

AMProject__Clean

Laura Gullicksen, Erich Gozebina, Daria Palitzsch

23/05/2025

TODO:

- Chainlink vs Link Descriptions
- Table Nummerierung
- Table Referenz in den Texten erwähnen
- Figure Captions / Table Captions im Chunk Header
- Begründung des Schrittes von LM zu Lasso (Varianzreduktion, Korrelationseffekte, etc.)

1. Introduction

#TODO: describe data choice -> Laura

2. Data & Descriptive Analysis

Data Aggregation and Strategy Frequency

The raw dataset provides CHAINLINK price data at *hourly frequency*. While such high-frequency data offers more granular insights, we chose to **aggregate the data to daily frequency** for the following reasons:

1. **Alignment with Trading Strategy:** Our core trading strategy is based on a **7-day momentum signal**, which inherently reflects **weekly price trends**. Applying such a signal at an hourly resolution would not be consistent with the strategy's time horizon.
2. **Noise Reduction:** Hourly crypto data can be highly volatile and noisy. Aggregating to daily returns reduces **microstructure noise**, **short-term reversals**, and **Whale-driven price spikes**, improving the signal-to-noise ratio.
3. **Practical Execution Perspective:** A strategy that rebalances daily is **more realistic to implement**, considering gas fees, latency and operational constraints on decentralized exchanges or CEX APIs.

4. **Interpretability and Robustness:** Daily returns are more interpretable and robust across backtests. Most financial and technical indicators (e.g., RSI, MACD, SMA) are commonly applied on daily charts.

Our strategy issues long/short signals based on the past 7-day log return of CHAINLINK, i.e.,

$$\text{Momentum}_t^{(7)} = \log \left(\frac{P_t}{P_{t-7}} \right)$$

This naturally assumes daily data, as each observation reflects the cumulative return over the previous seven days.

In summary, aggregating to daily frequency is a theoretically and practically sound choice. It ensures consistency between our signal construction, model estimation, and backtesting logic.

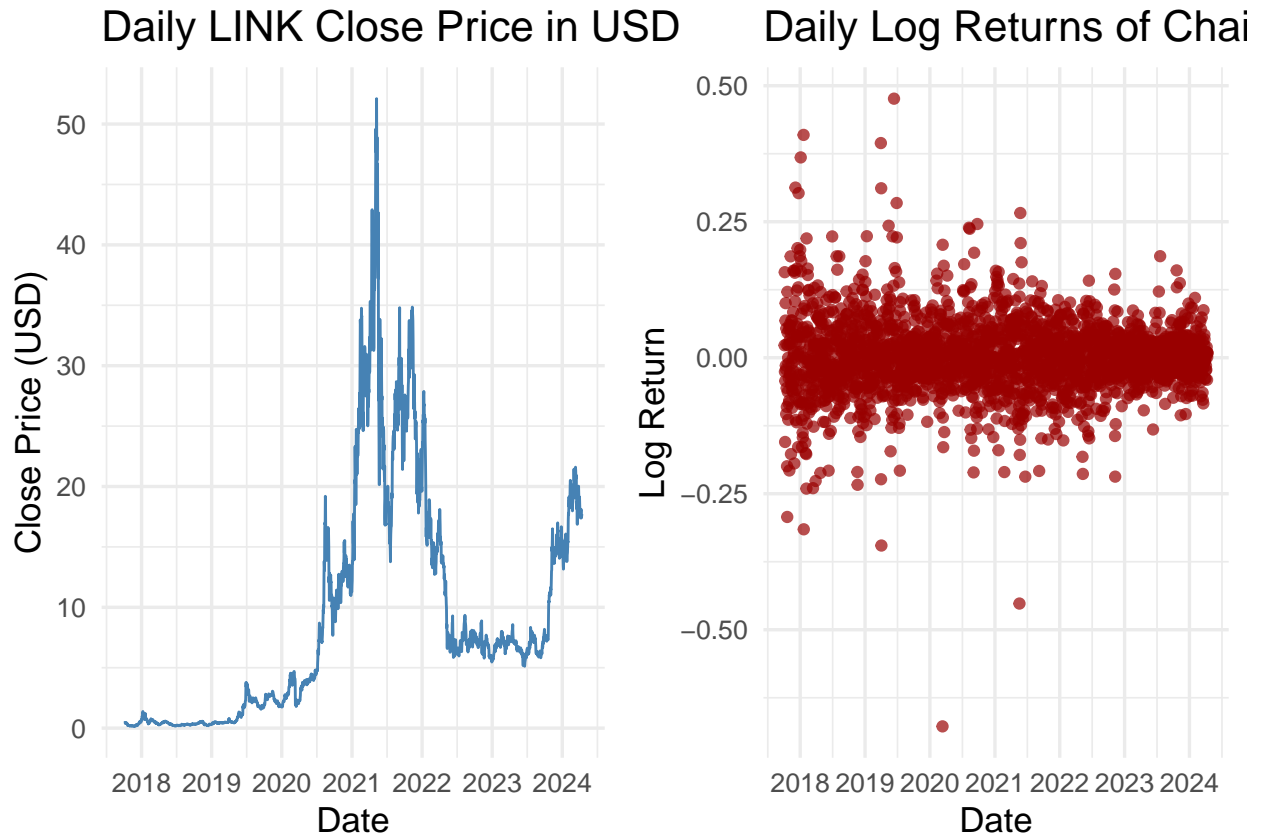


Figure 1: Close Price and Log Return of LINK

#DONE: Explain output and log choice -> Erich

Having a first glance on the LINK's price data, we see that there is not much movement until 2020, followed by sharp increases by factor ten in the consecutive nine months. The price reaches its peak of over 50\$/LINK in May 2021. In the second half of 2021 until April 2022 the data shows volatile behavior but with significant decrease on average. A trendless

period of relatively low volatility, starting in May 2022 and ending in October 2023, shows prices between 5 and 10\$/LINK. Finally, we observe another increase at the end of 2023 and beginning of 2024.

Moving away from non-stationary price data to stationary log returns, we observe a higher dispersion in the earlier years, suggesting a higher volatility then. Furthermore, most of points cluster around zero, indicating that there is no significant long term drift. Outliers, both positive and negative, imply extreme relative price movements, especially in the earlier period. Another important insight is the changing variance since the density of the points increases in the second half of the time window.

Why did we chose log returns over canonical (arithmetic) returns? Using log returns instead of canonical returns is a standard practice in financial econometrics and modeling.

$$\tilde{r}_t = \log(r_t + 1) = \log\left(\frac{p_t}{p_{t-1}}\right)$$

The underlying reason is the assumption that prices of an financial asset are log-normally distributed. This is reasonable since the log-normal distribution does not allow for negative values, which is also true for most asset prices (particularly for crypto currencies). Moreover, historical data provides evidence that the log-normal distribution gives a good fit for the prices of many financial assets. In reverse, since the logarithm function amplifies returns that are close to -1 more than positive returns, log-returns are distributed more symmetrically than canonical returns and indeed follow a normal distribution. Additionally, if returns are small, log returns approximate canonical returns very well. For x close to zero, it holds that

$$\log(x + 1) \approx x.$$

We can expect small returns since we shorten the considered time interval. Another important property is the additivity of log returns. It allows us to aggregate returns over multiple periods by summing up the pointwise log returns - a property that canonical returns miss. These properties make log returns more suitable for linear regression models, hypothesis testing, and machine learning regressors.

To better understand the characteristics of the Chainlink price and return series, we compute a set of descriptive statistics based on the daily close prices and the corresponding log returns. These statistics provide a first impression of the dataset's distribution, dispersion, and extreme values, and help assess whether further preprocessing or transformation steps are necessary before applying predictive models.

The summary statistics reveal that the mean daily log return of Chainlink is close to zero, while the standard deviation is relatively high, reflecting the well-known volatility of cryptocurrency markets. The minimum and maximum returns further highlight the presence of large price swings. The wide range between the minimum and maximum close prices illustrates the strong appreciation potential, but also the riskiness of the asset over the observation period.

Table 1: Summary Statistics for Chainlink Price and Returns

Statistic	Value
Number of Observations	2,383.0000
Mean Close Price	9.1484
Std. Dev. Close Price	9.5407
Minimum Close Price	0.1453
Maximum Close Price	52.1000
Mean Return	0.0016
Std. Dev. Return	0.0676
Minimum Return	-0.6776
Maximum Return	0.4762

To evaluate the temporal dependence structure of Chainlink’s daily log returns, we plot the autocorrelation function (ACF). The ACF helps determine whether past returns exhibit statistically significant correlation with future returns — a key consideration when assessing the potential for return predictability.

The autocorrelation function (ACF) of daily log returns shows no statistically significant linear dependence at any lag, indicating that past returns do not linearly predict future returns. This finding supports the weak-form Efficient Market Hypothesis (EMH). However, it does not rule out the presence of exploitable patterns captured by non-linear or directional indicators. Therefore, we proceed with a momentum-based trading strategy, leveraging the sign of multi-day past returns to generate long or short signals.

3. Standard Model

We develop a basic model that will serve as a starting point for an extended model. Our first approach to predict future returns is a simple linear regression. Why linear regression? First, this technique is the underlying mechanism of many advanced models that are often generalizations of the linear case. Therefore, it is a good fit for a starting point. In general, linear regression aims to identify linear relationships between input data and the target dimension. In our case the input data is price data, trading volume, market capitalization, and every predictor that is derived from those - returns for instance. The target dimension that we are going to predict is the return of the next day. The simplicity of linear mappings make results easy to interpret, whereby the model still remains powerful since many observed relationships are indeed of linear nature. Moreover, linear regression indicates the strength of those linear ties what makes it a helpful tool for decision making.

7-Day Momentum Signal Strategy

We define the 7-day momentum as the log return over the past 7 days:

$$\text{Momentum}_t = \log \left(\frac{P_t}{P_{t-7}} \right)$$

Autocorrelation of Daily Log Returns (LINK)

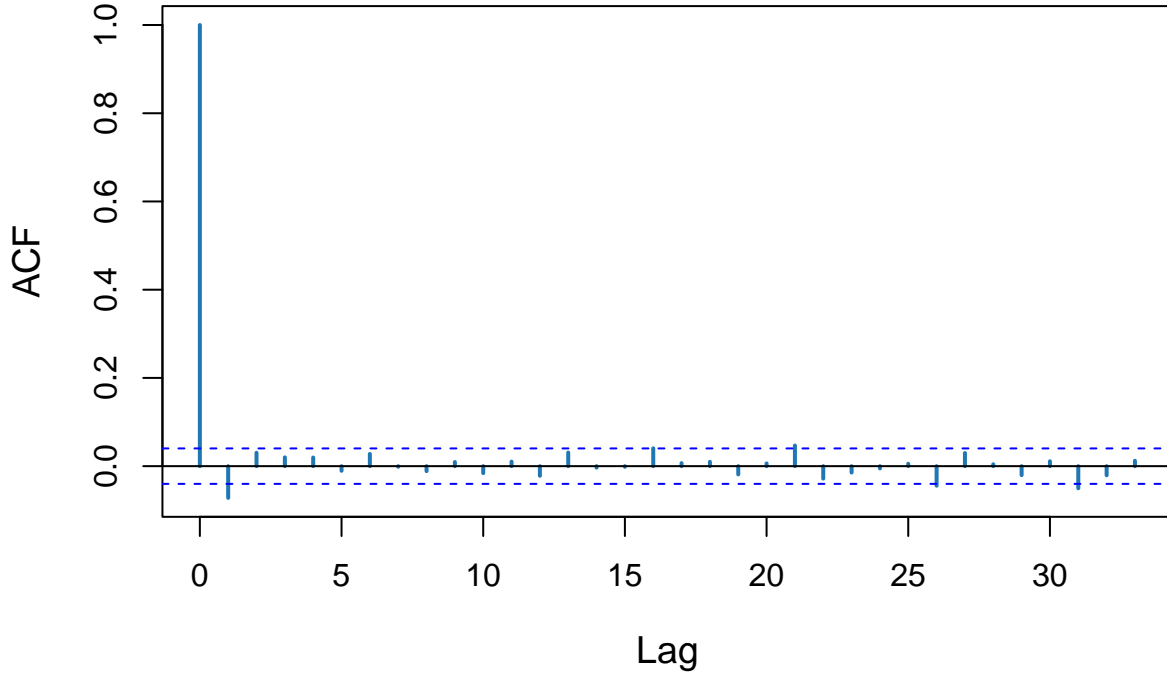


Figure 2: ACF of Daily Log Returns for Chainlink

The trading signal is then determined as:

$$\text{Signal}_t = \begin{cases} +1 & \text{if Momentum}_t > 0 \quad (\text{go long}) \\ -1 & \text{if Momentum}_t < 0 \quad (\text{go short}) \\ 0 & \text{otherwise (no position)} \end{cases}$$

The strategy return is computed as:

$$r_{t+1}^{\text{strategy}} = \text{Signal}_t \cdot r_{t+1}$$

where $r_{t+1} = \log\left(\frac{P_{t+1}}{P_t}\right)$ is the daily log return.

#TODO: insert standard model with momentum -> Erich **Regress target return on 7-day momentum**

Call: `lm(formula = target_return ~ momentum_7d, data = df_train_stdmodel)`

Residuals: Min 1Q Median 3Q Max -0.67856 -0.03704 -0.00046 0.03627 0.47476

Coefficients: Estimate Std. Error t value Pr(>|t|) (Intercept) 0.001334 0.001680 0.794 0.427 momentum_7d 0.002168 0.008816 0.246 0.806

Residual standard error: 0.07313 on 1898 degrees of freedom Multiple R-squared: 3.185e-05, Adjusted R-squared: -0.000495 F-statistic: 0.06046 on 1 and 1898 DF, p-value: 0.8058

Table 2: Regression Results: 7-Day Momentum Strategy

	<i>Dependent variable:</i>
	Strategy Return
Intercept	0.0022 (0.0088)
7-Day Momentum	0.0013 (0.0017)
Observations	1,900
R ²	0.00003
Adjusted R ²	−0.0005
Residual Std. Error	0.0731 (df = 1898)
F Statistic	0.0605 (df = 1; 1898)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

#TODO: explain result of standard 7 day momentum strategy -> Laura # Draft from Erich: The linear regression results in an estimation for the intercept of 0.001 having a p-value of 0.427 and a slope of 0.002 having a p-value of 0.806 Since the p-value for the slope clearly exceeds 0.05 (the threshold, which is commonly used for acceptance), we can not conclude a linear relation between the 7-day momentum and the future day return. Actually, this aligns with our finding in the ACF analysis, where no linear dependency between a lag of 7 days and the future day return was indicated. Thus, this result calls for an extended approach of future return prediction.

4. Extension

Extension of our OLS

To enhance the predictive power of the benchmark model, we extend it by incorporating a broader set of explanatory variables that capture not only short- and medium-term price dynamics, but also market sentiment, technical indicators, and inter-asset relationships. These include:

- Momentum indicators over 3, 7, and 14 days,
- Lagged daily returns (1-day and 2-day),
- A 7-day rolling volatility measure,
- Technical indicators such as the 14-day Relative Strength Index (RSI), MACD value and histogram, Simple Moving Average difference, and Average True Range (ATR),
- Day-of-week dummy variables to capture potential calendar effects,

- BTC-based predictors: daily BTC return, 7-day BTC momentum, and 7-day BTC volatility,
- ETH-based predictors: daily ETH return, 7-day ETH momentum, and 7-day ETH volatility,
- ETH trading volume: daily ETH volume return, 7-day ETH volume momentum, and 7-day ETH volume volatility,
- ETH market capitalization: daily ETH market capitalization return, 7-day ETH market capitalization momentum, and 7-day ETH market capitalization volatility,
- Ethereum gas fees: daily gas return, 7-day gas momentum, and 7-day gas volatility.

The extended predictive regression model is specified as:

$$r_{t+1} = \alpha + \sum_{h \in \{3,7,14\}} \beta_h \cdot \text{Momentum}_t^{(h)} + \gamma_1 \cdot r_t + \gamma_2 \cdot r_{t-1} + \delta \cdot \text{Volatility}_t^{(7)} + \sum_j \theta_j \cdot X_t^{(j)} + \varepsilon_{t+1}$$

where $X_t^{(j)}$ represents the set of technical indicators (RSI, MACD, ATR, SMA), weekday dummies, and BTC-based predictors.

$$\begin{aligned} r_{t+1} &:= \log \left(\frac{P_{t+1}}{P_t} \right) \quad (\text{one-day-ahead LINK return}) \\ \text{Momentum}_t^{(h)} &:= \log \left(\frac{P_t}{P_{t-h}} \right) \quad \text{for } h \in \{3, 7, 14\} \\ \text{Volatility}_t^{(7)} &:= \text{std} (r_{t-6}, \dots, r_t) \\ \text{BTC return}_t &:= \log \left(\frac{P_t^{\text{BTC}}}{P_{t-1}^{\text{BTC}}} \right) \\ \text{BTC Momentum}_t^{(7)} &:= \log \left(\frac{P_t^{\text{BTC}}}{P_{t-7}^{\text{BTC}}} \right) \\ \text{BTC Volatility}_t^{(7)} &:= \text{std} (r_{t-6}^{\text{BTC}}, \dots, r_t^{\text{BTC}}) \end{aligned}$$

The ETH-based predictors are constructed analogously to the BTC-based predictors. The parametrization of this model is estimated via Ordinary Least Squares (OLS) on the in-sample period. By incorporating this rich feature set, we aim to capture a range of return drivers including price trends, market overreaction, volatility clustering, inter-market dependencies, and behavioral biases tied to trading weekdays.

#DONE: add ethereum data -> Erich

Regress target return on all features

##

Call:

```

## lm(formula = target_return ~ momentum_3d + momentum_7d + momentum_14d +
##      return_lag1 + return_lag2 + volatility_7d + rsi_14 + sma_diff +
##      macd_val + macd_hist + atr_14 + monday + tuesday + wednesday +
##      thursday + friday + btc_return + btc_momentum_7d + btc_volatility_7d +
##      eth_return + eth_momentum_7d + eth_volatility_7d + eth_return_volume +
##      eth_momentum_7d_volume + eth_volatility_7d_volume + eth_return_marketcap +
##      eth_momentum_7d_marketcap + eth_volatility_7d_marketcap +
##      eth_return_gas + eth_momentum_7d_gas + eth_volatility_7d_gas,
##      data = df_train_extended)
##
## Residuals:
##      Min        1Q    Median        3Q        Max
## -0.35077 -0.03072 -0.00496  0.02357  0.46199
##
## Coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -2.915e-02  1.804e-02  -1.616   0.1062
## momentum_3d     -8.156e-02  3.387e-02  -2.408   0.0161 *
## momentum_7d     -1.558e-02  1.941e-02  -0.803   0.4223
## momentum_14d    -1.919e-02  1.745e-02  -1.099   0.2718
## return_lag1      7.957e-02  3.314e-02   2.401   0.0164 *
## return_lag2      9.199e-02  3.603e-02   2.553   0.0108 *
## volatility_7d    -7.095e-02  5.478e-02  -1.295   0.1954
## rsi_14           3.837e-04  3.351e-04   1.145   0.2522
## sma_diff        -1.010e-03  1.785e-03  -0.566   0.5717
## macd_val         3.820e-04  5.434e-04   0.703   0.4822
## macd_hist        1.085e-03  1.713e-03   0.633   0.5266
## atr_14          -2.003e-03  1.354e-03  -1.479   0.1393
## monday           6.314e-03  4.837e-03   1.305   0.1920
## tuesday          -9.696e-04  4.526e-03  -0.214   0.8304
## wednesday        1.387e-03  4.515e-03   0.307   0.7588
## thursday         4.120e-03  4.506e-03   0.914   0.3606
## friday           7.899e-03  4.472e-03   1.766   0.0775 .
## btc_return       -1.011e+00  3.740e-02 -27.033 <2e-16 ***
## btc_momentum_7d  -3.993e-03  1.426e-02  -0.280   0.7794
## btc_volatility_7d -1.576e-01  7.441e-02  -2.117   0.0344 *
## eth_return       -1.070e+00  7.790e+00  -0.137   0.8908
## eth_momentum_7d  -6.644e+00  4.127e+00  -1.610   0.1076
## eth_volatility_7d  2.317e+01  1.390e+01   1.667   0.0957 .
## eth_return_volume  1.642e-03  7.091e-03   0.232   0.8169
## eth_momentum_7d_volume  4.401e-03  4.984e-03   0.883   0.3773
## eth_volatility_7d_volume  2.616e-02  1.566e-02   1.670   0.0950 .
## eth_return_marketcap  1.088e+00  7.792e+00   0.140   0.8890
## eth_momentum_7d_marketcap  6.666e+00  4.128e+00   1.615   0.1065
## eth_volatility_7d_marketcap -2.298e+01  1.390e+01  -1.654   0.0983 .

```



```
## eth_return_gas          3.001e-03  5.264e-03   0.570   0.5686
## eth_momentum_7d_gas    -1.566e-03  3.571e-03  -0.439   0.6610
## eth_volatility_7d_gas   -3.109e-03  6.731e-03  -0.462   0.6442
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.05938 on 1847 degrees of freedom
## Multiple R-squared:  0.342, Adjusted R-squared:  0.331
## F-statistic: 30.97 on 31 and 1847 DF, p-value: < 2.2e-16

#TODO: generate nicer latex table output of the regression results -> Daria
#TODO: Description and interpretation of output -> Laura
```

Lasso Model

To prevent overfitting and perform automatic variable selection, we extend our linear modeling approach using the Lasso (Least Absolute Shrinkage and Selection Operator). The Lasso adds a penalty term to the standard OLS loss function, shrinking some coefficient estimates toward zero. This results in a sparse model that may improve predictive performance, particularly when dealing with multiple correlated predictors. Furthermore, reducing the number of relevant features allows for a better interpretation of simulation results.

The Lasso estimator is defined as the solution to the following optimization problem:

$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta_0, \beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

where:

- y_i is the target variable (e.g., one-day-ahead return),
- x_{ij} are the predictor variables,
- β_j are the coefficients,
- $\lambda \geq 0$ is the tuning parameter controlling the strength of the penalty.

As λ increases, more coefficients are shrunk toward zero. For $\lambda = 0$, the solution coincides with OLS.

We use 10-fold cross-validation to select the optimal λ that minimizes the mean squared prediction error on held-out data.

Optimal Lambda from Cross-Validation: $\lambda^* = 0.002780$

Given this optimal lambda, we now run the LASSO regression, including all variables from the previous OLS regression:

```
## 29 x 1 sparse Matrix of class "dgCMatrix"
##                                     s0
```

Table 3: Extended Regression Model: Predicting LINK Returns with Crypto Features

	<i>Dependent variable:</i>
	Target Return
Intercept	−0.0816** (0.0339)
Momentum (3d)	−0.0156 (0.0194)
Momentum (7d)	−0.0192 (0.0175)
Momentum (14d)	0.0796** (0.0331)
Lagged Return (1d)	0.0920** (0.0360)
Lagged Return (2d)	−0.0710 (0.0548)
Volatility (7d)	0.0004 (0.0003)
RSI (14)	−0.0010 (0.0018)
SMA Diff	0.0004 (0.0005)
MACD Value	0.0011 (0.0017)
MACD Histogram	−0.0020 (0.0014)
ATR (14)	0.0063 (0.0048)
Monday	−0.0010 (0.0045)
Tuesday	0.0014 (0.0045)
Wednesday	0.0041 (0.0045)

```

## (Intercept)                0.001414259
## momentum_3d                .
## momentum_7d                .
## momentum_14d               .
## return_lag1                .
## return_lag2                .
## volatility_7d               .
## rsi_14                     .
## sma_diff                   .
## macd_val                   .
## macd_hist                  .
## atr_14                     .
## monday                     .
## tuesday                    .
## wednesday                  .
## thursday                   .
## friday                      .
## btc_return                  -0.955207922
## btc_momentum_7d            .
## btc_volatility_7d           .
## eth_return_volume           .
## eth_momentum_7d_volume      0.001887347
## eth_volatility_7d_volume    .
## eth_return_marketcap        .
## eth_momentum_7d_marketcap   .
## eth_volatility_7d_marketcap .
## eth_return_gas              .
## eth_momentum_7d_gas         .
## eth_volatility_7d_gas       .

```

Table 4: Non-Zero Coefficients from LASSO Regression

Predictor	Coefficient
Intercept	0.001414
Bitcoin Daily Return	-0.955208
Ethereum 7-Day Volume Momentum	0.001887

The LASSO regression identified two non-zero predictors for explaining Chainlink (LINK) returns:

1. Bitcoin Daily Return (Coefficient: -0.9552): This variable has a large and negative coefficient, indicating that when Bitcoin's daily return increases by 1 unit (in our scaled units), the predicted return of our LINK-based trading strategy decreases by approximately 0.9552 units, all else equal. This suggests a strong inverse relationship

between BTC movements and our strategy, potentially due to hedging behavior or negative spillovers.

2. **Ethereum 7-Day Volume Momentum (Coefficient: 0.0019):** This predictor captures short-term trends in Ethereum’s trading volume. The positive but small coefficient implies that higher recent momentum in ETH trading volume is weakly associated with increased LINK returns, possibly due to spillover effects from rising market activity in related tokens.
3. **Intercept (Coefficient: 0.0014):** The intercept represents the model’s baseline prediction when all predictors are zero. Here, it suggests a small positive base return, though in practice this often has less interpretive value than the covariates.

5. Forecasting & Backtesting

In-Sample testing

To evaluate the performance of our predictive models, we begin by conducting in-sample (IS) testing. This involves fitting each model on a fixed training sample and evaluating how well the model explains historical variation in the data.

We assess in-sample performance using the following criteria:

- **Mean Squared Error (MSE):** Measures the average squared difference between predicted and actual returns.

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

- **Adjusted R^2 :** Indicates the proportion of variance explained by the model, adjusted for the number of predictors.

$$R_{\text{adj}}^2 = 1 - \frac{\text{RSS}/(n - p - 1)}{\text{TSS}/(n - 1)}$$

- **Directional Accuracy:** The fraction of times the predicted direction matches the actual direction of returns.

$$\text{Accuracy} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\text{sign}(\hat{y}_i) = \text{sign}(y_i))$$

These metrics are computed for all three models:

1. Benchmark (7-day momentum only),
2. Extended linear model with multiple features,
3. Lasso-regularized regression with automatic feature selection.

```
## # A tibble: 3 x 4
##   Model      MSE Directional_Accuracy  Adj_R2
##   <chr>      <dbl>                <dbl>    <dbl>
## 1 Benchmark 0.00527                0.505 -0.000533
## 2 Extended 0.00347                0.711  0.331
## 3 Lasso    0.00357                0.714  NA
```

#TODO: interpret results

Out-of-sample testing: # TODO: Verlauf des MSE im Vergleich von IS zu OOS

#TODO: review code, does not work at the moment

#TODO: Evaluate: o Sharpe ratio o Cumulative return o OOSR2 o Hit rate (how often you correctly predict direction)

Table 5: Table 5: Out-of-Sample Performance Metrics of LASSO Model

R-squared	MSE	Directional Accuracy	Sharpe Ratio (Actual)	Sharpe Ratio (Predicted)
0.5763	0.0016	0.9757	0.5055	0.8349

add transaction fees as extra path

6. Conclusion