

```
data = data.drop('Sub-product',axis=1)
data = data.drop('Sub-issue',axis=1)
data = data.drop('Consumer complaint narrative',axis=1)
data = data.drop('Company public response',axis=1)
data = data.drop('Tags',axis=1)
data = data.drop('Consumer consent provided?',axis=1)
data = data.drop('Submitted via',axis=1)
data = data.drop('Consumer disputed?',axis=1)
data = data.drop('Company response to consumer',axis=1)
data = data.drop('Date sent to company',axis=1)
data = data.drop('Timely response?',axis=1)
data = data.drop('Date received',axis=1)

data = data.dropna()

data['Product']=[re.sub('[^\w\s]+', '', s) for s in data['Product'].tolist()]
data['Issue']=[re.sub('[^\w\s]+', '', s) for s in data['Issue'].tolist()]
data['Company']=[re.sub('[^\w\s]+', '', s) for s in data['Company'].tolist()]
data['State']=[re.sub('[^\w\s]+', '', s) for s in data['State'].tolist()]
data['ZIP code']=[re.sub('[^\w\s]+', '', s) for s in data['ZIP code'].tolist()]
```

Verileri düzenlemek için kullandığım kodun bir parçası

Verileri düzenlerken birçok hatayla karşılaştım. Bu hataların en başında da boyut sorunu geliyordu. Csv dosyasının hepsini okumuyordu bu yüzden verilerin hepsini değil çok daha küçük bir kısmını alıyordu.

```
(1282355, 18)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1282355 entries, 0 to 1282354
Data columns (total 18 columns):
#   Column                                Non-Null Count  Dtype
---  ---
0   Date received                        1282355 non-null object
1   Product                             1282355 non-null object
2   Sub-product                         1047189 non-null object
3   Issue                               1282355 non-null object
4   Sub-issue                           751169 non-null object
5   Consumer complaint narrative        383564 non-null object
6   Company public response             449082 non-null object
7   Company                             1282355 non-null object
8   State                               1262955 non-null object
9   ZIP code                            1167057 non-null object
10  Tags                                175643 non-null object
11  Consumer consent provided?          690654 non-null object
12  Submitted via                       1282355 non-null object
13  Date sent to company                1282355 non-null object
14  Company response to consumer        1282348 non-null object
15  Timely response?                    1282355 non-null object
16  Consumer disputed?                  768501 non-null object
17  Complaint ID                        1282355 non-null int64
dtypes: int64(1), object(17)
memory usage: 176.1+ MB
```

Verilerin düzenlemeden önceki hali

Verileri düzenlemeden önce toplamda 18 sütun ve 1282355 tane de satır vardı. Bunların boyutu da 180 mb civarındaydı.

```
(1167023, 6)
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1167023 entries, 0 to 1282354
Data columns (total 6 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   Product               1167023 non-null object
1   Issue                 1167023 non-null object
2   Company               1167023 non-null object
3   State                 1167023 non-null object
4   ZIP code              1167023 non-null object
5   Complaint ID          1167023 non-null int64
dtypes: int64(1), object(5)
memory usage: 62.3+ MB
```

Verilerin düzenledikten sonraki hali

Verileri düzenledikten sonra 18 sütun sayısını istenen sütun sayısına yani 6 sütuna düşürdüm ve stopword ve noktalama işaretlerini de kaldırdıktan sonra total verinin boyutu 60 mb civarına kadar düştü.

```
def dongu(a, b, isim, Benzerlik):
    with time_it('thread'):
        count = 0
        for i in range(a, b):
            ilet = data[isim][i].split(" ")
            for k in range(a, b):
                count += 1
                klist = data[isim][k].split(" ")
                if (len(ilet) > len(klist)):
                    if(len(list(set(ilet) & set(klist)))/len(ilet)*100 > Benzerlik ):
                        oran_list.append(len(list(set(ilet) & set(klist)))/len(ilet)*100)
                        birinci_list.append(ilet)
                        ikinci_list.append(klist)
                    elif (len(ilet) < len(klist)):
                        if(len(list(set(ilet) & set(klist)))/len(klist)*100 > Benzerlik ):
                            oran_list.append(len(list(set(ilet) & set(klist)))/len(klist)*100)
                            birinci_list.append(ilet)
                            ikinci_list.append(klist)
                else:
                    if(len(list(set(ilet) & set(klist)))/len(ilet)*100 > Benzerlik ):
                        oran_list.append(len(list(set(ilet) & set(klist)))/len(ilet)*100)
                        birinci_list.append(ilet)
                        ikinci_list.append(klist)
```

Verileri karşılaştırmak için kullandığım arama metodu

Verileri karşılaştırmak için dataframe olan veri tipini list haline çevirmem gerekiyordu. Bu sorun beni en çok zorlayan sorunlardan biriydi. Aldığım hataları sürekli internete arayarak çözümlerini bulmaya çalışıyordum. Bu şekilde araştırmam bana çok fazla bilgi ve tecrübe kazandırmıştı. Bu şekilde devam ederek burada yaşadığım hataların bir çoğunu çözebildim.

```
def com_Multithread(thread_sayisi,isim,Benzerlik,id):
    with time_it('main Thread'):
        a = 100
        b = 0
        with ThreadPoolExecutor(thread_sayisi) as ex:
            for _ in range(10):
                b += a
                x = b-a
                if(id==0):
                    tanim = [x, b, isim, Benzerlik]
                    ex.submit(dongu, *tanim)
                else:
                    tanim = [x, b, isim, Benzerlik,id]
                    ex.submit(com_dongu, *tanim)

        raise_frame(window)
        BilgiEkrani()
```

Multithreading için kullandığım kod

Multithreading benim için çok fazla sorun oldu çünkü Multithreading python'da düzgün bir şekilde çalışmıyordu. thread sayısının işlem hızına hiçbir katkısı olmuyordu. Python'da bulunan GIL yüzünden threadler düzgün bir şekilde çalışmıyordu. Bu hata için uzun süre araştırdıktan sonra Multiprocessing buldum. Bu işlem hızını gerçekten artırıyordu. Ama hocamız Multithreading istediği için o şekilde bıraktım.

```
def BilgiEkranı():
    y = 1
    for i in range(len(zaman_list)):
        lb1 = Label(window, text=zaman_list[i])
        lb1.grid(column=3, row=y)
        y += 1

    lb1 = Label(window, text="Benzerlik Oranı")
    lb1.grid(column=0, row=0)
    lb1 = Label(window, text="Kayıt 1")
    lb1.grid(column=1, row=0)
    lb1 = Label(window, text="Kayıt 2")
    lb1.grid(column=2, row=0)
    lb1 = Label(window, text="Thread süreleri")
    lb1.grid(column=3, row=0)

    y = 2
    sayı= len(oran_list)
    for i in range(150):
        lb1 = Label(window, text=oran_list[i])
        lb2 = Label(window, text=birinci_list[i])
        lb3 = Label(window, text=ikinci_list[i])
        lb1.grid(column=0, row=y)
        lb2.grid(column=1, row=y)
        lb3.grid(column=2, row=y)
        y += 1
```

Arayüzün bilgi ekranını tasarlamak için kullandığım kod

```
def arayuz():
    lb1 = Label(window2, text="Kayıt: ")
    lb1.grid(column=0, row=0)
    kayit = Entry(window2)
    kayit.grid(column=1, row=0)
    lb2 = Label(window2, text="Thread Sayısı: ")
    lb2.grid(column=0, row=1)
    j = Entry(window2)
    j.grid(column=1, row=1)
    lb3 = Label(window2, text="Benzerlik Oranı: ")
    lb3.grid(column=0, row=2)
    k = Entry(window2)
    k.grid(column=1, row=2)
    lb3 = Label(window2, text="Complaint Id: ")
    lb3.grid(column=0, row=3)
    z = Entry(window2)
    z.grid(column=1, row=3)

    Button(window2, text="Onayla", command=lambda: con.Multithread(int(j.get()), kayit.get(), int(k.get()), int(z.get()))).grid(column=0, row=4)
```

Arayüzün giriş ekranını tasarlamak için kullandığım kod

Arayüz yapmak ile çok bilgim olmadığı için bütün projede almadığım kadar hatayı burada aldım. Ama bunları çok hızlı bir şekilde atılarak bir arayüz tasarlayabildim.

KOD BİLGİSİ(YALANCI KOD)

Kodu txt dosyası olarak ek olarak gönderdim.

SONUÇ

Sonuç olarak python ile Multithreading yaparak büyük veri üzerinde işlem yapabilmeyi ve bu verilerin nasıl işleneceğini öğrenmiş oldum. Bu işlediğim veriyi kullanıcıya düzgün bir şekilde sunmak için de basit de olsa bir arayüz tasarlamayı bu arayüzde girilen değerlerin kodun içinde çalışma mantığını anlamış oldum.

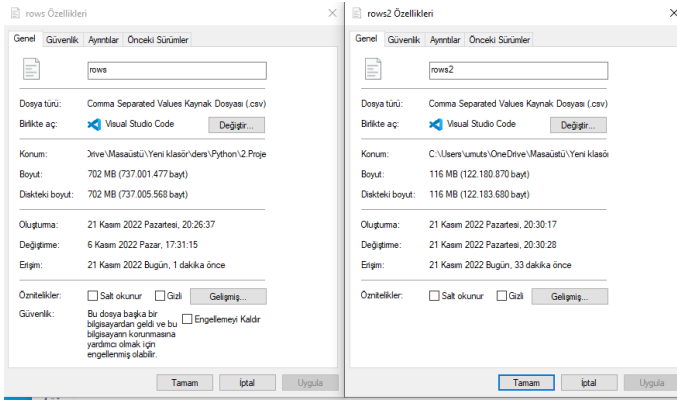
Aşağıda da projeme ait ekran görüntüleri:

Projenin giriş ekranı

Benzerlik Oranı	Kayıt 1	Kayıt 2	Thread süreleri
100.0	Checking savings account	Checking savings account	thread == 0.23400000000037835
100.0	Checking savings account	Checking savings account	thread == 0.25
100.0	Checking savings account	Checking savings account	thread == 0.218999999999914144
100.0	Checking savings account	Checking savings account	thread == 0.29700000000004802
100.0	Checking savings account	Checking savings account	thread == 0.36000000000005821
100.0	Checking savings account	Checking savings account	thread == 0.40600000000005956
100.0	Checking savings account	Checking savings account	thread == 0.35900000000037835
100.0	Checking savings account	Checking savings account	thread == 0.28100000000005856
100.0	Checking savings account	Checking savings account	thread == 0.32799999999995198
100.0	Checking savings account	Checking savings account	thread == 0.23399999999995396
100.0	Checking savings account	Checking savings account	main Thread == 1.07799999999995198
100.0	Checking savings account	Checking savings account	
100.0	Checking savings account	Checking savings account	
100.0	Checking savings account	Checking savings account	
100.0	Checking savings account	Checking savings account	
100.0	Checking savings account	Checking savings account	

Pronin arama yaptıktan sonraki çıktı ekranı

Çıktı ekranında giriş ekranında iste-nen verilerin gösterilmesini sağlıyorum.



Verinin ilk ve son halinin bir ekran görüntüsü

KAYNAKÇA

<https://www.kaggle.com/datasets/selener/consumer-complaint-database>

<https://mertmekatronik.com/thread-ve-multithread-nedir>

<https://www.w3schools.com/python/>

<https://www.geeksforgeeks.org/multithreading-python-set-1/>

<https://www.geeksforgeeks.org/python-programming-language/?ref=shm>

<https://docs.google.com/document/d/18CXhDb1ygxg-YXNBjNzfzZsDFosB5e6BfnXLlejd9l0/edit>

<https://github.com/tensorflow/swift/blob/master/docs/WhySwiftForTensorFlow.md>

<https://www.youtube.com/watch?v=gOtf-Yhn-k>

<https://docs.python.org/3/library/tkinter.html>