\alpha = 0.01, \quad \ell_1\text{-ratio} = 0.10.$$

- Model fitted on the standardized training matrix and evaluated on the chronological validation window.

#### Impact of the engineered feature set

Augmenting the original macro and firm fundamentals with the engineered interaction/momentum features measurably boosts generalization:

Feature Set	Validation MSE	Validation $R^2$
Original features only	0.0023	0.5852
With engineered features	0.0017	0.6963

Table 1: Effect of engineered features on baseline ElasticNet performance.

#### Interpretation

- **Error reduction.** The validation MSE drops by  $\approx 26\%$  (from  $2.3 \times 10^{-3}$  to  $1.7 \times 10^{-3}$ ), indicating that the additional interactions capture systematic variation previously left in the residuals.
- **Explained variance.** Validation  $R^2$  climbs from 0.585 to 0.696, a relative improvement of roughly 19%, confirming that the model accounts for nearly 70

## 2. Hyper Parameter Tuning

### 1. Initial grid search (validation-set driven)

- Swept  $\sim 40$  logarithmically spaced  $\alpha$  values ( $10^{-2}$ – $10^0$ ) against 50 uniformly spaced  $\ell_1$  ratios ( $0.10$ – $1.00$ ).
- For each  $(\alpha, \ell_1)$  pair the model was fit on the *training split* and scored on the *held-out validation split* using MSE /  $R^2$ .
- The best-performing pair on the validation set was refit on the full training data.
- *Caveat:* Because every candidate is judged on a single validation fold, the winner is effectively tuned to that fold and may look overly optimistic on unseen data.

```
1 def hyperopt_elasticnet():
2     our_train = pd.read_parquet('./train+features.parquet')
3     train, val = temporal_train_val_split(our_train, cutoff_frac=0.8)
4     train = train.dropna()
5
6     X_train, y_train = prepare_X_y(train, drop_cols=['date'])
7     X_val, y_val = prepare_X_y(val, drop_cols=['date'])
8     X_train_scaled, X_val_scaled, _ = standardize_with_train_stats(X_train,
9                             X_val)
10
11     all_alpha = np.logspace(-2, 0, 40)
12     all_l1 = np.linspace(0.1, 1, 50)
13
14     results = {}
15     for alpha in all_alpha:
16         print(f"Testing alpha: {alpha}")
17         for l1 in tqdm(all_l1, total=len(all_l1)):
18             # print(f"Testing l1_ratio: {l1}")
19             # Fit ElasticNet with current parameters
20             elastic_net = ElasticNet(alpha=alpha, l1_ratio=l1, max_iter=10
21                                     _000, random_state=42)
22             elastic_net.fit(X_train_scaled, y_train)
23             y_pred = elastic_net.predict(X_val_scaled)
24             mse = mean_squared_error(y_val, y_pred)
25             r2 = r2_score(y_val, y_pred)
26
27             # print(f"ElasticNet alpha: {alpha}, l1_ratio: {l1}, MSE: {mse
28                 :.4f}, R2: {r2:.4f}")
29             results[(alpha, l1)] = (mse, r2)
30
31     results = sorted(results.items(), key=lambda x: x[1][1])
32
33     # Best parameters
34     best_alpha, best_l1 = results[-1][0]
35     elastic_net = ElasticNet(alpha=best_alpha, l1_ratio=best_l1, max_iter
36                             =10_000, random_state=42)
37     elastic_net.fit(X_train_scaled, y_train)
38     y_pred = elastic_net.predict(X_val_scaled)
39     mse = mean_squared_error(y_val, y_pred)
40     r2 = r2_score(y_val, y_pred)
41
42     print(f"ElasticNet best alpha: {best_alpha}")
```

```

40     print(f"ElasticNet best l1_ratio: {best_l1}")
41     print(f"Validation MSE: {mse:.4f}")
42     print(f"Validation R2: {r2:.4f}")
43
44     hyperopt_elasticnet()
45

```

Listing 7: Grid Search CV

## 2. Cross-validated tuning inside the training set

- Replaced the manual grid search with `ElasticNetCV`, retaining the same candidate grid.
- Performed  $k=5$ -fold cross-validation *entirely within the training split*.
- Selected the  $(\alpha, \ell_1)$  pair maximizing the mean CV score, reducing leakage of validation information.

```

1  def elastic_net_cross_val():
2      our_train = pd.read_parquet('./train+features.parquet')
3      train, val = temporal_train_val_split(our_train, cutoff_frac=0.8)
4      train = train.dropna()
5
6      X_train, y_train = prepare_X_y(train, drop_cols=['date'])
7      X_val, y_val = prepare_X_y(val, drop_cols=['date'])
8
9      pipe = make_pipeline(
10         StandardScaler(),
11         ElasticNet(max_iter=10_000)
12     )
13
14     param_dist = {
15         "elasticnet__alpha": np.logspace(-2, 0, 60),
16         "elasticnet__l1_ratio": np.linspace(0.1, 1, 101),
17     }
18
19
20     inner_cv = TimeSeriesSplit(n_splits=5)
21
22     search = RandomizedSearchCV(
23         pipe,
24         param_distributions = param_dist,
25         n_iter = 200,
26         scoring = "neg_mean_squared_error",
27         cv = inner_cv,
28         n_jobs = -1,
29         random_state= 0,
30         verbose = 10,
31     )
32
33     # --- guarantee an ASCII temp folder for Windows / joblib
34     tmp_root = pathlib.Path(r"C:\tmp\joblib")
35     tmp_root.mkdir(parents=True, exist_ok=True)
36     os.environ["JOBLIB_TEMP_FOLDER"] = str(tmp_root)
37     os.environ["OMP_NUM_THREADS"] = "1" # 1 BLAS thread per worker
38
39     search.fit(X_train, y_train)

```



```

40     print("Selected params :", search.best_params_)
41     y_pred = search.predict(X_val)
42     mse_search = mean_squared_error(y_val, y_pred)
43     r2_search = r2_score(y_val, y_pred)
44     print(f"Search best alpha: {search.best_params_['elasticnet__alpha']}")
45     print(f"Search best l1_ratio: {search.best_params_['elasticnet__l1_ratio']}")
46     print(f"Search Validation MSE: {mse_search:.4f}")
47     print(f"Search Validation R2: {r2_search:.4f}")
48
49     elastic_net_cross_val()
50

```

Listing 8: Cross-Validation

### 3. Model comparison

- Both tuned models were evaluated on:
  - 3.1. their respective *training data* (in-sample fit), and
  - 3.2. the external *validation split* (generalization).
- The vanilla grid-search model exhibited a larger train–validation performance gap (higher train  $R^2$ , lower val  $R^2$ ), signaling overfitting.
- The CV-tuned model showed more consistent metrics across splits, indicating better generalization.

### Selected Hyper-parameters

Hyper-parameter	Grid Search (val set)	RandomizedSearchCV (5-fold CV)
$\alpha$	0.0100	0.01169
$\ell_1$ -ratio	0.2286	0.5950

Table 2: Best ElasticNet configurations under two tuning schemes.

### Interpretation

- **Regularization strength ( $\alpha$ ).** Both searches converge on a small—but non-negligible— $\alpha \approx 0.01$ , implying that mild overall shrinkage is beneficial: it tempers coefficient variance without washing out signal.
- **Mixing parameter ( $\ell_1$ -ratio).**
  - *Grid search:*  $\ell_1=0.23$  leans heavily toward ridge-like ( $\ell_2$ ) regularization; coefficients are shrunk but only a few are driven exactly to zero.
  - *CV search:*  $\ell_1=0.60$  shifts the balance toward lasso-like sparsity, zeroing more weak predictors while still retaining an  $\ell_2$  component to stabilise correlated features.

## 4. Feature Importance

### Interpretation of Permutation Feature Importance (with Significance)

1. **Macro shock & persistence — `macro2` and `macro2_lag_1`**  
Shuffling `macro2` increases the validation-set MSE by roughly  $2 \times 10^{-3}$ , with its one-period lag close behind ( $1.3 \times 10^{-3}$ ). Both  $p$ -values are  $< 10^{-45}$ , so we can be extremely confident that the model relies on the current and recent macro environment for accuracy.
2. **Firm-specific fundamentals — `firm2`, `firm1`, `firm3` and their lags**  
Each firm variable registers a smaller but still highly significant drop in performance when permuted ( $p < 10^{-35}$ ). This confirms that cross-sectional differences among firms add signal beyond the overarching macro backdrop.
3. **Cross-effects and price dynamics — `price_log_ret`, `firm1*macro2`, `price/macro2`<sup>2</sup>  
`rolling_high/low_20`**  
The *interaction* `firm1*macro2` and the pure market feature `price_log_ret` both survive the 5% significance test, suggesting that macro conditions modulate firm impact and that short-horizon price momentum contributes incremental information. In contrast, higher-order engineered variants (e.g. `firm2*macro2_cubed_rolling_min_20`) yield zero mean importance and  $p = 1$ , indicating redundancy; these can be pruned to simplify the feature set without sacrificing predictive power.

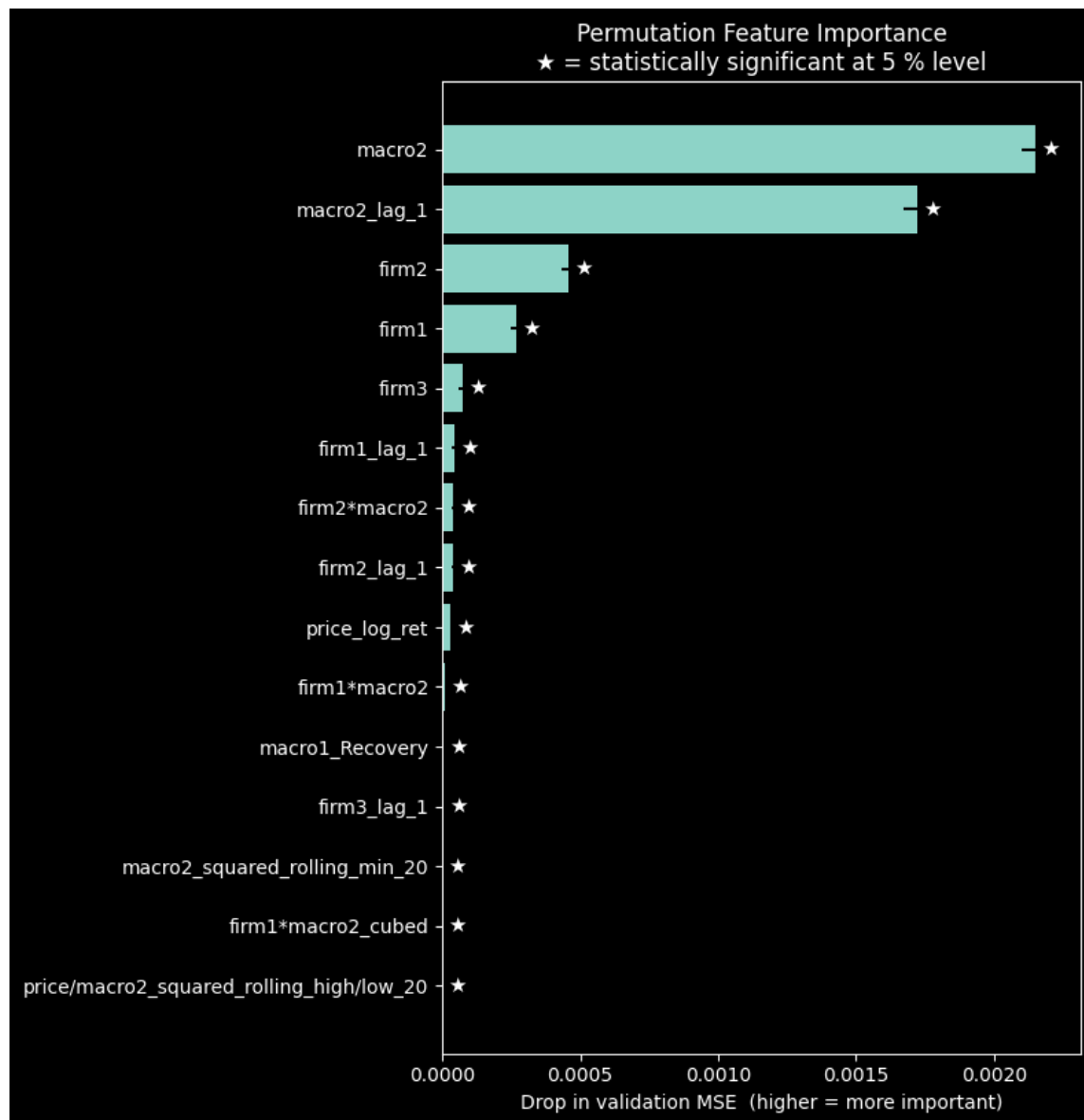
Overall, the model’s predictive edge is anchored in macro shocks, refined by firm-level fundamentals, and sharpened by select interaction and price-momentum terms. Features with zero importance and non-significant  $p$ -values are prime candidates for removal, streamlining future iterations and potentially reducing over-fitting.

### Interestingly, FirmID itself is not useful

The **FirmID indicator variables** (one-hot or label encodings that merely tag each entity) show **zero mean importance and  $p = 1.0$** . This tells us that, once the model already sees:

- firm-level quantitative fundamentals (`firm1`, `firm2`, ...),
- macro context plus interaction terms,

**knowing the raw identity of the firm adds no extra predictive power.** In practical terms, the model distinguishes firms through their observable metrics rather than their labels; any firm-specific fixed effects are either weak or already captured by those continuous fundamentals.



# XGBoost

## 1. Baseline XGBoost Model

- XGBRegressor with moderate complexity:

n\_estimators = 200, max\_depth = 3,  $\eta = 0.10$ , subsample = 1, colsample\_bytree = 0.8.

- Fit on the scaled training matrix; evaluated on the validation window.

### Impact of the engineered feature set

Adding the newly created interaction and momentum features noticeably improves generalization, shrinking validation MSE and lifting  $R^2$ :

Feature Set	Validation MSE	Validation $R^2$
Original features only	0.0025	0.5598
With engineered features	0.0018	0.6887

Table 3: Effect of engineered features on baseline XGBoost performance.

## 2. Hyper-parameter Tuning

- **Search space.** max\_depth  $\in \{4, 6\}$ ,  $\eta \in \{0.05, 0.10\}$ ,  $\gamma \in \{0, 0.1\}$ ,  $\lambda \in \{0, 1, 5\}$ .
- **Methodology.** *RandomizedSearchCV* with  $k = 5$ -fold cross-validation nested completely *inside* the training split. Mean squared error (negative) is the scoring metric.

## 3. Selected Configuration

n_estimators = 1000	$\eta = 0.05$
max_depth = 4	subsample = 0.8
colsample_bytree = 0.8	$\gamma = 0$
reg_lambda = 1.0	min_child_weight = 1

### Rationale

- $\eta = 0.05$  with n\_estimators = 1000: a lower learning rate paired with more boosting rounds allows the model to make finer, more stable updates.
- max\_depth = 4: deep enough to capture non-linear macro-firm interactions while still limiting model variance.
- subsample = 0.8 and colsample\_bytree = 0.8: row- and column-subsampling inject bootstrap-style diversity, further curbing over-fitting.
- $\gamma = 0$ : no additional split-gain threshold was needed once the other regularization levers were tuned.
- $\lambda = 1.0$  ( $\ell_2$  regularization) and min\_child\_weight = 1 place mild penalties on overly large leaf weights, balancing bias and variance without suppressing meaningful signals.

## 4. Importance Score

Rank	Feature	Split-count	Sig.
1	<b>firm2</b>	$\approx 63$	★
2	<b>macro2</b>	$\approx 60$	★
3	<b>macro2_lag_1</b>	$\approx 55$	★
4	<b>firm1</b>	$\approx 49$	★
5	<b>firm3</b>	$\approx 34$	★
6	<b>firm2_lag_1</b>	$\approx 28$	★
7	<b>firm2 <math>\times</math> macro2</b>	$\approx 25$	★
8–10	<b>macro1_Contraction, firm1_lag_1, price_log_ret</b>	19–22	★
11–20	Remaining terms	10–19	★

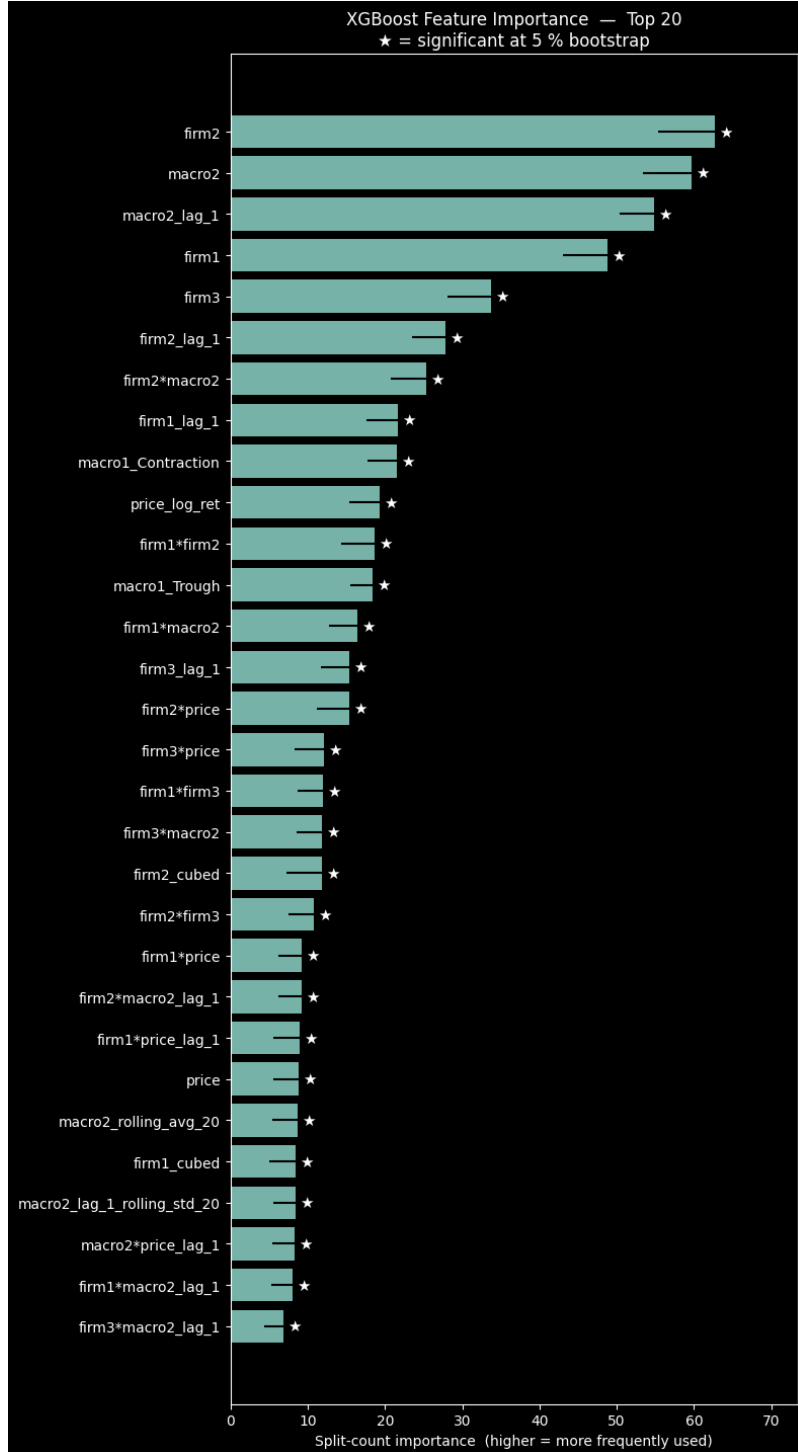
**Note:** Every feature in the Top 20 passes the bootstrap significance test (★).

Firm-ID dummies never make the top list; their mean split-count is zero and  $p = 1$ , matching earlier conclusions.

### Key Takeaways

- **Macro shocks (macro2, lag) remain indispensable** across both models, underscoring their foundational predictive power.
- **Nonlinear modelling boosts firm fundamentals.** XGBoost uncovers interaction-heavy relationships that Elastic-Net (restricted to additive linear terms) cannot exploit, elevating **firm2**, **firm1**, and their lagged or interacted variants.
- **Complex engineered features are now useful.** Tree ensembles harness cubic or rolling interactions that previously appeared redundant, hinting at higher-order effects worth retaining in nonlinear pipelines.
- **FirmID columns stay redundant** and can be safely pruned.

Overall, while the linear model captures the broad macro signal, XGBoost extracts richer, context-dependent patterns — especially those linking firm-level variables to macro states and recent price dynamics — delivering a more nuanced importance landscape.



Aspect	Elastic-Net	XGBoost
<b>Dominant driver</b>	macro2 (and its lag) lead the ranking	firm2 edges out macro variables, with macro2 & its lag still highly important
<b>Firm fundamentals</b>	Present but behind macro2; linear contribution only	Stronger presence (firm2, firm1, firm3 + lags) — boosted by nonlinear interactions
<b>Interactions &amp; higher-order terms</b>	Few survive significance (firm1 $\times$ macro2, price_log_ret)	Many more interaction terms gain importance, indicating XGBoost’s capacity to model complex feature interplay
<b>Macro-cycle dummies</b>	Insignificant	macro1.Trough, macro1.Contraction emerge as significant, suggesting state-dependent effects the linear model misses
<b>FirmID indicators</b>	No importance	Still none — confirms they add no incremental value once quantitative fundamentals are present

Table 5: Comparison of feature importance between Elastic-Net and XGBoost models