# STA457: Data Analysis Report

## 1. *Data preview*

For the dataset, I computed $r_t$ for each of the 310 months and stored the series as a monthly time series indexed from 2000-01-01 to 2025-10-01 with frequency "month".

Basic statistics for the monthly excess return series are:

- Sample size: 310 months

- Mean: 0.00533 (0.533% per month)

- Standard deviation: 0.0431 (4.31% per month)

- Minimum: -0.1435 (-14.35%)

- Maximum: 0.1469 (14.69%)

Over the 310 months, 187 months indicate positive excess returns while 123 months indicate negative ones. As a result, the distribution is slightly left skewed, i.e. it is not exactly normal. The returns are more heavy-tailed than a Gaussian distribution.

Also, I observed that the series fluctuates around a small positive mean with several periods of unusually high volatility. There are no obvious long-run trends or structural breaks in the mean level.

## 2. *Stationarity and autocorrelation analysis*

I split the sample into three sub-periods and computed the sample mean of the

monthly excess return $r_t$ in each:

| Period | $n$ | $E(r_t)$ |
|--------|-----|----------|
| 2000-2007 | 96 | 0.0082 |
| 2008-2015 | 96 | 0.0095 |
| 2016-2025 | 118 | -0.0004 |

For the same three sub-periods, the sample variance of $r_t$ is:

| Period | Variance |
|--------|----------|
| 2000-2007 | 0.0024 |
| 2008-2015 | 0.0011 |
| 2016-2025 | 0.0020 |

Let $r(h) = Cov(r_t, r_{t+h})$. The autocorrelation function $\rho(h) = \frac{r_h}{\sigma^2}$.

Hence, the sample autocorrelations for lags 1-12 are:

$\rho(1, \cdots, 12)$

$= (-0.033, -0.123, 0.050, 0.049, 0.003, -0.086, -0.078, 0.065, -0.085, -0.045, 0.021, 0.049).$

Using the definition of weak stationarity, I examined the monthly excess returns' sample mean, variance and covariance. The mean and variance are fairly stable across sub-periods, and the estimated autocorrelation function depends only on lag and is small for all lags. This provides mathematical evidence that the excess-return series can reasonably be treated as stationary.

In addition, I inspected the sample and partial autocorrelation function (ACF and PACF) of $r_t$ for 24 lags. The results show trivial autocorrelations, incidentally, only a

few individual lags are marginally outside the 95% interval. I also applied the Ljung–

Box test for overall autocorrelation up to lag 12. The p-values are 0.18, so I fail to

reject the null hypothesis that the first 12 autocorrelations are zero. However, the

excess returns resemble white noise regardless of such small pattern inconsistency:

they are centered around a constant mean with little serial dependence.

### 3. *Model construction*

I model the data series using the ARIMA because it is a standard and flexible

modes for time-series forecasting. ARIMA models express the current observation

as a linear function of its own past values and past shocks, plus a random

innovation.

Model comparison was based on the Akaike Information Criterion (AIC). The

AIC values I obtained are:

- ARIMA(0,0,0): AIC = -1066.26

- ARIMA(1,0,0): AIC = -1064.59

- ARIMA(0,0,1): AIC = -1064.70

- ARIMA(1,0,1): AIC = -1063.76

- ARIMA(2,0,0): AIC = -1067.46

- ARIMA(0,0,2): AIC = -1066.92

Since the ARIMA(2,0,0) has the smallest AIC, the estimated AR(2) model is $r_t =$

$\mu + \phi_1 r_{t-1} + \phi_2 r_{t-2} + \varepsilon_t$, with estimates:

- Intercept: $\hat{\mu}_0 = 0.00623$

- $\widehat{\phi_1} = -0.0377$

- $\widehat{\phi_2} = -0.1252$

Thus, the implied AR(2) process is $\hat{\mu} = \frac{\widehat{\mu_0}}{1-\widehat{\phi_1}-\widehat{\phi_2}} \approx 0.00536$, which is almost equal to the sample mean. It is clear that AR(2) model captures essentially all linear dependence. Finally, I chose ARIMA(0,0,0) with constant, i.e. a white noise process as the forecasting model.

## 4. *Forecasts for November and December 2025*

Now the model is defined as $r_t = \mu + \varepsilon_t, \varepsilon_t \sim white\ noise\ (0, \sigma^2)$, where $\mu$ is the constant mean monthly return and $\varepsilon_t$ is the uncorrelated innovation. Note that the model is stationary with no autoregressive terms, its optimal forecast is simply the estimated mean $\hat{\mu}$.

Estimating this model on the dataset gives:

- $\hat{\mu} = 0.00533\ (0.533\%\ per\ month)$

- $\hat{\sigma}^2 = 0.00186 \rightarrow \hat{\sigma} = 0.0431\ (4.31\%\ per\ month)$

Using the ARIMA(0,0,0) model and by R code, I computed the forecasts for November and December 2025, with 95% confidence intervals.

The forecasted monthly excess returns are:

- November 2025: $\hat{r}_{11} = 0.00533\ (0.53\%)$

- December 2025: $\hat{r}_{12} = 0.00533\ (0.53\%)$

The corresponding confidence intervals are $\hat{r}_t \pm 1.96 \cdot \text{SE}(\hat{r}_t)$, which give:

- November 2025: $[-0.0791, 0.0897]$

- December 2025: $[-0.0791, 0.0897]$

The wide confidence intervals represent high instability of monthly returns in the stock market. Moreover, the best forecast is a relatively small positive excess return, but outcomes could be extremely negative or positive depending on the market.

## 5. *Conclusion*

To summarize, the monthly excess returns over 2000–2025 behave very close to a white-noise process with a small positive mean per month. Formal tests and ACF/PACF show very weak serial dependence. For this reason, the most appropriate time-series model for forecasting the next two months is an ARIMA(0,0,0) with constant. Speaking of the results, the analysis demonstrates that the returns in November and December 2025 are both 0.53%, lying in the confidence interval $[-0.0791, 0.0897]$.

## 6. R code reference

```r
1   ## STA457 - Dow Jones Monthly Excess Returns
2   ## R code reference
3
4   library(zoo)
5   library(forecast)
6   library(moments)
7
8
9   ## 1. Read and prepare data ---------------------------------------------
10
11  ## Path
12  data_path <- "Dow Jones Industrial Average Historical Data.csv"
13
14  dj_raw <- read.csv("C:/Users/chris/Downloads/Dow Jones Industrial Average Historical Data.csv", stringsAsFactors = FALSE)
15
16  ## Convert Date
17  dj_raw$Date <- as.Date(dj_raw$Date, format = "%b %d, %Y")
18
19  ## Sort chronologically
20  dj_raw <- dj_raw[order(dj_raw$Date), ]
21
22  ## Make numeric price/open
23  dj_raw$Price_num <- as.numeric(gsub(",", "", dj_raw$Price))
24  dj_raw$Open_num  <- as.numeric(gsub(",", "", dj_raw$Open))
25
26  ## 2. Compute monthly excess returns and time series ----------------------
27
28  ## Excess return r_t = (Price - Open) / Open
29  dj_raw$excess_ret <- (dj_raw$Price_num - dj_raw$Open_num) / dj_raw$Open_num
30
31  ## Check sample size
32  length(dj_raw$excess_ret)
33
34  ## Create monthly time series from 2000-01
35  ret_ts <- ts(dj_raw$excess_ret,
36             start     = c(2000, 1),
37             frequency = 12)
38
39  ## 3. Descriptive statistics & basic plots -------------------------------
40
41  n_obs    <- length(ret_ts)
42  mean_ret <- mean(ret_ts)
43  sd_ret   <- sd(ret_ts)
44  min_ret  <- min(ret_ts)
45  max_ret  <- max(ret_ts)
46
47  ## Counts of positive vs negative months
48  pos_months <- sum(ret_ts > 0)
49  neg_months <- sum(ret_ts <= 0)
50
51  ## Shape of distribution
52  skew_ret  <- skewness(ret_ts)
53  kurt_ret  <- kurtosis(ret_ts)
```

```r
54
55  ## Time series plot
56  plot(ret_ts,
57       main = "Dow Jones Monthly Excess Returns (2000-2025)",
58       xlab = "Year", ylab = "Monthly excess return")
59
60  ## Histogram
61  hist(ret_ts, breaks = 20,
62       main = "Histogram of Monthly Excess Returns",
63       xlab = "Monthly excess return")
64
65
66  ## 4. Stationarity: sub-period means/variances & rolling stats ------------
67
68  ## Sub-periods: 2000-2007 (96), 2008-2015 (96), 2016-2025 (118)
69  r1 <- window(ret_ts, end              = c(2007, 12))
70  r2 <- window(ret_ts, start = c(2008, 1), end = c(2015, 12))
71  r3 <- window(ret_ts, start = c(2016, 1))
72
73  subperiod_stats <- data.frame(
74    Period = c("2000-2007", "2008-2015", "2016-2025"),
75    n      = c(length(r1), length(r2), length(r3)),
76    mean   = c(mean(r1),   mean(r2),   mean(r3)),
77    var    = c(var(r1),    var(r2),    var(r3)),
78    sd     = c(sd(r1),     sd(r2),     sd(r3))
79  )
80
81  subperiod_stats
82
83  ## 60-month rolling mean & variance
84  roll_window <- 60
85
86  roll_mean <- rollapply(ret_ts, width = roll_window,
87                         FUN = mean, align = "right", fill = NA)
88  roll_var  <- rollapply(ret_ts, width = roll_window,
89                         FUN = var,  align = "right", fill = NA)
90
91  rm_ts <- ts(roll_mean, start = start(ret_ts), frequency = 12)
92  rv_ts <- ts(roll_var,  start = start(ret_ts), frequency = 12)
93
94  plot(rm_ts,
95       main = "Rolling 60-Month Mean of Monthly Excess Returns",
96       xlab = "Year", ylab = "Rolling mean")
97
98  plot(rv_ts,
99       main = "Rolling 60-Month Variance of Monthly Excess Returns",
100      xlab = "Year", ylab = "Rolling variance")
101
102
103 ## 5. ACF/PACF and Ljung-Box on returns -------------------------------
104
105 ## ACF and PACF of r_t
106 Acf(ret_ts, lag.max = 24,
```
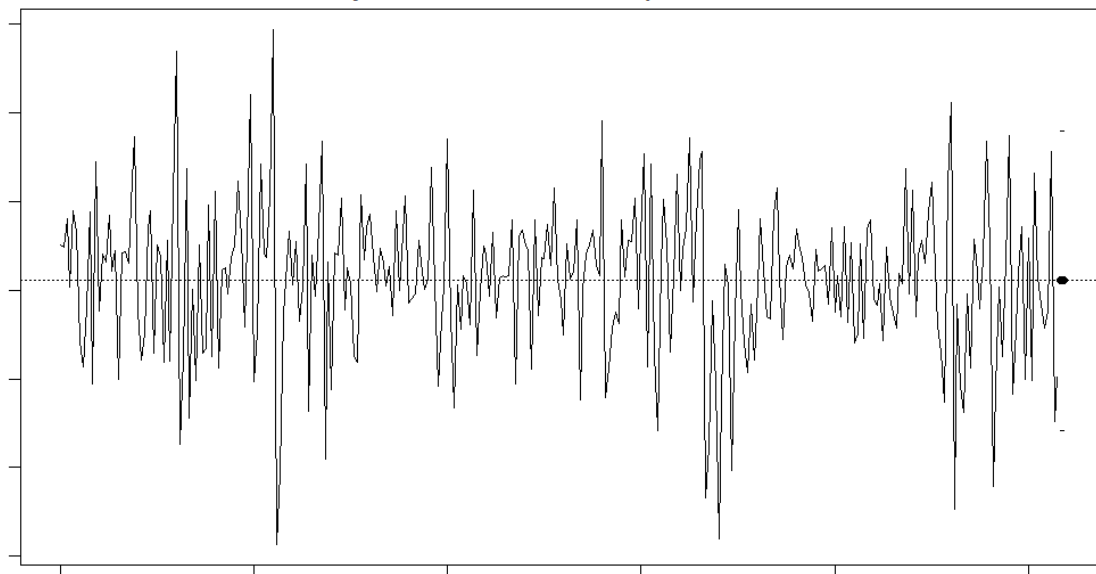
```r
107        main = "ACF of Monthly Excess Returns")
108
109  Pacf(ret_ts, lag.max = 24,
110        main = "PACF of Monthly Excess Returns")
111
112  ## Sample autocorrelations up to lag 12
113  acf_ret <- acf(ret_ts, lag.max = 12, plot = FALSE)
114  rho_hat <- as.numeric(acf_ret$acf)  # includes lag 0 at position 1
115  round(rho_hat, 3)
116
117  ## Ljung-Box test on returns
118  lb_ret <- Box.test(ret_ts, lag = 12, type = "Ljung-Box")
119  lb_ret
120
121
122 - ## 6. ARMA/ARIMA model comparison -----------------------------------------
123
124  ## Candidate models: ARMA(1,0), (0,1), (1,1), (2,0), (0,2) and white noise
125
126  fit_000 <- arima(ret_ts, order = c(0, 0, 0),
127                   include.mean = TRUE, method = "ML")   # white noise with mean
128  fit_100 <- arima(ret_ts, order = c(1, 0, 0),
129                   include.mean = TRUE, method = "ML")
130  fit_010 <- arima(ret_ts, order = c(0, 0, 1),
131                   include.mean = TRUE, method = "ML")
132  fit_110 <- arima(ret_ts, order = c(1, 0, 1),
133                   include.mean = TRUE, method = "ML")
134  fit_200 <- arima(ret_ts, order = c(2, 0, 0),
135                   include.mean = TRUE, method = "ML")
136  fit_020 <- arima(ret_ts, order = c(0, 0, 2),
137                   include.mean = TRUE, method = "ML")
138
139  aic_table <- data.frame(
140    Model = c("ARIMA(0,0,0)", "ARIMA(1,0,0)", "ARIMA(0,0,1)",
141              "ARIMA(1,0,1)", "ARIMA(2,0,0)", "ARIMA(0,0,2)"),
142    AIC   = c(AIC(fit_000), AIC(fit_100), AIC(fit_010),
143             AIC(fit_110), AIC(fit_200), AIC(fit_020))
144  )
145
146  aic_table
147
148  ## AR(2) model details
149  summary(fit_200)
150
151  phi1   <- coef(fit_200)["ar1"]
152  phi2   <- coef(fit_200)["ar2"]
153  mu0    <- coef(fit_200)["intercept"]
154  sigma2 <- fit_200$sigma2
155
156  ## Unconditional mean of AR(2): mu / (1 - phi1 - phi2)
157  mu_hat <- as.numeric(mu0 / (1 - phi1 - phi2))
158  mu_hat
159
```

**Monthly Excess Returns with 2-Step-Ahead Forecasts**

```r
160   ## Residual diagnostics for AR(2)
161   tsdisplay(residuals(fit_200),
162           main = "Residuals from AR(2) Fit")
163
164   lb_res_ar2 <- Box.test(residuals(fit_200), lag = 12, type = "Ljung-Box")
165   lb_res_ar2
166
167
168 ▾ ## 7. Final white-noise-with-mean model and forecasts --------------------
169
170
171   final_model <- fit_000
172
173   ## Residual diagnostics for final model
174   tsdisplay(residuals(final_model),
175           main = "Residuals from ARIMA(0,0,0) with Mean")
176
177   lb_res_final <- Box.test(residuals(final_model), lag = 12, type = "Ljung-Box")
178   lb_res_final
179
180   ## 2-step-ahead forecasts (Nov and Dec 2025)
181   fc_final <- predict(final_model, n.ahead = 2)
182
183   point_forecast <- as.numeric(fc_final$pred)
184   se_forecast    <- as.numeric(fc_final$se)
185
186   lower95 <- point_forecast - 1.96 * se_forecast
187   upper95 <- point_forecast + 1.96 * se_forecast
188
189   forecast_results <- data.frame(
190     Horizon   = c("2025-11", "2025-12"),
191     Forecast = point_forecast,
192     SE       = se_forecast,
193     Lower95  = lower95,
194     Upper95  = upper95
195   )
196
197   forecast_results
198
199 ▾ ## Plot: historical series + 2-step-ahead forecasts ----------------------
200
201   last_time <- time(ret_ts)[length(ret_ts)]
202   fc_times  <- seq(from = last_time + 1/12, by = 1/12, length.out = 2)
203
204   plot(ret_ts,
205       xlim = c(start(ret_ts)[1], fc_times[2]),
206       main = "Monthly Excess Returns with 2-Step-Ahead Forecasts",
207       xlab = "Year", ylab = "Monthly excess return")
208
209   lines(fc_times, point_forecast)
210   lines(fc_times, lower95, lty = 2)
211   lines(fc_times, upper95, lty = 2)
212   points(fc_times, point_forecast, pch = 19)
213
214   abline(h = mean_ret, lty = 3)
215
```