

```
In [1]: 1 import numpy as np
        2 import pandas as pd
        3 import matplotlib.pyplot as plt
```

executed in 2.30s, finished 14:05:22 2021-11-08

```
In [2]: 1 data = pd.read_csv('./data.csv')
```

executed in 268ms, finished 14:05:36 2021-11-08

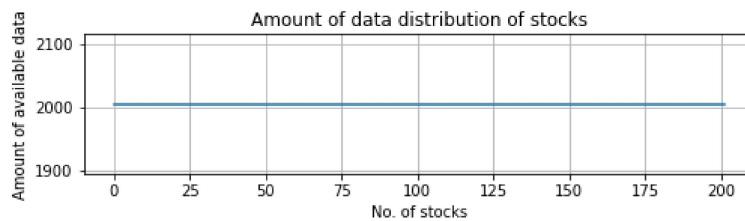
▼ EDA of raw_data : 'data.csv'

```
In [37]: 1 # Filtering out the stocks with inadequate data
        2
        3 tickcounts= data.ticker.value_counts()
        4
        5 availstocks = list(tickcounts[:202].index)# Since the dataset is already small. We could confidently throw away some
        6 print('Number of available stocks:%d'%len(availstocks))
        7 print('='*40)
        8 plt.figure(figsize = (8,1.7))
        9 plt.plot(tickcounts.values)
       10 plt.ylabel('Amount of available data')
       11 plt.xlabel('No. of stocks')
       12 plt.title('Amount of data distribution of stocks')
       13 plt.grid()
       14 plt.show()
       15
       16 data = data[data['ticker'].isin(availstocks)]
       17
       18 # Check if there are missing values in the dataframe
       19
       20 print('The number of missing value = %d'%data.isnull().sum().sum())
       21 print('='*40)
       22
```

executed in 260ms, finished 14:20:02 2021-11-08

Number of available stocks:202

=====



The number of missing value = 0

=====

▼ Idea-generating for strategy

Since I only get price and vols, I should firstly generate some basic technical factors like moving average of prices and vols

- ☐ Let's consider the 'last' as the 'last tick'

```
In [46]: 1 _name = availstocks[0]
        2 df = data[data['ticker']==_name]
        3 df = df[['last', 'volume']]
        4 df.index = np.arange(df.shape[0])
```

executed in 32ms, finished 14:24:40 2021-11-08

▼ Price factor

```
In [84]: 1 df['ma5'] = df['last'].rolling(window = 5).mean()
2 df['ma10'] = df['last'].rolling(window = 10).mean()
3 df['ma20'] = df['last'].rolling(window = 20).mean()
4
5 df['long_sig'] = (df['ma5']>df['ma10'])&(df['ma10']>df['ma20'])
6 df['short_sig'] = (df['ma5']<df['ma10'])&(df['ma10']<df['ma20'])
7 df.tail(4)
```

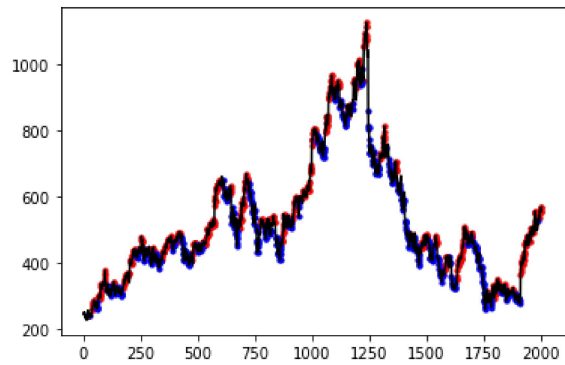
executed in 22ms, finished 14:51:18 2021-11-08

Out[84]:

	last	volume	ma5	ma10	ma20	long_sig	short_sig
2001	552.0	1963700	555.8	549.1	533.6	True	False
2002	554.0	1450400	555.4	550.3	536.0	True	False
2003	565.0	1850200	557.6	553.9	538.0	True	False
2004	567.0	2451300	560.2	556.9	540.8	True	False

```
In [77]: 1 long_sig = list(df[df['long_sig']==1].index)
2 long_price = [df['last'].iloc[idx] for idx in long_sig]
3 short_sig = list(df[df['short_sig']==1].index)
4 short_price = [df['last'].iloc[idx] for idx in short_sig]
5
6 plt.scatter(long_sig, long_price,s = 10, color = 'r')
7 plt.scatter(short_sig, short_price,s = 10, color = 'blue')
8 plt.plot(df['last'].values,color = 'black')
9 plt.show()
```

executed in 199ms, finished 14:44:13 2021-11-08



In the graph above, we observe that the

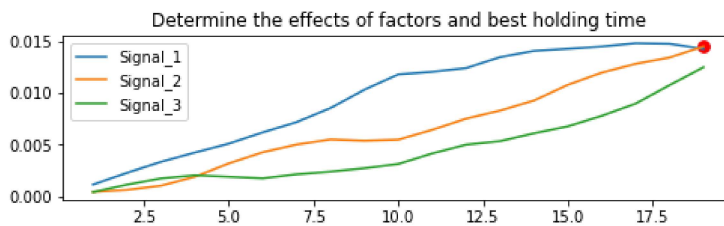
- ☐ red points appear when the stock rise.
- ☐ blue points appear when the stock drops.

```

In [176]: 1 # Let's check if this factor is useful
2
3 # Dividing df into 3 part by long_sig, no_sig, short_sig
4
5 no_sig = list(df[(df['long_sig']==0) & (df['short_sig']==0)].index)
6 # Calculate the returns individually
7
8     # indices = long_sig#####
9     # length = 15####
10
11 def calret(indices, length):
12     newlist = [i for i in indices if (i<2005-length-1) ]
13     new2list = [i+length for i in indices if (i<2005-length-1) ]
14     df1 = df.loc[newlist]['last'].values
15     df2 = df.loc[new2list]['last'].values
16     ret = (df2-df1)/df1
17     return ret
18
19 # Here we use the return mean as an indicator
20 hold_len = list(range(1,20))
21 hold_relong = [calret(long_sig,t).mean() for t in hold_len]
22 hold_retno = [calret(no_sig,t).mean() for t in hold_len]
23 hold_retsort = [calret(short_sig,t).mean() for t in hold_len]
24
25 plt.figure(figsize= (8,2))
26 plt.title('Determine the effects of factors and best holding time')
27 plt.plot(hold_len, hold_relong, label = 'Signal_1')
28 plt.plot(hold_len, hold_retno, label = 'Signal_2')
29 plt.plot(hold_len, hold_retsort, label = 'Signal_3')
30 plt.scatter(hold_len[-1], hold_retno[-1], s = 60, color= 'r')
31 plt.legend()
32 plt.show()

```

executed in 283ms, finished 15:38:00 2021-11-08



1 The reason I stopped searching is because the red dot above. I will explain it in further version :)

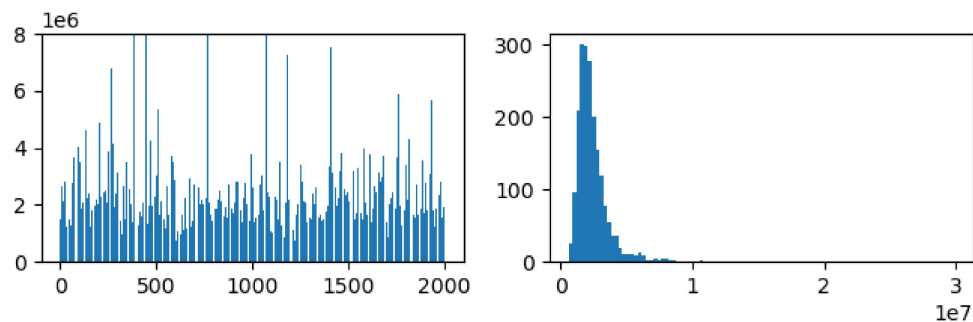
volume factor

```

In [137]: 1 plt.figure(figsize =(8,2), dpi = 100)
2 plt.subplot(121)
3 plt.bar(range(len(df['volume'])), df['volume'])
4 plt.ylim(0, 0.8e7)
5
6 plt.subplot(122)
7 plt.hist(df['volume'], bins = 100)
8
9 plt.show()
10
11 # df['volume'].describe()

```

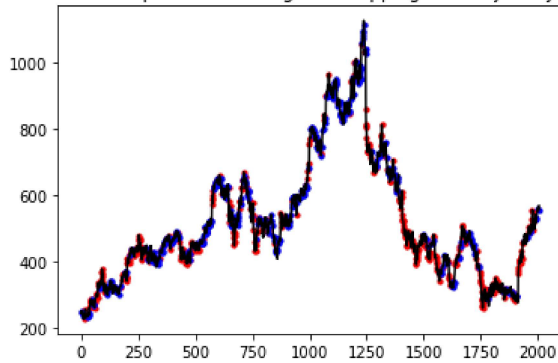
executed in 2.97s, finished 15:13:43 2021-11-08



```
In [186]: 1 # volin_sig = df[  
2  
3 from statsmodels.tsa.stattools import adfuller  
4 adfuller(df['volume'])# Stationary  
5 # Hence we could use the past dataset in future  
6 import scipy.stats as stats  
7 import statsmodels.api as sm  
8  
9  
10 df['hi'] = (df['volume'] > df['volume'].describe()['75%'])  
11 df['lo'] = (df['volume'] < df['volume'].describe()['25%'])  
12  
13 hi_sig = list(df[df['hi']==1].index)  
14 lo_sig = list(df[df['lo']==1].index)  
15  
16 hi_price = [df['last'].iloc[idx] for idx in hi_sig]  
17 lo_price = [df['last'].iloc[idx] for idx in lo_sig]  
18  
19 plt.scatter(hi_sig, hi_price, s = 10, color = 'r')  
20 plt.scatter(lo_sig, lo_price, s = 10, color = 'blue')  
21 plt.plot(df['last'].values, color = 'black')  
22 plt.title('This factor did not separate the rising and dropping so well just by appearance')  
23 plt.show()
```

executed in 322ms, finished 15:43:19 2021-11-08

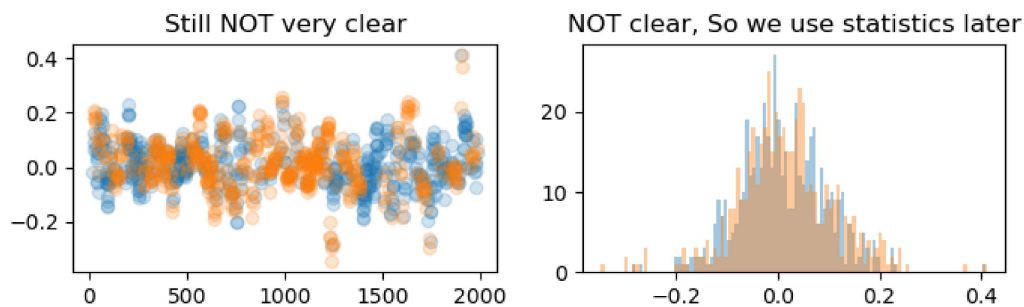
This factor did not separate the rising and dropping so well just by appearance



I intended to directly show, not really clear.

```
In [188]: 1 # 15days return
2
3 tmpdf = df.copy()
4 tmpdf['15last'] = tmpdf['last'].shift(-15)
5 tmpdf['ret15'] = (tmpdf['15last']-tmpdf['last'])/tmpdf['last']
6
7 tmpdf = tmpdf.dropna()
8
9 hidf = tmpdf[tmpdf['hi']==1]
10 lodf = tmpdf[tmpdf['lo']==1]
11
12 plt.figure(figsize = (8,2), dpi = 100)
13 plt.subplot(121)
14 plt.title('Still NOT very clear')
15 plt.scatter(list(hidf.index), hidf['ret15'].values, alpha = 0.2)
16 plt.scatter(list(lodf.index), lodf['ret15'].values, alpha = 0.2)
17 plt.subplot(122)
18 plt.title('NOT clear, So we use statistics later')
19 plt.hist(hidf['ret15'].values, alpha = 0.4, bins = 100)
20 plt.hist(lodf['ret15'].values, alpha = 0.4, bins = 100)
21 plt.show()
```

executed in 502ms, finished 15:43:54 2021-11-08



▼ Use t-stats to tell whether my volume factor could take effect

```
In [175]: 1 from scipy import stats
2 print(stats.ttest_ind(hidf['ret15'].values, lodf['ret15'].values))
3 print('====Hence, we cannot use these pair of factors, because p>0.05====')
```

executed in 8ms, finished 15:37:25 2021-11-08

Ttest_indResult(statistic=-1.0943294732577258, pvalue=0.2740766809301102)
====Hence, we cannot use these pair of factors, because p>0.05=====

```
In [184]: 1 print(stats.ttest_ind(hold_retlong, hold_retshort))
2 print(stats.ttest_ind(hold_retlong, hold_retno))
3 print('====Hence, Keep the factor, because p<0.1=====')
```

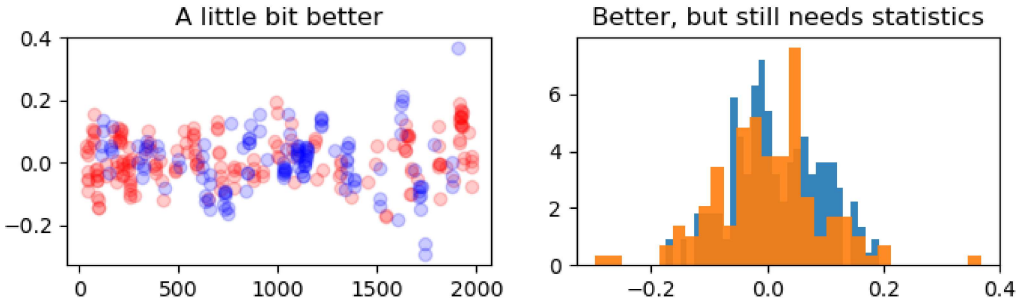
executed in 15ms, finished 15:41:03 2021-11-08

Ttest_indResult(statistic=3.8286724512489636, pvalue=0.0004954811015490555)
Ttest_indResult(statistic=2.0062038115433154, pvalue=0.0523925377251018)
====Hence, Keep the factor, because p<0.1=====

▼ The time left not too much, so I try a combination

```
In [216]: 1 tmpdf = tmpdf.dropna()
2
3 longhi = list(tmpdf[(tmpdf['ma5']>tmpdf['ma10'])&(tmpdf['ma10']>tmpdf['ma20'])&(tmpdf['hi']==1)].index)
4 nolo = list(tmpdf[(tmpdf['ma5']<tmpdf['ma10'])&(tmpdf['ma10']<tmpdf['ma20'])&(tmpdf['lo']==1)].index)
5 longhired = tmpdf[(tmpdf['ma5']>tmpdf['ma10'])&(tmpdf['ma10']>tmpdf['ma20'])&(tmpdf['hi']==1)][ret15']
6 nolored = tmpdf[(tmpdf['ma5']<tmpdf['ma10'])&(tmpdf['ma10']<tmpdf['ma20'])&(tmpdf['lo']==1)][ret15']
7
8 plt.figure(figsize = (8,2), dpi = 100)
9 plt.subplot(121)
10 plt.title('A little bit better')
11 plt.scatter(longhi, longhired.values, alpha = 0.2, c = 'r')
12 plt.scatter(nolo, nolored.values, alpha = 0.2, c = 'b')
13 plt.subplot(122)
14 plt.title('Better, but still needs statistics')
15 plt.hist(longhired.values, alpha = 0.9, bins = 30, density = True )
16 plt.hist(nolored.values, alpha = 0.9, bins = 30, density = True )
17 plt.show()
```

executed in 337ms, finished 16:22:55 2021-11-08



```
In [196]: 1 print(stats.ttest_ind(longhired, nolored))
2 print('====Hence, Keep the factor, because p<0.1====')
3 print('The combination works')
```

executed in 20ms, finished 15:53:09 2021-11-08

Ttest_indResult(statistic=-4.531104161898143, pvalue=8.386041356579832e-06)
====Hence, Keep the factor, because p<0.1====
The combination works

```
In [ ]:
```

1

▼ No.1 Strategy

Item	Details
Targets	202 stocks in available stocks
Buy points	Signal_1,hi_1, position is not full
Holding length	15 days
stop loss	-10% for every deal

```

In [305]: 1 class backtest1:
2         def __init__(self, name, stoploss=0.1, holding_length=15):
3             # Generating the price matrix
4
5             tdf = data[data['ticker']==name].copy()
6             tdf = tdf[['date', 'last', 'volume']]
7             tdf.index = np.arange(tdf.shape[0])
8
9             tdf['ma5'] = tdf['last'].rolling(window = 5).mean()
10            tdf['ma10'] = tdf['last'].rolling(window = 10).mean()
11            tdf['ma20'] = tdf['last'].rolling(window = 20).mean()
12            tdf['long_sig'] = (tdf['ma5']>tdf['ma10'])&(tdf['ma10']>tdf['ma20'])
13            tdf['short_sig'] = (tdf['ma5']<tdf['ma10'])&(tdf['ma10']<tdf['ma20'])
14
15            self.hi = tdf['volume'].iloc[:1500].describe()['75%']
16            self.lo = tdf['volume'].iloc[:1500].describe()['25%']
17            tdf['hi'] = (tdf['volume']>self.hi)
18            tdf['lo'] = (tdf['volume']<self.lo)
19
20            tdf['15last'] = tdf['last'].shift(-holding_length)
21            tdf['ret15'] = (tdf['15last']-tdf['last'])/tdf['last']
22            self.tdf = tdf.iloc[1500:]
23            self.stoploss = stoploss
24            self.holding_length = holding_length
25
26        def backtesting(self):
27            self.pos = 0
28            self.cap = 1
29            self.earn = []
30            self.net = [1]
31            tdf = self.tdf
32            nexttradeday = 0
33            buyp = tdf['last'].iloc[0]
34            for idx in range(tdf.shape[0]-self.holding_length-1):
35                if idx == nexttradeday:
36                    m = (tdf['last'].iloc[idx]-buyp)/buyp
37                    self.net.append((1+m)*self.net[-1])
38                    # print(self.net[-1], 'asdf', buyp, tdf['last'].iloc[idx])
39                    self.pos = 0
40                if (self.pos==1):
41                    m = (tdf['last'].iloc[idx]-buyp)/buyp
42                    self.net.append((1+m)*self.net[-1])
43                    # print(self.net[-1], tdf['date'].iloc[idx], buyp, tdf['last'].iloc[idx], '====')
44                    continue
45                if (not tdf['hi'].iloc[idx]) or (not tdf['long_sig'].iloc[idx]):
46                    self.net.append(self.net[-1])
47                    continue
48
49                self.pos = 1
50                buyp = tdf['last'].iloc[idx]
51                nexttradeday = idx+self.holding_length
52                self.earn.append(tdf['ret15'].iloc[idx])
53                self.cap *= (1+tdf['ret15'].iloc[idx])
54                # print(self.cap, tdf['date'].iloc[idx], tdf['last'].iloc[idx], tdf['15last'].iloc[idx], tdf['ret15'].iloc[idx])
55
56                self.net.append(self.net[-1])
57                self.cap = self.net[-1]
58            return
59

```

executed in 23ms, finished 17:07:44 2021-11-08

```

In [321]: 1 ts = []
2         for i in range(50):
3             test = backtest1(availstocks[i], holding_length = 10)
4             test.backtesting()
5             t = test.net
6             t = t[:480]
7             ts.append(t)
8
9         ts = np.array(ts)
10        # plt.plot(range(len(t)), t)
11

```

executed in 1.81s, finished 17:10:43 2021-11-08

```

In [322]: 1 len(t)

```

executed in 13ms, finished 17:10:43 2021-11-08

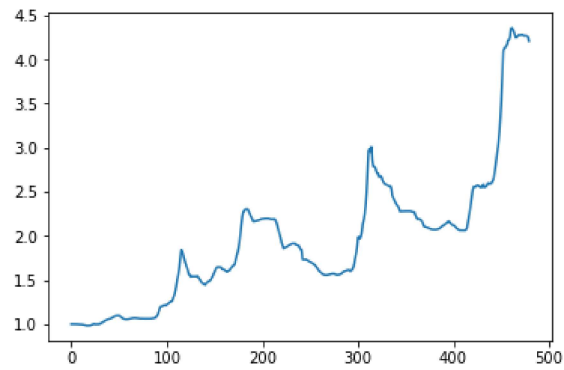
Out[322]: 480

In [323]:

```
1 tsmean = ts.mean(axis = 0)
2 plt.plot(tsmean)
3
```

executed in 161ms, finished 17:11:02 2021-11-08

Out[323]: <matplotlib.lines.Line2D at 0x233a7362550>



In [324]:

```
1 def MaxDrawdown(return_list):
2     maxcum = np.zeros(len(return_list))
3     b = return_list[0]
4     for i in range(0, len((return_list))):
5         if (return_list[i]>b):
6             b = return_list[i]
7             maxcum[i] = b
8     i = np.argmax((maxcum-return_list)/maxcum)
9     if i == 0:
10         return 0
11     j = np.argmax(return_list[:i])
12
13     return (return_list[j]-return_list[i])/return_list[j]
```

executed in 16ms, finished 17:12:25 2021-11-08

In [331]:

```
1 test.tdf.iloc[:480]
```

executed in 25ms, finished 17:14:00 2021-11-08

Out[331]:

	date	last	volume	ma5	ma10	ma20	long_sig	short_sig	hi	lo	15last	ret15
1500	2019-02-19	4075.9183	2444400	4000.59030	3945.02654	3935.330855	True	False	False	False	4129.9904	0.013266
1501	2019-02-20	4099.2253	2712700	4027.06700	3958.54457	3941.297435	True	False	False	False	4116.0063	0.004094
1502	2019-02-21	4115.0740	2789200	4061.00190	3976.91044	3951.505875	True	False	False	False	4068.4601	-0.011328
1503	2019-02-22	4116.9386	1784500	4102.20858	3999.65802	3961.061725	True	False	False	True	4013.4557	-0.025136
1504	2019-02-25	4150.5006	2510500	4111.53136	4035.08458	3971.176940	True	False	False	False	4017.1848	-0.032120
...
1975	2021-02-04	4060.0000	2023600	3948.40000	3973.40000	3907.200000	False	False	False	True	4117.0000	0.014039
1976	2021-02-05	4060.0000	1963200	3986.80000	3969.40000	3938.600000	True	False	False	True	4199.0000	0.034236
1977	2021-02-08	4141.0000	2062900	4033.80000	3973.50000	3963.800000	True	False	False	True	4176.0000	0.008452
1978	2021-02-09	4143.0000	1600400	4080.80000	3992.70000	3988.450000	True	False	False	True	4314.0000	0.041274
1979	2021-02-10	4325.0000	3488400	4145.80000	4033.30000	4019.200000	True	False	True	False	4161.0000	-0.037919

480 rows × 12 columns

In [343]:

```
1 print('MaxDraw:%f'%MaxDrawdown(tsmean))
2 print('Sharpe:', (tsmean[-1]**0.5-1.0)/np.std(tsmean))
```

executed in 6ms, finished 17:16:46 2021-11-08

MaxDraw: 0.325306

Sharpe: 1.340815908280916

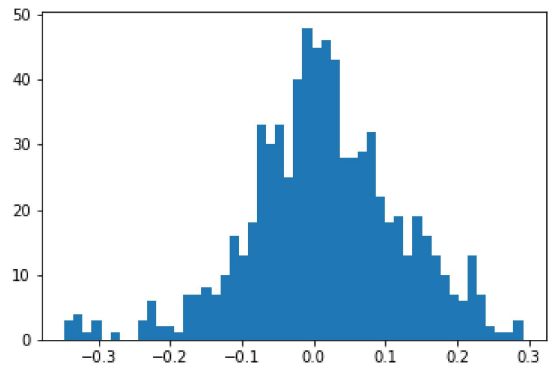
- ▼ Further research: Try to look into the fails of in the first factor.
- ▼ Other supporting materials that are not considered as main

materials

In [111]:

```
1 plt.hist(ret, bins =50)
2 plt.show()
3 ret.mean()
```

executed in 207ms, finished 14:56:10 2021-11-08



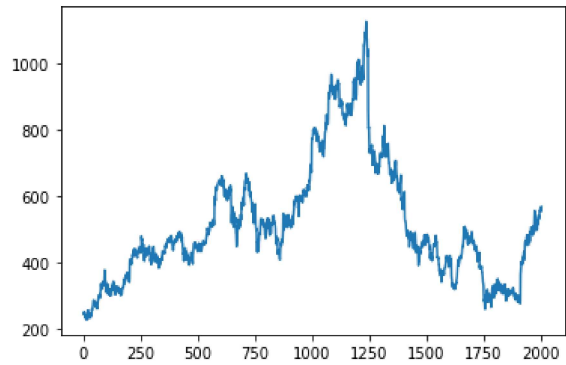
Out[111]: 0.014159445419081963

In [35]:

```
1 prices = data[data['ticker']=='availstocks[0']]['last']
2 plt.plot(prices.values)
```

executed in 180ms, finished 14:16:22 2021-11-08

Out[35]: [<matplotlib.lines.Line2D at 0x2339d66b790>]



In []:

```
1
```

In []:

```
1
```