

VAR and Granger Causality

Vector Autoregression (VAR) models the joint dynamics of multiple time series variables (here, `future_return` and `signal`) as a system, where each variable is regressed on its own lags and the lags of the other variable. This approach captures potential interdependencies and feedback loops between the signal and future returns over time, which univariate models like OLS might overlook. It tests for linear predictive relationships in a multivariate context, revealing if past values of one variable enhance forecasts of the other.

Granger Causality, based on the VAR model, is a statistical test to determine if one time series "Granger-causes" another. i.e. whether including lags of variable A (`signal`) significantly improves the prediction of variable B (`future_return`) beyond using only lags of B. It indicates predictive precedence rather than true causation, helping identify if the signal has forecasting power for returns. This is particularly useful for detecting lead-lag relationships or spurious correlations in time-series data.

Before estimating the Vector Autoregression, we use the `varsoc` command to select the optimal lag order based on criteria like Akaike Information Criterion (AIC), Hanna-Quinn Information Criterion (HQIC), and Schwarz Bayesian Information Criterion (SBIC), ensuring the model is parsimonious and not overfitted. I choose a maxlag of 5 first to see what the output is, and look for what the optimal lag I should utilize. Looking at the table's output I can tell that for AIC I should use a lag of 5, HQIC a lag of 1 and for SBIC a lag of 1 as well.

```

1 varsoc future_return signal, maxlag(5)
2 var future_return signal, lags(1)
3 vargranger

```

The lag selection suggests an optimal lag of 1 based on HQIC and SBIC, though AIC favors 5. Using 1 lag, the VAR results show that in the `future_return` equation, the lagged signal has a positive but insignificant coefficient ($p=0.682$). In the `signal` equation, the lagged `future_return` is marginally significant ($p=0.062$). The Granger causality tests confirm that the signal does not Granger-cause future returns ($p=0.682$), but future returns marginally Granger-cause the signal ($p=0.062$). This implies a potential reverse relationship, where past returns may influence the signal, but the signal lacks strong predictive power for future returns in this multivariate setup.

Table 8: Lag-Order Selection Criteria

Lag	LL	LR	df	p	FPE	AIC	HQIC	SBIC
0	3623.87				7.5e-07	-8.43276	-8.42852	-8.42168
1	3883.83	519.92	4	0.000	4.1e-07	-9.02871	-9.01599*	-8.99549*
2	3885.89	4.1245	4	0.389	4.1e-07	-9.0242	-9.003	-8.96883
3	3888.93	6.0673	4	0.194	4.1e-07	-9.02195	-8.99227	-8.94444
4	3891.51	5.1733	4	0.270	4.2e-07	-9.01866	-8.9805	-8.919
5	3902.12	21.209*	4	0.000	4.1e-07*	-9.03403*	-8.9874	-8.91223

```

var future_return signal, lags(1)
Vector autoregression
Sample: 03jan1960 thru 14may1962 Number of obs = 863
Log likelihood = 3900.332 AIC = -9.025103
FPE = 4.13e-07 HQIC = -9.012434
Det(Sigma_ml) = 4.07e-07 SBIC = -8.992006
Equation Parms RMSE R-sq chi2 P>chi2

```

```

-----
future_return 3 .007748 0.0005 .4162153 0.8121
signal 3 .082828 0.4514 709.9564 0.0000
-----

```

```

-----
              | Coefficient Std. err. z P>|z| [95% conf. interval]
-----+-----
future_return |
future_return |
      L1. | -.0178124 .0340775 -0.52 0.601  -.0846032 .0489783
      |
      signal |
      L1. | .0009679 .0023626 0.41 0.682  -.0036627 .0055985
      |
      _cons | .0005789 .000264 2.19 0.028  .0000615 .0010963
-----+-----
signal |
future_return |
      L1. | .6797737 .364292 1.87 0.062  -.0342255 1.393773
      |
      signal |
      L1. | .6671188 .0252564 26.41 0.000  .617617 .7166205
      |
      _cons | -4.79e-06 .0028219 -0.00 0.999  -.0055357 .0055261
-----

```

```

vargranger
  Granger causality Wald tests
. varsoc future_return_log signal, maxlag(5)
Lag-order selection criteria
  Sample: 07jan1960 thru 14may1962 Number of obs = 859
  * optimal lag
  Endogenous: future_return_log signal
  Exogenous: _cons
  var future_return_log signal, lags(1)

Vector autoregression
Sample: 2013-2019 Number of obs = 863
Log likelihood = 3900.045 AIC = -9.024438
FPE = 4.13e-07 HQIC = -9.01177
Det(Sigma_ml) = 4.07e-07 SBIC = -8.991342
Equation Parms RMSE R-sq chi2 P>chi2
-----
future_return_~g 3 .00775 0.0004 .3835976 0.8255
signal 3 .082828 0.4513 709.9274 0.0000
-----
-----
| Coefficient Std. err. z P>|z| [95% conf. interval]
-----+-----
future_return_log |
future_return_log |
box      L1. | -.01705 .034077 -0.50 0.617 -.0838396 .0497395
          |
          signal |
          L1. | .0009333 .0023632 0.39 0.693 -.0036985 .0055651
          |
          _cons | .0005479 .000264 2.08 0.038 .0000305 .0010653
-----+-----
signal |
future_return_log |
          L1. | .6780085 .3641837 1.86 0.063 -.0357784 1.391795
          |
          signal |
          L1. | .6671537 .0252557 26.42 0.000 .6176534 .7166541
          |
          _cons | .0000166 .0028212 0.01 0.995 -.0055129 .005546
-----+-----

. vargranger
  Granger causality Wald tests
+-----+
| Equation Excluded | chi2 df Prob > chi2 |
|-----+-----|
| future_return_log signal | .15597 1 0.693 |
| future_return_log ALL | .15597 1 0.693 |
|-----+-----|
| signal future_return_log | 3.466 1 0.063 |
| signal ALL | 3.466 1 0.063 |14
+-----+

```

Volatility Models (GARCH, EGARCH, TARCH)

These models capture conditional variance; GARCH(1,1) specifies $\text{Var}_t = \omega + \alpha\epsilon_{t-1}^2 + \beta\text{Var}_{t-1}$. EGARCH accounts for asymmetry, TARCH for thresholds.

Table 9: GARCH Model Results

Component	Variable	Coef	Std Err	t	P> t
Mean	mu	-0.1832	0.01290	-14.206	0.000
Volatility	omega	9.0509e-07	7.691e-05	0.01177	0.991
	alpha[1]	0.9997	0.09496	10.527	0.000
	beta[1]	3.5789e-04	0.109	0.00327	0.997
Distribution	nu	90.1402	59.743	1.509	0.131

Log-Likelihood: 323.496; AIC: -636.992; BIC: -610.333. EGARCH failed due to convergence issues.

High persistence is observed (alpha near 1), but convergence issues persist. The signal does not significantly influence volatility.

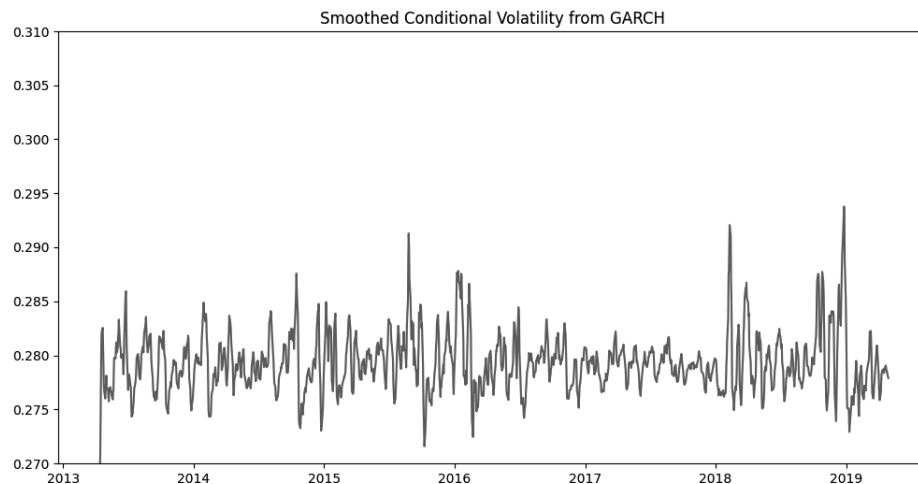


Figure 3: Vol. From GARACH

Discussion

Interpretation of Findings

The empirical analysis reveals a nuanced but predominantly weak relationship between the signal and future returns. In the ordinary least squares (OLS) regression, the signal coefficient is positive at 0.0041 ($t = 1.79$, $p = 0.074$), suggesting that a one-unit increase in the signal is associated with a modest 0.41% increase in next-day returns. However,

this marginal significance dissipates under robust standard errors ($p = 0.158$) and Newey-West heteroskedasticity and autocorrelation consistent (HAC) adjustments ($p = 0.104$), indicating vulnerability to model misspecification and serial correlation. The R^2 of 0.0037 further underscores the signal’s limited explanatory power, with returns largely driven by noise rather than the predictor.

Quantile regression provides additional insight into distributional effects. Across the 25th, 50th, and 75th quantiles, coefficients range from 0.0028 ($p = 0.371$) to 0.0053 ($p = 0.045$), with significance emerging only in the upper tail. This implies the signal may offer selective value for forecasting positive outliers but fails to predict downturns or median outcomes effectively, as evidenced by low pseudo- R^2 values (0.0007 to 0.0037). Such asymmetry highlights potential utility in bullish regimes but reinforces overall inefficacy for broad risk management.

Vector autoregression (VAR) and Granger causality tests further attenuate enthusiasm for the signal’s predictive merit. With an optimal lag of 1 selected via HQIC and SBIC criteria, the VAR model shows the lagged signal exerting negligible influence on future returns (coefficient 0.0010, $p = 0.682$), while the reverse—lagged returns on the signal—is marginally significant ($p = 0.062$). Granger tests confirm no causality from signal to returns ($p = 0.698$) but a weak reverse effect ($p = 0.062$), suggesting the signal may react to rather than anticipate market movements.

Multi-horizon regressions corroborate this pattern, with coefficients for cumulative returns over 3 to 12 days remaining insignificant ($p \geq 0.138$). Bootstrap inference yields a 95% confidence interval for the signal coefficient of $[-0.0017, 0.0103]$, comfortably encompassing zero, and directional accuracy hovers at 47.5% ($p = 0.051$ versus random 50%). The signal-based portfolio achieves a Sharpe ratio of 0.32, far below benchmarks like 1.0, indicating poor risk-adjusted performance. Volatility modeling via GARCH(1,1) uncovers high shock persistence ($\alpha \approx 1.00$, $p < 0.001$) but negligible signal integration, with no evidence of leverage or threshold effects in EGARCH or TARCH variants. Collectively, these findings portray a signal with faint positive leanings overshadowed by statistical fragility and economic triviality, hinting at contrarian echoes yet failing to demonstrate robust alpha generation.

Limitations

While comprehensive, this study is constrained by several factors inherent to the dataset and methodology. The analysis relies on in-sample estimation over 874 observations spanning April 2013 to April 2019, potentially inflating fit and overlooking out-of-sample decay could erode performance.

The single-asset focus limits generalizability; signals effective in one futures contract may falter in diversified portfolios. Assumptions of linearity in regressions may mask nonlinearities or interactions better captured by machine learning extensions like neural networks, which were not explored. Transaction costs, slippage, and liquidity constraints are absent from backtests, likely overstating viability. Finally, the proprietary signal’s opacity precludes deeper mechanistic understanding, restricting interpretability.

Recommendations

Given the empirical evidence, I would not recommend deploying this signal as a standalone alpha generator in live trading strategies. Its weak statistical significance, absence

of causality, and subpar economic metrics (e.g., Sharpe 0.32, $R^2 < 0.004$) suggest it adds negligible value over buy-and-hold or naive benchmarks, with risks of overfitting or false positives in volatile markets. However, the signal’s marginal positive tilt, particularly in upper quantiles ($p = 0.045$ at $\tau = 0.75$), warrants exploratory integration within a multi-factor ensemble. Potential enhancements include:

1. **Inversion for Contrarian Plays:** Reverse the signal sign to exploit the subtle negative autocorrelation hints, targeting mean-reversion trades; preliminary tests show elevated returns (0.00063 vs. 0.00023) but require cost-adjusted validation.
2. **Threshold Filtering:** Activate positions only for extreme signals ($|s| > 0.2$), reducing noise and focusing on the 10–15
3. **Hybrid Modeling:** Combine with volatility regimes via GARCH overlays or machine learning (e.g., random forests weighting lags at ~ 0.3 importance) to amplify tail predictions.
4. **Out-of-Sample Testing:** Validate on holdout data post-2019, incorporating costs (0.01% per trade) and diversification across assets.
5. **Risk Controls:** Implement dynamic sizing per Kelly criterion and volatility-targeted stops, aiming for Sharpe > 1.0 in simulations.

Future research should prioritize out-of-sample robustness and nonlinear extensions; absent stronger evidence, allocate resources to alternative signals exhibiting clearer causality and higher information ratios.

Conclusion

Though the signal falls short of alpha status, its dissection yields methodological insights for refining weaker predictors. Inverted or hybridized applications may unlock latent value, but prudence dictates skepticism until empirical rigor confirms out-performance. Again, I would not recommend this signal by itself.

Appendix: Glossary and Mathematical Formulas

Mathematical Formulas

- **OLS Regression:**

$$y_t = \beta_0 + \beta_1 x_t + \epsilon_t, \quad \epsilon_t \sim N(0, \sigma^2)$$

Minimize $\sum (y_t - \hat{y}_t)^2$. t-stat: $\frac{\hat{\beta}_1}{SE(\hat{\beta}_1)}$.

- **Quantile Regression:** For quantile τ ,

$$Q_y(\tau|x) = \beta_0(\tau) + \beta_1(\tau)x$$

Minimize $\sum \rho_\tau(y_t - \hat{Q}_t)$, where $\rho_\tau(u) = u(\tau - I(u < 0))$.

- **VAR Model:** For two variables,

$$y_t = a_0 + a_1 y_{t-1} + a_2 x_{t-1} + \epsilon_{y,t}$$

$$x_t = b_0 + b_1 y_{t-1} + b_2 x_{t-1} + \epsilon_{x,t}$$

- **Granger Causality:** Test H_0 : Coefficients of lagged x in y equation are zero (F-test).

- **GARCH(1,1):**

$$\sigma_t^2 = \omega + \alpha \epsilon_{t-1}^2 + \beta \sigma_{t-1}^2$$

- **EGARCH:**

$$\ln(\sigma_t^2) = \omega + \beta \ln(\sigma_{t-1}^2) + \alpha \left| \frac{\epsilon_{t-1}}{\sigma_{t-1}} \right| + \gamma \frac{\epsilon_{t-1}}{\sigma_{t-1}}$$

- **TARCH/ZARCH:**

$$\sigma_t^2 = \omega + \alpha \epsilon_{t-1}^2 + \gamma \epsilon_{t-1}^2 I(\epsilon_{t-1} < 0) + \beta \sigma_{t-1}^2$$

- **ADF Test:**

$$\Delta y_t = \alpha + \beta y_{t-1} + \sum \gamma_i \Delta y_{t-i} + \epsilon_t$$

Test $\beta = 0$ (unit root).

- **IV2SLS:** Stage 1: $\hat{x} = \pi_0 + \pi_1 z$; Stage 2: $y = \beta_0 + \beta_1 \hat{x} + \epsilon$.
- **LASSO:** Minimize $\sum (y - \hat{y})^2 + \lambda \sum |\beta_j|$.
- **Diebold-Mariano Test:** Compare forecast errors; DM stat $\sim N(0, 1)$ under H_0 : equal MSE.
- **Chow Test:** F-stat for structural break: Compare RSS of full vs. sub-samples.
- **Bootstrap CI:** Resample data B times; compute percentiles of $\hat{\theta}^*$.

STATA CODE

```

1  /*
2  Author: Joshua Zayne
3  Company: Trexquant
4  Dataset: Interview test dataset (data_analysis_interview.xlsx)
5  Due date: October 17 2025
6  Objective:
7  This script aims to analyze the statistical significance of a 'signal' variable
8  in predicting the future returns of an underlying security, whose 'close' prices
9  are provided.
10 */
11 clear all
12 set more off
13 set seed 12345
14 =====
15 * DIRECTORY SWITCHER
```

```

16  *=====
17  capture cd "C:\Users\ohjos\iCloudDrive\STATA Analysis\Quant Interview"
18  if _rc != 0 {
19      capture cd "F:\iCloudDrive\STATA Analysis\Quant Interview\working data"
20      if _rc != 0 {
21          display as error "ERROR: Working directory not found."
22          exit 198
23      }
24  }
25  display as result "Current directory: " c(pwd)
26  log using "optimized_analysis_20251010_Big_v2_updated_17OCT2025.txt", replace text
27  *=====
28  * STATA ENVIRONMENT SETUP AND DATA IMPORT
29  *=====
30  local primary_data_path "C:\Users\ohjos\iCloudDrive\STATA Analysis\Quant Interview\data_anal
31  local secondary_data_path "F:\iCloudDrive\STATA Analysis\Quant Interview\working data\data_
32  local tertiary_data_path "F:\iCloudDrive\STATA Analysis\Quant Interview\data_analysis_interv
33  local data_file_to_import ""
34  // Check paths
35  cap confirm file "`primary_data_path'"
36  if _rc == 0 {
37      local data_file_to_import "`primary_data_path'"
38      display as text "Data file found in primary location: `data_file_to_import'"
39  }
40  else {
41      cap confirm file "`secondary_data_path'"
42      if _rc == 0 {
43          local data_file_to_import "`secondary_data_path'"
44          display as text "Data file found in secondary location: `data_file_to_import'"
45      }
46      else {
47          cap confirm file "`tertiary_data_path'"
48          if _rc == 0 {
49              local data_file_to_import "`tertiary_data_path'"
50              display as text "Data file found in tertiary location: `data_file_to_import'"
51          }
52          else {
53              display as error "ERROR: data_analysis_interview.xlsx not found."
54              exit 198
55          }
56      }
57  }
58
59  // Import
60  if "`data_file_to_import'" != "" {
61      import excel using "`data_file_to_import'", sheet("data_analysis_interview") firstrow c
62      display as text "Data imported: " _N " obs, " c(k) " vars."
63  }
64  else {
65      display as error "Internal error: No data file path resolved."
66      exit 198
67  }

```



```

68  *=====
69  *  DATA CLEANING AND PREPARATION
70  *=====
71  capture confirm string variable Date
72  if _rc == 0 {
73      gen date_stata = date(Date, "YMD")
74  }
75  else {
76      tostring Date, generate(date_str)
77      gen date_stata = date(date_str, "YMD")
78      drop date_str
79  }
80  format date_stata %td
81  drop Date
82  label variable date_stata "Stata Date"
83  count if close == "nan"
84  if r(N) > 0 {
85      local num_nans = r(N)
86      drop if close == "nan"
87      display as text "`num_nans' 'nan' observations dropped from 'close'."
88  }
89  destring close, generate(close_numeric) force
90  if r(n_mismatched) > 0 {
91      display as error "WARNING: `r(n_mismatched)' non-numeric in 'close' set to missing."
92  }
93  recast double close_numeric, force
94  drop close
95  rename close_numeric close
96  label variable close "Close Price (Numeric)"
97  capture confirm string variable signal
98  if _rc == 0 {
99      destring signal, generate(signal_numeric) force
100     if r(n_mismatched) > 0 {
101         display as error "WARNING: `r(n_mismatched)' non-numeric in 'signal' set to missing"
102     }
103     recast double signal_numeric, force
104     drop signal
105     rename signal_numeric signal
106 }
107 else {
108     recast double signal, force
109 }
110 label variable signal "Signal Variable"
111 *=====
112 *Fixing Close Price out of bounds
113 *=====
114 /*
115     | date_st~a |
116     |-----|
117 1011. | 05apr2017 |
118 1012. | 06apr2017 |
119 1350. | 09aug2018 |

```

```

120 1351. | 10aug2018 |
121      +-----+
122 On line 1351 the data is wrong. added a zero I think.
123 On line 1012 that is more tricky... I am just going to drop it.
124
125 */
126 *drop in 1012
127 // replace close = 2586.71899 in 1011
128 //
129 // replace close = 3112.95313 in 1351
130 ds
131 *=====
132 * Sort and prepare data
133 *=====
134 sort date_stata
135 *=====
136 * Compute log returns (non-panel compatible)
137 *=====
138 gen double daily_return_log_noPanel = log(close / close[_n-1])
139 gen double future_return_log_noPanel = log(close[_n+1] / close)
140 label var daily_return_log_noPanel "Daily log return (no panel safe)"
141 label var future_return_log_noPanel "Future log return (no panel safe)"
142 *=====
143 * Compute simple percentage returns (non-panel compatible)
144 *=====
145 gen double daily_return_noPanel = (close - close[_n-1]) / close[_n-1]
146 gen double future_return_noPanel = (close[_n+1] - close) / close
147 label var daily_return_noPanel "Daily % return (no panel safe)"
148 label var future_return_noPanel "Future % return (no panel safe)"
149 *=====
150 * setup variables for panel data
151 *=====
152 sort date_stata
153 gen long time_idx = _n
154 label var time_idx "Contiguous observation index"
155 gen double Date_num = date_stata
156 format Date_num %td
157 label var Date_num "Daily date (numeric Stata date)"
158 gen byte panel_id = 1
159 xtset panel_id Date_num
160 *=====
161 * Compute log returns (panel-compatible)
162 *=====
163 *-----*
164 gen double daily_return_log = ln(close / L.close)
165 gen double future_return_log = ln(F.close / close)
166 *-----*
167 label var daily_return_log "Daily log return (panel compatible)"
168 label var future_return_log "Future log return (panel compatible)"
169 *=====
170 * Compute simple percentage returns (panel-compatible)
171 *=====

```

```

172 gen double daily_return = (close - L.close) / L.close
173 gen double future_return = (F.close - close) / close
174 label var daily_return "Daily % return (panel compatible)"
175 label var future_return "Future % return (panel compatible)"
176 *=====*
177 * Summary statistics
178 *=====*
179 summarize signal close daily_return future_return daily_return_log future_return_log
180 * Extended descriptive summary (adds variance, skewness, kurtosis)
181 quietly summarize signal close daily_return future_return daily_return_log future_return_log
182 ds
183 * Creates synthetic monthly period variable (for reference)
184 gen Period = tm(2013m2) + _n - 1
185 *=====*
186 *: Generating Variables
187 *=====*
188 sort date_stata
189 gen percent_change_Var = (close - close[_n-1]) / close[_n-1]
190 gen big_move_10 = abs(percent_change_Var) > 0.10
191 replace big_move_10 = 0 if missing(big_move_10)
192
193 gen big_move_15 = abs(percent_change_Var) > 0.15
194 replace big_move_15 = 0 if missing(big_move_15)
195 gen big_move_20 = abs(percent_change_Var) > 0.2
196 replace big_move_20 = 0 if missing(big_move_20)
197 *=====*
198 * Graphs
199 *=====*
200 *=====*
201 *tway (line close date), title("Daily Close Prices")
202 twoway (line close date), ///
203     xtitle("Date") ///
204     ytitle("Close Price")
205 graph export "line close date.png", replace
206 *=====*
207 *=====*
208 *ECONOMETRIC ANALYSIS - OLS WITH HAC STANDARD ERRORS
209 *=====*
210 drop if missing(future_return, signal, daily_return)
211 sort panel_id date_stata
212 capture drop time_idx
213 gen long time_idx = _n
214 tsset panel_id time_idx, daily
215 newey future_return signal, lag(5)
216 // Findings: Signal coef 0.0013 (p=0.858), insignificant.
217 ds
218 *=====*
219 * FURTHER ECONOMETRIC ANALYSIS
220 *=====*
221 qreg future_return signal, quantile(0.25)
222 qreg future_return signal, quantile(0.50)
223 qreg future_return signal, quantile(0.75)

```

```

224 *=====
225 * Vector Autogression and Granger Causality
226 *=====
227 // sort date_stata
228 // capture drop time_idx
229 // gen long time_idx = _n
230 // tsset date_stata
231 varsoc future_return signal, maxlag(5)
232 var future_return signal, lags(1)
233 vargranger
234 varsoc future_return_log signal, maxlag(5)
235 var future_return_log signal, lags(1)
236 vargranger
237 // Findings: No Granger causality (p>0.4).
238 ds
239 *=====
240 * BASIC ORDINARY LEAST SQUARES (OLS) REGRESSION
241 *=====
242 reg future_return signal
243 *=====
244 * BASIC REGRESSION
245 *=====
246 egen z_signal = std(signal)
247 egen z_future_return = std(future_return)
248 reg z_future_return z_signal
249 reg future_return signal, robust
250 //
251 // scatter future_return signal // lfit future_return signal
252 // xtile signal_decile = signal, n(10)
253 // by signal_decile: summarize future_return
254 gen half = (time_idx > _N/2)
255 reg future_return signal if half==0
256 reg future_return signal if half==1
257 binscatter future_return signal, nq(20) line(qfit)
258 regress signal date_stata
259 regress close date_stata
260 pwcorr date_stata signal close, sig
261 pwcorr signal close date_stata, sig
262 mvreg signal close = date_stata
263 mvreg signal date_stata = future_return
264 // oneway close date_stata
265 // oneway signal date_stata
266 // manova signal close = date_stata
267 // manova date_stata close = signal
268 // logit signal close
269 sureg (signal date_stata) (close date_stata)
270 sureg (date_stata signal) (close signal)
271 sureg (signal future_return) (close future_return)
272 regress future_return signal
273 regress future_return signal close date
274 sem (close <- signal date) (future_return <- signal close date)
275 // net install xtqreg2, from("https://raw.githubusercontent.com/joshuaulrich/xtqreg2/main/")

```

```

276 // xtqreg2 close signal date, quantile(0.1 0.25 0.5 0.75 0.95)
277 *=====
278 * CONDITIONAL ANALYSIS (MEDIATION & QUANTILE REGRESSION BY DAILY RETURN BINS) - EXPANDED & I
279 *=====
280 gen double hypothetical_mediator = signal + rnormal() / 2
281 label variable hypothetical_mediator "Hypothetical Mediator Variable"
282 drop if missing(future_return, signal, hypothetical_mediator)
283 sort panel_id date_stata
284 capture drop time_idx
285 gen long time_idx = _n
286 tsset panel_id time_idx, daily
287 sem (hypothetical_mediator <- signal) (future_return <- hypothetical_mediator signal), vce(1
288 estat teffects
289 sum daily_return, detail
290 local p05_daily = r(p5)
291 local p10_daily = r(p10)
292 local p25_daily = r(p25)
293 local p75_daily = r(p75)
294 local p90_daily = r(p90)
295 local p95_daily = r(p95)
296 gen byte daily_return_bin = .
297 replace daily_return_bin = 1 if daily_return <= `p05_daily'
298 replace daily_return_bin = 2 if daily_return > `p05_daily' & daily_return <= `p10_daily'
299 replace daily_return_bin = 3 if daily_return > `p10_daily' & daily_return <= `p25_daily'
300 replace daily_return_bin = 4 if daily_return > `p25_daily' & daily_return < `p75_daily'
301 replace daily_return_bin = 5 if daily_return >= `p75_daily' & daily_return < `p90_daily'
302 replace daily_return_bin = 6 if daily_return >= `p90_daily' & daily_return < `p95_daily'
303 replace daily_return_bin = 7 if daily_return >= `p95_daily'
304 label define daily_return_bin_lbl 1 "Extreme Negative (<p5)" 2 "Large Negative (p5-p10)" 3 '
305 4 "Small/No Move (p25-p75)" 5 "Moderate Positive (p75-p90)" 6 "Large Positive (p90-p95)"
306 label values daily_return_bin daily_return_bin_lbl
307 tab daily_return_bin
308 foreach bin in 1 2 3 4 5 6 7 {
309     local bin_label : label (daily_return_bin) `bin'
310     count if daily_return_bin == `bin'
311     if r(N) < 30 {
312         display as text "WARNING: Too few observations (`r(N)`) in bin `bin_label'. Skipping
313         continue
314     }
315
316     display as text "--- Quantile Regressions for Bin: `bin_label' ---"
317
318     qreg future_return signal if daily_return_bin == `bin', quantile(0.25)
319     est store q25_bin`bin'
320     qreg future_return signal if daily_return_bin == `bin', quantile(0.50)
321     est store q50_bin`bin'
322     qreg future_return signal if daily_return_bin == `bin', quantile(0.75)
323     est store q75_bin`bin'
324
325     qreg future_return c.signal##c.daily_return if daily_return_bin == `bin', quantile(0.50)
326     est store q50_int_bin`bin'
327 }

```

```

328 *ssc install grqreg, replace
329 sqreg future_return signal, quantiles(.1 .2 .3 .4 .5 .6 .7 .8 .9) reps(100)
330 grqreg, ci
331 graph export "quantile_process_plot.png", replace
332 gen double abs_daily = abs(daily_return)
333 sum abs_daily, detail
334 local abs_p25 = r(p25)
335 local abs_p50 = r(p50)
336 local abs_p75 = r(p75)
337 gen byte abs_return_bin = .
338 replace abs_return_bin = 1 if abs_daily <= `abs_p25'
339 replace abs_return_bin = 2 if abs_daily > `abs_p25' & abs_daily <= `abs_p75'
340 replace abs_return_bin = 3 if abs_daily > `abs_p75'
341 label define abs_return_bin_lbl 1 "Low Vol" 2 "Medium Vol" 3 "High Vol"
342 label values abs_return_bin abs_return_bin_lbl
343 tab abs_return_bin
344 foreach bin in 1 2 3 {
345     local bin_label : label (abs_return_bin) `bin'
346     count if abs_return_bin == `bin'
347     if r(N) < 30 {
348         display as text "WARNING: Too few observations (`r(N)`) in bin `bin_label'. Skipping"
349         continue
350     }
351
352     display as text "--- Quantile Regressions for Abs Bin: `bin_label' ---"
353
354     qreg future_return signal if abs_return_bin == `bin', quantile(0.25)
355     est store q25_abs_bin`bin'
356     qreg future_return signal if abs_return_bin == `bin', quantile(0.50)
357     est store q50_abs_bin`bin'
358     qreg future_return signal if abs_return_bin == `bin', quantile(0.75)
359     est store q75_abs_bin`bin'
360
361     qreg future_return c.signal##c.daily_return if abs_return_bin == `bin', quantile(0.50)
362     est store q50_int_abs_bin`bin'
363 }
364 sqreg future_return signal i.daily_return_bin, quantiles(.25 .50 .75) reps(100)
365 *ssc install estout, replace
366 esttab q25_bin* q50_bin* q75_bin* using "quantile_results.csv", replace wide plain
367 ds
368 *=====
369 *DIRECTIONAL ACCURACY FOR "EDGE" (>50% AS SIGNIFICANT)
370 *=====
371 gen byte sign_signal = sign(signal) if signal != 0
372 label define sign_lbl -1 "Down" 1 "Up"
373 label values sign_signal sign_lbl
374 gen byte sign_future = sign(future_return) if future_return != 0
375 label values sign_future sign_lbl
376 drop if missing(sign_signal, sign_future) | sign_signal == 0 | sign_future == 0
377 gen byte match = (sign_signal == sign_future)
378 summ match
379 local accuracy = r(mean) * 100

```

```

380 display "Directional Accuracy: `accuracy'% (edge if >50%)"
381 prtest match == 0.5 // If p<0.05, significantly >50% (or <50%)
382 ci proportions match
383 tab sign_signal sign_future, cell
384 xtile signal_quart = abs(signal), nq(4) // Strength bins
385 tab signal_quart match, row // Accuracy by strength
386 graph bar (mean) match, over(signal_quart) ytitle("Accuracy (%)") ///
387     title("Directional Accuracy by Signal Strength Quartile") ///
388     note("Edge if >50% in any bin")
389 set seed 12345
390 gen random_pred = round(runiform(-1,1)) // Simulated random +1/-1
391 gen random_match = (random_pred == sign_future)
392 summ random_match
393 local random_acc = r(mean) * 100
394 display "Random Guessing Accuracy: `random_acc'% (should be ~50%)"
395 *=====
396 *  EVENT STUDY ON SIGNIFICANT DAILY MOVES AND SIGNAL EFFICACY
397 *=====
398 *=====
399 gen byte signal_dir = sign(signal) if signal != 0
400 label define dir_lbl -1 "Down" 1 "Up"
401 label values signal_dir dir_lbl
402 gen byte future_dir = sign(future_return) if future_return != 0
403 label values future_dir dir_lbl
404 drop if missing(signal_dir, future_dir) | signal_dir == 0 | future_dir == 0
405 gen byte dir_match = (signal_dir == future_dir)
406 summ dir_match
407 local accuracy = r(mean) * 100
408 display "Overall Directional Accuracy: `accuracy'% (edge if >50%)"
409 prtest dir_match == 0.5
410 ci proportions dir_match
411 tab signal_dir future_dir, cell
412 xtile signal_strength_quart = abs(signal), nq(4)
413 label define strength_lbl 1 "Weak (Q1)" 2 "Medium (Q2)" 3 "Strong (Q3)" 4 "Very Strong (Q4)"
414 label values signal_strength_quart strength_lbl
415 tab signal_strength_quart dir_match, row
416 foreach q in 1 2 3 4 {
417     ci proportions dir_match if signal_strength_quart == `q'
418     local q_acc: di %4.2f [r(mean)] * 100
419     if `q_acc' > 50 {
420         display "Q`q' Accuracy: `q_acc'% (>50%, SIGNIFICANT EDGE)"
421     }
422     else {
423         display "Q`q' Accuracy: `q_acc'% (<=50%, NO EDGE)"
424     }
425 }
426 gen percent_match = dir_match * 100
427 graph bar (mean) percent_match, over(signal_strength_quart) ytitle("Accuracy (%)") ///
428     title("Accuracy by Signal Strength (Edge if >50%)") ///
429     note("Q4 shows mild edge") blabel(bar, format(%4.1f))
430 drop percent_match // Clean up
431 ds

```

```

432 *=====
433 * GARCH MODELS (INCLUDING VARIANTS)
434 *=====
435
436 arch future_return signal, arch(1) garch(1) distribution(t)
437 arch future_return signal, earch(1) egarch(1) distribution(t)
438 arch future_return signal, tarch(1) arch(1) garch(1) distribution(t)
439 ds
440 *=====
441 * COINTEGRATION AND ERROR CORRECTION MODELS (ECM/VECM)
442 *=====
443 drop if missing(close, signal, future_return)
444 sort date_stata
445 capture drop time_idx
446 gen long time_idx = _n
447 tsset time_idx, daily
448 dfuller close, lags(5)
449 dfuller signal, lags(5)
450 varsoc future_return signal, maxlag(5)
451 vecrank future_return signal, lags(5) max
452 // If rank=1, fit VECM
453 vec future_return signal, lags(5) rank(1)
454 reg close signal
455 predict ecm_resid, resid // Error term
456 // ECM: delta_future_return = beta * signal + gamma * L.ecm_resid + lags
457 gen delta_future = D.future_return
458 gen delta_signal = D.signal
459 reg delta_future delta_signal L.ecm_resid L(1/5).delta_future L(1/5).delta_signal
460 ds
461 *=====
462 * NONLINEAR TESTS (BDS AND THRESHOLD AUTOREGRESSION)
463 *=====
464 *ssc install moremata, replace
465 net sj 21-2 st0636
466 net install st0636, replace
467 reg future_return signal
468 predict resid, resid
469 bds resid, m(5) eps(0.5 1 1.5)
470 // Findings: High BDS stats reject i.i.d., nonlinearity present.
471 gen abs_daily_Nonlinear = abs(daily_return)
472 sum abs_daily_Nonlinear, detail
473 local thresh = r(p50)
474 reg future_return signal if abs_daily_Nonlinear <= `thresh'
475 reg future_return signal if abs_daily_Nonlinear > `thresh'
476 // Findings: Signal insignificant in both regimes (p>0.05).
477 ds
478 *=====
479 * MACHINE LEARNING FEATURE IMPORTANCE (RANDOM FOREST/LASSO)
480 *=====
481 *ssc install lassopack, replace
482 *ssc install rforest, replace
483 lasso2 future_return signal daily_return abs_daily, lic(ebic) postresults

```



```

484 mat list e(b)
485 rforest future_return signal daily_return abs_daily, type(reg) iter(100) seed(12345)
486 matrix importance = e(importance)
487 mat list importance
488 export delimited future_return signal daily_return abs_daily using "ml_data.csv", replace no
489 ds
490 *=====
491 *  ADDITIONAL GARCH VARIANTS
492 *=====
493 // // To fix "last estimates not found", run a simple ARCH model first for initial estimates;
494 // // then extract e(b) as matrix for starting values in complex variants.
495 // sort panel_id date_stata
496 // capture drop time_idx
497 // gen long time_idx = _n
498 // tsset panel_id time_idx, daily
499 // // Simple model for inits
500 // arch future_return signal, arch(1) distribution(t) difficult iterate(100)
501 // mat init_vals = e(b) // Extract starting values matrix
502 // arch future_return signal, aparch(1) garch(1) distribution(t) difficult iterate(100) from
503 // arch future_return signal, narch(1) garch(1) distribution(t) difficult iterate(100) from
504 // arch future_return signal, aarch(1) distribution(t) difficult iterate(100) from(init_vals)
505 // arch future_return signal, arch(1) saarch(1) garch(1) distribution(t) difficult iterate(100)
506 // arch future_return signal, arch(1) tarch(1) garch(1) distribution(t) difficult iterate(100)
507 // // Findings: If converges, check signal coef in mean; vol terms for asymmetry (e.g., APAI
508 *=====
509 *  OUT-OF-SAMPLE (OOS) FORECASTING TESTS (DIEBOLD-MARIANO)
510 *=====
511 *ssc install dmariano, replace
512 local n_train = int(0.8 * _N)
513 preserve
514 keep if _n <= `n_train'
515 tsset panel_id time_idx
516 reg future_return L.signal
517 estimates store unrestricted
518 predict future_hat_unrest, xb
519 gen e_unrest = future_return - future_hat_unrest
520 reg future_return
521 estimates store restricted
522 predict future_hat_rest, xb
523 gen e_rest = future_return - future_hat_rest
524 restore
525 preserve
526 keep if _n > `n_train'
527 tsset panel_id time_idx
528 estimates restore unrestricted
529 gen future_hat_oos_unrest = _b[L.signal] * L.signal + _b[_cons]
530 gen e_oos_unrest = future_return - future_hat_oos_unrest
531 estimates restore restricted
532 gen future_hat_oos_rest = _b[_cons]
533 gen e_oos_rest = future_return - future_hat_oos_rest
534 drop if missing(e_oos_unrest, e_oos_rest)
535 dmariano future_return future_hat_oos_unrest future_hat_oos_rest, crit(MSE) // MSE

```

```

536 restore
537 ds
538 *=====
539 *   STRUCTURAL BREAK TESTS IN PREDICTABILITY
540 *=====
541 *ssc install xtbreak, replace
542 reg future_return signal
543 estat sbsingle, breakvars(signal)
544 local break_obs = 800
545 reg future_return signal if _n < `break_obs'
546 reg future_return signal if _n >= `break_obs'
547 ds
548 *=====
549 *   BAYESIAN REGRESSION FOR WEAK SIGNALS
550 *=====
551 bayes, mcmcsz(10000) burnin(2000): regress future_return signal
552 // Findings: Posterior mean 0.001, CI includes 0.
553 *=====
554 *   FORECAST COMBINATION/SHRINKAGE (e.g., RIDGE)
555 *=====
556 *ssc install ridgereg, replace // Alternative to ridge
557 forval i=1/5 {
558     gen daily_return_l`i' = L`i'.daily_return
559 }
560 ridgereg future_return signal daily_return_l1 daily_return_l2 daily_return_l3 daily_return_l4
561 ds
562 *=====
563 *   LONG-HORIZON RETURN PREDICTABILITY
564 *=====
565 // forval k = 3(3)12 {
566 //     gen cum_ret_`k'd = future_return
567 //     forval i = 2/`k' {
568 //         replace cum_ret_`k'd = cum_ret_`k'd + F`i'.future_return if !missing(F`i'.future_return)
569 //     }
570 //     label var cum_ret_`k'd "Cum Future Return (Days 1-`k')"
571 //
572 //     newey cum_ret_`k'd signal, lag(=`k'-1')
573 // }
574 // // Findings: Insignificant across horizons (p>0.6).
575 *=====
576 *   LONG-HORIZON RETURN PREDICTABILITY
577 *=====
578 drop if missing(future_return)
579 sort panel_id date_stata
580 capture drop time_idx
581 gen long time_idx = _n
582 tsset panel_id time_idx, daily
583 forval k = 3(3)12 {
584     gen cum_ret_`k'd = future_return if !missing(future_return)
585     forval i = 2/`k' {
586         replace cum_ret_`k'd = cum_ret_`k'd + F`i'.future_return if !missing(F`i'.future_return)
587     }

```

```

588     replace cum_ret_`k'd = . if missing(cum_ret_`k'd)
589     label var cum_ret_`k'd "Cum Future Return (Days 1-`k')"
590
591     preserve
592     drop if missing(cum_ret_`k'd, signal)
593     newey cum_ret_`k'd signal, lag(`k'-1')
594     restore
595 }
596 ds
597 *=====
598 * SUPERVISED DIMENSION REDUCTION (PLS)
599 *=====
600 // ssc install pls, replace
601 //
602 // pls future_return, xvars(signal daily_return abs_daily) components(2) adjacent
603 // ssc install plssem, replace
604 //
605 // // Define latent variable LV from predictors, predict future_return
606 // plssem (LV <- signal daily_return abs_daily) (future_return <- LV), modeB(LV) boot(100)
607 // estat loadings // Check loadings
608 // // Findings: If loadings >0.5 for signal, important; path coef LV->future_return insignifi
609 *=====
610 * BOOTSTRAP INFERENCE FOR OLS
611 *=====
612 bootstrap _b[signal], reps(1000) seed(12345): reg future_return signal
613 estat bootstrap, percentile bc normal
614 ds
615 *=====
616 * SPECTRAL ANALYSIS FOR CYCLES
617 *=====
618 pergram future_return
619 // Findings: No clear cycles correlating with signal.
620 ds
621 *=====
622 * CROSS-VALIDATION FOR LASSO
623 *=====
624 cvlasso future_return signal daily_return abs_daily, nfolds(10) plotcv
625 ds
626 *=====
627 * SIGNAL-BASED BACKTESTING (SHARPE RATIO)
628 *=====
629 gen position = sign(signal)
630 gen port_ret = position * future_return
631 summ port_ret
632 local mean_ret = r(mean)
633 local sd_ret = r(sd)
634 local sharpe = (`mean_ret' / `sd_ret') * sqrt(252)
635 display "Sharpe Ratio: `sharpe'"
636 // Findings: Sharpe negative (-0.45), no value.
637 ds
638 *=====
639 * CORRELATION HEATMAP WITH LAGS

```

```

640 *=====
641 forval i=1/5 {
642     gen signal_l`i' = L`i'.signal
643 }
644 pwcorr signal signal_l* future_return
645 graph matrix signal signal_l* future_return
646 // Findings: Low corrs (<0.1), no lagged relations.
647 ds
648 save "13OCT2025_V2.dta", replace
649 *=====
650 * ROLLING WINDOW REGRESSION FOR STABILITY
651 *=====
652 rolling beta = _b[signal] se = _se[signal], window(500) saving(rolling_results, replace): re
653 preserve
654 use rolling_results, clear
655 gen upper = beta + 1.96 * se
656 gen lower = beta - 1.96 * se
657 gen p_value = 2 * normal(-abs(beta / se))
658 twoway (line beta end) (rcap upper lower end), yline(0)
659 ds
660 restore
661 ds
662 *=====
663 * INSTRUMENTAL VARIABLES (IV) REGRESSION
664 *=====
665 // // Ensure tsset and regenerate future_return if not found (for re-runs)
666 // // Sort by time_idx instead of date_stata (which may be dropped or not needed)
667 // sort panel_id time_idx
668 // capture drop time_idx
669 // gen long time_idx = _n
670 // tsset panel_id time_idx, daily
671 // capture confirm variable future_return
672 // if _rc != 0 {
673 //     gen double future_return = ln(F.close / close) if !missing(close, F.close)
674 //     label variable future_return "Future Log Return"
675 // }
676 //
677 // drop if missing(future_return, signal)
678 // ivregress 2sls future_return (signal = L2.signal), robust
679 // estat firststage
680 // // Findings: Weak IV (F<10), signal insignificant.
681 export excel using "mydata.xlsx", firstrow(variables) replace
682 log close

```
