

# Data Management

## Contents

<b>Raw Data</b>	<b>2</b>
<b>Cleanup</b>	<b>2</b>
<b>Building Additional Variables</b>	<b>3</b>
Volatility By Hour . . . . .	3
Social Media Post Count . . . . .	3
Adding Dummy for Social Media Post . . . . .	3
Sentiment Analysis . . . . .	4
Adding . . . . .	5
<b>Data Save</b>	<b>5</b>
<b>Merging All Data</b>	<b>5</b>
First Merge . . . . .	5
<b>Using Alpha Vantage API</b>	<b>5</b>
<b>Tutorials</b>	<b>6</b>
<b>Symbols Explanation</b>	<b>6</b>

## Raw Data

```
# 1. Political

#truthsocial
raw_truths <- read.csv(here("data/political_data", "truths_new.csv"))

#twitter
raw_tweets <- read.csv(here("data/political_data", "tweets.csv"))

# 2. Financial

#S&P500
data_loader(symbol="SPY")

#STOXX50
data_loader(symbol="VGK")

#CSI 300 (China)
data_loader(symbol="ASHR")
```

## Cleanup

```
#MAKE FUNCTIONS TO CLEANUP EASY

# 1. Tweets
tweets = raw_tweets

#only keep original Tweets
tweets <- tweets %>% filter(isRetweet != "t")
tokens <- tokens(tweets$text)
dfm <- dfm(tokens)

#cleanup
tweets = data.frame(tweets$date, tweets$text)
colnames(tweets) = c("timestamp", "tweet_text")
tweets$timestamp = as.POSIXct(tweets$timestamp, format = "%Y-%m-%d %H:%M:%S")
second(tweets$timestamp) = 0

# 2. Truths
truthsbackup <- truths_processor(raw_truths)
truths = truthsbackup

#cleanup
truths <- truths %>% filter(media != 1)
truths = data.frame(truths$date_time_parsed, truths$post)
colnames(truths) = c("timestamp", "truths_text")
truths$timestamp = as.POSIXct(truths$timestamp, format = "%Y-%m-%d %H:%M:%S")
```

```

second(truths$timestamp) = 0

# Merging social media data since it is not overlapping
names(truths)[names(truths) == 'truths_text'] <- 'tweet_text'
social = rbind(tweets,truths)
social <- social[order(social$timestamp, decreasing=F), ]

# 3. Financial

#remove index
SPY = raw_SPY[-1]
VGK = raw_VGK[-1]
ASHR = raw_ASHR[-1]

```

## Building Additional Variables

### Volatility By Hour

```

SPY = r.vol_hourly(SPY,merge=T)
VGK = r.vol_hourly(VGK,merge=T)
ASHR = r.vol_hourly(ASHR,merge=T)

```

### Social Media Post Count

```

#convert to datatable
social = as.data.table(social)

#count by hour
tweet_count = social[, .N, by=.(year(timestamp), month(timestamp),
                                day(timestamp), hour(timestamp))]

#fix timestamp by hour
tweet_count$timestamp = as.POSIXct(sprintf("%04d-%02d-%02d %02d:00:00",
                                           tweet_count$year, tweet_count$month, tweet_count$day,
                                           tweet_count$hour), format = "%Y-%m-%d %H:00:00")

#remove useless columns and reorder by oldest first
tweet_count = dplyr::select(tweet_count, timestamp, N)
tweet_count = tweet_count[ order(tweet_count$timestamp , decreasing = F ),]

```

### Adding Dummy for Social Media Post

```

#using post count we create dummy
social_hourly = tweet_count %>% mutate(dummy = if_else(N > 0, 1, 0))

```

## Sentiment Analysis

```
#sentiment analysis on post text
nrc_scores <- get_nrc_sentiment(social$tweet_text)

#add to main social dataframe
social <- bind_cols(social, nrc_scores)

#aggregate by hour
sent_hour <- social %>%
  mutate(timestamp = floor_date(timestamp, unit = "hour")) %>%
  group_by(timestamp) %>%
  summarise(across(anger:positive, sum), .groups = 'drop')
sent_hour = as.data.frame(sent_hour)

#get proportion of sentiment for each post
social <- social %>%
  mutate(total_sentiment = anger + anticipation + disgust + fear +
         joy + sadness + surprise + trust,
         total_posneg = positive + negative,
         prop_anger = anger / total_sentiment,
         prop_anticipation = anticipation / total_sentiment,
         prop_disgust = disgust / total_sentiment,
         prop_fear = fear / total_sentiment,
         prop_joy = joy / total_sentiment,
         prop_sadness = sadness / total_sentiment,
         prop_surprise = surprise / total_sentiment,
         prop_trust = trust / total_sentiment,
         prop_negative = negative / total_posneg,
         prop_positive = positive / total_posneg)
social[is.na(social)] <- 0

#same but hourly
sent_hour <- sent_hour %>%
  mutate(total_sentiment = anger + anticipation + disgust + fear +
         joy + sadness + surprise + trust,
         total_posneg = positive + negative,
         prop_anger = anger / total_sentiment,
         prop_anticipation = anticipation / total_sentiment,
         prop_disgust = disgust / total_sentiment,
         prop_fear = fear / total_sentiment,
         prop_joy = joy / total_sentiment,
         prop_sadness = sadness / total_sentiment,
         prop_surprise = surprise / total_sentiment,
         prop_trust = trust / total_sentiment,
         prop_negative = negative / total_posneg,
         prop_positive = positive / total_posneg)
sent_hour[is.na(sent_hour)] <- 0
social_hourly <- left_join(social_hourly, sent_hour, by="timestamp")
```

## Adding

```
# 1. Add count for tweets ()  
  
# 2. Sentiments ()  
  
# 3. Dummy Tweet ()  
  
# 4. Dummy Important Word  
  
# 5. Dummy Emotional Word
```

```
#is there a point?
```

## Data Save

```
#financial  
write.csv(SPY, here("data/mothership/SPY.csv"), row.names=F)  
write.csv(VGK, here("data/mothership/VGK.csv"), row.names=F)  
write.csv(ASHR, here("data/mothership/ASHR.csv"), row.names=F)  
  
#social media  
write.csv(social, here("data/mothership/social.csv"), row.names=F)  
write.csv(social_hourly, here("data/mothership/socialhourly.csv"), row.names=F)
```

## Merging All Data

### First Merge

```
#run script to load and merge all data  
rm(list=ls())  
source(here("helperfunctions/fulldata_loader.R"))  
  
#i think there's no point to this chunk just complicates things
```

## Using Alpha Vantage API

```
library(alphavantage)  
  
av_api_key(Sys.getenv("ALPHAVANTAGE_API_KEY"))  
  
#for past month  
data=av_get("ASHR", av_fun = "TIME_SERIES_INTRADAY", interval = "1min",  
            adjusted="false", extended_hours="false", outputsize = "full")
```

```

#for a particular month
data2=av_get("SPY", av_fun = "TIME_SERIES_INTRADAY", interval = "1min",
            adjusted="false", extended_hours="false",
            month="2025-04", outputsize = "full") #create loop for more

write.csv(data,"~/ASHR.csv", row.names = T) #saves to documents
write.csv(data2,"~/SPY-2025-04.csv", row.names = T)

library(alphavantage)

av_api_key(Sys.getenv("ALPHAVANTAGE_API_KEY"))

year = "2022"
months = c("01","02","03","04","05","06","07","08","09","10","11","12")
market = "SPY"

for (t in 1:length(months)) {
  date = paste(year, months[t], sep="-")
  dataloop = av_get(market, av_fun="TIME_SERIES_INTRADAY",interval="1min",
                    adjusted="false", extended_hours="false",
                    month=date, outputsize="full")
  filename = paste(market,date,sep="-")
  filename = paste("~/", filename, sep="")
  filename = paste(filename,".csv",sep="")
  write.csv(dataloop,filename)
}

```

## Tutorials

Manual smp500: [https://cafim.sssup.it/~giulio/other/alpha\\_vantage/index.html#orgaaf54ef](https://cafim.sssup.it/~giulio/other/alpha_vantage/index.html#orgaaf54ef)

AlphaVantageR Tutorial: [https://github.com/business-science/alphavantager/blob/master/man/av\\_get.Rd](https://github.com/business-science/alphavantager/blob/master/man/av_get.Rd)

Intra-Day Analysis: <https://arxiv.org/html/2406.17198v1>

## Symbols Explanation

- ONEQ = NASDAQ Composite
- SPY = S&P500
- SMI = Swiss Market Index
- VTHR = Russell 3000 (US)
- VTI = CRSP US Total Market Index
- VGK = Euro Stoxx 50
- ASHR = basically china