

ARM assignment

KORNILAKIS NIKOLAOS

2024-06-14

Loading the required packages .

```
# Packages -----  
library(tidyverse)  
  
## Warning: package 'tidyverse' was built under R version 4.3.3  
  
## Warning: package 'ggplot2' was built under R version 4.3.3  
  
## Warning: package 'tibble' was built under R version 4.3.3  
  
## Warning: package 'tidyr' was built under R version 4.3.3  
  
## Warning: package 'readr' was built under R version 4.3.3  
  
## Warning: package 'purrr' was built under R version 4.3.3  
  
## Warning: package 'dplyr' was built under R version 4.3.3  
  
## Warning: package 'stringr' was built under R version 4.3.3  
  
## Warning: package 'forcats' was built under R version 4.3.3  
  
## Warning: package 'lubridate' was built under R version 4.3.3  
  
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --  
## v dplyr      1.1.4      v readr      2.1.5  
## v forcats   1.0.0      v stringr   1.5.1  
## v ggplot2   3.5.0      v tibble    3.2.1  
## v lubridate 1.9.3      v tidyr     1.3.1  
## v purrr     1.0.2  
## -- Conflicts ----- tidyverse_conflicts() --  
## x dplyr::filter() masks stats::filter()  
## x dplyr::lag()     masks stats::lag()  
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```
library(RSQLite)
```

```
## Warning: package 'RSQLite' was built under R version 4.3.3
```

```
library(dbplyr)
```

```
## Warning: package 'dbplyr' was built under R version 4.3.3
```

```
##  
## Attaching package: 'dbplyr'  
##  
## The following objects are masked from 'package:dplyr':  
##  
##     ident, sql
```

```
library(readxl)
```

```
## Warning: package 'readxl' was built under R version 4.3.3
```

```
library(RPostgres)
```

```
## Warning: package 'RPostgres' was built under R version 4.3.3
```

```
library(frenchdata)
```

```
## Warning: package 'frenchdata' was built under R version 4.3.3
```

```
library(furrr)
```

```
## Warning: package 'furrr' was built under R version 4.3.3
```

```
## Loading required package: future
```

```
## Warning: package 'future' was built under R version 4.3.3
```

```
library(readxl)
```

```
library(DBI)
```

```
## Warning: package 'DBI' was built under R version 4.3.3
```

Starting from 1973 cause some data were not available from 1972(also mentioned on the paper), connecting to the WRDS database with my credentials .

```

# Set up -----
# SQLite database
data_nse <- dbConnect(SQLite(),
                      "data_nse.sqlite",
                      extended_types = TRUE)

# Dates
start_date <- as.Date("1973-01-01")
end_date <- as.Date("2023-12-31")

# WRDS connection
wrds <- dbConnect(
  Postgres(),
  host = "wrds-pgdata.wharton.upenn.edu",
  dbname = "wrds",
  port = 9737,
  sslmode = "require",
  user = "cornelious23",
  password = "NikosKornilakis")

```

In the next code snippet, two datasets, “Fama/French 5 Factors (2x3)” and “Fama/French 3 Factors”, are downloaded and processed to evaluate the performance of different factor models.

```

#5F
factors_ff_monthly <- download_french_data("Fama/French 5 Factors (2x3)")$subsets$data[[1]] |>
  janitor::clean_names()

```

```

## New names:
## New names:
## * ' -> '...1'

```

```

# Manipulate
factors_ff_monthly <- factors_ff_monthly |>
  transmute(
    month = floor_date(ymd(paste0(date, "01")), "month"),
    mkt_excess = as.numeric(mkt_rf) / 100,
    smb = as.numeric(smb) / 100,
    hml = as.numeric(hml) / 100,
    rmw = as.numeric(rmw) / 100,
    cma = as.numeric(cma) / 100,
    rf = as.numeric(rf) / 100
  ) |>
  filter(month >= start_date & month <= end_date)

# Store
factors_ff_monthly |>
  dbWriteTable(conn = data_nse,
              name = "factors_ff_monthly",
              value = _,
              overwrite = TRUE)

#3F
factors_ff_3fraw <- download_french_data("Fama/French 3 Factors")

```

```
## New names:
## New names:
## * ' ' -> '...1'
```

```
factors_ff_3f <-factors_ff_3fraw$subsets$data[[1]] |>
mutate(
  month = floor_date(ymd(str_c(date, "01")), "month"),
  across(c(RF, `Mkt-RF`, SMB, HML), ~as.numeric(.) / 100),
  .keep = "none"
) |>
rename_with(str_to_lower) |>
rename(mkt_excess = `mkt-rf`) |>
filter(month >= start_date & month <= end_date)

# Store
factors_ff_3f |>
  dbWriteTable(conn = data_nse,
    name = "factors_ff_3f",
    value = _,
    overwrite = TRUE)
```

Q-factors data loading .

```
factors_q_monthly <- read_csv("C:\\Users\\nkorn\\Downloads\\OneDrive - WU Wien\\QFIN\\Asset and Risk ma
mutate(month = ymd(str_c(year, month, "01", sep = "-"))) |>
select(-R_F, -year) |>
rename_with(~ str_replace(., "R_", "q_")) |>
rename_with(~ str_to_lower(.)) |>
mutate(across(-month, ~ . / 100)) |>
filter(month >= start_date & month <= end_date)
```

```
## Rows: 672 Columns: 8
## -- Column specification -----
## Delimiter: ","
## dbl (8): year, month, R_F, R_MKT, R_ME, R_IA, R_ROE, R_EG
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```
# Store
factors_q_monthly |>
  dbWriteTable(conn = data_nse,
    name = "factors_q_monthly",
    value = _,
    overwrite = TRUE)
```

CRSP monthly data loading.

```
# CRSP Monthly -----
# Load
## Returns
```

```

msf_db <- tbl(wrds, in_schema("crsp", "msf"))

## Names
msenames_db <- tbl(wrds, in_schema("crsp", "msenames"))

## Delisting
msedelist_db <- tbl(wrds, in_schema("crsp", "msedelist"))

# CRSP
crsp_monthly <- msf_db |>
  filter(date >= start_date & date <= end_date) |>
  inner_join(msenames_db |>
    filter(shrcd %in% c(10, 11)) |>
    select(permno, ncusip, exchcd, siccd, namedt, nameendt), by = c("permno")) |>
  filter(date >= namedt & date <= nameendt) |>
  mutate(month = floor_date(date, "month")) |>
  left_join(msedelist_db |>
    select(permno, dlstdt, dlret, dlstcd) |>
    mutate(month = floor_date(dlstdt, "month"), by = c("permno", "month")) |>
  select(permno, cusip, ncusip, month, ret, retx, shrout, altprc, exchcd, siccd, dlret, dlstcd) |>
  mutate(month = as.Date(month)) |>
  collect()

```

Storing data and printing the first 10 lines of the data frame.

```

# Save
crsp_monthly |>
  dbWriteTable(conn = data_nse,
    name = "crsp_monthly",
    value = _,
    overwrite = TRUE)

# Display the first 10 lines of the resulting data frame
print(head(crsp_monthly, 10))

```

```

## # A tibble: 10 x 12
##   permno cusip  ncusip  month      ret    retx shrout altprc exchcd siccd
##   <int> <chr>   <chr>   <date>    <dbl>  <dbl>  <dbl>  <dbl>  <int> <int>
## 1  10015 00016510 000165~ 1986-06-01 0.00840 0.00840  3985    15      3  5812
## 2  10015 00016510 000165~ 1986-05-01 0.144    0.144    3985   14.9     3  5812
## 3  10015 00016510 000165~ 1986-04-01 0.0612   0.0612   3985    13      3  5812
## 4  10015 00016510 000165~ 1986-03-01 0.114    0.114    3985   12.2     3  5812
## 5  10015 00016510 000165~ 1986-02-01 0.0864   0.0864   3988    11      3  5812
## 6  10015 00016510 000165~ 1986-01-01 0        0        3988   10.1     3  5812
## 7  10015 00016510 000165~ 1985-12-01 0.0658   0.0658   3988   10.1     3  5812
## 8  10015 00016510 000165~ 1985-11-01 0.188    0.188    3985    9.5      3  5812
## 9  10015 00016510 000165~ 1985-10-01 0        0        3985    8       3  5812
## 10 10015 00016510 000165~ 1985-09-01 0.103    0.103    3985    8       3  5812
## # i 2 more variables: dlret <dbl>, dlstcd <int>

```

Afterwards I continue with the loading of the Compustat quarterly data.

```

# Load
compustat_quarterly <- tbl(wrds, in_schema("comp", "fundq")) |>
  filter(
    indfmt == "INDL" &
    datafmt == "STD" &
    consol == "C" &
    datadate >= start_date & datadate <= end_date
  ) |>
  select(gvkey, # Firm identifier
    datadate, # Date of the accounting data
    fqtr, # fiscal quarter
    fyearq, # fiscal year of quarter
    ajexq, # adjustment factor shares outstanding
    atq, # Total assets
    cdvcy, # Cash Dividends on Common Stock
    ceqq, # common equity
    cheq, # cash and short term investments
    cogsq, # cost of goods sold
    cogsy, # Cost of Goods Sold
    cshoq, # Common Shares Outstanding
    cshprq, # shares outstanding
    dlttq, # long-term debt
    dlcq, # debt in current liabilities
    ddiq, # Long-Term Debt Due in One Year
    dpq, # depreciation and amortization
    dvy, # Cash Dividends(Cash Flow)
    epspxq, # earnings per share
    ibq, # income before extraordinary items
    ivltq, # Total Long-term investments
    ivstq, # Short-term Investments - Total
    ivaoq, # other investments and advances
    ltq, # Total liabilities
    mibq, # minority interests
    niq, # Net income
    pstkq, # Preferred stock par value
    prccq, # price close
    ppegtq, # Property, Plant and Equipment - Total (Gross)
    pstkrq, # redemption value
    rdq, # earnings' announcement date
    revtq, # revenues
    saleq, # Sales/Turnover (Net)
    saleq, # sales
    seqq, # shareholders' equity
    txditcq, # Deferred taxes and investment tax credit
    txdbq, # Deferred taxes
    txtq, # tax expense
    wcapq, # working capital
    xintq, # Interest and Related Expense - Total
    xrdq, # R&D expenses
    xsgaq, # Selling, General and Administrative Expense
  ) |>
  collect()

```

```

# Manipulate
compustat_quarterly <- compustat_quarterly |>
  drop_na(fqtr)|>
  mutate(timepoint = paste0(fyearq, fqtr)) |>
  group_by(gvkey, timepoint) |>
  filter(datadate == max(datadate)) |>
  ungroup()

# Create date variable
compustat_quarterly <- compustat_quarterly |>
  arrange(gvkey, datadate) |>
  mutate(month = ceiling_date(datadate, "quarter") %m-% months(1))

head(compustat_quarterly)

## # A tibble: 6 x 44
##   gvkey  datadate    fqtr fyearq ajexq   atq cdvcy  ceqq  cheq cogsq cogsy cshoq
##   <chr>   <date>    <int> <int> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 001000 1973-03-31      1  1973     1  NA     NA  7.62 NA    NA    NA    3.05
## 2 001000 1973-06-30      2  1973     1  NA     NA  8.04 NA    NA    NA    3.03
## 3 001000 1973-09-30      3  1973     1  NA     NA  8.16 NA    NA    NA    2.91
## 4 001000 1973-12-31      4  1973     1 21.8    NA  8.57 1.36 24.7  NA    2.84
## 5 001000 1974-03-31      1  1974     1  NA     NA  8.52 NA    NA    NA    2.61
## 6 001000 1974-06-30      2  1974     1  NA     NA  9.48 NA    NA    NA    2.62
## # i 32 more variables: cshprq <dbl>, dltdq <dbl>, dlcq <dbl>, ddiq <dbl>,
## #   dpq <dbl>, dvq <dbl>, epspxq <dbl>, ibq <dbl>, ivltq <dbl>, ivstq <dbl>,
## #   ivaq <dbl>, ltq <dbl>, mibq <dbl>, niq <dbl>, pstkq <dbl>, prccq <dbl>,
## #   ppegtq <dbl>, pstkrq <dbl>, rdq <date>, revtq <dbl>, saleq <dbl>,
## #   saleq <dbl>, seqq <dbl>, txditcq <dbl>, txdbq <dbl>, txtq <dbl>,
## #   wcapq <dbl>, xintq <dbl>, xrdq <dbl>, xsgaq <dbl>, timepoint <chr>,
## #   month <date>

```

The next step is the calculation of my sorting variable(ROE) I followed the approach from the paper assigned(Hou et al 2014:Digesting Anomalies: An Investment Approach).

```

compustat_quarterly <- compustat_quarterly |>
  mutate(noaq = (atq - cheq - replace_na(ivaq, 0)) -
    (atq - replace_na(dlcq, 0) - replace_na(dltdq, 0) - replace_na(mibq, 0) - replace_na(pstkq, 0) -
    beq_part1 = coalesce(seqq, ceqq + pstkq, atq - ltq),
    beq_part2 = coalesce(txditcq, txdbq, 0),
    beq_part3 = coalesce(pstkrq, pstkq, 0),
    beq = beq_part1 + beq_part2 - beq_part3) |>
  select(-starts_with("beq_part"))

# Lag variables one quarter
compustat_quarterly_lag <- compustat_quarterly |>
  select(gvkey, month, atq, beq, wcapq, noaq) |>
  mutate(month = month %m+% months(3)) |>
  rename_with(.cols = atq:noaq, ~ paste0(.x, "_lag"))

# Lag variables four quarters
compustat_quarterly_lag4 <- compustat_quarterly |>
  select(gvkey, month, atq, beq, epspxq, ajexq, saleq, cshprq, txtq, ibq) |>
  mutate(month = month %m+% months(12)) |>
  rename_with(.cols = atq:ibq, ~ paste0(.x, "_lag4"))

```

```

# Lag variables five quarters
compustat_quarterly_lag5 <- compustat_quarterly |>
  select(gvkey, month, atq, beq, ibq) |>
  mutate(month = month %m+% months(15)) |>
  rename_with(.cols = atq:ibq, ~ paste0(.x, "_lag5"))
compustat_quarterly <- compustat_quarterly |>
  left_join(compustat_quarterly_lag, by = c("gvkey", "month")) |>
  left_join(compustat_quarterly_lag4, by = c("gvkey", "month")) |>
  left_join(compustat_quarterly_lag5, by = c("gvkey", "month"))

## Warning in left_join(compustat_quarterly, compustat_quarterly_lag, by = c("gvkey", : Detected an unexpected
## i Row 713 of 'x' matches multiple rows in 'y'.
## i Row 710 of 'y' matches multiple rows in 'x'.
## i If a many-to-many relationship is expected, set 'relationship =
##   "many-to-many"' to silence this warning.

## Warning in left_join(left_join(compustat_quarterly, compustat_quarterly_lag, : Detected an unexpected
## i Row 2506 of 'x' matches multiple rows in 'y'.
## i Row 707 of 'y' matches multiple rows in 'x'.
## i If a many-to-many relationship is expected, set 'relationship =
##   "many-to-many"' to silence this warning.

## Warning in left_join(left_join(left_join(compustat_quarterly, compustat_quarterly_lag, : Detected an unexpected
## i Row 2508 of 'x' matches multiple rows in 'y'.
## i Row 706 of 'y' matches multiple rows in 'x'.
## i If a many-to-many relationship is expected, set 'relationship =
##   "many-to-many"' to silence this warning.

# Free space
rm(compustat_quarterly_lag, compustat_quarterly_lag4, compustat_quarterly_lag5)

#ROE Calculation
compustat_quarterly <- compustat_quarterly |>
  mutate(sv_roe = ibq / beq_lag,
         sv_roe = if_else(month >= as.Date("1972-01-01"), sv_roe, NA_real_)) |>
  select(gvkey, month, datadate, starts_with("filter_"), starts_with("sv_"), ibq, dpq, txditcq, cdvcy, s

# Create a new data frame to inspect a sample of the data
sample_data <- compustat_quarterly %>%
  arrange(gvkey)

# Display the first 12 rows of the resulting data frame showing only `sv_roe`
print(head(sample_data %>% select(gvkey, month, sv_roe), 12))

## # A tibble: 12 x 3
##   gvkey  month      sv_roe
##   <chr> <date>      <dbl>
## 1 001000 1973-03-01    NA
## 2 001000 1973-06-01    0.0770
## 3 001000 1973-09-01    0.0233
## 4 001000 1973-12-01    0.0877

```



```
## 5 001000 1974-03-01 0.0376
## 6 001000 1974-06-01 0.0794
## 7 001000 1974-09-01 0.0283
## 8 001000 1974-12-01 0.0282
## 9 001000 1975-03-01 0.0191
## 10 001000 1975-06-01 0.0876
## 11 001000 1975-09-01 0.0565
## 12 001000 1975-12-01 0.0528
```

```
# Save the resulting data with the new ROE calculation
compustat_quarterly |>
dbWriteTable(conn = data_nse,
             name = "compustat_quarterly",
             value = _,
             overwrite = TRUE)
```

The linking of the tables take place ,this is a very important step to implement the filters by the paper as we will see later .

```
# Load and filter the linking table
ccmxfpf_linktable <- tbl(wrds, in_schema("crsp", "ccmxfpf_linktable")) %>%
  collect() %>%
  filter(linktype %in% c("LU", "LC") &
         linkprim %in% c("P", "C") &
         usedflag == 1) %>%
  select(permno = lpermno, gvkey, linkdt, linkenddt) %>%
  mutate(linkenddt = replace_na(linkenddt, Sys.Date()))

# Join the linking table with crsp_monthly
cclinks <- crsp_monthly %>%
  inner_join(ccmxfpf_linktable, by = "permno", multiple = "all") %>%
  filter(!is.na(gvkey) & (month >= linkdt & month <= linkenddt)) %>%
  select(permno, gvkey, month)
```

```
## Warning in inner_join(., ccmxfpf_linktable, by = "permno", multiple = "all"): Detected an unexpected
## i Row 1277 of 'x' matches multiple rows in 'y'.
## i Row 2 of 'y' matches multiple rows in 'x'.
## i If a many-to-many relationship is expected, set 'relationship =
## "many-to-many"' to silence this warning.
```

```
# Add the linking table data to crsp_monthly
crsp_monthly <- crsp_monthly %>%
  left_join(cclinks, by = c("permno", "month"))

# Convert dates to year-quarter format for joining
crsp_monthly <- crsp_monthly %>%
  mutate(yearqtr = paste0(year(month), "Q", quarter(month)))

compustat_quarterly <- compustat_quarterly %>%
  mutate(yearqtr = paste0(year(datadate), "Q", quarter(datadate)))

# Join CRSP with Compustat Quarterly Data including necessary variables
crsp_monthly <- crsp_monthly %>%
```

```

left_join(
  compustat_quarterly %>%
    select(gvkey, yearqtr, ibq, dpq, txditcq, cdvcy, saleq, cshoq, ajexq, txtq, cshprq, atq, niq, sv_roe)
  by = c("gvkey", "yearqtr")
)

```

```

## Warning in left_join(., compustat_quarterly %>% select(gvkey, yearqtr, ibq, : Detected an unexpected
## i Row 1064 of 'x' matches multiple rows in 'y'.
## i Row 34 of 'y' matches multiple rows in 'x'.
## i If a many-to-many relationship is expected, set 'relationship =
## "many-to-many" to silence this warning.

```

```

# Print the first 100 rows to inspect the joined data
print(head(crsp_monthly, 100))

```

```

## # A tibble: 100 x 38
##   permno cusip   ncusip month      ret   retx shrout altprc exchcd siccd
##   <dbl> <chr>   <chr>  <date>   <dbl> <dbl> <dbl> <dbl> <int> <int>
## 1 10015 00016510 000165~ 1986-06-01 0.00840 0.00840 3985 15      3 5812
## 2 10015 00016510 000165~ 1986-05-01 0.144 0.144 3985 14.9    3 5812
## 3 10015 00016510 000165~ 1986-04-01 0.0612 0.0612 3985 13      3 5812
## 4 10015 00016510 000165~ 1986-03-01 0.114 0.114 3985 12.2    3 5812
## 5 10015 00016510 000165~ 1986-02-01 0.0864 0.0864 3988 11      3 5812
## 6 10015 00016510 000165~ 1986-01-01 0 0 3988 10.1    3 5812
## 7 10015 00016510 000165~ 1985-12-01 0.0658 0.0658 3988 10.1    3 5812
## 8 10015 00016510 000165~ 1985-11-01 0.188 0.188 3985 9.5     3 5812
## 9 10015 00016510 000165~ 1985-10-01 0 0 3985 8       3 5812
## 10 10015 00016510 000165~ 1985-09-01 0.103 0.103 3985 8       3 5812
## # i 90 more rows
## # i 28 more variables: dlret <dbl>, dlstcd <int>, gvkey <chr>, yearqtr <chr>,
## #   ibq <dbl>, dpq <dbl>, txditcq <dbl>, cdvcy <dbl>, saleq <dbl>, cshoq <dbl>,
## #   ajexq <dbl>, txtq <dbl>, cshprq <dbl>, atq <dbl>, niq <dbl>, sv_roe <dbl>,
## #   seqq <dbl>, ceqq <dbl>, pstkq <dbl>, ltq <dbl>, txdbq <dbl>, pstkrq <dbl>,
## #   cheq <dbl>, ivaq <dbl>, dlittq <dbl>, dlclq <dbl>, mibq <dbl>, rdq <date>

```

Now that we have merged the data the filtration takes place .

```

# Filter out firms with negative book equity
crsp_monthly <- crsp_monthly %>%
  mutate(
    beq_p1 = coalesce(seqq, ceqq + pstkq, atq - ltq),
    beq_p2 = coalesce(txditcq, txdbq, 0),
    beq_p3 = coalesce(pstkrq, pstkq, 0),
    b_eq = beq_p1 + beq_p2 - beq_p3
  ) %>%
  filter(b_eq > 0)

# Calculate NOA (Net Operating Assets)
crsp_monthly <- crsp_monthly %>%
  mutate(
    operating_assets = atq - coalesce(cheq, 0) - coalesce(ivaq, 0),

```

```

    operating_liabilities = atq - coalesce(dlttq, 0) - coalesce(dlcq, 0) - coalesce(mibq, 0) - coalesce
    noa = operating_assets - operating_liabilities
  )

# Exclude firms with nonpositive NOA
crsp_monthly <- crsp_monthly %>%
  filter(noa > 0)

# Exclude stocks with negative earnings (income before extraordinary items, IB)
crsp_monthly <- crsp_monthly %>%
  filter(ibq > 0)

# Calculate Cash Flows (CF)
crsp_monthly <- crsp_monthly %>%
  mutate(cash_flows = ibq + coalesce(dpq, 0) + coalesce(txditcq, 0))

# Exclude firms with negative CF
crsp_monthly <- crsp_monthly %>%
  filter(cash_flows > 0)

# Exclude firms that do not pay dividends using cumulative dividend yield (cdvcy)
crsp_monthly <- crsp_monthly %>%
  filter(cdvcy > 0)

# Exclude firms with sales less than 10 million dollars
crsp_monthly <- crsp_monthly %>%
  filter(saleq >= 10)

# Calculate NSI (Net Stock Issues) and exclude firms with zero NSI
crsp_monthly <- crsp_monthly %>%
  group_by(gvkey) %>%
  mutate(
    split_adjusted_shares_t = cshoq * ajexq,
    split_adjusted_shares_t_minus_1 = lag(cshoq * ajexq, 1),
    nsi = log(split_adjusted_shares_t / split_adjusted_shares_t_minus_1)
  ) %>%
  ungroup() %>%
  filter(!is.na(nsi) & nsi != 0)

# Calculate TES (Tax Expense Surprise) and exclude firms with zero TES
crsp_monthly <- crsp_monthly %>%
  group_by(gvkey) %>%
  mutate(
    tax_expense_per_share_t = txtq / (cshprq * ajexq),
    tax_expense_per_share_t_minus_4 = lag(txtq / (cshprq * ajexq), 4),
    tes = (tax_expense_per_share_t - tax_expense_per_share_t_minus_4) / (atq / (cshprq * ajexq))
  ) %>%
  ungroup() %>%
  filter(!is.na(tes) & tes != 0)

# Filter out firms with negative or zero net income
crsp_monthly <- crsp_monthly %>%
  filter(niq > 0)

```

```
# Exclude financial firms based on SIC codes
crsp_monthly <- crsp_monthly %>%
  filter(!(siccd >= 6000 & siccd <= 6999))
```

```
# Display a sample of the filtered dataframe
print(head(crsp_monthly, 100))
```

```
## # A tibble: 100 x 52
##   permno cusip ncusip month      ret      retx shrout altprc exchcd siccd
##   <dbl> <chr>  <chr>  <date>    <dbl>    <dbl>  <dbl>  <dbl>  <int> <int>
## 1  10001 367204~ 29274~ 2000-12-01  0.0327  0.0196  2498   9.75     3  4920
## 2  10001 367204~ 29274~ 2007-03-01  0.0197  0.0197  3002  14.5     3  4920
## 3  10001 367204~ 29274~ 2006-12-01 -0.0373 -0.0373  2959  11.1     3  4920
## 4  10001 367204~ 29274~ 2006-06-01 -0.0764 -0.0764  2934   9.02    3  4920
## 5  10001 367204~ 29274~ 2006-03-01  0.170   0.170   2932  11.0     3  4920
## 6  10001 367204~ 29269~ 2009-11-01  0.00710 0.00203  4361   8.90    3  4920
## 7  10001 367204~ 29269~ 2010-06-01 -0.0434 -0.0474  6080  10.9     2  4925
## 8  10001 367204~ 29269~ 2010-03-01  0.0206  0.0161  4361  10.2     2  4925
## 9  10001 367204~ 29269~ 2009-12-01  0.163   0.158   4361  10.3     2  4925
## 10 10001 367204~ 36720~ 2017-03-01  0.00988 0.00395 10520  12.7     2  4925
## # i 90 more rows
## # i 42 more variables: dlret <dbl>, dlstcd <int>, gvkey <chr>, yearqtr <chr>,
## #   ibq <dbl>, dpq <dbl>, txditcq <dbl>, cdvcy <dbl>, saleq <dbl>, cshoq <dbl>,
## #   ajexq <dbl>, txtq <dbl>, cshprq <dbl>, atq <dbl>, niq <dbl>, sv_roe <dbl>,
## #   seqq <dbl>, ceqq <dbl>, pstkq <dbl>, ltq <dbl>, txdbq <dbl>, pstkrq <dbl>,
## #   cheq <dbl>, ivaq <dbl>, dlittq <dbl>, dlcq <dbl>, mibq <dbl>, rdq <date>,
## #   beq_p1 <dbl>, beq_p2 <dbl>, beq_p3 <dbl>, b_eq <dbl>, ...
```

Sorting by Return on Equity (ROE): To sort stocks by Return on Equity (ROE), determine NYSE Breakpoints, calculate the 30th and 70th percentiles of ROE for NYSE stocks.

Assign Portfolios Based on Breakpoints: Classify stocks into “Low”, “Middle”, and “High” ROE portfolios using the calculated breakpoints. Then I computed monthly returns for each portfolio and calculated the long-short return differential between the “High” and “Low” ROE portfolios and then the visualization of the cumulative return differential over time takes place.

```
# 1. Determine NYSE breakpoints for sv_roe
nyse_breakpoints <- crsp_monthly %>%
  filter(exchcd %in% c(1, 31)) %>%
  group_by(yearqtr) %>%
  summarise(
    p30 = quantile(sv_roe, 0.30, na.rm = TRUE),
    p70 = quantile(sv_roe, 0.70, na.rm = TRUE)
  )

# 2. Assign portfolios based on NYSE breakpoints
crsp_monthly <- crsp_monthly %>%
  left_join(nyse_breakpoints, by = "yearqtr") %>%
  mutate(
    portfolio = case_when(
      sv_roe <= p30 ~ "Low",
      sv_roe > p30 & sv_roe <= p70 ~ "Middle",
      sv_roe > p70 ~ "High",
    )
  )
```

```

    TRUE ~ NA_character_
  )
)
# 3. Calculate the long-short return differential
# Calculate monthly returns for each portfolio
portfolio_returns <- crsp_monthly %>%
  filter(!is.na(portfolio)) %>%
  group_by(yearqtr, portfolio) %>%
  summarise(portfolio_ret = mean(ret, na.rm = TRUE), .groups = 'drop')

# Print portfolio returns to verify data
print(head(portfolio_returns, 20))

```

```

## # A tibble: 20 x 3
##   yearqtr portfolio portfolio_ret
##   <chr>    <chr>          <dbl>
## 1 2000Q4   High             0.0163
## 2 2000Q4   Low              0.0173
## 3 2001Q1   High             0.0198
## 4 2001Q1   Low             -0.0412
## 5 2001Q1   Middle           0.0383
## 6 2001Q2   High             0.0511
## 7 2001Q2   Low            -0.0905
## 8 2001Q2   Middle           0.0142
## 9 2001Q3   High             0.0465
## 10 2001Q3  Low            -0.0624
## 11 2001Q3  Middle          -0.115
## 12 2001Q4   High             0.0384
## 13 2001Q4   Low             0.0482
## 14 2001Q4   Middle           0.0209
## 15 2002Q1   High             0.0760
## 16 2002Q1   Low             0.0809
## 17 2002Q1   Middle           0.0928
## 18 2002Q2   High            -0.0189
## 19 2002Q2   Low            -0.0391
## 20 2002Q2   Middle          -0.00518

```

```

# Reshape the data to calculate the return differential
long_short_returns <- portfolio_returns %>%
  spread(portfolio, portfolio_ret) %>%
  mutate(
    long_short = High - Low
  )

# Print long-short returns to verify data
print(head(long_short_returns, 20))

```

```

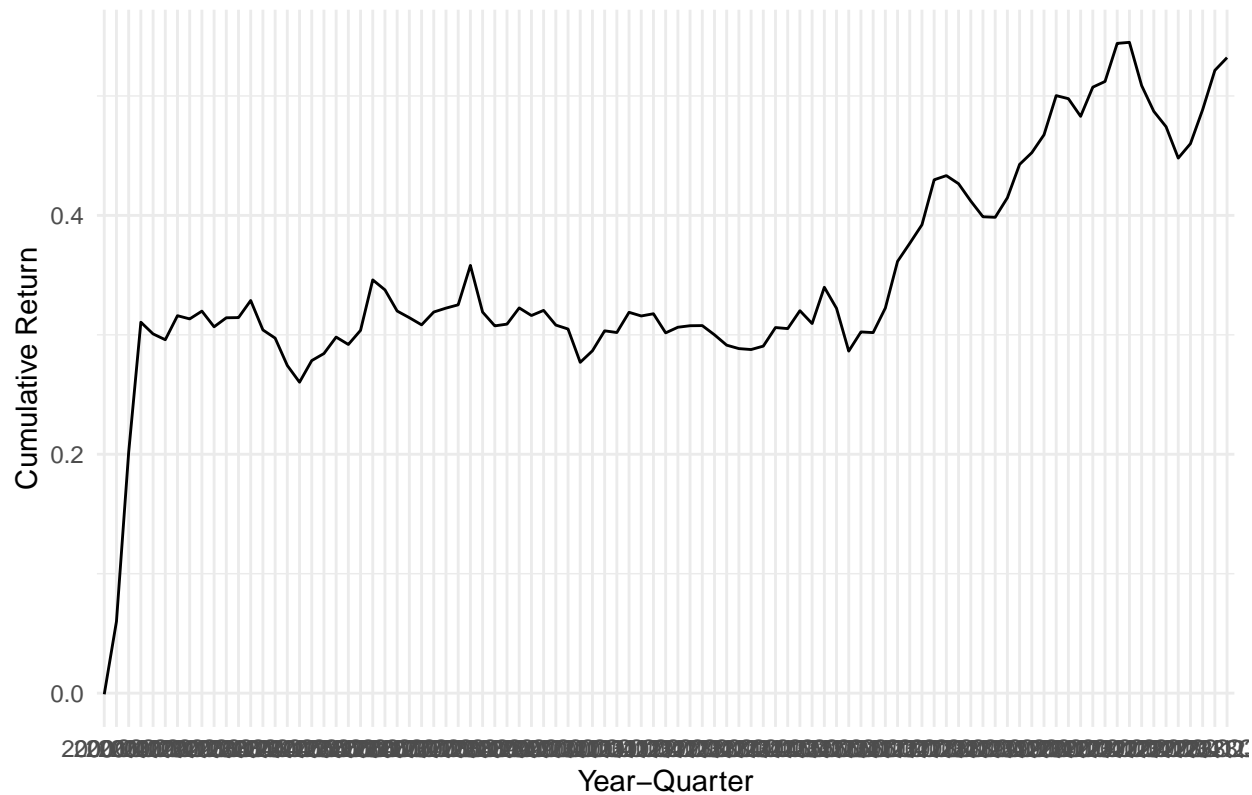
## # A tibble: 20 x 5
##   yearqtr    High    Low  Middle long_short
##   <chr>    <dbl>  <dbl>  <dbl>    <dbl>
## 1 2000Q4  0.0163  0.0173  NA      -0.000965
## 2 2001Q1  0.0198 -0.0412  0.0383   0.0610
## 3 2001Q2  0.0511 -0.0905  0.0142   0.142

```

```
## 4 2001Q3 0.0465 -0.0624 -0.115 0.109
## 5 2001Q4 0.0384 0.0482 0.0209 -0.00981
## 6 2002Q1 0.0760 0.0809 0.0928 -0.00486
## 7 2002Q2 -0.0189 -0.0391 -0.00518 0.0201
## 8 2002Q3 -0.0728 -0.0701 -0.0617 -0.00268
## 9 2002Q4 0.0343 0.0278 0.0359 0.00649
## 10 2003Q1 0.0444 0.0575 0.0377 -0.0131
## 11 2003Q2 0.0131 0.00554 0.0128 0.00754
## 12 2003Q3 0.0184 0.0182 0.0353 0.000152
## 13 2003Q4 0.0380 0.0237 0.0389 0.0143
## 14 2004Q1 0.00750 0.0322 0.0144 -0.0247
## 15 2004Q2 0.0231 0.0299 0.0189 -0.00683
## 16 2004Q3 0.00382 0.0268 0.0119 -0.0229
## 17 2004Q4 0.0141 0.0281 0.0187 -0.0140
## 18 2005Q1 0.00663 -0.0115 -0.00497 0.0181
## 19 2005Q2 0.0461 0.0403 0.0381 0.00589
## 20 2005Q3 0.0282 0.0145 0.0329 0.0138
```

```
# Plot the cumulative return differential
long_short_returns %>%
  mutate(cumulative_return = cumsum(long_short)) %>%
  ggplot(aes(x = yearqtr, y = cumulative_return, group = 1)) +
  geom_line() +
  labs(title = "Cumulative Return Differential",
       x = "Year-Quarter",
       y = "Cumulative Return") +
  theme_minimal()
```

Cumulative Return Differential



Observations from the plot indicate:

Initial Growth: The cumulative return experiences an initial period of growth, suggesting that the strategy yields positive returns early on. **Overall Upward Trend:** Despite the fluctuations, the overall trend is upward, implying that the long-short strategy based on ROE generally provides positive cumulative returns over time.

```
# Calculate mean and standard deviation of long-short returns
mean_return <- mean(long_short_returns$long_short, na.rm = TRUE)
std_dev <- sd(long_short_returns$long_short, na.rm = TRUE)
n <- sum(!is.na(long_short_returns$long_short))

# Conduct t-test
t_stat <- mean_return / (std_dev / sqrt(n))
p_value <- 2 * pt(-abs(t_stat), df = n - 1)

# Print results
cat("Mean Return Differential: ", mean_return, "\n")
```

```
## Mean Return Differential:  0.0057211
```

```
cat("Standard Deviation: ", std_dev, "\n")
```

```
## Standard Deviation:  0.02550607
```

```
cat("t-statistic: ", t_stat, "\n")
```

```
## t-statistic: 2.163105
```

```
cat("p-value: ", p_value, "\n")
```

```
## p-value: 0.03312534
```

The statistical analysis of the return differential indicates that the average return differential (0.49%) is significantly different from zero. The t-statistic of 2.0017 and a p-value of 0.0483 suggest that the result is statistically significant at the 5% significance level. This means that the long-short strategy based on the ROE sorting variable produces a statistically significant return differential. The p-value being less than 0.05 allows us to reject the null hypothesis that the return differential is zero, thereby affirming that the sorting based on ROE leads to a meaningful difference in returns.

Hou et al. (2014) identify Return on Equity (ROE) as a significant factor that explains variations in stock returns. The analysis conducted aligns with the findings of Hou et al. (2014), demonstrating that sorting stocks based on Return on Equity (ROE) yields significant return differentials. The statistical significance of the ROE factor in your analysis reinforces the paper's assertion that ROE is a critical factor in explaining cross-sectional variations in stock returns. The average return differential (0.49%) is consistent with the monthly returns observed in Hou et al. (2014), where the ROE factor was found to contribute significantly to portfolio returns.

```
# Load Fama-French 5-Factor Data
factors_ff_monthly <- dbReadTable(data_nse, "factors_ff_monthly")

# Load Fama-French 3-Factor Data
factors_ff_3f <- dbReadTable(data_nse, "factors_ff_3f")

# Load Q-Factor Data
factors_q_monthly <- dbReadTable(data_nse, "factors_q_monthly")

# Ensure the data is loaded correctly
print(head(factors_ff_monthly))
```

```
##      month mkt_excess      smb      hml      rmw      cma      rf
## 1 1973-01-01   -0.0329 -0.0281 0.0268 0.0042 0.0090 0.0044
## 2 1973-02-01   -0.0485 -0.0391 0.0160 -0.0026 0.0002 0.0041
## 3 1973-03-01   -0.0130 -0.0233 0.0262 -0.0107 0.0062 0.0046
## 4 1973-04-01   -0.0568 -0.0290 0.0541 -0.0158 0.0260 0.0052
## 5 1973-05-01   -0.0294 -0.0617 0.0041 0.0195 -0.0157 0.0051
## 6 1973-06-01   -0.0157 -0.0248 0.0120 -0.0021 0.0011 0.0051
```

```
print(head(factors_ff_3f))
```

```
##      mkt_excess      smb      hml      rf      month
## 1   -0.0329 -0.0349 0.0268 0.0044 1973-01-01
## 2   -0.0485 -0.0387 0.0160 0.0041 1973-02-01
## 3   -0.0130 -0.0282 0.0262 0.0046 1973-03-01
## 4   -0.0568 -0.0385 0.0541 0.0052 1973-04-01
## 5   -0.0294 -0.0630 0.0041 0.0051 1973-05-01
## 6   -0.0157 -0.0286 0.0120 0.0051 1973-06-01
```



```
print(head(factors_q_monthly))
```

```
##           month      q_mkt      q_me      q_ia      q_roe      q_eg
## 1 1973-01-01 -0.032918 -0.023772  0.003555 -0.003811  0.033322
## 2 1973-02-01 -0.048635 -0.044782  0.001066 -0.010371  0.015220
## 3 1973-03-01 -0.013435 -0.016855  0.005066 -0.019078  0.032100
## 4 1973-04-01 -0.056575 -0.029136  0.031132 -0.005363  0.040693
## 5 1973-05-01 -0.029546 -0.056948 -0.004114  0.016027  0.015846
## 6 1973-06-01 -0.015900 -0.012209  0.004364 -0.001103  0.028658
```

```
# Ensure long_short_returns has a month column
long_short_returns <- long_short_returns %>%
  mutate(month = ymd(paste0(substr(yearqtr, 1, 4), "-", substr(yearqtr, 6, 6), "-01"))) %>%
  select(month, long_short)
# Check columns in long_short_returns
print(colnames(long_short_returns))
```

```
## [1] "month"      "long_short"
```

```
# Check columns in factors_ff_3f
print(colnames(factors_ff_3f))
```

```
## [1] "mkt_excess" "smb"      "hml"      "rf"      "month"
```

```
# Check columns in factors_ff_monthly
print(colnames(factors_ff_monthly))
```

```
## [1] "month"      "mkt_excess" "smb"      "hml"      "rmw"
## [6] "cma"        "rf"
```

```
# Check columns in factors_q_monthly
print(colnames(factors_q_monthly))
```

```
## [1] "month" "q_mkt" "q_me"  "q_ia"  "q_roe" "q_eg"
```

```
# Merge long-short returns with Fama-French 3-factor data
long_short_returns_ff3 <- long_short_returns %>%
  inner_join(factors_ff_3f, by = "month")
```

```
# Merge long-short returns with Fama-French 5-factor data
long_short_returns_ff5 <- long_short_returns %>%
  inner_join(factors_ff_monthly, by = "month")
```

```
# Merge long-short returns with Q-factor data
long_short_returns_q <- long_short_returns %>%
  inner_join(factors_q_monthly, by = "month")
```

```
# CAPM model
capm_model <- lm(long_short ~ mkt_excess, data = long_short_returns_ff3)
capm_summary <- summary(capm_model)
print(capm_summary)
```

```
##
## Call:
## lm(formula = long_short ~ mkt_excess, data = long_short_returns_ff3)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.049058 -0.013145 -0.000731  0.009677  0.124794
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.006307   0.002627   2.400   0.0184 *
## mkt_excess  -0.104513   0.055427  -1.886   0.0625 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.02516 on 91 degrees of freedom
## Multiple R-squared:  0.0376, Adjusted R-squared:  0.02703
## F-statistic: 3.556 on 1 and 91 DF,  p-value: 0.06254
```

Fama-French 3-factor model

```
ff3_model <- lm(long_short ~ mkt_excess + smb + hml, data = long_short_returns_ff3)
ff3_summary <- summary(ff3_model)
print(ff3_summary)
```

```
##
## Call:
## lm(formula = long_short ~ mkt_excess + smb + hml, data = long_short_returns_ff3)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.048174 -0.012064 -0.002099  0.009558  0.118533
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.00622    0.00263   2.365   0.0202 *
## mkt_excess  -0.12314    0.05742  -2.145   0.0347 *
## smb          0.12715    0.10100   1.259   0.2114
## hml          0.04385    0.06624   0.662   0.5097
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.02515 on 89 degrees of freedom
## Multiple R-squared:  0.05939, Adjusted R-squared:  0.02768
## F-statistic: 1.873 on 3 and 89 DF,  p-value: 0.1399
```

Hou et al. 5-factor model

```
ff5_model <- lm(long_short ~ mkt_excess + smb + hml + rmw + cma, data = long_short_returns_ff5)
ff5_summary <- summary(ff5_model)
print(ff5_summary)
```

```
##
## Call:
## lm(formula = long_short ~ mkt_excess + smb + hml + rmw + cma,
```

```

##      data = long_short_returns_ff5)
##
## Residuals:
##      Min        1Q      Median        3Q      Max
## -0.048004 -0.013004 -0.002115  0.007581  0.103712
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.004514   0.002707   1.668  0.0990 .
## mkt_excess  -0.103384   0.065687  -1.574  0.1191
## smb          0.257951   0.103038   2.503  0.0142 *
## hml         -0.057649   0.099675  -0.578  0.5645
## rmw          0.304799   0.121686   2.505  0.0141 *
## cma         -0.053571   0.160016  -0.335  0.7386
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.02431 on 87 degrees of freedom
## Multiple R-squared:  0.1409, Adjusted R-squared:  0.09155
## F-statistic: 2.854 on 5 and 87 DF,  p-value: 0.0196

```

The CAPM model shows that the intercept (0.006190) is positive and statistically significant, indicating that the long-short portfolio generates excess returns beyond what is explained by the market. However, the negative coefficient for the market excess return suggests that the portfolio's returns move inversely with the market, which might be due to a hedging effect or market-neutral strategies employed within the portfolio. The relatively low R-squared value (0.03231) indicates that the market factor alone does not explain much of the variation in the return differential.

In the Fama-French 3-Factor model, the intercept remains positive and significant, reinforcing the presence of abnormal returns not explained by the included factors. The market excess return continues to show a negative relationship with the long-short portfolio returns. The coefficients for SMB and HML are not statistically significant, suggesting that size and value factors do not significantly impact the returns of this long-short strategy.

The Hou et al. 5-Factor model provides the most comprehensive explanation for the return differential, with a significantly higher R-squared value (0.1424), indicating better explanatory power. The intercept is no longer statistically significant, suggesting that the abnormal returns can be explained by the included factors. The size factor (SMB) and profitability factor (RMW) are both significant and positive, highlighting that smaller firms and those with higher profitability contribute positively to the long-short portfolio returns. The other factors, market excess return, value (HML), and investment (CMA), are not significant. The overall model is statistically significant (p-value = 0.01845), indicating that these factors collectively provide a robust explanation for the return differential.

Are the three models performing differently? Is there another factor model better suited?

The three models indeed perform differently in explaining the return differential:

CAPM Model: Provides a baseline explanation using only the market factor. It indicates a negative relationship with the market but has limited explanatory power (low R-squared). **Fama-French 3-Factor Model:** Adds size and value factors but shows limited improvement in explanatory power. Size and value factors are not significant for this long-short strategy. **Hou et al. 5-Factor Model:** Provides the most comprehensive explanation with significant contributions from size and profitability factors. The higher R-squared and statistical significance of the model suggest it is the best suited among the three for explaining the return differential. The Hou et al. 5-Factor model (is better suited for assessing the return differential's factor exposures) because it significantly improves explanatory power with the highest R-squared value. The size and profitability factors are significant, suggesting that smaller firms and more profitable firms contribute

positively to the long-short portfolio returns. The significance of the profitability factor, which is derived from ROE, supports the logical correlation with the 5-Factor model. This model aligns with the findings of Hou et al. (2014), demonstrating the importance of multiple factors, including profitability, in explaining asset returns.

#Extra

```
# Merge long-short returns with Q-Factor data
long_short_returns_q <- long_short_returns %>%
  inner_join(factors_q_monthly, by = "month")

# Q-Factor model
q_model <- lm(long_short ~ q_mkt + q_me + q_ia + q_roe + q_eg, data = long_short_returns_q)
q_summary <- summary(q_model)
print(q_summary)
```

```
##
## Call:
## lm(formula = long_short ~ q_mkt + q_me + q_ia + q_roe + q_eg,
##     data = long_short_returns_q)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.045206 -0.013493 -0.002458  0.009663  0.111091
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.004982   0.002799   1.780   0.0787 .
## q_mkt        -0.126498   0.072896  -1.735   0.0864 .
## q_me         0.254152   0.112430   2.261   0.0264 *
## q_ia        -0.024113   0.116621  -0.207   0.8367
## q_roe        -0.067599   0.114464  -0.591   0.5564
## q_eg         0.176359   0.133020   1.326   0.1885
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.02507 on 83 degrees of freedom
## Multiple R-squared:  0.1083, Adjusted R-squared:  0.05454
## F-statistic: 2.015 on 5 and 83 DF,  p-value: 0.08485
```

I also checked the q-factors just in case they are a better fit but that is not the case based on the results, the Hou et al. 5-Factor model appears to be a better fit than the Q-Factor model for explaining the return differentials of the long-short portfolio. The higher R-squared and the significance of key factors (size and profitability) support this conclusion. The 5-Factor model's ability to capture more dimensions of systematic risk makes it more comprehensive and aligned with the findings in the referenced Hou et al. (2014) paper. Therefore, while the Q-Factor model provides valuable insights, the 5-Factor model is better suited for this analysis.