

Live Case: S&P500 (1 of 3)

Aug 10, 2023.

S&P 500

The S&P 500, also called the Standard & Poor's 500, is a stock market index that tracks the performance of 500 major publicly traded companies listed on U.S. stock exchanges. It serves as a widely accepted benchmark for assessing the overall health and performance of the U.S. stock market.

S&P Dow Jones Indices, a division of S&P Global, is responsible for maintaining the index. The selection of companies included in the S&P 500 is determined by a committee, considering factors such as market capitalization, liquidity, and industry representation.

The S&P is a float-weighted index, meaning the market capitalizations of the companies in the index are adjusted by the number of shares available for public trading. <https://www.investopedia.com/terms/s/sp500.asp>

The performance of the S&P 500 is frequently used to gauge the broader stock market and is commonly referenced by investors, analysts, and financial media. It provides a snapshot of how large-cap U.S. stocks are faring and is considered a reliable indicator of overall market sentiment.

Typically, the S&P 500 index consists of 500 stocks. However, in reality, there are actually 503 stocks included. This discrepancy arises because three of the listed companies have multiple share classes, and each class is considered a separate stock that needs to be included in the index.

Among these 503 stocks, Apple, the technology giant, holds the top position with a market capitalization of \$2.35 billion. Following Apple, Microsoft and Amazon.com rank as the second and third largest stocks in the S&P 500, respectively. The next positions are held by Nvidia Corp, Tesla, Berkshire Hathaway, and two classes of shares from Google's parent company, Alphabet..

S&P 500 Data - Preliminary Analysis

We will analyze a real-world, recent dataset containing information about the S&P500 stocks. The dataset is located in a Google Sheet

The data is disorganized and challenging to understand. We will review the data and proceed in a step-by-step manner.

Read the S&P500 data from a Google Sheet into a tibble dataframe.

1. The complete URL is
`https://docs.google.com/spreadsheets/d/11ahk9uWxBkDqrhNm7qYmiTwrlSC53N1zvXYfv7ttOCM/`
2. The Google Sheet ID is: `11ahk9uWxBkDqrhNm7qYmiTwrlSC53N1zvXYfv7ttOCM`. We can use the function `gsheet2tbl` in package `gsheet` to read the Google Sheet into a tibble or dataframe, as demonstrated in the following code.

```
# Read S&P500 stock data present in a Google Sheet.
library(gsheet)
prefix <- "https://docs.google.com/spreadsheets/d/"
sheetID <- "11ahk9uWxBkDqrhNm7qYmiTwrlSC53N1zvXYfv7ttOCM"
url500 <- paste(prefix,sheetID) # Form the URL to connect to
sp500 <- gsheet2tbl(url500) # Read it into a tibble called sp500
```

No encoding supplied: defaulting to UTF-8.

Review the data

1. We want to understand the different data columns and their data structure. For this purpose, we run the `str()` function.

```
str(sp500)
```

```
Classes 'tbl_df', 'tbl' and 'data.frame':  503 obs. of  36 variables:
 $ Date           : chr  "8/12/2023" "8/12/2023" "8/12/2023" "8/12/2023" ...
 $ Stock          : chr  "A" "AAL" "AAP" "AAPL" ...
 $ Description    : chr  "Agilent Technologies, Inc." "American Airl..."
 $ Sector         : chr  "Health Technology" "Transportation" "Retail..."
 $ Industry       : chr  "Medical Specialties" "Airlines" "Specialty..."
```

```

$ Market.Capitalization      : num  3.70e+10 1.03e+10 4.16e+09 3.01e+12 2.63e+1
$ Price                      : num  125.2 15.8 69.9 191.2 148.9 ...
$ X52.Week.Low              : num  113.3 11.7 63.6 124.2 131 ...
$ X52.Week.High             : num  160 19.1 212 198 168 195 116 84.8 328 553 .
$ Return.on.Equity..TTM.     : num  24.8 NA 14.6 146 NA 295 NA NA 30.7 33.7 ..
$ Return.on.Assets..TTM.     : num  12.7 3.9 3.35 27.6 NA 2.86 NA NA 14.9 17.9
$ Return.on.Invested.Capital..TTM. : num  16.51 8.01 6.17 57.18 NA ...
$ Gross.Margin..TTM.        : num  54.1 23.8 43.8 43.2 87 ...
$ Operating.Margin..TTM.     : num  23.78 9.39 5.63 29.16 40.28 ...
$ Net.Margin..TTM.          : num  19.19 4.98 3.61 24.49 15.46 ...
$ Price.to.Earnings.Ratio..TTM. : num  27.6 4.29 10.41 32.48 30.54 ...
$ Price.to.Book..FY.        : num  6.97 NA 1.55 60.15 15.27 ...
$ Enterprise.Value.EBITDA..TTM. : num  20 5.58 8.78 24.9 NA 12.3 NA NA 17.6 34.2 .
$ EBITDA..TTM.              : num  1.97e+09 7.16e+09 9.21e+08 1.24e+11 NA ...
$ EPS.Diluted..TTM.         : num  4.54 3.69 6.72 5.89 4.88 ...
$ EBITDA..TTM.YoY.Growth.    : num  10.52 1074.1 -16 -5.36 NA ...
$ EBITDA..Quarterly.YoY.Growth. : num  8.2 72.2 -39.01 -4.58 NA ...
$ EPS.Diluted..TTM.YoY.Growth. : num  9.17 NA -25.21 -4.33 -31.01 ...
$ EPS.Diluted..Quarterly.YoY.Growth. : num  11.69944 156.4148 -68.3683 -0.00656 122.310
$ Price.to.Free.Cash.Flow..TTM. : num  31.44 7.12 NA 31.08 NA ...
$ Free.Cash.Flow..TTM.YoY.Growth. : num  11.81 NA -100.23 -7.85 NA ...
$ Free.Cash.Flow..Quarterly.YoY.Growth. : num  55.7078 -10.2542 -176.135 -0.0312 NA ...
$ Debt.to.Equity.Ratio..MRQ. : num  0.473 NA 1.582 1.763 NA ...
$ Current.Ratio..MRQ.       : num  2.37 0.749 1.244 0.94 NA ...
$ Quick.Ratio..MRQ.         : num  1.708 0.656 0.238 0.878 NA ...
$ Dividend.Yield.Forward    : num  0.705 NA 1.436 0.498 3.963 ...
$ Dividends.per.share..Annual.YoY.Growth. : num  8.25 NA 84.62 5.88 7.53 ...
$ Price.to.Sales..FY.       : num  5.487 0.212 0.381 7.915 4.56 ...
$ Revenue..TTM.YoY.Growth.   : num  7.86 29.909 1.415 -0.254 -2.312 ...
$ Revenue..Quarterly.YoY.Growth. : num  6.85 4.72 1.29 -2.51 -4.92 ...
$ Technical.Rating           : chr  "Strong Sell" "Strong Sell" "Buy" "Sell" ..

```

2. The `str(sp500)` output provides valuable insights into the structure and data types of the columns in the `sp500` tibble. Let's delve into the details.
3. The output reveals that `sp500` is a tibble with dimensions $[503 \times 36]$. This means it consists of 503 rows, each representing a specific S&P500 stock, and 36 columns containing information about each stock.
4. Here is a preliminary breakdown of the information associated with each column:
 - The columns labeled `Date`, `Stock`, `Description`, `Sector`, and `Industry` are character columns. They respectively represent the date, stock ticker symbol, description, sector, and industry of each S&P500 stock.

- Columns such as `Market.Capitalization`, `Price`, `X52.Week.Low`, `X52.Week.High`, and other numeric columns contain diverse financial metrics and stock prices related to the S&P500 stocks.
 - The column labeled `Technical.Rating` is a character column that assigns a technical rating to each stock.
5. By examining the `str(sp500)` output, we gain a preliminary understanding of the data types and column names present in the `sp500` tibble, enabling us to grasp the structure of the dataset.

Rename Data Columns

1. The names of the data columns are lengthy and confusing.
2. We will rename the data columns to make it easier to work with the data, using the `rename_with()` function.

```
# Define a mapping of new column names
new_names <- c(
  "Date", "Stock", "StockName", "Sector", "Industry",
  "MarketCap", "Price", "Low52Wk", "High52Wk",
  "ROE", "ROA", "ROIC", "GrossMargin",
  "OperatingMargin", "NetMargin", "PE",
  "PB", "EVEBITDA", "EBITDA", "EPS",
  "EBITDA_YOY", "EBITDA_QYOY", "EPS_YOY",
  "EPS_QYOY", "PFCF", "FCF",
  "FCF_QYOY", "DebtToEquity", "CurrentRatio",
  "QuickRatio", "DividendYield",
  "DividendsPerShare_YOY", "PS",
  "Revenue_YOY", "Revenue_QYOY", "Rating"
)
# Rename the columns using the new_names vector
sp500 <- sp500 %>%
  rename_with(~ new_names, everything())
```

This code is designed to rename the columns of the `sp500` tibble using a predefined mapping of new column names. Let's go through the code step by step:

1. A vector named `new_names` is created, which contains the desired new names for each column in the `sp500` tibble. Each element in the `new_names` vector corresponds to a specific column in the `sp500` tibble and represents the desired new name for that column.

2. The `%>%` operator, often referred to as the pipe operator, is used to pass the `sp500` tibble to the subsequent operation in a more readable and concise manner.
3. The `rename_with()` function from the `dplyr` package is applied to the `sp500` tibble. This function allows us to rename columns based on a specified function or formula.
4. In this case, a formula `~ new_names` is used as the first argument of `rename_with()`. This formula indicates that the new names for the columns should be sourced from the `new_names` vector.
5. The second argument, `everything()`, specifies that the renaming should be applied to all columns in the `sp500` tibble.
6. Finally, the resulting tibble with the renamed columns is assigned back to the `sp500` variable, effectively updating the tibble with the new column names.
7. We could also use the following code to rename the columns.

```
# Rename the columns using the new_names vector
colnames(sp500) <- new_names
```

In essence, the code uses the `new_names` vector as a mapping to assign new column names to the `sp500` tibble, ensuring that each column is given the desired new name specified in `new_names`.

Review the data again after renaming columns

1. We review the column names again after renaming them, using the `colnames()` function can help.

```
colnames(sp500)
```

[1] "Date"	"Stock"	"StockName"
[4] "Sector"	"Industry"	"MarketCap"
[7] "Price"	"Low52Wk"	"High52Wk"
[10] "ROE"	"ROA"	"ROIC"
[13] "GrossMargin"	"OperatingMargin"	"NetMargin"
[16] "PE"	"PB"	"EVEBITDA"
[19] "EBITDA"	"EPS"	"EBITDA_YOY"
[22] "EBITDA_QYOY"	"EPS_YOY"	"EPS_QYOY"
[25] "PFCF"	"FCF"	"FCF_QYOY"
[28] "DebtToEquity"	"CurrentRatio"	"QuickRatio"
[31] "DividendYield"	"DividendsPerShare_YOY"	"PS"

[34] "Revenue_YOY"

"Revenue_QYOY"

"Rating"

Understand the Data Columns

1. The complete data has 36 columns. Our goal is to gain a deeper understanding of what the data columns mean.
2. We reorganize the column names into eight tables, labeled Table 1a, 1b.. 1h.
 - a. The column names described in Table 1a. concern basic **Company Information** of each stock.

Table 1a: Data Columns giving basic Company Information	
ColumnName	Description
Date	Date (e.g. "7/15/2023")
Stock	Stock Ticker (e.g. AAL)
StockName	Name of the company (e.g "American Airlines Group, Inc.")
Sector	Sector the stock belongs to (e.g. "Transportation")
Industry	Industry the stock belongs to (e.g "Airlines")
MarketCap	Market capitalization of the company
Price	Recent Stock Price

- b. The column names described in Table 1b. are related to **Technical Analysis** of each stock, including the 52-Week High and Low prices.

Table 1b: Data Columns related to Pricing and Technical Analysis	
ColumnName	Description
Low52Wk	52-Week Low Price
High52Wk	52-Week High Price
Rating	Technical Rating

- c. The column names described in Table 1c. are related to the **Profitability** of each stock.

Table 1c: Data Columns related to Profitability	
ColumnName	Description
ROE	Return on Equity
ROA	Return on Assets
ROIC	Return on Invested Capital
GrossMargin	Gross Profit Margin

Table 1c: Data Columns related to Profitability	
ColumnName	Description
OperatingMargin	Operating Profit Margin
NetMargin	Net Profit Margin

The column names described in Table 1d are related to the **Earnings** of each stock.

Table 1d: Data Columns related to Earnings	
ColumnName	Description
PE	Price-to-Earnings Ratio
PB	Price-to-Book Ratio
EVEBITDA	Enterprise Value to EBITDA Ratio
EBITDA	EBITDA
EPS	Earnings per Share
EBITDA_YOY	EBITDA Year-over-Year Growth
EBITDA_QYOY	EBITDA Quarterly Year-over-Year Growth
EPS_YOY	EPS Year-over-Year Growth
EPS_QYOY	EPS Quarterly Year-over-Year Growth

The column names described in Table 1e are related to the **Free Cash Flow** of each stock.

Table 1e: Data Columns related to Free Cash Flow	
ColumnName	Description
PFCF	Price-to-Free Cash Flow
FCF	Free Cash Flow
FCF_QYOY	Free Cash Flow Quarterly Year-over-Year Growth

The column names described in Table 1f concern the **Liquidity** of each stock.

Table 1f: Data Columns related to Liquidity	
ColumnName	Description
DebtToEquity	Debt-to-Equity Ratio
CurrentRatio	Current Ratio
QuickRatio	Quick Ratio

The column names described in Table 1g are related to the **Revenue** of each stock.

Table 1g: Data Columns related to Revenue	
ColumnName	Description
PS	Price-to-Sales Ratio
Revenue_YOY	Revenue Year-over-Year Growth
Revenue_QYOY	Revenue Quarterly Year-over-Year Growth

The column names described in Table 1h are related to the **Dividends** of each stock.

Table 1h: Data Columns related to Dividends	
ColumnName	Description
DividendYield	Dividend Yield
DividendsPerShare_YOY	Annual Dividends per Share Year-over-Year Growth

Remove Rows containing no data or Null values

1. The following code checks if the “Stock” column in the sp500 dataframe contains any null or blank values. If there are null or blank values present, it removes the corresponding rows from the sp500 dataframe, resulting in a filtered dataframe without null or blank values in the “Stock” column.

```
# Check for blank or null values in the "Stock" column
hasNull <- any(sp500$Stock == "" | is.null(sp500$Stock))
if (hasNull) {
  # Remove rows with null or blank values from the dataframe tibble
  sp500 <- sp500[!(is.null(sp500$Stock) | sp500$Stock == ""), ]
}
```

Here’s an alternate code using `dplyr` to achieve the same result:

```
library(dplyr)
# Check for blank or null values in the "Stock" column
hasNull <- any(sp500 %>% pull(Stock) == "" | is.null(sp500 %>% pull(Stock)))
if (hasNull) {
  # Remove rows with null or blank values from the dataframe tibble
  sp500 <- sp500 %>% filter(!(is.null(Stock) | Stock == ""))
}
```



```
# View the filtered dataframe
nrow(sp500)
```

```
[1] 503
```

Thus, we have 502 stocks of the S&P500 in our dataset.

S&P500 Sector

The S&P500 shares are divided into multiple Sectors. Each stock belongs to a unique sector. Thus, it makes sense to model Sector as a `factor()` variable.

```
sp500$Sector <- as.factor(sp500$Sector)
```

It makes sense to convert Sector to a factor variable, since there are 19 distinct Sectors in the S&P500 and each stock belongs to a unique sector. We confirm that Sector is now modelled as a factor variable, by running the `str()` function.

```
str(sp500$Sector)
```

```
Factor w/ 19 levels "Commercial Services",...: 11 18 16 7 11 6 11 9 17 17 ...
```

Now that Sectors is a factor variable, we can use the `levels()` function to review the different levels it can take.

```
levels(sp500$Sector)
```

```
[1] "Commercial Services"    "Communications"        "Consumer Durables"
[4] "Consumer Non-Durables"  "Consumer Services"      "Distribution Services"
[7] "Electronic Technology"  "Energy Minerals"        "Finance"
[10] "Health Services"        "Health Technology"      "Industrial Services"
[13] "Non-Energy Minerals"    "Process Industries"     "Producer Manufacturing"
[16] "Retail Trade"           "Technology Services"    "Transportation"
[19] "Utilities"
```

The `table()` function allows us to count how many stocks are part of each sector.

```
table(sp500$Sector)
```

Commercial Services	Communications	Consumer Durables
13	3	12
Consumer Non-Durables	Consumer Services	Distribution Services
31	29	9
Electronic Technology	Energy Minerals	Finance
49	16	92
Health Services	Health Technology	Industrial Services
12	47	9
Non-Energy Minerals	Process Industries	Producer Manufacturing
7	24	31
Retail Trade	Technology Services	Transportation
23	50	15
Utilities		
31		

Thus, we can see how many stocks are part of each one of the 19 sectors.

We can sum them to confirm that they add up to 502.

```
sum(table(sp500$Sector))
```

```
[1] 503
```

This completes our review of the Sector variable.

Stock Ratings

In the data, the S&P500 shares have Technical Ratings such as {Buy, Sell, ..}. Since each Stock has a unique Technical Rating, it makes sense to model the data column Rating as a `factor()` variable.

```
sp500$Rating <- as.factor(sp500$Rating)
```

We confirm that Rating is now modelled as a factor variable, by running the `str()` function.

```
str(sp500$Rating)
```

Factor w/ 5 levels "Buy","Neutral",...: 5 5 1 3 1 3 3 3 5 5 ...

We can use the `levels()` function to review the different levels it can take.

```
levels(sp500$Rating)
```

```
[1] "Buy"          "Neutral"      "Sell"         "Strong Buy"   "Strong Sell"
```

The `table()` function allows us to count how many stocks have each Rating.

```
table(sp500$Rating)
```

Buy	Neutral	Sell	Strong Buy	Strong Sell
130	85	203	13	72

Thus, we can see how many stocks have ratings ranging from “Strong Sell” to “Strong Buy”. This completes our review of Technical Rating.

6. Low52WkPerc: Create a new column to track Share Prices relative to their 52 Week Low, as described in the chapter Live Case: S&P500 (1 of 3).

```
sp500 <- sp500 %>% mutate(Low52WkPerc = round((Price - Low52Wk)*100 / Low52Wk,2))
colnames(sp500)
```

```
[1] "Date"          "Stock"          "StockName"
[4] "Sector"        "Industry"       "MarketCap"
[7] "Price"         "Low52Wk"       "High52Wk"
[10] "ROE"           "ROA"           "ROIC"
[13] "GrossMargin"   "OperatingMargin" "NetMargin"
[16] "PE"            "PB"            "EVEBITDA"
[19] "EBITDA"        "EPS"           "EBITDA_YOY"
[22] "EBITDA_QYOY"   "EPS_YOY"       "EPS_QYOY"
[25] "PFCF"          "FCF"           "FCF_QYOY"
[28] "DebtToEquity"  "CurrentRatio"   "QuickRatio"
[31] "DividendYield" "DividendsPerShare_YOY" "PS"
[34] "Revenue_YOY"   "Revenue_QYOY"   "Rating"
[37] "Low52WkPerc"
```

7. Creating a new column `MarketCapBillions = MarketCap/1000,000,000`

```
#sp500 <- sp500 %>% mutate(MarketCapBillions = MarketCap/ 10000000000)
#colnames(sp500)
```

Where are we now?

We believe this dataset of S&P500 shares is now ready for further analysis. We end this stage of our analysis in this chapter, by running the `str()` function to review the data columns.

```
#str(sp500)
```

Summary of Chapter 6 – Exploring S&P500 Data

Chapter 6 embarks on an exploration of the S&P500, a significant stock market index encompassing 500 major publicly traded companies in the U.S. The chapter introduces the index's role as a benchmark for assessing the overall health and performance of the U.S. stock market, maintained by S&P Dow Jones Indices.

Part 1 of the chapter delves into a real-world dataset containing information about S&P500 stocks. The data is loaded into a tibble using the R package `gsheet`, and its structure is examined using the `str()` function. To facilitate data management, column names are renamed using the `rename_with()` function from `dplyr`, and a detailed breakdown of column information is presented across eight tables.

Part 2 addresses data quality, ensuring a cleaner dataset by removing rows with null or blank values in the “Stock” column. Additionally, the “Sector” and “Rating” columns are transformed into factor variables, reflecting the distinct sectors and technical ratings each stock holds. The distribution of sectors and ratings is analyzed using various functions. After data preparation, the dataset is considered ready for further analysis.

Chapter 6 skillfully guides readers through the intricacies of exploring S&P500 data, employing practical examples and R code to foster a deeper understanding of the dataset's structure and content. Further exploration is encouraged with a wealth of references for continued learning and analysis.