Search WRDS     🔍

📊 Get Data ▾    📈 Analytics ▾    🎓 Classroom ▾    💻 Research ▾    🎧 Support ▾

# Fama-French Factors (Python)

Based on the SAS based research application, this Python code replicates Fama and French's (1993) methodology to construct size and value factors.

## Background

This set of Python code is written based on the original SAS code that replicates the Fama French risk factors SMB and HML. Please refer to the original Fama-French Factors page for detailed discussion on methodology. The flow of the code and the dataset naming convention mimics the SAS code for easy cross-reference.

Top of Section

## Technical Requirement

This code makes use of the Python `wrds` module, a freely-available Python package that WRDS has developed that allows for easy connectivity to our data from within your Python program. The `wrds` module is available both on our high-performance computing cluster, the WRDS Cloud, and on your local workstation. Please see the below documentation for instructions on setting up your Python environment from either location to be able to connect to WRDS. Once set up, you can then run the Python code presented on this page.

- PYTHON: On the WRDS Cloud
- PYTHON: From Your Computer (Jupyter/Spyder)

Top of Section

## Fama French Factors Python Sample Code

📊 Get Data ▾    📈 Analytics ▾    🎓 Classroom ▾    💻 Research ▾    🎧 Support ▾

```
  1    ##########################################
```

```python
# Fama French Factors
# April 2018
# Qingyi (Freda) Song Drechsler
#########################################

import pandas as pd
import numpy as np
import datetime as dt
import wrds
import psycopg2
import matplotlib.pyplot as plt
from dateutil.relativedelta import *
from pandas.tseries.offsets import *
from scipy import stats

###################
# Connect to WRDS #
###################
conn=wrds.Connection()

###################
# Compustat Block #
###################
comp = conn.raw_sql("""
                    select gvkey, datadate, at, pstkl, txditc,
                    pstkrv, seq, pstk
                    from comp.funda
                    where indfmt='INDL'
                    and datafmt='STD'
                    and popsrc='D'
                    and consol='C'
                    and datadate >= '01/01/1959'
                    """)

comp['datadate']=pd.to_datetime(comp['datadate']) #convert datadate to date fmt
comp['year']=comp['datadate'].dt.year

# create preferrerd stock
comp['ps']=np.where(comp['pstkrv'].isnull(), comp['pstkl'], comp['pstkrv'])
comp['ps']=np.where(comp['ps'].isnull(),comp['pstk'], comp['ps'])
comp['ps']=np.where(comp['ps'].isnull(),0,comp['ps'])

comp['txditc']=comp['txditc'].fillna(0)

# create book equity
comp['be']=comp['seq']+comp['txditc']-comp['ps']
comp['be']=np.where(comp['be']>0, comp['be'], np.nan)

# number of years in Compustat
comp=comp.sort_values(by=['gvkey','datadate'])
comp['count']=comp.groupby(['gvkey']).cumcount()

comp=comp[['gvkey','datadate','year','be','count']]

###################
# CRSP Block      #
###################
# sql similar to crspmerge macro
crsp_m = conn.raw_sql("""
                    select a.permno, a.permco, a.date, b.shrcd, b.exchcd,
                    a.ret, a.retx, a.shrout, a.prc
                    from crsp.msf as a
                    left join crsp.msenames as b
```

```python
                            on a.permno=b.permno
                            and b.namedt<=a.date
                            and a.date<=b.nameendt
                            where a.date between '01/01/1959' and '12/31/2017'
                            and b.exchcd between 1 and 3
                            """)

# change variable format to int
crsp_m[['permco','permno','shrcd','exchcd']]=crsp_m[['permco','permno','shrcd','exchcd']].ast


# Line up date to be end of month
crsp_m['date']=pd.to_datetime(crsp_m['date'])
crsp_m['jdate']=crsp_m['date']+MonthEnd(0)

# add delisting return
dlret = conn.raw_sql("""
                        select permno, dlret, dlstdt
                        from crsp.msedelist
                        """)
dlret.permno=dlret.permno.astype(int)
dlret['dlstdt']=pd.to_datetime(dlret['dlstdt'])
dlret['jdate']=dlret['dlstdt']+MonthEnd(0)

crsp = pd.merge(crsp_m, dlret, how='left',on=['permno','jdate'])
crsp['dlret']=crsp['dlret'].fillna(0)
crsp['ret']=crsp['ret'].fillna(0)
crsp['retadj']=(1+crsp['ret'])*(1+crsp['dlret'])-1
crsp['me']=crsp['prc'].abs()*crsp['shrout'] # calculate market equity
crsp=crsp.drop(['dlret','dlstdt','prc','shrout'], axis=1)
crsp=crsp.sort_values(by=['jdate','permco','me'])

### Aggregate Market Cap ###
# sum of me across different permno belonging to same permco a given date
crsp_summe = crsp.groupby(['jdate','permco'])['me'].sum().reset_index()
# largest mktcap within a permco/date
crsp_maxme = crsp.groupby(['jdate','permco'])['me'].max().reset_index()
# join by jdate/maxme to find the permno
crsp1=pd.merge(crsp, crsp_maxme, how='inner', on=['jdate','permco','me'])
# drop me column and replace with the sum me
crsp1=crsp1.drop(['me'], axis=1)
# join with sum of me to get the correct market cap info
crsp2=pd.merge(crsp1, crsp_summe, how='inner', on=['jdate','permco'])
# sort by permno and date and also drop duplicates
crsp2=crsp2.sort_values(by=['permno','jdate']).drop_duplicates()

# keep December market cap
crsp2['year']=crsp2['jdate'].dt.year
crsp2['month']=crsp2['jdate'].dt.month
decme=crsp2[crsp2['month']==12]
decme=decme[['permno','date','jdate','me','year']].rename(columns={'me':'dec_me'})

### July to June dates
crsp2['ffdate']=crsp2['jdate']+MonthEnd(-6)
crsp2['ffyear']=crsp2['ffdate'].dt.year
crsp2['ffmonth']=crsp2['ffdate'].dt.month
crsp2['1+retx']=1+crsp2['retx']
crsp2=crsp2.sort_values(by=['permno','date'])

# cumret by stock
crsp2['cumretx']=crsp2.groupby(['permno','ffyear'])['1+retx'].cumprod()
# lag cumret
crsp2['lcumretx']=crsp2.groupby(['permno'])['cumretx'].shift(1)
```

```python
128
129    # lag market cap
130    crsp2['lme']=crsp2.groupby(['permno'])['me'].shift(1)
131
132    # if first permno then use me/(1+retx) to replace the missing value
133    crsp2['count']=crsp2.groupby(['permno']).cumcount()
134    crsp2['lme']=np.where(crsp2['count']==0, crsp2['me']/crsp2['1+retx'], crsp2['lme'])
135
136    # baseline me
137    mebase=crsp2[crsp2['ffmonth']==1][['permno','ffyear', 'lme']].rename(columns=
138    {'lme':'mebase'})
139
140    # merge result back together
141    crsp3=pd.merge(crsp2, mebase, how='left', on=['permno','ffyear'])
142    crsp3['wt']=np.where(crsp3['ffmonth']==1, crsp3['lme'], crsp3['mebase']*crsp3['lcumretx'])
143
144    decme['year']=decme['year']+1
145    decme=decme[['permno','year','dec_me']]
146
147    # Info as of June
148    crsp3_jun = crsp3[crsp3['month']==6]
149
150    crsp_jun = pd.merge(crsp3_jun, decme, how='inner', on=['permno','year'])
151    crsp_jun=crsp_jun[['permno','date', 'jdate',
152    'shrcd','exchcd','retadj','me','wt','cumretx','mebase','lme','dec_me']]
153    crsp_jun=crsp_jun.sort_values(by=['permno','jdate']).drop_duplicates()
154
155    #######################
156    # CCM Block           #
157    #######################
158    ccm=conn.raw_sql("""
159                    select gvkey, lpermno as permno, linktype, linkprim,
160                    linkdt, linkenddt
161                    from crsp.ccmxpf_linktable
162                    where substr(linktype,1,1)='L'
163                    and (linkprim ='C' or linkprim='P')
164                    """)
165
166    ccm['linkdt']=pd.to_datetime(ccm['linkdt'])
167    ccm['linkenddt']=pd.to_datetime(ccm['linkenddt'])
168    # if linkenddt is missing then set to today date
169    ccm['linkenddt']=ccm['linkenddt'].fillna(pd.to_datetime('today'))
170
171    ccm1=pd.merge(comp[['gvkey','datadate','be', 'count']],ccm,how='left',on=['gvkey'])
172    ccm1['yearend']=ccm1['datadate']+YearEnd(0)
173    ccm1['jdate']=ccm1['yearend']+MonthEnd(6)
174
175    # set link date bounds
176    ccm2=ccm1[(ccm1['jdate']>=ccm1['linkdt'])&(ccm1['jdate']<=ccm1['linkenddt'])]
177    ccm2=ccm2[['gvkey','permno','datadate','yearend', 'jdate','be', 'count']]
178
179    # link comp and crsp
180    ccm_jun=pd.merge(crsp_jun, ccm2, how='inner', on=['permno', 'jdate'])
181    ccm_jun['beme']=ccm_jun['be']*1000/ccm_jun['dec_me']
182
183    # select NYSE stocks for bucket breakdown
184    # exchcd = 1 and positive beme and positive me and shrcd in (10,11) and at least 2 years
185    in comp
186    nyse=ccm_jun[(ccm_jun['exchcd']==1) & (ccm_jun['beme']>0) & (ccm_jun['me']>0) &
187    (ccm_jun['count']>1) & ((ccm_jun['shrcd']==10) | (ccm_jun['shrcd']==11))]
188    # size breakdown
189    nyse_sz=nyse.groupby(['jdate'])['me'].median().to_frame().reset_index().rename(columns=
190    {'me':'sizemedn'})
```

```python
191   # beme breakdown
192   nyse_bm=nyse.groupby(['jdate'])['beme'].describe(percentiles=[0.3, 0.7]).reset_index()
193   nyse_bm=nyse_bm[['jdate','30%','70%']].rename(columns={'30%':'bm30', '70%':'bm70'})
194
195   nyse_breaks = pd.merge(nyse_sz, nyse_bm, how='inner', on=['jdate'])
196   # join back size and beme breakdown
197   ccm1_jun = pd.merge(ccm_jun, nyse_breaks, how='left', on=['jdate'])
198
199
200   # function to assign sz and bm bucket
201   def sz_bucket(row):
202       if row['me']==np.nan:
203           value=''
204       elif row['me']<=row['sizemedn']:
205           value='S'
206       else:
207           value='B'
208       return value
209
210   def bm_bucket(row):
211       if 0<=row['beme']<=row['bm30']:
212           value = 'L'
213       elif row['beme']<=row['bm70']:
214           value='M'
215       elif row['beme']>row['bm70']:
216           value='H'
217       else:
218           value=''
219       return value
220
221   # assign size portfolio
222   ccm1_jun['szport']=np.where((ccm1_jun['beme']>0)&(ccm1_jun['me']>0)&
223   (ccm1_jun['count']>=1), ccm1_jun.apply(sz_bucket, axis=1), '')
224   # assign book-to-market portfolio
225   ccm1_jun['bmport']=np.where((ccm1_jun['beme']>0)&(ccm1_jun['me']>0)&
226   (ccm1_jun['count']>=1), ccm1_jun.apply(bm_bucket, axis=1), '')
227   # create positivebmeme and nonmissport variable
228   ccm1_jun['posbm']=np.where((ccm1_jun['beme']>0)&(ccm1_jun['me']>0)&(ccm1_jun['count']>=1),
229   1, 0)
230   ccm1_jun['nonmissport']=np.where((ccm1_jun['bmport']!=''), 1, 0)
231
232   # store portfolio assignment as of June
233   june=ccm1_jun[['permno','date', 'jdate', 'bmport','szport','posbm','nonmissport']]
234   june['ffyear']=june['jdate'].dt.year
235
236   # merge back with monthly records
237   crsp3 =
238   crsp3[['date','permno','shrcd','exchcd','retadj','me','wt','cumretx','ffyear','jdate']]
239   ccm3=pd.merge(crsp3,
240           june[['permno','ffyear','szport','bmport','posbm','nonmissport']], how='left', on=
241   ['permno','ffyear'])
242
243   # keeping only records that meet the criteria
244   ccm4=ccm3[(ccm3['wt']>0)& (ccm3['posbm']==1) & (ccm3['nonmissport']==1) &
245           ((ccm3['shrcd']==10) | (ccm3['shrcd']==11))]
246
247   ###########################
248   # Form Fama French Factors #
249   ###########################
250
251   # function to calculate value weighted return
252   def wavg(group, avg_name, weight_name):
253       d = group[avg_name]
```

```
254        w = group[weight_name]
255        try:
256            return (d * w).sum() / w.sum()
257        except ZeroDivisionError:
258            return np.nan
259
260    # value-weigthed return
261    vwret=ccm4.groupby(['jdate','szport','bmport']).apply(wavg,
262    'retadj','wt').to_frame().reset_index().rename(columns={0: 'vwret'})
263    vwret['sbport']=vwret['szport']+vwret['bmport']
264
265    # firm count
266    vwret_n=ccm4.groupby(['jdate','szport','bmport'])
267    ['retadj'].count().reset_index().rename(columns={'retadj':'n_firms'})
268    vwret_n['sbport']=vwret_n['szport']+vwret_n['bmport']
269
270    # tranpose
271    ff_factors=vwret.pivot(index='jdate', columns='sbport', values='vwret').reset_index()
272    ff_nfirms=vwret_n.pivot(index='jdate', columns='sbport', values='n_firms').reset_index()
273
274    # create SMB and HML factors
275    ff_factors['WH']=(ff_factors['BH']+ff_factors['SH'])/2
276    ff_factors['WL']=(ff_factors['BL']+ff_factors['SL'])/2
277    ff_factors['WHML'] = ff_factors['WH']-ff_factors['WL']
278
279    ff_factors['WB']=(ff_factors['BL']+ff_factors['BM']+ff_factors['BH'])/3
280    ff_factors['WS']=(ff_factors['SL']+ff_factors['SM']+ff_factors['SH'])/3
       ff_factors['WSMB'] = ff_factors['WS']-ff_factors['WB']
       ff_factors=ff_factors.rename(columns={'jdate':'date'})

       # n firm count
       ff_nfirms['H']=ff_nfirms['SH']+ff_nfirms['BH']
       ff_nfirms['L']=ff_nfirms['SL']+ff_nfirms['BL']
       ff_nfirms['HML']=ff_nfirms['H']+ff_nfirms['L']

       ff_nfirms['B']=ff_nfirms['BL']+ff_nfirms['BM']+ff_nfirms['BH']
       ff_nfirms['S']=ff_nfirms['SL']+ff_nfirms['SM']+ff_nfirms['SH']
       ff_nfirms['SMB']=ff_nfirms['B']+ff_nfirms['S']
       ff_nfirms['TOTAL']=ff_nfirms['SMB']
       ff_nfirms=ff_nfirms.rename(columns={'jdate':'date'})
```

Top of Section

## Outcome Discussion

We compare the output from the Python code with the one from the Fama-French Factors database, and the results are very close.

```
1    ###################
2    # Compare With FF #
3    ###################
4    _ff = conn.get_table(library='ff', table='factors_monthly')
5    _ff=_ff[['date','smb','hml']]
6    _ff['date']=_ff['date']+MonthEnd(0)
7
8    _ffcomp = pd.merge(_ff, ff_factors[['date','WSMB','WHML']], how='inner', on=['date'])
9    _ffcomp70=_ffcomp[_ffcomp['date']>='01/01/1970']
10   print(stats.pearsonr(_ffcomp70['smb'], _ffcomp70['WSMB']))
11   print(stats.pearsonr(_ffcomp70['hml'], _ffcomp70['WHML']))
```
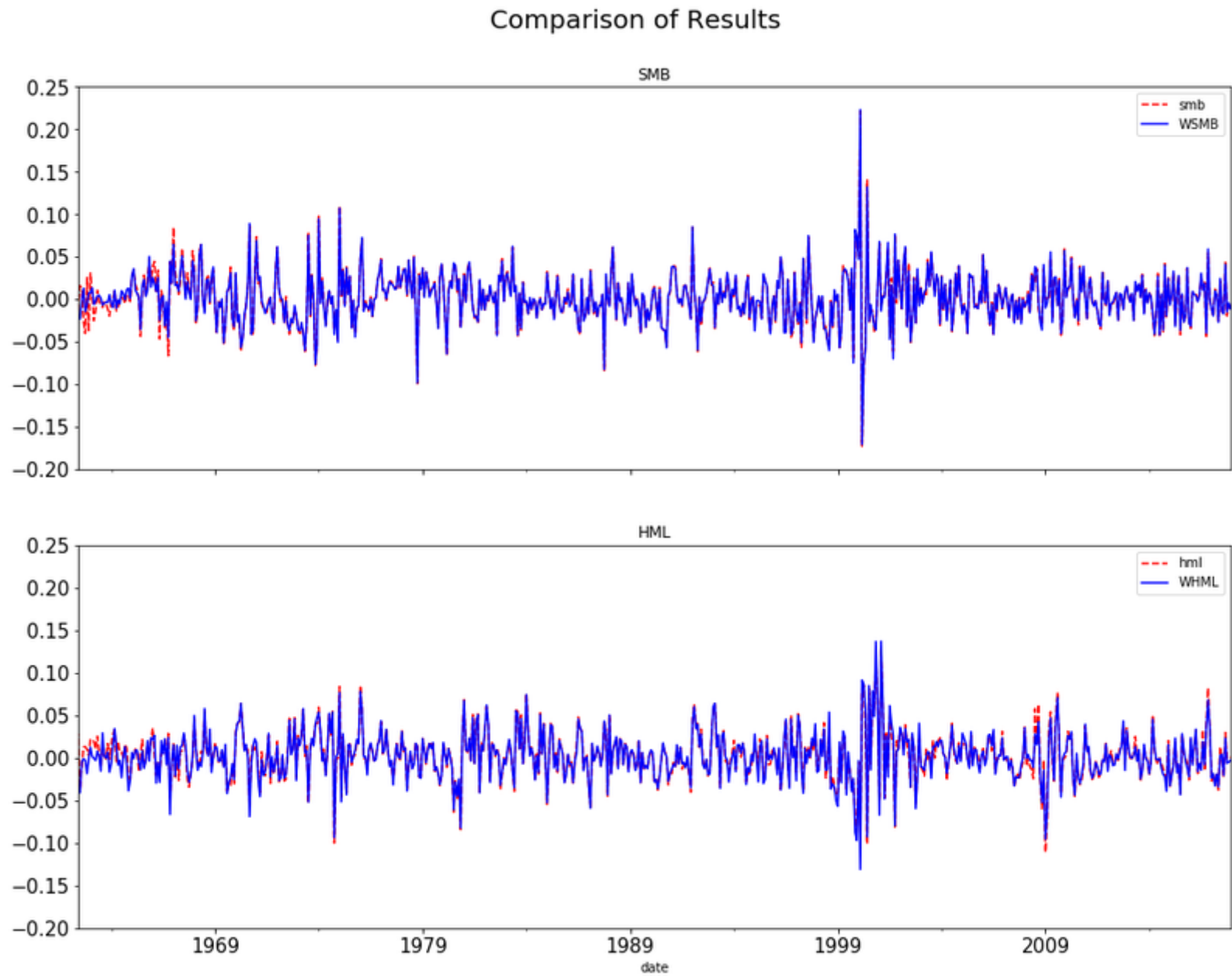
The table below lists the correlation between the Python series of SMB and HML with the Fama French factors respectively, for the sample period of 1970 onwards.

| | Correlation |
|---|---|

| | **Correlation** |
|---|---|
| | **Correlation** |
| **SMB** | 99.6% |
| **HML** | 98.1% |

Including the earlier sample slightly decreases the correlation as Fama French adopted sources other than CRSP and Compustat for the earlier sample when calculating the original series.

Lastly, we present the figure below which compares the entire series of factors. Solid blue line represents the risk factor generated from the Python code, and dash red line represents the original data series from Fama French library.



Comparison of Results

Top of Section

Top

# Table of Contents

Account Preferences

Info / Support Request

Privacy Policy

Sample Data

Conference Calendar

Impactful Research

Account Preferences

Info / Support Request

Privacy Policy

Sample Data

Conference Calendar

Impactful Research