

This file is used when we met discontinuities in the part of data crawling, to restart the crawling process from the last checkpoint.

In []:

```
1 # Run this document if you raised internet error when crawling the data.
2 # Note that if no temporary data has been saved, then no need to run this document.
3
4 import pandas as pd
5 import akshare as ak
6 import re
7 import datetime
8 import time
```

In []:

```
1 startDate, endDate = "19950101", "20211231"
2
3 if int(startDate) > int(endDate) or int(endDate) > datetime.date.today().year * 10000 + \
4     datetime.date.today().month * 100 + datetime.date.today().day:
5     print("Invalid Time Interval")
6     quit()
```

In []:

```
1 def ParseDate(date) -> datetime.date:
2     date = list(map(int, re.findall(pattern="[0-9]+", string=str(date))))
3     return datetime.date(date[0], date[1], date[2])
```

In []:

```
1 close_df = pd.read_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_temp_data\\close_temp.csv",
2                         encoding="gbk", index_col="trade_date")
3 return_df = pd.read_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_temp_data\\return_temp.csv",
4                          encoding="gbk", index_col="trade_date")
5 BM_df = pd.read_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_temp_data\\BM_temp.csv",
6                     encoding="gbk", index_col="trade_date")
7 MV_df = pd.read_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_temp_data\\MV_temp.csv",
8                     encoding="gbk", index_col="trade_date")
9
10 close_df.index = pd.Series(map(ParseDate, close_df.index), name="trade_date")
11 return_df.index = pd.Series(map(ParseDate, return_df.index), name="trade_date")
12 BM_df.index = pd.Series(map(ParseDate, BM_df.index), name="trade_date")
13 MV_df.index = pd.Series(map(ParseDate, MV_df.index), name="trade_date")
```

In []:

```
1 calender = pd.read_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_data\\calender.csv")
2 calender = pd.Series(map(lambda x: datetime.date(int(x[0:4]), int(x[5:7]), int(x[8:])),
3                         calender["trade_date"].values), index=list(calender.iloc[:, 0]), name=
```



In []:

```
1 exceptionLst = pd.read_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_temp_data\\exceptionLst.csv")
2 exceptionLst = list(map(lambda x: (int(re.findall("[0-9]+", x)[0]), re.findall("\\' ([0-9]+)"
3                             exceptionLst.iloc[:, 0]))
```

In []:

```
1 def GetCodeLst(fromWhat: str) -> list:
2     if fromWhat == "HS300":
3         hs300_Stocks = pd.read_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_data\\hs300_stocks.csv",
4                                     encoding="gbk").set_index("code")
5         return list(map(lambda x: re.search(pattern="[0-9]+", string=x).group(),
6                             list(hs300_Stocks.index)))
7
8     elif fromWhat == "A":
9         return list(ak.stock_info_sh_name_code(indicator="主板A股")["代码"]) + \
10                list(ak.stock_info_sh_name_code(indicator="科创板")["代码"]) + \
11                list(ak.stock_info_sz_name_code(indicator="A股列表")["A股代码"])
12
13
```

In []:

```
1 codeLst = GetCodeLst(fromWhat="A")
```

In []:

```
1 def ParseDate(date: str) -> datetime.date:
2     date = list(map(int, re.findall(pattern="[0-9]+", string=str(date))))
3     return datetime.date(date[0], date[1], date[2])
```

Restart from the last checkpoint.

Noted that a manual input of checkpoint (value of what_now) is required.

In []:

```
1 # what_now should be the value of last checkpoint + 1
2 what_now, failure, maximum_failure_allowed, length = 1201, 0, 3, len(codeLst)
3 print(f"Data of {length} stocks in total need to be collected, waiting.....")
4 while what_now < length:
5     code = codeLst[what_now]
6
7     try:
8         this_stock_hist_daily = ak.stock_zh_a_hist(symbol=code, period="daily",
9                                                     start_date=startDate, end_date=endDate,
10                                                     adjust="hfq")["日期", "收盘", "涨跌幅"].se
11
12         this_stock_BMMV_daily = \
13             ak.stock_a_lg_indicator(symbol=code)["trade_date", "pb", "total_mv"].set_index("t
14
15         this_stock_hist_daily.index = map(ParseDate, list(this_stock_hist_daily.index))
16         this_stock_BMMV_daily.index = map(ParseDate, list(this_stock_BMMV_daily.index))
17
18         failure = 0
19
20         try:
21             close_df[code] = this_stock_hist_daily["收盘"]
22             return_df[code] = this_stock_hist_daily["涨跌幅"]
23             BM_df[code] = 1 / this_stock_BMMV_daily["pb"]
24             MV_df[code] = this_stock_BMMV_daily["total_mv"]
25
26             print(f"{what_now}/{length}. Data collected and merged for code: {code}")
27
28         except:
29             print(f"{what_now}/{length}. Met an unknown error when merging data of code: {code}")
30             exceptionLst.append((what_now, code))
31
32         if what_now % 100 == 0:
33             close_df.to_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_temp_data\\close_temp.csv",
34                             index=True, header=True)
35             return_df.to_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_temp_data\\return_temp.c
36                               index=True, header=True)
37             BM_df.to_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_temp_data\\BM_temp.csv", ind
38                           header=True)
39             MV_df.to_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_temp_data\\MV_temp.csv", ind
40                           header=True)
41
42             if len(exceptionLst) != 0:
43                 pd.Series(exceptionLst).to_csv(
44                     "C:\\Users\\tianj\\Project 1\\data\\HS300_temp_data\\calender.csv",
45                     index=False, header=True)
46
47             print(f"Temporary file is saved at: {code}. Position is: {what_now}")
48
49         if what_now % 30 == 0:
50             print("Resuming in 45 seconds.....")
51             time.sleep(45)
52
53         what_now += 1
54
55     except:
56         failure += 1
57         print(f"{what_now}/{length}. Problem encountered at code: {code}. Failure = {failure}")
58
59         if failure > maximum_failure_allowed:
```

```
60         break
61     else:
62         print(f"Retrying in {60 * failure} seconds.....")
63         time.sleep(60 * failure)
64         continue
65
66
67 else:
68     close_df.to_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_data\\close.csv", index=True, header=True)
69     return_df.to_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_data\\return.csv", index=True, header=True)
70     BM_df.to_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_data\\BM.csv", index=True, header=True)
71     MV_df.to_csv("C:\\Users\\tianj\\Project 1\\data\\HS300_data\\MV.csv", index=True, header=True)
72
73     if len(exceptionLst) == 0:
74         print("All data are collected and merged successfully")
75     else:
76         print("exceptionLst is not empty: failed to merge some data")
77
```