Mugs Data Engineer internship

ASSIGNMENT

NAME- Khushi Tolani

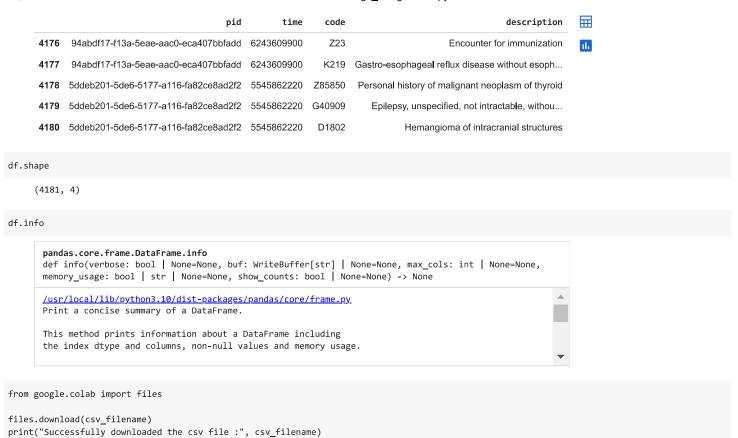
EMAIL ID- kkhushitolanii@gmail.com

PHONE NO.- 8318094789

```
!wget https://physionet.org/static/published-projects/mimic-iv-fhir-demo/mimic-iv-clinical-database-demo-on-fhir-2.0.zip
--2024-05-07 22:20:54-- <a href="https://physionet.org/static/published-projects/mimic-iv-fhir-demo/mimic-iv-clinical-database-demo-on-fhir-2.0.">https://physionet.org/static/published-projects/mimic-iv-fhir-demo/mimic-iv-clinical-database-demo-on-fhir-2.0.</a>
             Resolving physionet.org (physionet.org)... 18.13.52.205
             Connecting to physionet.org (physionet.org) 18.13.52.205 :443... connected.
             HTTP request sent, awaiting response... 200 OK
             Length: 121863908 (116M) [application/zip]
             Saving to: 'mimic-iv-clinical-database-demo-on-fhir-2.0.zip'
            mimic-iv-clinical-d 100%[=======>] 116.22M 269KB/s
                                                                                                                                                                                     in 7m 25s
             2024-05-07 22:28:20 (267 KB/s) - 'mimic-iv-clinical-database-demo-on-fhir-2.0.zip' saved [121863908/121863908]
!unzip mimic-iv-clinical-database-demo-on-fhir-2.0.zip -d mimic_fhir_dataset
             Archive: mimic-iv-clinical-database-demo-on-fhir-2.0.zip
                  inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/LICENSE.txt
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/SHA256SUMS.txt
                 inflating: \verb|mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/Condition.ndjson| and the following of the following 
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/Encounter.ndjson
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/EncounterICU.ndjson
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/Location.ndjson
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/MedicationAdministration.ndjson
                  inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/MedicationAdministrationICU.ndjson
                 inflating: \verb|mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/Medication Dispense.ndjson | and the property of the 
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/MedicationRequest.ndjson
                  inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/ObservationChartevents.ndjson
                 inflating: mimic fhir dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/ObservationDatetimeevents.ndjson
                  inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/ObservationLabevents.ndjson
                  inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/ObservationMicroOrg.ndjson
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/ObservationMicroSusc.ndjson
                 inflating: \verb|mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/0bservation \verb|MicroTest.nd|| some analysis of the property of the pro
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/ObservationOutputevents.ndjson
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/Organization.ndjson
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/Patient.ndjson
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/Procedure.ndjson
                  inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/ProcedureICU.ndjson
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/Specimen.ndjson
                 inflating: mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/SpecimenLab.ndjson
import json
import pandas as pd
import os
def read_ndjson(file_path):
          data = []
          with open(file path, 'r') as file:
                    for line in file:
                             data.append(json.loads(line))
          return data
dataset_dir = '/content/mimic_fhir_dataset/mimic-iv-clinical-database-demo-on-fhir-2.0/mimic-fhir/'
patients = read_ndjson(dataset_dir + 'Patient.ndjson')
conditions = read_ndjson(dataset_dir + 'Condition.ndjson')
encounters = read_ndjson(dataset_dir + 'Encounter.ndjson')
encounters_icu = read_ndjson(dataset_dir + 'EncounterICU.ndjson')
```

```
from datetime import datetime
import csv
def get_condition_timestamp(condition, encounters_data):
    encounter\_id = condition.get('encounter', \{\}).get('reference', '').split('/')[-1]\\
    for encounter in encounters_data:
       if encounter['id'] == encounter_id:
           start_time = encounter.get('period', {}).get('start')
           if start_time:
               timestamp = datetime.fromisoformat(start_time.replace('Z', '+00:00')).timestamp()
               return int(timestamp)
    return None
patient_conditions = {}
for condition in conditions:
   patient_id = condition.get('subject', {}).get('reference', '').split('/')[-1]
    if patient_id not in patient_conditions:
       patient_conditions[patient_id] = []
    condition_timestamp = get_condition_timestamp(condition, encounters)
    if condition_timestamp is not None:
       condition_data = {
           'pid': patient_id,
           'time': condition_timestamp,
            \label{lem:condition} $$ 'description': condition.get('code', {}).get('coding', [{}])[0].get('display', '') $$
       }
       patient_conditions[patient_id].append(condition_data)
print("Step 2: Patient Conditions:")
for patient_id, conditions in patient_conditions.items():
    print(f"Patient ID: {patient_id}")
    for condition in conditions:
       print(condition)
```

```
Patient ID: 94abdf17-f13a-5eae-aac0-eca407bbfadd
     {'pid': '94abdf17-f13a-5eae-aac0-eca407bbfadd',
                                                           'time': 6243609900, 'code': 'I214', 'description': 'Non-ST elevation (NSTEMI) myocarc
                                                          'time': 6243609900, 'code': 'Z7982', 'description': 'Long term (current) use of aspir
     {'pid': '94abdf17-f13a-5eae-aac0-eca407bbfadd',
                                                                                'code': 'Z7902', 'description': 'Long term (current) use of antit 'code': 'I7102', 'description': 'Dissection of abdominal aorta'}
     {'pid': '94abdf17-f13a-5eae-aac0-eca407bbfadd',
                                                          'time': 6243609900,
     {'pid': '94abdf17-f13a-5eae-aac0-eca407bbfadd',
                                                          'time': 6243609900,
     {'pid': '94abdf17-f13a-5eae-aac0-eca407bbfadd', 'time': 6243609900, 'code': 'Z23', 'description': 'Encounter for immunization'}
     {'pid': '94abdf17-f13a-5eae-aac0-eca407bbfadd'
                                                          'time': 6243609900,
                                                                                'code': 'K219', 'description': 'Gastro-esophageal reflux disease
     Patient ID: 5ddeb201-5de6-5177-a116-fa82ce8ad2f2
     {'pid': '5ddeb201-5de6-5177-a116-fa82ce8ad2f2',
                                                           'time': 5545862220, 'code': 'Z85850', 'description': 'Personal history of malignant r
     {'pid': '5ddeb201-5de6-5177-a116-fa82ce8ad2f2',
                                                          'time': 5545862220, 'code': 'G40909', 'description': 'Epilepsy, unspecified, not intr
'time': 5545862220, 'code': 'D1802', 'description': 'Hemangioma of intracranial struc
     {'pid': '5ddeb201-5de6-5177-a116-fa82ce8ad2f2',
                                                          'time': 5545862220,
csv_filename = 'patient_conditions.csv'
csv_columns = ['pid', 'time', 'code', 'description']
with open(csv_filename, 'w', newline='') as csvfile:
    writer = csv.DictWriter(csvfile, fieldnames=csv_columns)
    writer.writeheader()
    for patient_id, conditions in patient_conditions.items():
        for condition in conditions:
             writer.writerow(condition)
print("CSV File Generated:", csv_filename)
     CSV File Generated: patient_conditions.csv
import pandas as pd
df = pd.read_csv(csv_filename)
df
                                               pid
                                                           time
                                                                    code
                                                                                                           description
                                                                                                                           0
            b410dd44-7d65-56f9-974f-2751e8aa80e2 5616018000
                                                                   Z8546
                                                                            Personal history of malignant neoplasm of pros...
                                                                                                                           11.
             b410dd44-7d65-56f9-974f-2751e8aa80e2 5387190060
                                                                   V1072
                                                                                      Personal history of hodgkin's disease
        1
             b410dd44-7d65-56f9-974f-2751e8aa80e2 5639393940
                                                                   Z7902
                                                                             Long term (current) use of antithrombotics/ant...
        2
             b410dd44-7d65-56f9-974f-2751e8aa80e2 5426582400
        3
                                                                   49390
                                                                                      Asthma, unspecified type, unspecified
             b410dd44-7d65-56f9-974f-2751e8aa80e2 5426582400
                                                                    2724
                                                                                      Other and unspecified hyperlipidemia
        4
      4176
              94abdf17-f13a-5eae-aac0-eca407bbfadd 6243609900
                                                                     723
                                                                                               Encounter for immunization
      4177
              94abdf17-f13a-5eae-aac0-eca407bbfadd 6243609900
                                                                    K219
                                                                          Gastro-esophageal reflux disease without esoph...
             5ddeb201-5de6-5177-a116-fa82ce8ad2f2 5545862220
                                                                           Personal history of malignant neoplasm of thyroid
      4178
                                                                  Z85850
             5ddeb201-5de6-5177-a116-fa82ce8ad2f2 5545862220
                                                                  G40909
      4179
                                                                              Epilepsy, unspecified, not intractable, withou...
      4180 5ddeb201-5de6-5177-a116-fa82ce8ad2f2 5545862220
                                                                  D1802
                                                                                     Hemangioma of intracranial structures
     4181 rows × 4 columns
               Generate code with df
                                          View recommended plots
 Next steps:
df.head()
                                                                                                                     ▦
                                            pid
                                                        time
                                                               code
                                                                                                      description
      0 b410dd44-7d65-56f9-974f-2751e8aa80e2 5616018000
                                                              Z8546
                                                                      Personal history of malignant neoplasm of pros...
         b410dd44-7d65-56f9-974f-2751e8aa80e2
                                                 5387190060
                                                              V1072
                                                                                Personal history of hodgkin's disease
      2 b410dd44-7d65-56f9-974f-2751e8aa80e2
                                                 5639393940
                                                             77902
                                                                       Long term (current) use of antithrombotics/ant...
         b410dd44-7d65-56f9-974f-2751e8aa80e2
                                                 5426582400
                                                              49390
                                                                                Asthma, unspecified type, unspecified
      4 b410dd44-7d65-56f9-974f-2751e8aa80e2 5426582400
                                                               2724
                                                                                Other and unspecified hyperlipidemia
 Next steps:
               Generate code with df
                                          View recommended plots
df.tail()
```



Successfully downloaded the csv file : patient_conditions.csv