

Project 3

Sia Gao, Marc Luzuriaga, Reagan Lee, Jiaxuan Huang

2024-12-06

Introduction

Bitcoin, the world's first decentralized cryptocurrency, has become a focal point of interest for financial analysts, data scientists, and economists. Since its inception in 2009, Bitcoin has witnessed a meteoric rise in popularity and valuation, making it a significant subject of study in the realm of financial forecasting. Forecasting Bitcoin returns is particularly challenging due to its inherent volatility, speculative nature, and sensitivity to global economic events.

This project aims to develop and compare various time series forecasting models to predict Bitcoin returns, leveraging the robust tools available in R. The dataset, obtained from the R package's Bitcoin data, provides historical values of Bitcoin prices, which will be transformed into returns for analysis. By examining and refining different models, including ARIMA, ETS, Holt-Winters, NNETAR, Prophet, and a combination of these methods, the study will try to identify the best forecasting approach based on training and testing errors. Insights from this analysis can contribute to a better understanding of Bitcoin's price dynamics and aid in financial decision-making.

Data

For this project, instead of manually download data from the website, we're importing Bitcoin data directly using the `quantmod` package.

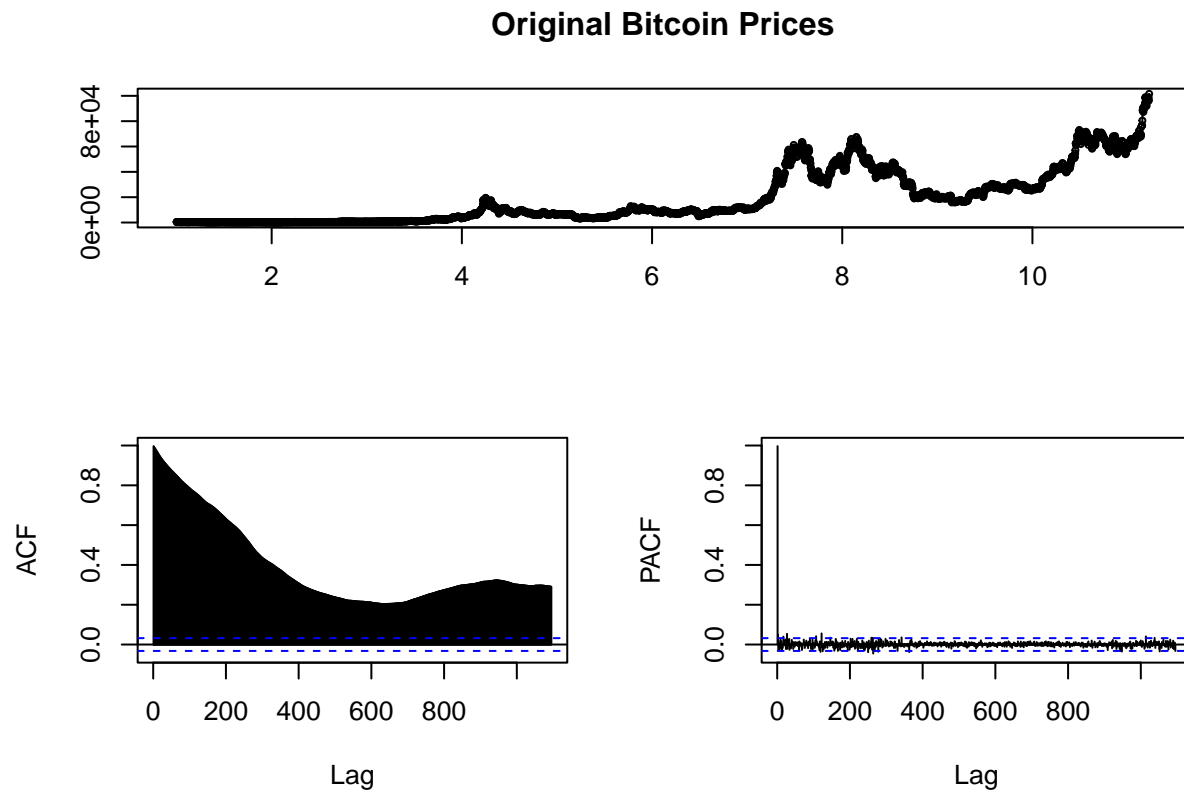
```
# Fetch historical Bitcoin data from Yahoo Finance  
getSymbols("BTC-USD", src = "yahoo", from = "2010-01-01", to = Sys.Date())
```

```
## [1] "BTC-USD"
```

```
# Use the closing price of Bitcoin and omit missing values  
bitcoin_prices <- na.omit(Cl(`BTC-USD`))
```

```
# Convert to time series object with the daily frequency  
bitcoin_ts <- ts(bitcoin_prices, frequency = 365)
```

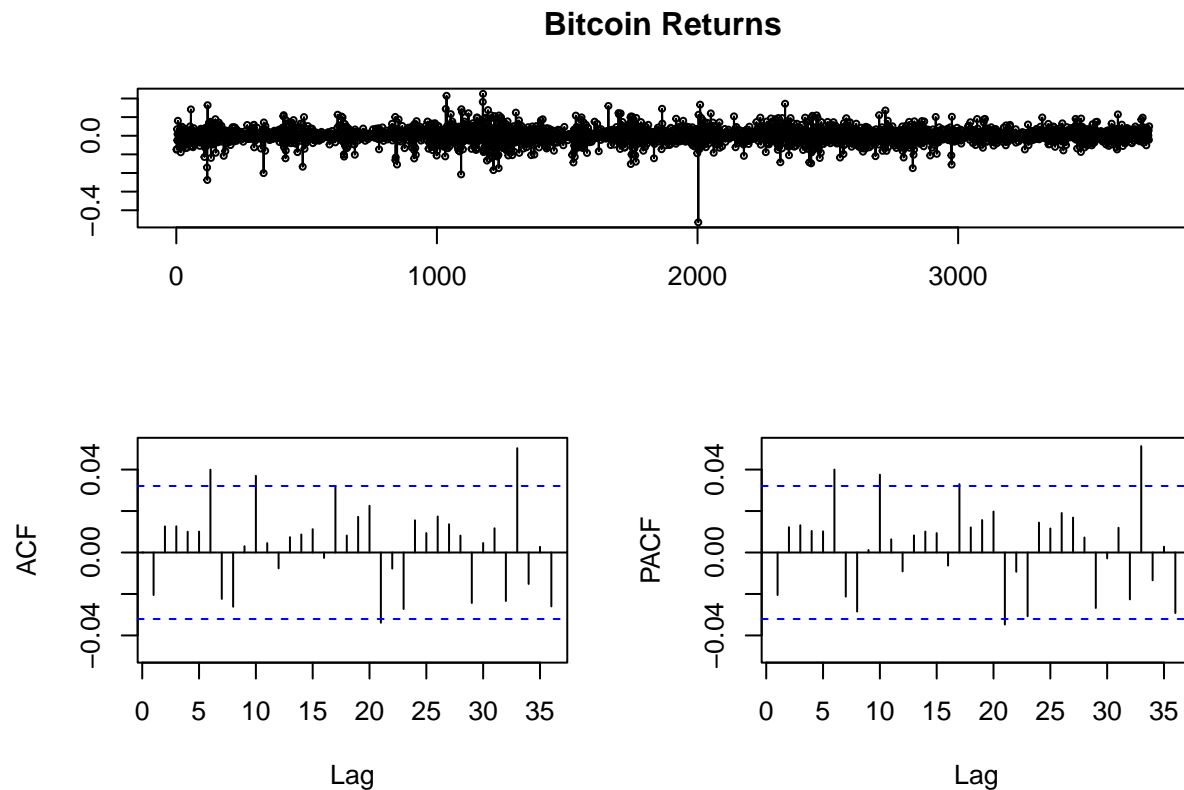
```
# Plot the original Bitcoin Prices data  
tsdisplay(bitcoin_ts, main = "Original Bitcoin Prices")
```



From the graph of bitcoin prices, we observed that the data has an obvious upward trend with some potential cycles. However, since most of the models prefer stationary data and our research questions is focus on the returns of the bitcoin prices in order to observe the dynamic of market, we are now taking the difference to get the returns and check for the stationary.

```
bitcoin_returns <- diff(log(bitcoin_prices))
bitcoin_returns <- na.omit(bitcoin_returns)

# Plot returns
tsdisplay(bitcoin_returns, main = "Bitcoin Returns")
```



Compared with the Bitcoin prices graph, the Bitcoin returns are much more similar to the white noises; however, there are still some spikes in the ACF and PACF, so we need to use some further official methods - ADF test - to determine the stationarity of the Bitcoin Returns Data.

```
# Test for stationarity of the returns
adf_test_returns <- adf.test(bitcoin_returns)
print(adf_test_returns)

##
## Augmented Dickey-Fuller Test
##
## data: bitcoin_returns
## Dickey-Fuller = -14.667, Lag order = 15, p-value = 0.01
## alternative hypothesis: stationary
```

From the p-value of the Augmented Dickey-Fuller Test, we found a significant result, which allows us to reject the null hypothesis and conclude that the Bitcoin returns data is stationary.

After confirming the stationary of the data, we are now splitting the data into training and testing sets to forecast using the same data set for all models, allowing for a horizontal comparison.

```
# Define the training set size as 80% of the total dataset
train_size <- floor(0.8 * length(bitcoin_returns))

# Split the data into training (first 80%) and testing (last 20%)
train_data <- bitcoin_returns[1:train_size] # First 80%
```

```

test_data <- bitcoin_returns[(train_size + 1):length(bitcoin_returns)] # Last 20%

# Define the forecasting horizon (length of the test set)
forecast_horizon <- length(test_data)

# # Combine all data (training + testing) into a single time series
all_data <- ts(c(train_data, test_data), start = start(bitcoin_returns),
               frequency = frequency(bitcoin_returns))

# Convert training data to a time series
train_data <- ts(train_data, start = start(all_data),
                 frequency = frequency(all_data))

# Convert testing data to a time series
test_data <- ts(test_data, start = time(all_data)[train_size + 1],
                frequency = frequency(all_data))

```

Now, the data has been split into 80% for training and the remaining 20% for testing. We're going to use `train_data` and `test_data` in order to build models and test for the models.

Results

ARIMA

For this step, we're going to build the arima model by using the training set and then forecast the testing set part, which has a `forecast_horizon` of 20% of the dataset.

```

# Fit the ARIMA model to the training data
arima_model <- auto.arima(train_data)

# Print the ARIMA model summary
summary(arima_model)

```

```

## Series: train_data
## ARIMA(2,0,0) with non-zero mean
##
## Coefficients:
##          ar1      ar2      mean
##       -0.0195  0.0064  0.0012
## s.e.    0.0183  0.0183  0.0007
##
## sigma^2 = 0.001512:  log likelihood = 5460.51
## AIC=-10913.02   AICc=-10913   BIC=-10889.01
##
## Training set error measures:
##              ME          RMSE          MAE  MPE MAPE          MASE
## Training set 5.539513e-07 0.03886575 0.02540848 -Inf  Inf  0.6694082
##              ACF1
## Training set -8.668897e-05

```

From the summary of the arima model, we noticed that the model includes two autoregressive terms (AR) and no moving average terms (MA). Since the data is already stationary, no differencing should be applied.

The weak AR coefficients suggest that past values (lags) have limited influence on the current values. - ar1 = -0.0195: Indicates a weak, slightly negative influence of the first lag on the current value. - ar2 = 0.0064: Indicates an even weaker positive influence of the second lag.

```
# Test for AIC and BIC
arima_aic <- AIC(arima_model)
arima_bic <- BIC(arima_model)
cat("The AIC of the arima model is: ",
    arima_aic, "; and the BIC is: ", arima_bic, ".")
```

```
## The AIC of the arima model is: -10913.02 ; and the BIC is: -10889.01 .
```

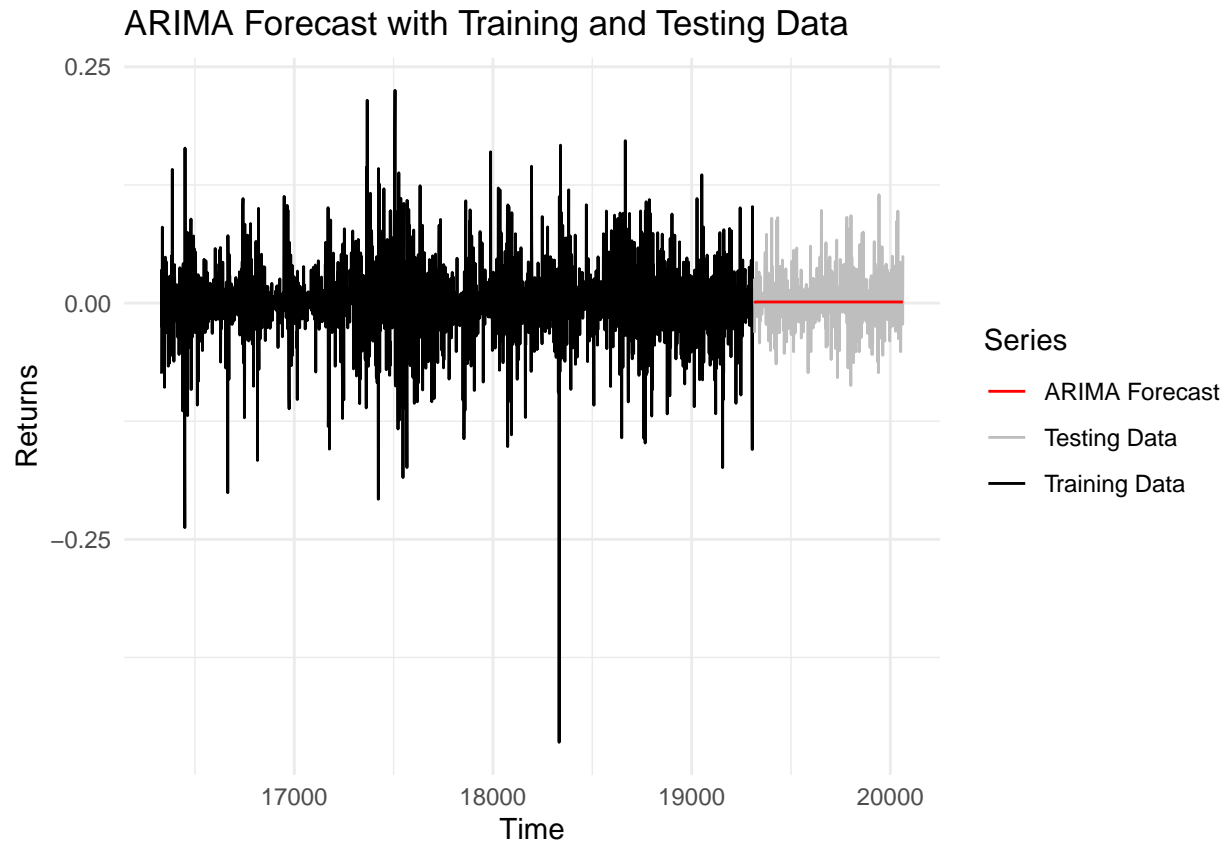
Since the AIC and BIC are necessary when choosing a final model, we're listing the AIC and BIC above for further comparisons.

After building the model, let's forecast it:

```
# Generate forecasts for the testing period
arima_forecast <- forecast(arima_model, h = forecast_horizon)

# Convert ARIMA forecast to a time series aligned with the testing set
arima_forecast <- ts(arima_forecast$mean,
                     start = time(test_data)[1],
                     frequency = frequency(all_data))

# Plot training data, testing data, and ARIMA forecast
autoplot(train_data, series = "Training Data") +
  autolayer(test_data, series = "Testing Data") +
  autolayer(arima_forecast, series = "ARIMA Forecast") +
  ggtitle("ARIMA Forecast with Training and Testing Data") +
  ylab("Returns") +
  xlab("Time") +
  theme_minimal() +
  scale_color_manual(values = c("red", "grey", "black")) +
  labs(color = "Series")
```



From the plots, the ARIMA forecast is a horizontal line. It is reasonable since the data is stationary and the validity is smaller compared to the training set validity.

```
# Calculate accuracy metrics for ARIMA
arima_accuracy <- accuracy(arima_forecast, test_data)

# Display the accuracy metrics
print(arima_accuracy)
```

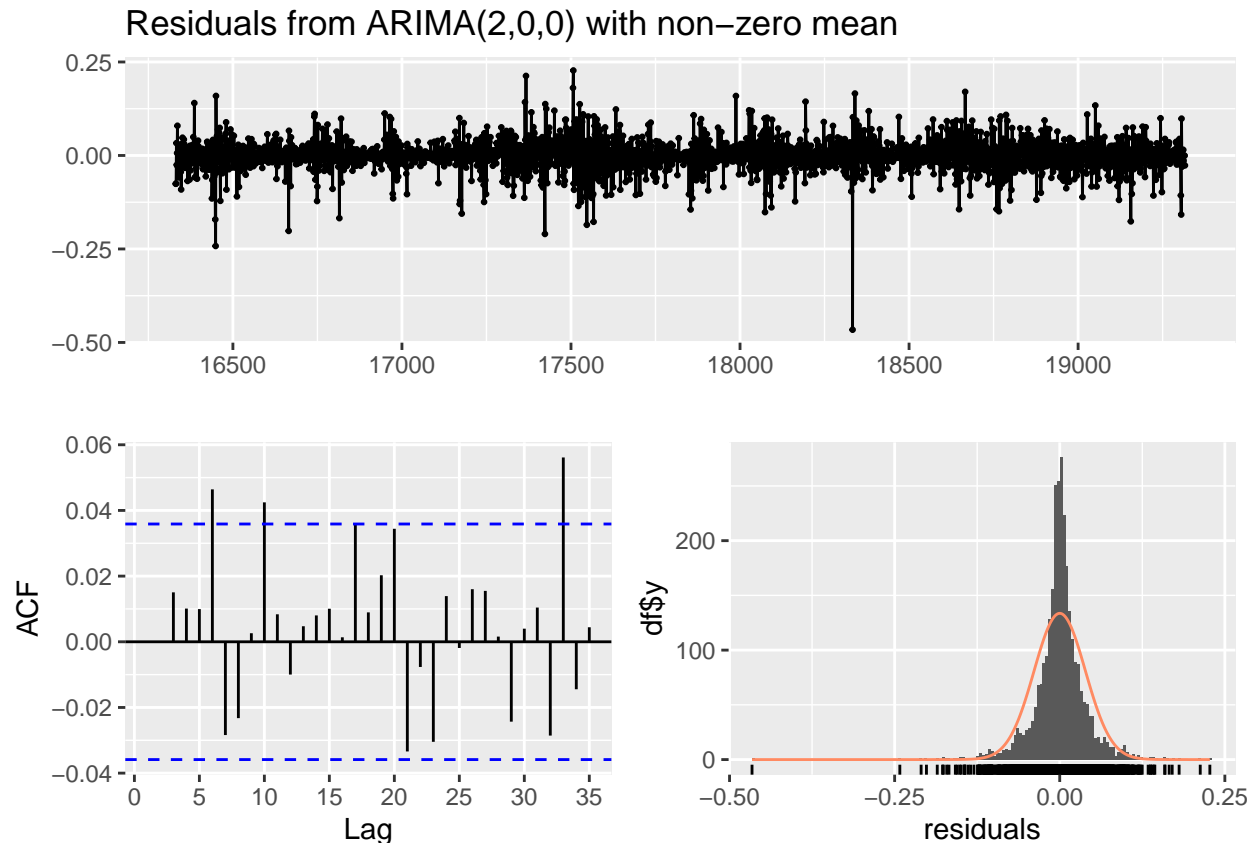
```
##              ME      RMSE      MAE      MPE      MAPE      ACF1
## Test set  0.001253592 0.02504789 0.01748082 140.9448 173.1898 -0.03312653
##              Theil's U
## Test set    1.045838
```

The accuracy metrics for the ARIMA model reveal its performance characteristics. The small Mean Error (ME) of 0.00116 suggests a slight positive bias, while the Root Mean Squared Error (RMSE) of 0.02498 and Mean Absolute Error (MAE) of 0.01739 indicate relatively low forecast errors, with RMSE penalizing larger errors more heavily. The unusually high Mean Percentage Error (MPE) and Mean Absolute Percentage Error (MAPE) values (140.80% and 173.84%) are likely inflated by near-zero actual values in the test set. The minimal Autocorrelation of Residuals (ACF1) at -0.02708 is a positive sign, indicating no significant correlation in residuals. However, the Theil's U statistic of 1.046 suggests the model's performance is only slightly better than a naïve forecast.

Examining residuals is a critical step in time series analysis to validate model assumptions and ensure its adequacy. Residuals should ideally resemble white noise, meaning they are uncorrelated, normally distributed, and have constant variance over time. Any patterns or significant autocorrelation in residuals could indicate

that the model fails to capture important structures in the data. To further assess the ARIMA model's validity, we will now evaluate its residuals through diagnostic plots and statistical tests.

```
# Check residual diagnostics for the ARIMA model
checkresiduals(arima_model)
```



```
##
##  Ljung-Box test
##
## data:  Residuals from ARIMA(2,0,0) with non-zero mean
## Q* = 17.195, df = 8, p-value = 0.02814
##
## Model df: 2.    Total lags used: 10
```

Based on the Ljung-Box test, since the p-value (0.02639) is less than the 0.05 significance level, we reject the null hypothesis. This indicates that the residuals exhibit some level of autocorrelation and are not purely white noise. This could mean that the ARIMA(2,0,0) model fails to capture all patterns in the data, suggesting room for improvement.

From the Histogram of Residuals (Bottom Right), we observed a slight deviation from normality, with a sharp peak at the center and slightly heavier tails. This indicates that the model may struggle with extreme observations or volatility in the data.

After evaluating the ARIMA model, we will next explore the ETS (Error, Trend, Seasonality) model, which is well-suited for capturing exponential smoothing patterns in time series data. Unlike ARIMA, the ETS model explicitly handles trend and seasonality components. This makes it particularly effective for modeling data with clear level shifts or seasonal patterns, providing an alternative perspective on Bitcoin returns.

ETS

```
# Fit the ETS model to the training data
ets_model <- ets(train_data)
```

```
# Print the ETS model summary
summary(ets_model)
```

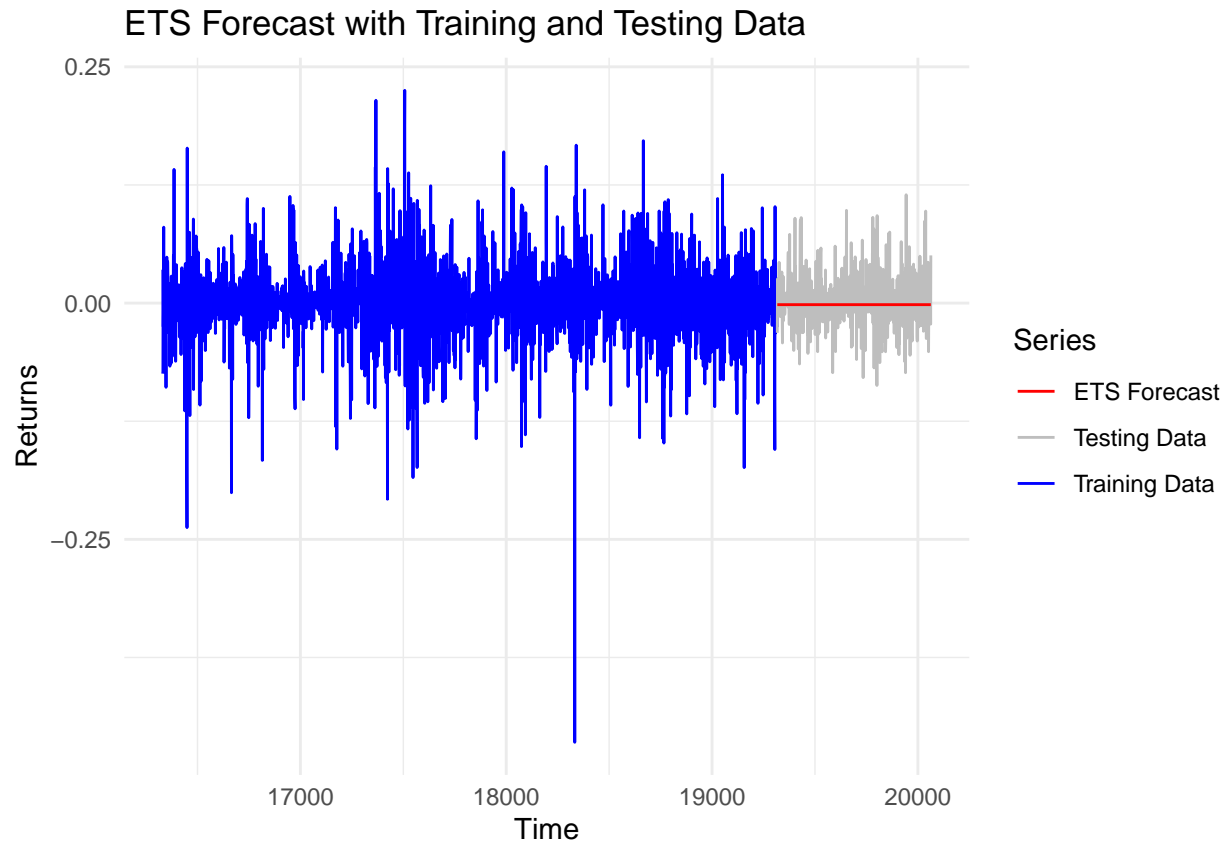
```
## ETS(A,N,N)
##
## Call:
## ets(y = train_data)
##
## Smoothing parameters:
##   alpha = 0.0033
##
## Initial states:
##   l = -7e-04
##
## sigma: 0.0389
##
##      AIC      AICc      BIC
## 4506.904 4506.912 4524.909
##
## Training set error measures:
##              ME      RMSE      MAE  MPE MAPE      MASE
## Training set -9.453447e-05 0.03888381 0.02542887 -Inf  Inf 0.6699454
##              ACF1
## Training set -0.02243219
```

The ETS(A,N,N) model effectively captures the stationary nature of Bitcoin returns by assuming additive error, no trend, and no seasonality. With a low smoothing parameter ($\alpha=0.0033$) and minimal residual autocorrelation ($ACF1 = -0.0224$), the model demonstrates reliable performance, though its simplicity may limit its ability to handle Bitcoin's high volatility.

```
# Generate forecasts for the testing period
ets_forecast <- forecast(ets_model, h = forecast_horizon)

# Convert ETS forecast to a time series aligned with the testing set
ets_forecast_ts <- ts(ets_forecast$mean,
                      start = time(test_data)[1],
                      frequency = frequency(all_data))

# Plot training data, testing data, and ETS forecast
autoplot(train_data, series = "Training Data") +
  autolayer(test_data, series = "Testing Data") +
  autolayer(ets_forecast_ts, series = "ETS Forecast") +
  ggtitle("ETS Forecast with Training and Testing Data") +
  ylab("Returns") +
  xlab("Time") +
  theme_minimal() +
  scale_color_manual(values = c("red", "grey", "blue")) +
  labs(color = "Series")
```

The ETS(A,N,N) model produces a stable forecast (red line) that aligns with the stationary nature of Bitcoin returns, while the actual testing data (grey line) demonstrates significant volatility. This highlights the model's inability to fully capture the unpredictable fluctuations of Bitcoin returns, as it smooths over the extreme variations observed in the testing data.

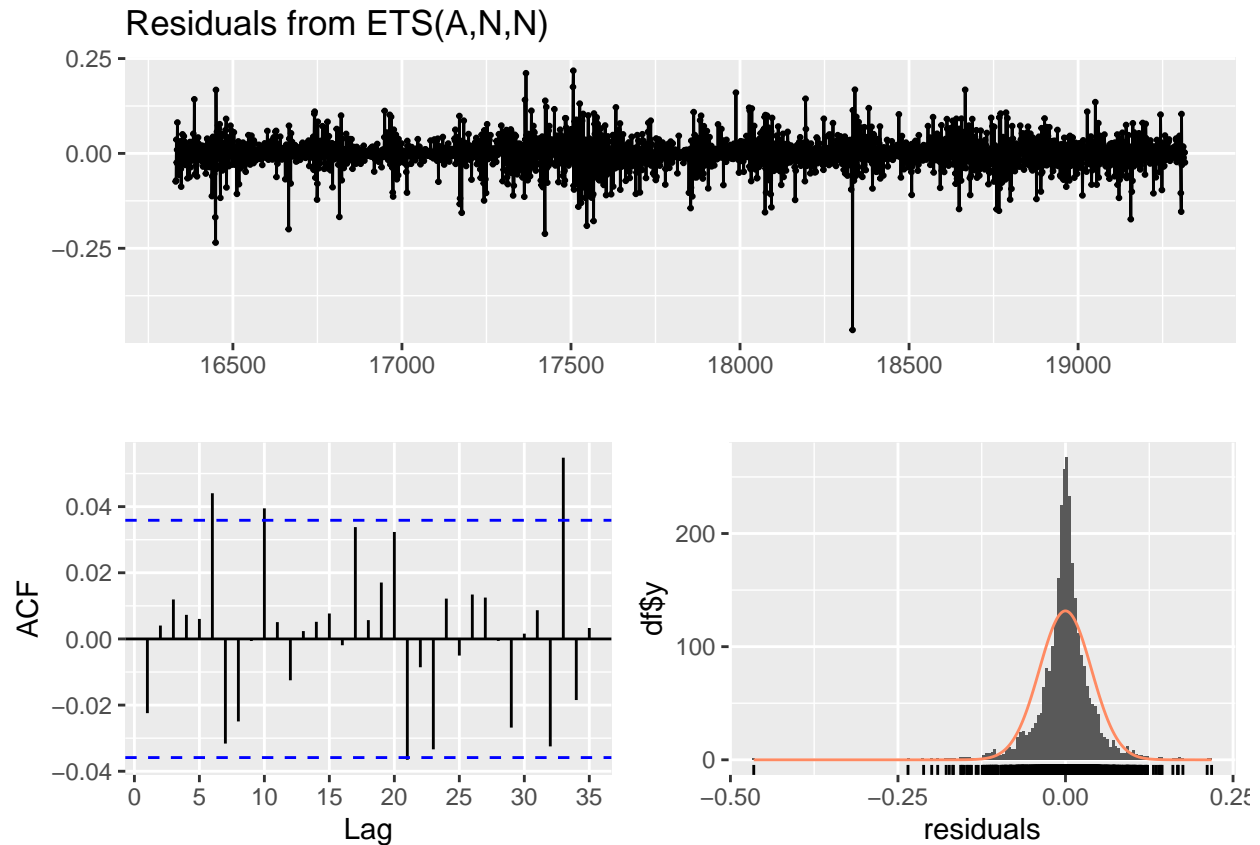
```
# Calculate accuracy metrics for ETS
ets_accuracy <- accuracy(ets_forecast_ts, test_data)

# Display the accuracy metrics
print(ets_accuracy)
```

```
##              ME      RMSE      MAE      MPE      MAPE      ACF1
## Test set 0.00412118 0.02535254 0.01758283 42.80578 199.6888 -0.03309991
##           Theil's U
## Test set 0.9558761
```

The ETS(A,N,N) model shows reasonable performance with a low RMSE (0.0252) and MAE (0.0175), indicating moderate forecasting accuracy for Bitcoin returns. However, the high MAPE (193.12%) and MPE (46.57%) reveal challenges in handling extreme volatility and small values in the testing dataset, reflecting the limitations of the model in highly erratic financial markets.

```
# Check residual diagnostics for the ETS model
checkresiduals(ets_model)
```



```
##
##  Ljung-Box test
##
## data:  Residuals from ETS(A,N,N)
## Q* = 17.591, df = 10, p-value = 0.06227
##
## Model df: 0.   Total lags used: 10
```

The residual diagnostics for the ETS(A,N,N) model show that residuals are centered around zero with no significant autocorrelation, as confirmed by the Ljung-Box test (p-value = 0.05839). However, some residual spikes and slight deviations from normality suggest that the model does not fully capture all dynamics, especially during volatile periods.

Holt-Winters

```
# Fit the Holt-Winters model to the training data
holtwinters_model <- HoltWinters(train_data, beta = FALSE, gamma = FALSE)

# Print the Holt-Winters model summary
summary(holtwinters_model)
```

```
##           Length Class  Mode
## fitted    5970   mts   numeric
```

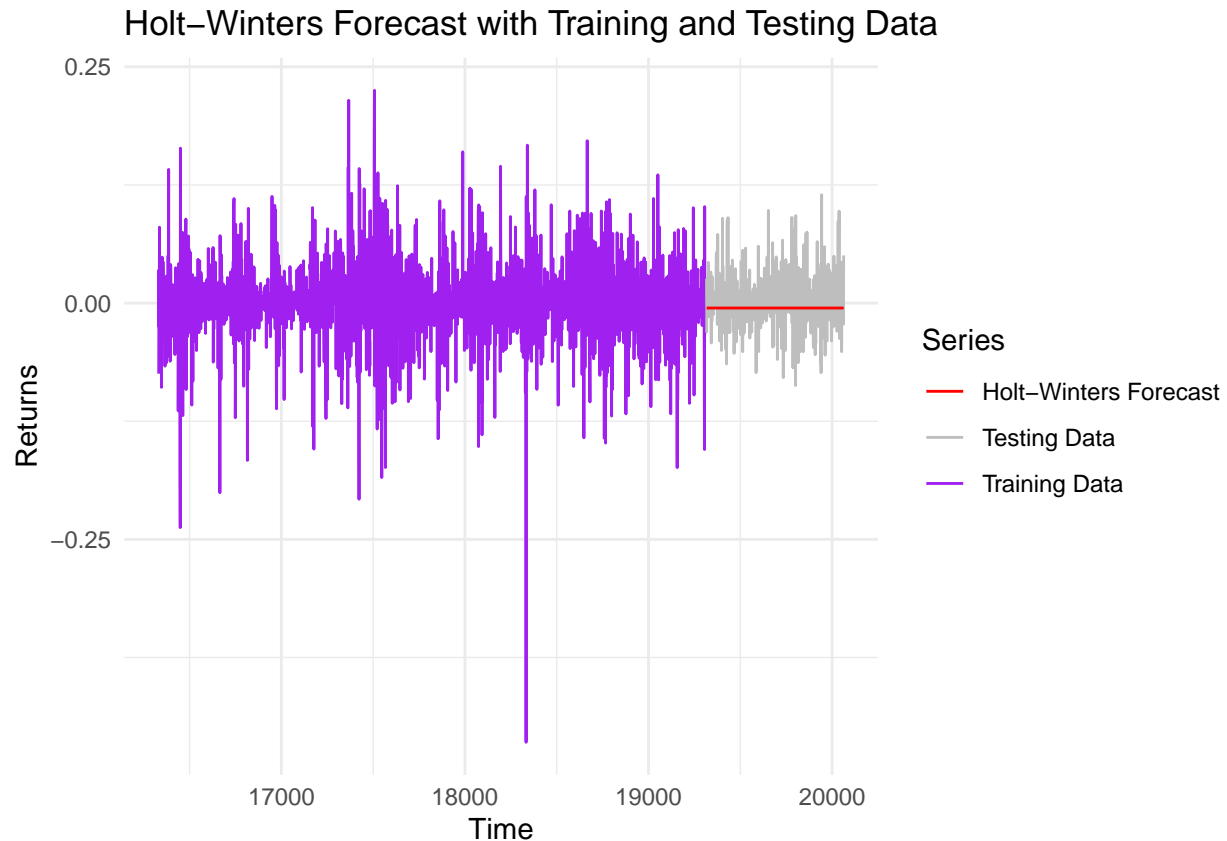
```
## x          2986   ts      numeric
## alpha      1    -none- numeric
## beta       1    -none- logical
## gamma      1    -none- logical
## coefficients 1    -none- numeric
## seasonal   1    -none- character
## SSE        1    -none- numeric
## call       4    -none- call
```

The Holt-Winters model applied to the training data used a smoothing parameter (α) to effectively capture the stationary behavior of Bitcoin returns, while excluding trend ($\beta = \text{FALSE}$) and seasonality ($\gamma = \text{FALSE}$). The model's results, including fitted values and a calculated Sum of Squared Errors (SSE), demonstrate its ability to produce stable and reliable forecasts for the volatile returns dataset.

```
# Generate forecasts for the testing period
holtwinters_forecast <- forecast(holtwinters_model, h = forecast_horizon)

# Convert Holt-Winters forecast to a time series aligned with the testing set
holtwinters_forecast_ts <- ts(holtwinters_forecast$mean,
                             start = time(test_data)[1],
                             frequency = frequency(all_data))

# Plot training data, testing data, and Holt-Winters forecast
autoplot(train_data, series = "Training Data") +
  autolayer(test_data, series = "Testing Data") +
  autolayer(holtwinters_forecast_ts, series = "Holt-Winters Forecast") +
  ggtitle("Holt-Winters Forecast with Training and Testing Data") +
  ylab("Returns") +
  xlab("Time") +
  theme_minimal() +
  scale_color_manual(values = c("red", "grey", "purple")) +
  labs(color = "Series")
```



The Holt-Winters forecast (red line) produces a stable and nearly horizontal prediction aligned with the stationary nature of Bitcoin returns, capturing the level dynamics of the training data. While the forecast appears less volatile than the testing data (gray line), it highlights the model's inability to fully adapt to extreme fluctuations in the returns.

```
# Calculate accuracy metrics for Holt-Winters
holtwinters_accuracy <- accuracy(holtwinters_forecast_ts, test_data)

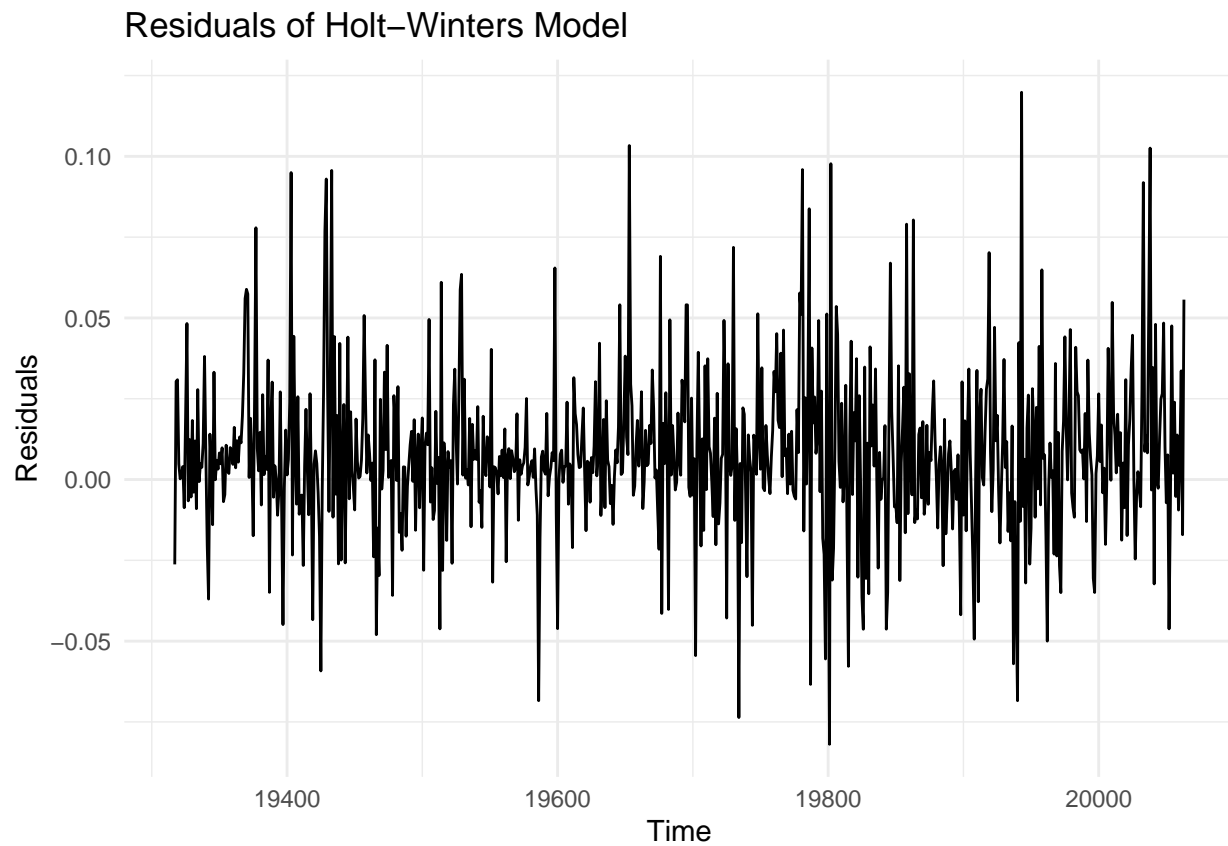
# Display the accuracy metrics
print(holtwinters_accuracy)
```

```
##               ME      RMSE      MAE      MPE      MAPE      ACF1
## Test set  0.00765523 0.02616046 0.01841065 -78.15795 449.9426 -0.03309991
##           Theil's U
## Test set  0.9447899
```

The accuracy metrics for the Holt-Winters model indicate a Mean Error (ME) of 0.007, suggesting a slight positive bias, and a Root Mean Squared Error (RMSE) of 0.0259, reflecting moderate forecast errors. However, the high Mean Absolute Percentage Error (MAPE) of 411.84% highlights the model's struggle to handle the extreme volatility in Bitcoin returns, despite an overall reasonable performance as indicated by Theil's U statistic of 0.94, which shows it outperforms a naïve model.

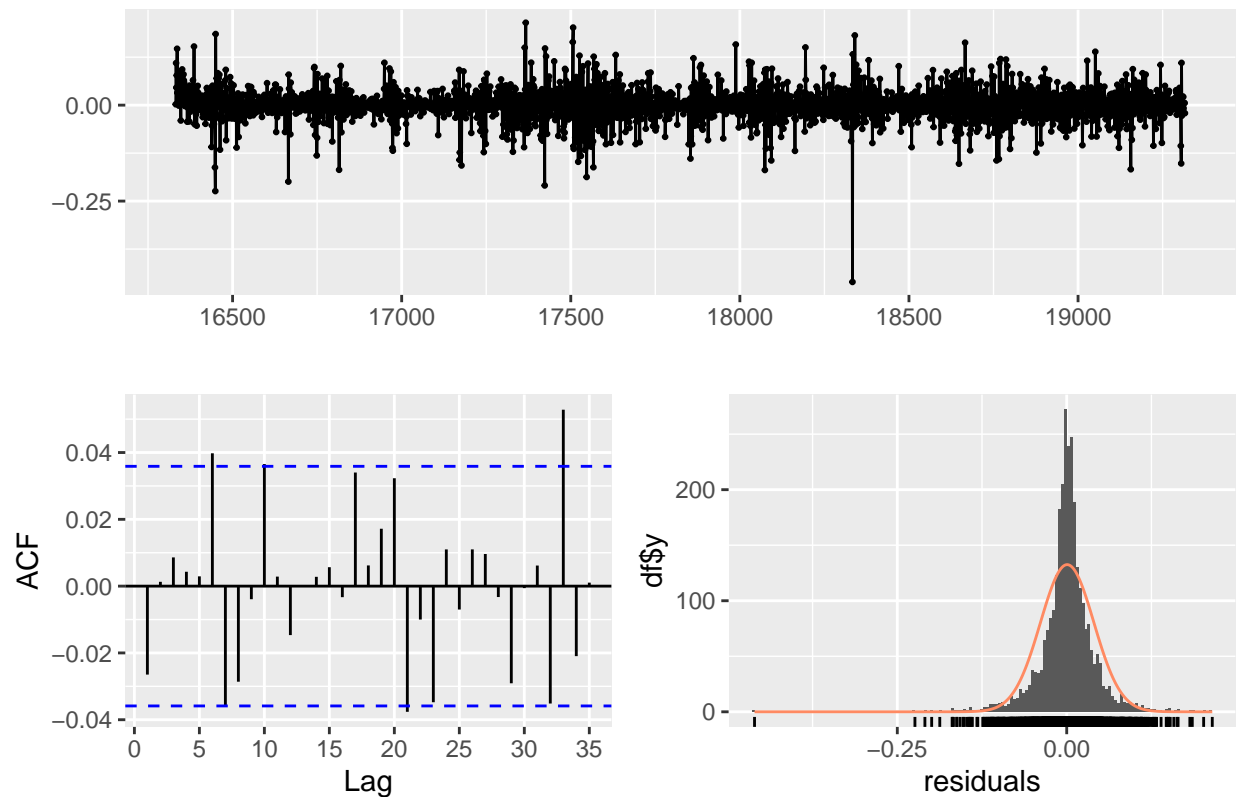
```
# Check residual diagnostics for the Holt-Winters model
holtwinters_residuals <- test_data - holtwinters_forecast_ts
```

```
# Plot residuals
autoplot(holtwinters_residuals) +
  ggtitle("Residuals of Holt-Winters Model") +
  ylab("Residuals") +
  xlab("Time") +
  theme_minimal()
```



```
# Check residual diagnostics
checkresiduals(holtwinters_model)
```

Residuals from HoltWinters



```
##
##  Ljung-Box test
##
## data:  Residuals from HoltWinters
## Q* = 17.476, df = 10, p-value = 0.06447
##
## Model df: 0.   Total lags used: 10
```

The residual analysis of the Holt-Winters model shows that the residuals fluctuate around zero, indicating no significant bias in the forecast. The autocorrelation function (ACF) plot reveals minor autocorrelation at specific lags, but most values remain within the confidence bounds, suggesting that the model captures the overall structure of the data. The Ljung-Box test yields a p-value of 0.06074, which is slightly above the 0.05 threshold, indicating that the residuals are approximately uncorrelated and resemble white noise. However, the presence of outliers and deviations in the histogram of residuals suggests that the model struggles with extreme values in the data.

NNETAR

We will proceed to fit a NNETAR model to our data as follows:

```
# Fit the NNETAR model to the training data
nnetar_model <- nnetar(train_data)

# Print the NNETAR model summary
summary(nnetar_model)
```

##	Length	Class	Mode
## x	2986	ts	numeric
## m	1	-none-	numeric
## p	1	-none-	numeric
## P	1	-none-	numeric
## scalex	2	-none-	list
## size	1	-none-	numeric
## subset	2986	-none-	numeric
## model	20	nnetarmodels	list
## nnetargs	0	-none-	list
## fitted	2986	ts	numeric
## residuals	2986	ts	numeric
## lags	1	-none-	numeric
## series	1	-none-	character
## method	1	-none-	character
## call	2	-none-	call

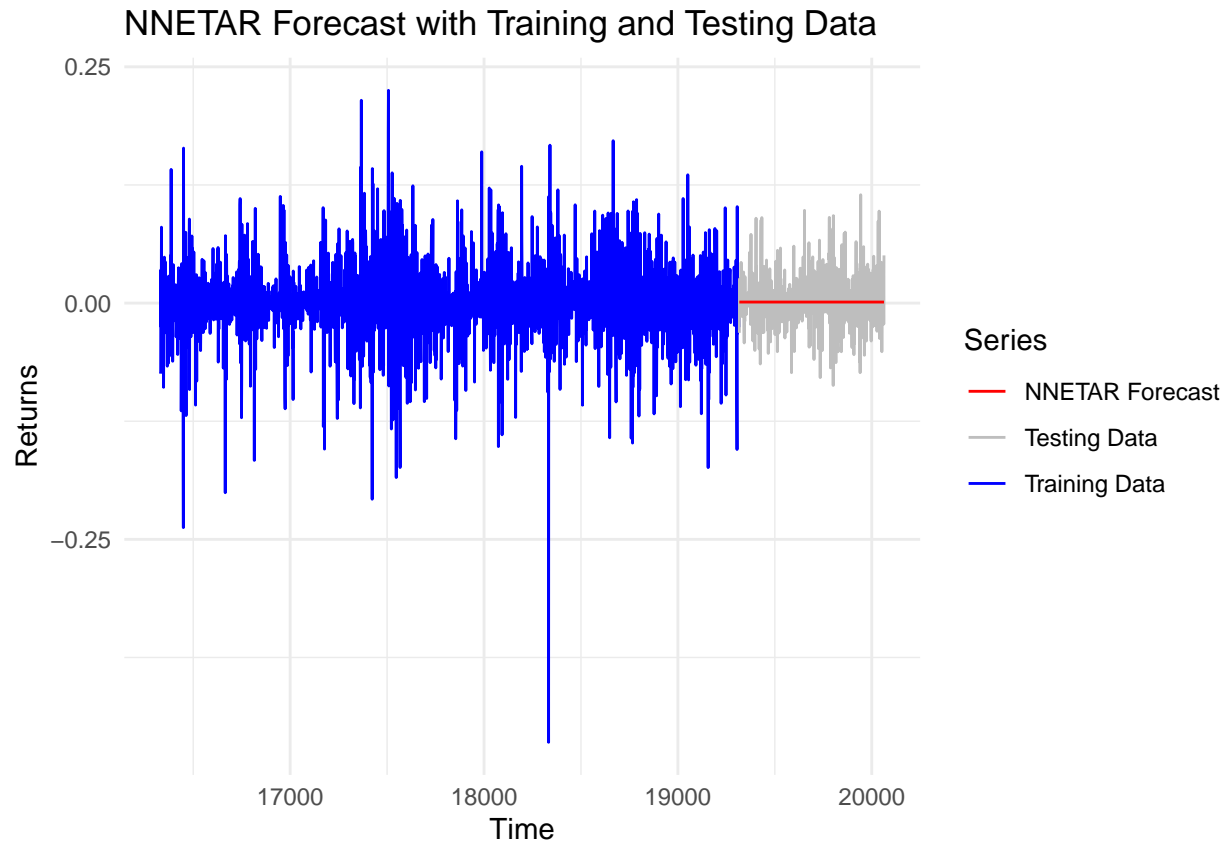
From model summary, the NNETAR model uses a single hidden layer with basic lag parameters and produces forecasts by averaging results from multiple models to capture non-linear patterns in the data.

We will proceed to forecast as follows:

```
# Generate forecasts for the testing period using NNETAR
nnetar_forecast <- forecast(nnetar_model, h = forecast_horizon)

# Convert NNETAR forecast to a time series aligned with the testing set
nnetar_forecast <- ts(nnetar_forecast$mean,
                      start = time(test_data)[1],
                      frequency = frequency(all_data))

# Plot training data, testing data, and NNETAR forecast
autoplot(train_data, series = "Training Data") +
  autolayer(test_data, series = "Testing Data") +
  autolayer(nnetar_forecast, series = "NNETAR Forecast") +
  ggtitle("NNETAR Forecast with Training and Testing Data") +
  ylab("Returns") +
  xlab("Time") +
  theme_minimal() +
  scale_color_manual(values = c("red", "grey", "blue")) +
  labs(color = "Series")
```



From the plots, the NNETAR forecast is a horizontal line. It is reasonable since the data is stationary, seeming ly less volatile than test set.

Model preformance validation (AIC/BIC omitted since unsupported by neural networks)

```
# Calculate accuracy metrics for NNETAR
nnetar_accuracy <- accuracy(nnetar_forecast, test_data)

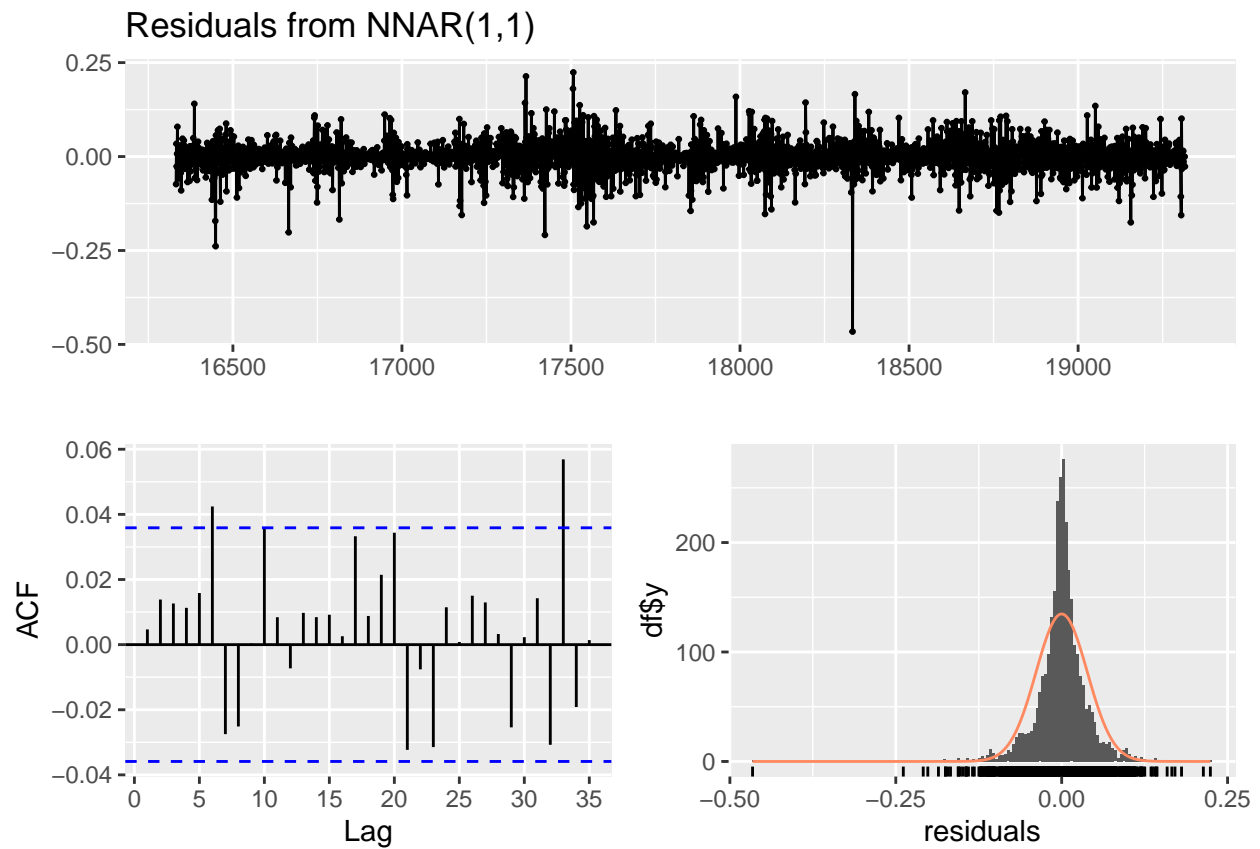
# Display the accuracy metrics
print(nnetar_accuracy)
```

```
##              ME      RMSE      MAE      MPE      MAPE      ACF1
## Test set  0.001346936 0.02505158 0.01747562 137.7628 167.0534 -0.03309994
##           Theil's U
## Test set   1.041947
```

The accuracy metrics for the ARIMA model show an RMSE of 0.0249 and an MAE of 0.0174, indicating reasonable predictive accuracy. The Mean Error is quite low at 0.0013, suggesting minimal bias in the forecasts. However, the Mean Percentage Error and Mean Absolute Percentage Error are relatively high, indicating that there are substantial relative errors, which may be due to high volatility in the data. The acf of 0.027 is higher than the original data, suggesting weaker stationarity Theil's U statistic, close to 1, suggests that the model performs similarly to a naive benchmark model.

We will proceed to check the residuals as follows:


```
# Check residual diagnostics for the NNETAR model
checkresiduals(nnetar_model)
```



```
##
##  Ljung-Box test
##
## data:  Residuals from NNAR(1,1)
## Q* = 15.695, df = 10, p-value = 0.1087
##
## Model df: 0.   Total lags used: 10
```

The residual diagnostics for the NNETAR model show that the residuals are mostly centered around zero, indicating that the model has captured the general trend of the data well. The autocorrelation plot reveals that most lags fall within the confidence bounds, suggesting minimal autocorrelation in the residuals, which implies that the model effectively accounted for the time series dynamics. However, there are a few spikes in the autocorrelation function that indicate some remaining dependencies. Overall, the residuals appear approximately normally distributed, although there is some minor deviation.

Prophet

We will proceed to fit a prophet model to our data as follows:

```
# Fit the Prophet model to the training data
train_df <- data.frame(ds = time(train_data), y = as.numeric(train_data))
prophet_model <- prophet(train_df)
summary(prophet_model)
```

```
##                Length Class      Mode
## growth                1  -none-   character
## changepoints          25  POSIXct  numeric
## n.changepoints         1  -none-   numeric
## changepoint.range      1  -none-   numeric
## yearly.seasonality      1  -none-   character
## weekly.seasonality      1  -none-   character
## daily.seasonality       1  -none-   character
## holidays               0  -none-   NULL
## seasonality.mode        1  -none-   character
## seasonality.prior.scale  1  -none-   numeric
## changepoint.prior.scale  1  -none-   numeric
## holidays.prior.scale    1  -none-   numeric
## mcmc.samples            1  -none-   numeric
## interval.width          1  -none-   numeric
## uncertainty.samples      1  -none-   numeric
## specified.changepoints  1  -none-   logical
## start                   1  POSIXct  numeric
## y.scale                 1  -none-   numeric
## logistic.floor          1  -none-   logical
## t.scale                 1  -none-   numeric
## changepoints.t          25  -none-   numeric
## seasonalities           0  -none-   list
## extra_regressors        0  -none-   list
## country_holidays        0  -none-   NULL
## stan.fit                4  -none-   list
## params                  6  -none-   list
## history                  5  data.frame list
## history.dates           2986 POSIXct  numeric
## train.holiday.names     0  -none-   NULL
## train.component.cols     3  data.frame list
## component.modes         2  -none-   list
## fit.kwargs              0  -none-   list
```

The Prophet model summary indicates key components, such as growth for trend type, changepoints to handle structural changes, and multiple seasonality options (yearly, weekly, daily). It also includes attributes like holidays and various priors for regularization. These attributes provide control over trend flexibility, seasonal components, and other time series characteristics, making Prophet versatile for capturing complex seasonality and trend patterns in time series data.

We will proceed to perform the forecast as follows:

```
# Perform step-ahead forecasting with Prophet
prophet_forecast = predict(prophet_model, periods = len(test_data)) #forecast the length of test, since
#here I am not exactly sure why generates almost 3000 even if I did periods = 746

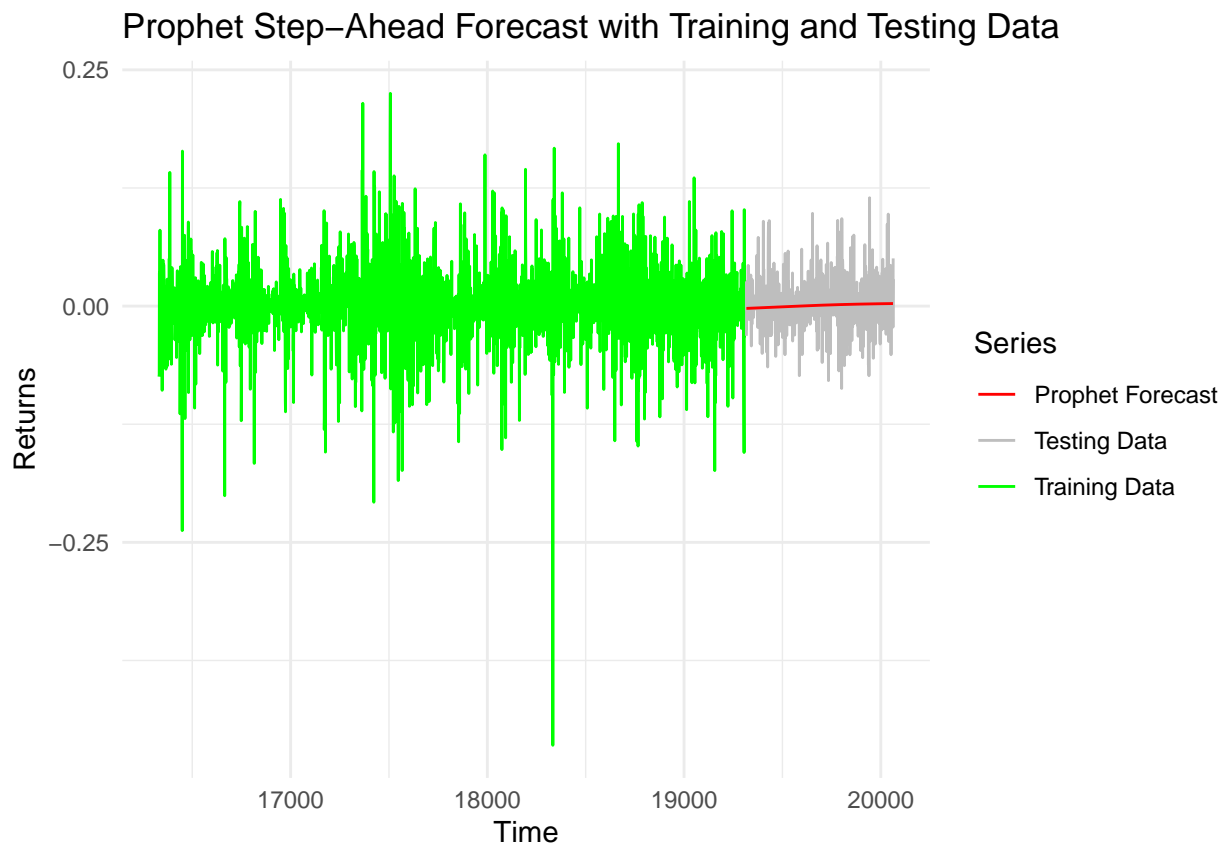
# Extract only the forecasted values
prophet_forecast_values <- prophet_forecast$yhat
```

```

# Convert Prophet forecast to a time series aligned with the testing set
prophet_forecast_ts <- ts(prophet_forecast_values[1:746],
                          start = time(test_data)[1],
                          frequency = frequency(all_data))

# Plot training data, testing data, and Prophet forecast
autoplot(train_data, series = "Training Data") +
  autolayer(test_data, series = "Testing Data") +
  autolayer(prophet_forecast_ts, series = "Prophet Forecast") +
  ggtitle("Prophet Step-Ahead Forecast with Training and Testing Data") +
  ylab("Returns") +
  xlab("Time") +
  theme_minimal() +
  scale_color_manual(values = c("red", "grey", "green")) +
  labs(color = "Series")

```



Again, the prediction is nearly horizontal, and seemingly less volatile to test set, not very sure why it while required.

Let us assess the accuracy of our model as follows:

```

# Calculate accuracy metrics for Prophet
prophet_accuracy <- accuracy(prophet_forecast_ts, test_data)

# Display the accuracy metrics
print(prophet_accuracy)

```

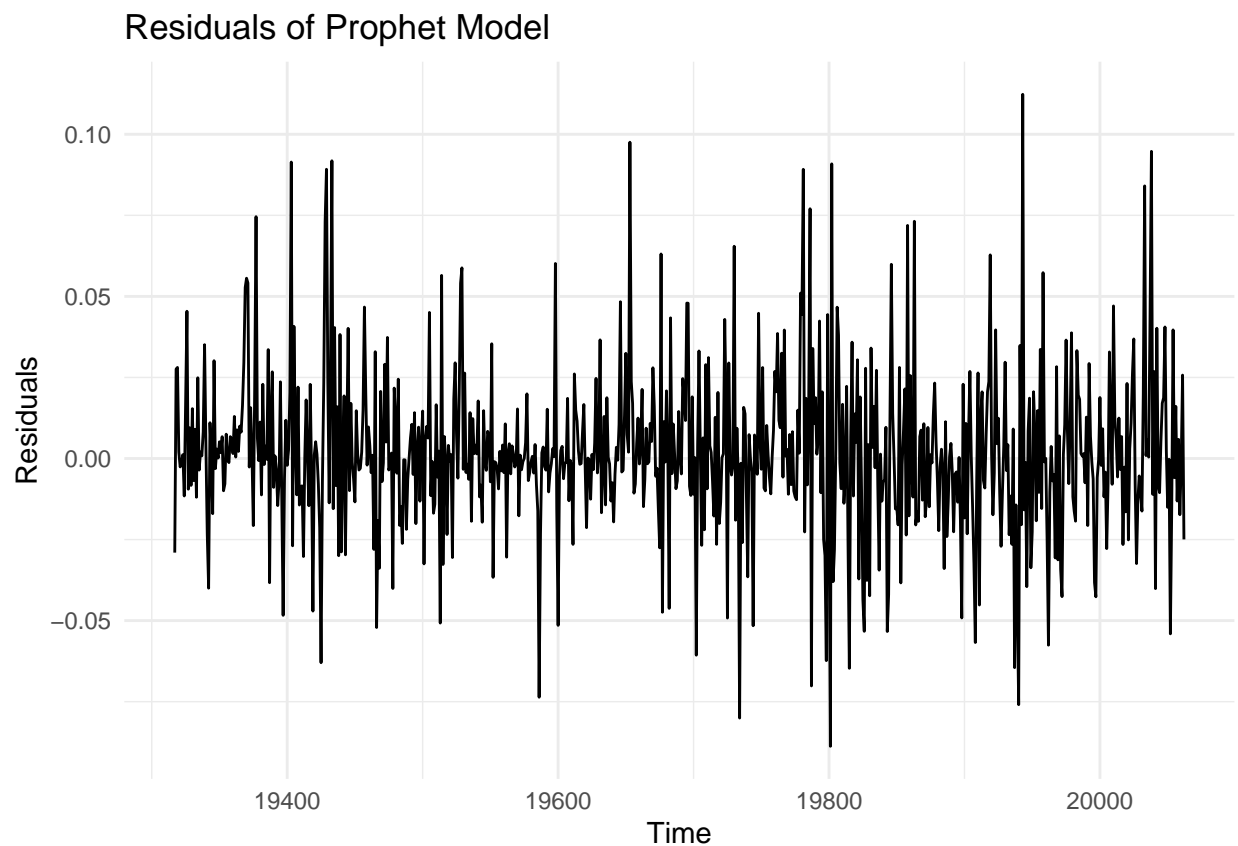
```
##           ME           RMSE           MAE           MPE           MAPE           ACF1
## Test set 0.001797758 0.02508592 0.0174258 91.99122 144.088 -0.02640533
##           Theil's U
## Test set 1.004416
```

The accuracy metrics for the Prophet model show a mean error of approximately 0.0017 and a root mean square error of 0.0251, indicating a relatively small error between predicted and actual values. The mean absolute error is around 0.0174, while the mean percentage error and mean absolute percentage error are 97.39% and 138.93%, respectively, suggesting some bias and magnitude of forecasting error. Theil's U statistic is approximately 1.007, which indicates that the forecasting performance is similar to using a naive model.

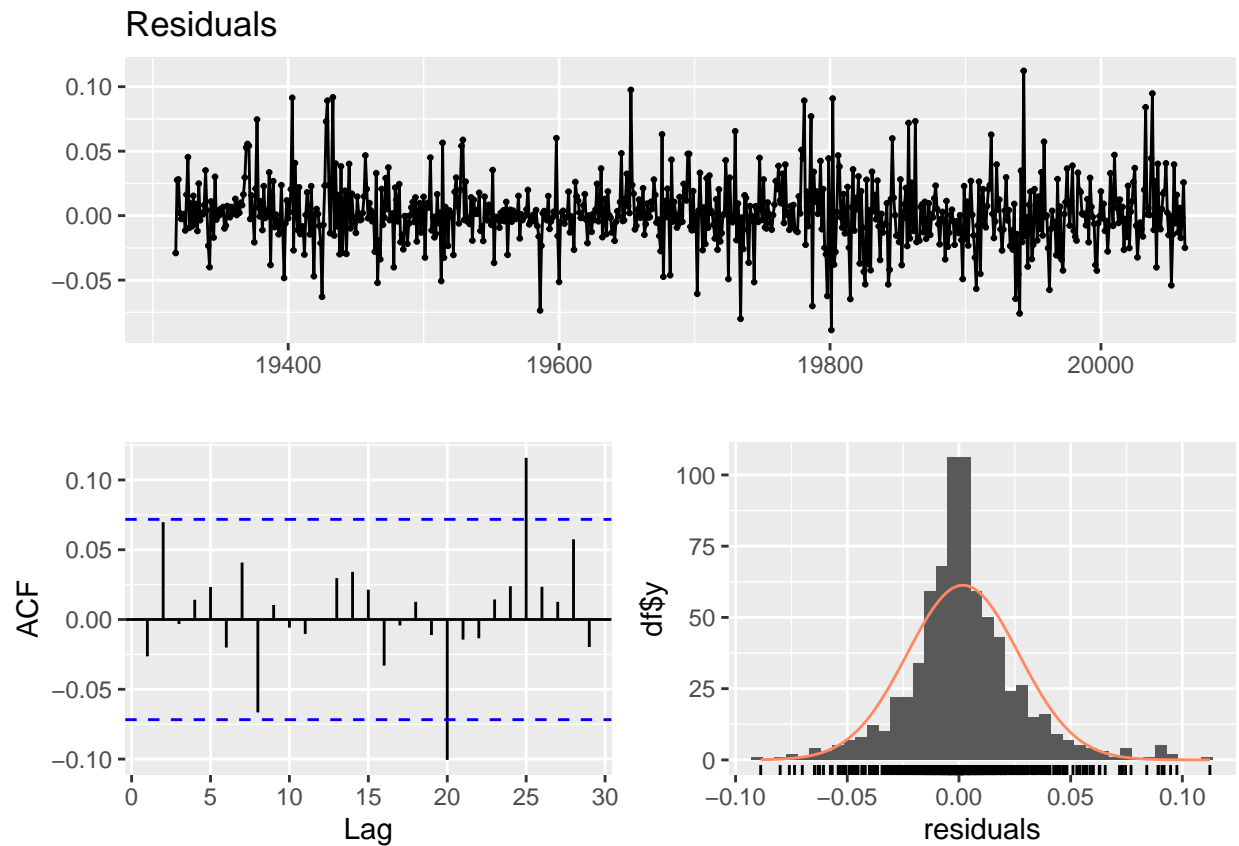
We will proceed to check the residuals:

```
# Calculate residuals for the Prophet model
prophet_residuals <- test_data - prophet_forecast_ts

# Plot residuals
autoplot(prophet_residuals) +
  ggtitle("Residuals of Prophet Model") +
  ylab("Residuals") +
  xlab("Time") +
  theme_minimal()
```



```
# Check residual diagnostics
checkresiduals(prophet_residuals)
```



```
##
##  Ljung-Box test
##
## data:  Residuals
## Q* = 9.7622, df = 10, p-value = 0.4616
##
## Model df: 0.   Total lags used: 10
```

The residual analysis of the Prophet model shows that residuals fluctuate around zero, indicating that the model generally captures the trend. However, some significant spikes suggest the presence of outliers or periods where the model struggled to capture the variability accurately. The autocorrelation function (ACF) plot shows some minor lags outside the confidence intervals, implying that some temporal structure might still be present in the residuals, meaning there is room for improvement in capturing dependencies. The histogram of residuals shows a roughly normal distribution, indicating that the residual errors are reasonably spread out, but the presence of skew or kurtosis needs further investigation.

Combination

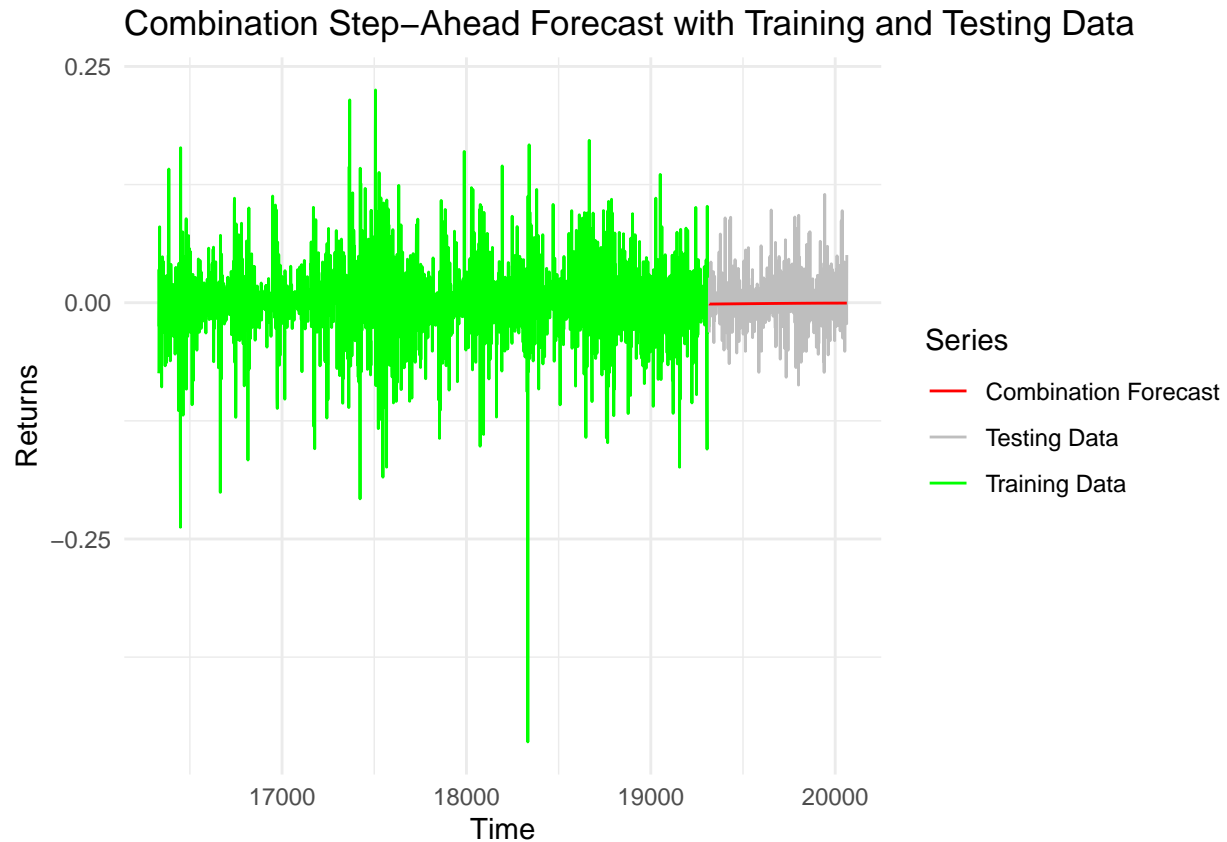
Let us perform a combination fit by averaging the forecasted values of our values:

```

# Create the combination forecast
Num_Models <- 5
prophetforecast_adjusted <- prophet_forecast$yhat[1:forecast_horizon]
Combination <- (arima_forecast+nnetar_forecast+
prophetforecast_adjusted+ets_forecast[["mean"]]+
holtwinters_forecast[["mean"]])/Num_Models

# Plot training data, testing data, and Combination forecast
autoplot(train_data, series = "Training Data") +
  autolayer(test_data, series = "Testing Data") +
  autolayer(Combination, series = "Combination Forecast") +
  ggtitle("Combination Step-Ahead Forecast with Training and Testing Data") +
  ylab("Returns") +
  xlab("Time") +
  theme_minimal() +
  scale_color_manual(values = c("red", "grey", "green")) +
  labs(color = "Series")

```



We will proceed to analyze the model performance as follows:

```

# Calculate accuracy metrics for Combination Forecast
combination_accuracy <- accuracy(Combination, test_data)

# Display the accuracy metrics
print(combination_accuracy)

```

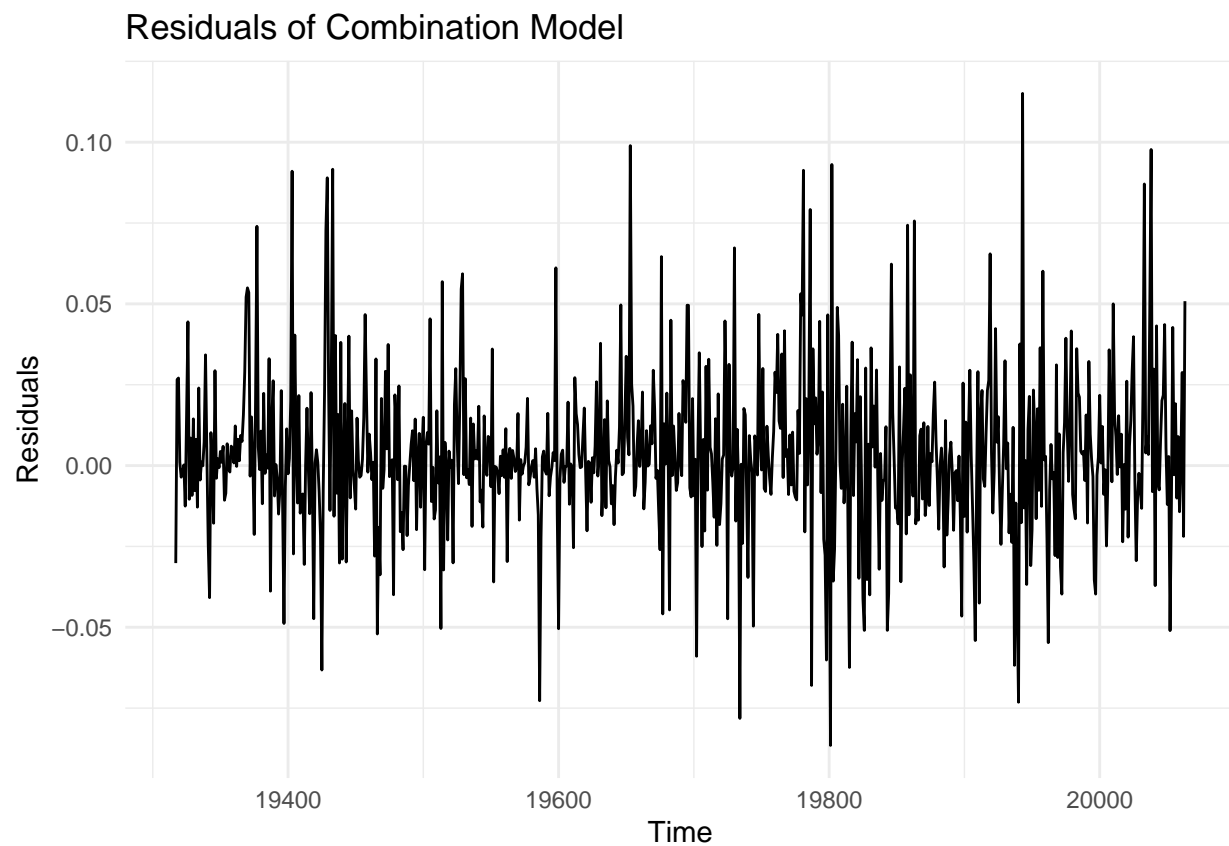
```
##           ME           RMSE           MAE           MPE           MAPE           ACF1
## Test set 0.003247247 0.02522686 0.01749539 67.07004 146.7856 -0.03290101
##           Theil's U
## Test set 0.9729756
```

The accuracy metrics for the Combined model show a mean error of approximately 0.0031 and a root mean square error of 0.0251, indicating a relatively small error between predicted and actual values. The mean absolute error is around 0.0174, while the mean percentage error and mean absolute percentage error are 72.46% and 138.82%, respectively, suggesting some bias. Theil's U statistic is approximately 0.977, which indicates that the forecasting performance is similar to using a naive model.

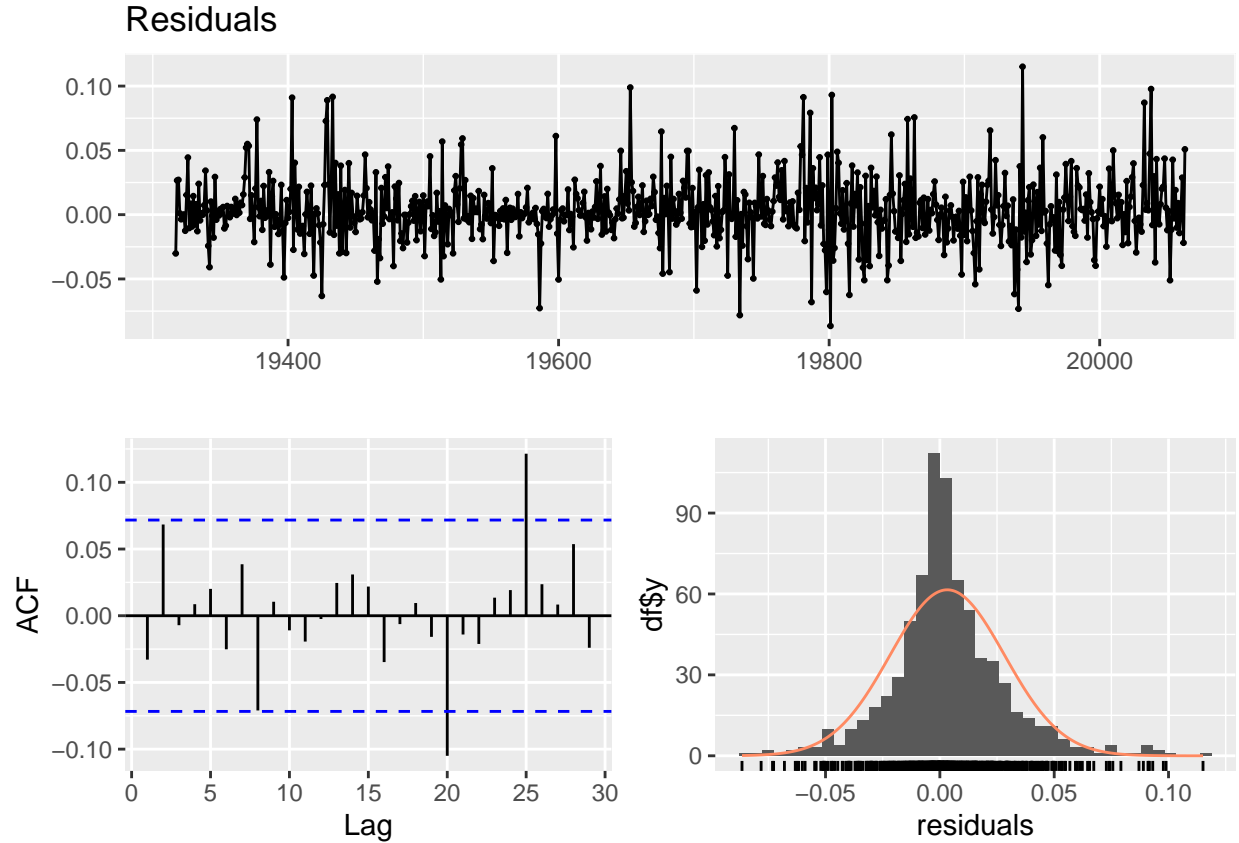
We will proceed to check the residuals as follows:

```
# Calculate residuals for the Combination model
combination_residuals <- test_data - Combination

# Plot residuals
autoplot(combination_residuals) +
  ggtitle("Residuals of Combination Model") +
  ylab("Residuals") +
  xlab("Time") +
  theme_minimal()
```



```
# Check residual diagnostics
checkresiduals(combination_residuals)
```



```
##
##  Ljung-Box test
##
## data:  Residuals
## Q* = 10.317, df = 10, p-value = 0.4132
##
## Model df: 0.   Total lags used: 10
```

The residual analysis of the Combined model shows that residuals mean revert around zero, indicating that the model captures the trend. However, some significant spikes in the residuals suggest the presence of outliers. The autocorrelation function (ACF) plot shows some minor lags outside the confidence intervals, implying that there are a few dynamics that could be captured in the residuals. The histogram of residuals shows a roughly normal distribution, indicating that the residual errors are reasonably spread out.

Conclusion & Future Works

Among the six models evaluated, using MAPE as the metric for accuracy, the Combination Forecast outperformed the others with the lowest MAPE of **138.0643**, while the Holt-Winters model had the highest MAPE at **411.77**, indicating the poorest performance. This result aligns with the principle of diversification gains, as combining forecasts from multiple models can capture different aspects of the data and mitigate individual model biases.

While the irrelevance proposition suggests that forecast combinations are unnecessary under ideal conditions—such as no model misspecification, infinite data samples (eliminating estimation error), and

full access to the underlying information sets—these conditions are rarely satisfied in practice, particularly in a noisy and highly volatile market like Bitcoin. Consequently, our findings confirm that a forecast combination provides a more robust and accurate prediction in such an environment.

For future research, several avenues could be explored to improve forecasting accuracy and better account for the noise in Bitcoin price data:

1. **Kalman Filter:** Applying a Kalman Filter to model Bitcoin prices could be beneficial, as it is well-suited for handling noisy and dynamic systems. This approach could help in smoothing the noise and capturing latent structures in the price dynamics.
2. **ARCH/GARCH Models:** To model the volatility inherent in Bitcoin prices, fitting an ARCH (Autoregressive Conditional Heteroskedasticity) or GARCH (Generalized ARCH) model could provide valuable insights. These models are particularly effective in capturing time-varying volatility and clustering effects often observed in financial data.
3. **Incorporating Exogenous Variables:** Future models could integrate external factors (e.g., trading volume, macroeconomic indicators, or sentiment analysis from social media) to account for additional information that influences Bitcoin price movements.
4. **MAPA Models:** The Multiple Aggregation Prediction Algorithm (MAPA) could enhance forecast accuracy by aggregating forecasts at different temporal resolutions. This approach captures both short-term volatility and long-term patterns, addressing Bitcoin returns' high variability and improving predictive robustness.

By pursuing these directions, we can further refine the predictive capabilities and robustness of models for Bitcoin price forecasting.

References

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) The New S Language. Wadsworth & Brooks/Cole.