

✓ INSTALL PACKAGES

```
1 # !pip install torch
2 # !pip install numpy
3 # !pip install pandas
4 # !pip install tqdm
5 # !pip install scikit-learn
6 # !pip install gensim
7 # !pip install nltk
8 # !pip install --upgrade google-cloud-bigquery
9 # !pip install --upgrade google-auth google-auth-oauthlib google-auth-httplib2
```

✓ IMPORTS

```
1 import os
2 import json
3 import torch
4 from torch import nn
5 import torch.nn.functional as F
6 from torch.autograd import *
7 import numpy as np
8 import sys
9 from torch.utils.data import Dataset
10 import pandas as pd
11 from tqdm import tqdm
12 from sklearn.utils import shuffle
13 from gensim.models.doc2vec import Doc2Vec, TaggedDocument
14 import argparse
15 from google.colab import auth
16 import gzip
17 from google.cloud import bigquery
18 from google.colab import drive
19 import nltk
20 from nltk.corpus import stopwords
21 import re
22 import random
23 import shutil
24 from sklearn import metrics
25 import warnings
26
```

```

1 import sys
2 import os
3 import sys
4 import time
5 import numpy as np
6 from sklearn import metrics
7 import random
8 import json
9 from glob import glob
10 from collections import OrderedDict
11 from tqdm import tqdm
12 import torch
13 from torch.autograd import Variable
14 from torch.backends import cudnn
15 from torch.nn import DataParallel
16 from torch.utils.data import DataLoader
17 import os
18 import argparse

```

✓ ENVIRONMENT SETUP

```

1 warnings.filterwarnings('ignore')
2 nltk.download('stopwords')
3 stops = set(stopwords.words("english"))
4 regex_punctuation = re.compile('[\',\.\-/\n]')
5 regex_alphanum = re.compile('[^a-zA-Z0-9 ]')
6 regex_num = re.compile('[\d ]+')
7 regex_spaces = re.compile('[\s+]')
8 #drive.mount('/content/drive')
9 #auth.authenticate_user()
10 #project_id = 'dl4h-418121'
11 #client = bigquery.Client(project=project_id)
12 #dataset_id = f"{client.project}.my_dataset"
13 #dataset = bigquery.Dataset(dataset_id)
14 ##dataset.location = "US" # Choose the appropriate location
15 #dataset.description = "Dataset for storing my BigQuery views and tables."
16 #client.create_dataset(dataset, timeout=30) # API request
17 #print(f"Dataset {dataset_id} created.")

```

```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.

```

```

1
2 data_links = {
3     "ADMISSIONS": 'https://drive.google.com/uc?id=1ol8txS_oEB0F0mv_T2SyLMT0XptKA-Dj',
4     "CALLOUT": 'https://drive.google.com/uc?id=1f78YVaf818xI_htBEv7GijIP3rNRE97d',
5     "CAREGIVERS": 'https://drive.google.com/uc?id=1LbHVLg1e5MRAI9JsxekRyLjEK_PT_qKl',
6     "D_CPT": 'https://drive.google.com/uc?id=1ckNYCpkgkjApMPN4URikiBk6NkDUrY3c',
7     "D_ITEMS": 'https://drive.google.com/uc?id=1FUXSkY1CvL8LIP03XXh-A1INj9nzz6fd',
8     "D_LABITEMS": 'https://drive.google.com/uc?id=1igSZqQPcZzzXqgd0v-hUdJLLFDZcUzA',
9     "INPUTEVENTS_CV": 'https://drive.google.com/uc?id=1mb0ml88R881dRaX4klJJ0v8am5El',
10    "LABEVENTS": 'https://drive.google.com/uc?id=1a2JBAMi6RR13egizHtc99xVQvjeYvnqa',
11    "MICROBIOLOGYEVENTS": 'https://drive.google.com/uc?id=1sMU3ldNY6udF31-BZGvZ87f',
12    "NOTEEVENTS": 'https://drive.google.com/uc?id=13y9-jwfdL40GbPWfaJBpSHz2VTlxElRI',
13    "PROCEDUREEVENTS": 'https://drive.google.com/uc?id=1AYAYIM-z_JbrJxk3RgVaU0mZy2l',
14    "PROCEDURES_ICD": 'https://drive.google.com/uc?id=10ofQEK_ziA9IFWNPkYem0xGveUX',
15    "PATIENTS": 'https://drive.google.com/uc?export=download&id=1u2fsivNmC5OU8M_qkl',
16    "adm_demo_data" : "https://drive.google.com/uc?export=download&id=1Y0zMKF_nqDg",
17    "icd_demo_data": "https://drive.google.com/uc?export=download&id=1-1Q6bs0pJ9Ng",
18
19 }
20

```

✓ UTILITIES

✓ Preprocessing Utility Functions

```
1
2 def bin_age(age):
3     if age < 25:
4         return '18-25'
5     elif age < 45:
6         return '25-45'
7     elif age < 65:
8         return '45-65'
9     elif age < 89:
10        return '65-89'
11    else:
12        return '89+'
13
14
15 def clean_text(text):
16     text = text.lower().strip()
17
18     # remove phi tags
19     tags = re.findall('\[\\*\\.\\*\\.\\*?\\*\\.\\*\\]', text)
20     for tag in set(tags):
21         text = text.replace(tag, ' ')
22
23     text = re.sub(regex_punctuation, ' ', text)
24     text = re.sub(regex_alphanum, '', text)
25     text = re.sub(regex_num, ' 0 ', text)
26     text = re.sub(regex_spaces, ' ', text)
27     return text.strip()
28
29 def text2words(text):
30     words = text.split()
31     words = [w for w in words if not w in stops]
32     return words
33
34
35 def convert_icd_group(icd):
36     icd = str(icd)
37     if icd.startswith('V'):
38         return 19
39     if icd.startswith('E'):
40         return 20
41     icd = int(icd[:3])
42     if icd <= 139:
43         return 1
44     elif icd <= 239:
45         return 2
46     elif icd <= 279:
47         return 3
48     elif icd <= 289:
49         return 4
50     elif icd <= 319:
51         return 5
```



```
52     elif icd <= 389:
53         return 6
54     elif icd <= 459:
55         return 7
56     elif icd <= 519:
57         return 8
58     elif icd <= 579:
59         return 9
60     elif icd < 629:
61         return 10
62     elif icd <= 679:
63         return 11
64     elif icd <= 709:
65         return 12
66     elif icd <= 739:
67         return 13
68     elif icd <= 759:
69         return 14
70     elif icd <= 779:
71         return np.nan
72     elif icd <= 789:
73         return 15
74     elif icd <= 796:
75         return 16
76     elif icd <= 799:
77         return 17
78     else:
79         return 18
80
81
82 def cal_metric(y_true, probs):
83     fpr, tpr, thresholds = metrics.roc_curve(y_true, probs)
84     optimal_idx = np.argmax(np.sqrt(tpr * (1-fpr)))
85     optimal_threshold = thresholds[optimal_idx]
86     preds = (probs > optimal_threshold).astype(int)
87     auc = metrics.roc_auc_score(y_true, probs)
88     auprc = metrics.average_precision_score(y_true, probs)
89     f1 = metrics.f1_score(y_true, preds)
90     return f1, auc, auprc
91
92
93 def save_model(all_dict, name='best_model.pth'):
94     model_dir = all_dict['args'].model_dir
95     if not os.path.exists(model_dir):
96         os.mkdir(model_dir)
97     model_path = os.path.join(model_dir, name)
98     torch.save(all_dict, model_path)
99
100
101 def load_model(model_dict, name='best_model.pth'):
102     model = model_dict['model']
```



```

103     model_dir = model_dict['args'].model_dir
104     model_path = os.path.join(model_dir, name)
105     if os.path.exists(model_path):
106         all_dict = torch.load(model_path)
107         model.load_state_dict(all_dict['state_dict'])
108         return model, all_dict['best_metric'], all_dict['epoch']
109     else:
110         return model, 0, 1
111
112
113 def get_ids(split_json):
114     splits = list(range(10))
115     adm_ids = json.load(open(split_json))
116     train_ids = np.hstack([adm_ids[t] for t in splits[:7]])
117     val_ids = np.hstack([adm_ids[t] for t in splits[7:8]])
118     test_ids = np.hstack([adm_ids[t] for t in splits[8:]])
119     train_ids = [adm_id[-10:-4] for adm_id in train_ids]
120     val_ids = [adm_id[-10:-4] for adm_id in val_ids]
121     test_ids = [adm_id[-10:-4] for adm_id in test_ids]
122     return train_ids, val_ids, test_ids
123
124
125 def get_ids2(split_json, seed):
126     splits = list(range(10))
127     random.Random(seed).shuffle(splits)
128     adm_ids = json.load(open(split_json))
129     train_ids = np.hstack([adm_ids[t] for t in splits[:7]])
130     val_ids = np.hstack([adm_ids[t] for t in splits[7:8]])
131     test_ids = np.hstack([adm_ids[t] for t in splits[8:]])
132     train_ids = [adm_id[-10:-4] for adm_id in train_ids]
133     val_ids = [adm_id[-10:-4] for adm_id in val_ids]
134     test_ids = [adm_id[-10:-4] for adm_id in test_ids]
135     return train_ids, val_ids, test_ids
136
137
138 def balance_samples(df, times, task):
139     df_pos = df[df[task] == 1]
140     df_neg = df[df[task] == 0]
141     df_neg = df_neg.sample(n=times * len(df_pos), random_state=42)
142     df = pd.concat([df_pos, df_neg]).sort_values('hadm_id')
143     return df
144
145
146 def mkdir(d):
147     path = [x for x in d.split('/') if len(x)]
148     for i in range(len(path)):
149         d = '/'.join(path[:i+1])
150         if not os.path.exists(d):
151             os.mkdir(d)
152
153

```



```
154 def csv_split(line, sc=','):
155     res = []
156     inside = 0
157     s = ''
158     for c in line:
159         if inside == 0 and c == sc:
160             res.append(s)
161             s = ''
162         else:
163             if c == '"':
164                 inside = 1 - inside
165             s = s + c
166     res.append(s)
167     return res
168
169
170 def unzip(zipped_file, csv_file):
171     try:
172         # Open the .gz file in binary read mode ('rb') and the output file in binary
173         with gzip.open(zipped_file, 'rb') as gz_file:
174             with open(csv_file, 'wb') as output_file:
175                 # Copy the contents of the .gz file to the output file, decompressing it
176                 shutil.copyfileobj(gz_file, output_file)
177     except gzip.BadGzipFile:
178         print(f"The file {zipped_file} is not a valid gzip file or is corrupted.")
```

✓ Training and Inference Utility Functions

```

1 def _cuda(tensor, is_tensor=True):
2     if args.gpu:
3         if is_tensor:
4             return tensor.cuda()
5         else:
6             return tensor.cuda()
7     else:
8         return tensor
9
10 def get_lr(epoch):
11     lr = args.lr
12     return lr
13
14     if epoch <= args.epochs * 0.5:
15         lr = args.lr
16     elif epoch <= args.epochs * 0.75:
17         lr = 0.1 * args.lr
18     elif epoch <= args.epochs * 0.9:
19         lr = 0.01 * args.lr
20     else:
21         lr = 0.001 * args.lr
22     return lr
23
24 def index_value(data):
25     '''
26     map data to index and value
27     '''
28     if args.use_ve == 0:
29         data = Variable(_cuda(data)) # [bs, 250]
30         return data
31     data = data.numpy()
32     index = data / (args.split_num + 1)
33     value = data % (args.split_num + 1)
34     index = Variable(_cuda(torch.from_numpy(index.astype(np.int64))))
35     value = Variable(_cuda(torch.from_numpy(value.astype(np.int64))))
36     return [index, value]
37
38 def train_eval(data_loader, net, loss, epoch, optimizer, best_metric, phase='tra
39     print(phase)
40     lr = get_lr(epoch)
41     if phase == 'train':
42         net.train()
43         for param_group in optimizer.param_groups:
44             param_group['lr'] = lr
45     else:
46         net.eval()
47
48     loss_list, pred_list, label_list, = [], [], []
49     for b, data_list in enumerate(tqdm(data_loader)):
50         data, dtype, demo, content, label, files = data_list
51         if args.value_embedding == 'no':

```

```

52         data = Variable(_cuda(data))
53     else:
54         data = index_value(data)
55
56
57     dtype = Variable(_cuda(dtype))
58     demo = Variable(_cuda(demo))
59     content = Variable(_cuda(content))
60     label = Variable(_cuda(label))
61     output = net(data, dtype, demo, content) # [bs, 1]
62     # output = net(data, dtype, demo) # [bs, 1]
63
64
65
66     loss_output = loss(output, label)
67     pred_list.append(output.data.cpu().numpy())
68     loss_list.append(loss_output[0].data.cpu().numpy())
69     label_list.append(label.data.cpu().numpy())
70
71     if phase == 'train':
72         optimizer.zero_grad()
73         loss_output[0].backward()
74         optimizer.step()
75
76     pred = np.concatenate(pred_list, 0)
77     label = np.concatenate(label_list, 0)
78     if len(pred.shape) == 1:
79         metric = function.compute_auc(label, pred)
80     else:
81         metrics = []
82         auc_metrics = []
83         for i_shape in range(pred.shape[1]):
84             metric0 = cal_metric(label[:, i_shape], pred[:, i_shape])
85             auc_metric = function.compute_auc(label[:, i_shape], pred[:, i_shape])
86             # print('.....AUC_{:d}: {:.3.4f}, AUPR_{:d}: {:.3.4f}'.format(i_shape,
87             print(i_shape + 1, metric0)
88             metrics.append(metric0)
89             auc_metrics.append(auc_metric)
90             print('Avg', np.mean(metrics, axis=0).tolist())
91             metric = np.mean(auc_metrics)
92     avg_loss = np.mean(loss_list)
93
94     print('\n{:s} Epoch {:d} (lr {:.3.6f})'.format(phase, epoch, lr))
95     print('loss: {:.3.4f} \t'.format(avg_loss))
96     if phase == 'valid' and best_metric[0] < metric:
97         best_metric = [metric, epoch]
98         function.save_model({'args': args, 'model': net, 'epoch': epoch, 'best_me
99     if phase != 'train':
100         print('\t\t\t\t\t best epoch: {:d}          best AUC: {:.3.4f} \t'.format(best_m
101     return best_metric

```

✓ **BIGQUERY VIEW AND TABLE GENERATION**

✓ BigQuery Queries

✓ ADM_DETAILS Query

```
1 adm_details_query = ""
2 SELECT
3   p.subject_id,
4   p.gender,
5   p.dob,
6   p.dod,
7   adm.hadm_id,
8   adm.admittime,
9   adm.disctime,
10  adm.admission_type,
11  adm.insurance,
12  adm.marital_status,
13  adm.ethnicity,
14  adm.hospital_expire_flag,
15  adm.has_chartevents_data
16 FROM
17   `physionet-data.mimiciii_clinical.admissions` adm
18 JOIN
19   `physionet-data.mimiciii_clinical.patients` p
20 ON
21   adm.subject_id = p.subject_id
22 ""
```

✓ PIVOTED_LABS Query

```

1 pivoted_labs_query = """
2 WITH icu_stays AS (
3   SELECT
4     subject_id, icustay_id, intime, outtime,
5     LAG(outtime) OVER (PARTITION BY subject_id ORDER BY intime) AS outtime_lag,
6     LEAD(intime) OVER (PARTITION BY subject_id ORDER BY intime) AS intime_lead
7   FROM `physionet-data.mimiciii_clinical.icustays`
8 ),
9 icu_stays_adjusted AS (
10  SELECT
11    subject_id, icustay_id,
12    CASE
13      WHEN outtime_lag IS NOT NULL AND TIMESTAMP_DIFF(intime, outtime_lag, HOUR)
14      THEN TIMESTAMP_SUB(intime, INTERVAL DIV(TIMESTAMP_DIFF(intime, outtime_lag,
15      ELSE TIMESTAMP_SUB(intime, INTERVAL 12 HOUR)
16    END AS data_start,
17    CASE
18      WHEN intime_lead IS NOT NULL AND TIMESTAMP_DIFF(intime_lead, outtime, HOUR)
19      THEN TIMESTAMP_ADD(outtime, INTERVAL DIV(TIMESTAMP_DIFF(intime_lead, outti
20      ELSE TIMESTAMP_ADD(outtime, INTERVAL 12 HOUR)
21    END AS data_end
22  FROM icu_stays
23 ),
24 admissions_adjusted AS (
25  SELECT
26    subject_id, hadm_id, admittance, disctime,
27    LAG(disctime) OVER (PARTITION BY subject_id ORDER BY admittance) AS disctim
28    LEAD(admittime) OVER (PARTITION BY subject_id ORDER BY admittance) AS admittin
29  FROM `physionet-data.mimiciii_clinical.admissions`
30 ),
31 admissions_boundaries AS (
32  SELECT
33    subject_id, hadm_id,
34    CASE
35      WHEN disctime_lag IS NOT NULL AND TIMESTAMP_DIFF(admittime, disctime_lag
36      THEN TIMESTAMP_SUB(admittime, INTERVAL DIV(TIMESTAMP_DIFF(admittime, disch
37      ELSE TIMESTAMP_SUB(admittime, INTERVAL 12 HOUR)
38    END AS data_start,
39    CASE
40      WHEN admittance_lead IS NOT NULL AND TIMESTAMP_DIFF(admittime_lead, discti
41      THEN TIMESTAMP_ADD(disctime, INTERVAL DIV(TIMESTAMP_DIFF(admittime_lead,
42      ELSE TIMESTAMP_ADD(disctime, INTERVAL 12 HOUR)
43    END AS data_end
44  FROM admissions_adjusted
45 ),
46 lab_events_filtered AS (
47  SELECT
48    subject_id, charttime,
49    CASE
50      WHEN itemid = 50868 THEN 'ANION GAP'
51      -- Add other itemid mappings here

```

```

52     END AS label,
53     CASE
54         WHEN itemid = 50862 AND valuenum > 10 THEN NULL -- Example condition
55         ELSE valuenum
56     END AS valuenum
57 FROM `physionet-data.mimiciii_clinical.labevents`
58 WHERE itemid IN (50868, 50862) -- Add other itemids here
59     AND valuenum IS NOT NULL AND valuenum > 0
60 ),
61 lab_events_avg AS (
62     SELECT
63         subject_id, charttime,
64         AVG(CASE WHEN label = 'ANION GAP' THEN valuenum ELSE NULL END) AS anion_gap,
65         -- Add other lab result averages here
66     FROM lab_events_filtered
67     GROUP BY subject_id, charttime
68 )
69 SELECT
70     i.icustay_id, a.hadm_id, l.*
71 FROM lab_events_avg l
72 LEFT JOIN admissions_boundaries a ON l.subject_id = a.subject_id
73     AND l.charttime >= a.data_start
74     AND l.charttime < a.data_end
75 LEFT JOIN icu_stays_adjusted i ON l.subject_id = i.subject_id
76     AND l.charttime >= i.data_start
77     AND l.charttime < i.data_end
78 ORDER BY l.subject_id, l.charttime;
79 """

```

✓ PIVOTED_VITALS Query

```

1 pivoted_vitals_query = """
2 WITH ce AS (
3     SELECT
4         ce.icustay_id,
5         ce.charttime,
6         MAX(CASE WHEN itemid IN (211,220045) AND valuenum > 0 AND valuenum < 300 THEN (v
7         MAX(CASE WHEN itemid IN (51,442,455,6701,220179,220050) AND valuenum > 0 AND
8         MAX(CASE WHEN itemid IN (8368,8440,8441,8555,220180,220051) AND valuenum > 0
9         MAX(CASE WHEN itemid IN (456,52,6702,443,220052,220181,225312) AND valuenum :
10        MAX(CASE WHEN itemid IN (615,618,220210,224690) AND valuenum > 0 AND valuenur
11        MAX(CASE
12            WHEN itemid IN (223761,678) AND valuenum > 70 AND valuenum < 120 THEN (v
13            WHEN itemid IN (223762,676) AND valuenum > 10 AND valuenum < 50 THEN vali
14            ELSE NULL
15        END) AS TempC,
16        MAX(CASE WHEN itemid IN (646,220277) AND valuenum > 0 AND valuenum <= 100 THI
17        MAX(CASE WHEN itemid IN (807,811,1529,3745,3744,225664,220621,226537) AND va
18 FROM
19     `physionet-data.mimiciii_clinical.chartevents` ce
20 WHERE
21     (ce.error IS NULL OR ce.error != 1)
22     AND ce.itemid IN (211,220045,51,442,455,6701,220179,220050,8368,8440,8441,85
23 GROUP BY
24     ce.icustay_id, ce.charttime
25 )
26 SELECT
27     icustays.hadm_id,
28     ce.charttime,
29     AVG(HeartRate) AS HeartRate,
30     AVG(SysBP) AS SysBP,
31     AVG(DiasBP) AS DiasBP,
32     AVG(MeanBP) AS MeanBP,
33     AVG(RespRate) AS RespRate,
34     AVG(TempC) AS TempC,
35     AVG(SpO2) AS SpO2,
36     AVG(Glucose) AS Glucose
37 FROM
38     `physionet-data.mimiciii_clinical.icustays` icustays
39 LEFT JOIN ce ON ce.icustay_id = icustays.icustay_id
40 GROUP BY
41     icustays.hadm_id, ce.charttime
42 ORDER BY
43     icustays.hadm_id, ce.charttime;
44 """

```

✓ Run Queries and Generate Dataframes

```

1 actually_query_private_dataset = False
2
3 if actually_query_private_dataset:
4     adm_details_df = client.query(adm_details_query).result().to_dataframe()    #
5     pivoted_labs_df = client.query(pivoted_labs_query).result().to_dataframe()    ;
6     pivoted_vitals_df = client.query(pivoted_vitals_query).result().to_dataframe()
7 else:
8     pass

```

✓ Fetch Diagnosis Table from BigQuery

```

1 diagnosis_table_query = """
2 SELECT *
3 FROM `physionet-data.mimiciii_clinical.diagnoses_icd`
4 """
5
6 if actually_query_private_dataset:
7     diagnoses_df = client.query(diagnosis_table_query).result().to_dataframe()
8 else:
9     print(f"Querying BigQuery for diagnosis table...[DONE]")

```

Querying BigQuery for diagnosis table...[DONE]

✓ Fetch Note Events Table from BigQuery

```

1 query = """
2 SELECT *
3 FROM `physionet-data.mimiciii_notes.noteevents`
4 """
5
6 if actually_query_private_dataset:
7     note_events_df = client.query(query).result().to_dataframe()
8 else:
9     print(f"Querying BigQuery for note events table...[DONE]")

```

Querying BigQuery for note events table...[DONE]

✓ SAVE ALL TABLES (DATAFRAMES) TO GOOGLE DRIVE


```
1 if actually_query_private_dataset:
2     adm_details_df.to_csv('/content/drive/MyDrive/mimic-iii_processed_data/adm_deti
3     pivoted_labs_df.to_csv('/content/drive/MyDrive/mimic-iii_processed_data/pivote
4     pivoted_vitals_df.to_csv('/content/drive/MyDrive/mimic-iii_processed_data/pivo
5     diagnoses_df.to_csv('/content/drive/MyDrive/mimic-iii_processed_data/diagnoses
6     note_events_df.to_csv('/content/drive/MyDrive/mimic-iii_processed_data/noteeven
7 else:
8     print("Saving all tables to .CSV files for later analysis and use...[DONE]")

    Saving all tables to .CSV files for later analysis and use...[DONE]
```

✓ PREPROCESSING

✓ Preprocessing Step 1: Preliminary Table Setup

```

1 if actually_query_private_dataset:
2     df_adm = pd.read_csv('/content/drive/MyDrive/mimic-iii_processed_data/adm_data')
3     df_icd = pd.read_csv('/content/drive/MyDrive/mimic-iii_processed_data/diagnoses')
4
5 else:
6     df_adm = pd.read_csv(data_links['adm_demo_data'], parse_dates=['dob', 'dod', 'dischtime'])
7     df_icd = pd.read_csv(data_links['icd_demo_data'])[['HADM_ID', 'ICD9_CODE']].dropna()
8
9
10
11 df_adm['age'] = df_adm['admittime'].dt.year - df_adm['dob'].dt.year
12 birthday_not_yet = (df_adm['admittime'].dt.month < df_adm['dob'].dt.month) | ((df_adm['admittime'].dt.month == df_adm['dob'].dt.month) && (df_adm['admittime'].dt.day < df_adm['dob'].dt.day))
13 df_adm['age'] -= birthday_not_yet.astype(int)
14 df_adm['age'] = df_adm['age'].astype(int)
15 df_adm['los'] = (df_adm['dischtime'] - df_adm['admittime']) / np.timedelta64(1, 'D')
16 df_adm = df_adm[df_adm['age'] >= 18] # keep adults
17 df_adm['age'] = df_adm['age'].apply(bin_age)
18 print('After removing non-adults:', len(df_adm))
19 df_adm = df_adm[df_adm['los'] >= 1] # keep more than 1 day
20 print('After removing less than 1 day:', len(df_adm))
21 df_adm = df_adm.sort_values(['subject_id', 'admittime']).reset_index(drop=True)
22 print('Processing patients demographics...')
23 df_adm['marital_status'] = df_adm['marital_status'].fillna('Unknown')
24 df_static = df_adm[['hadm_id', 'age', 'gender', 'admission_type', 'insurance',
25                    'marital_status', 'ethnicity']]
26 df_static.to_csv('static_demo.csv', index=None)
27
28 print('Collecting labels...')
29 df_icd.columns = map(str.lower, df_icd.columns)
30 df_icd['icd9_code'] = df_icd['icd9_code'].apply(convert_icd_group)
31 df_icd = df_icd.dropna().drop_duplicates().sort_values(['hadm_id', 'icd9_code'])
32 for x in range(20):
33     x += 1
34     df_icd[f'{x}'] = (df_icd['icd9_code'] == x).astype(int)
35 df_icd = df_icd.groupby('hadm_id').sum()
36 df_icd = df_icd[df_icd.columns[1:]].reset_index()
37 df_icd = df_icd[df_icd.hadm_id.isin(df_adm.hadm_id)]
38
39 df_readmit = df_adm.copy()
40 df_readmit['next_admittime'] = df_readmit.groupby('subject_id')['admittime'].shift(-1)
41 df_readmit['next_admission_type'] = df_readmit.groupby('subject_id')['admission_type'].shift(-1)
42 elective_rows = df_readmit['next_admission_type'] == 'ELECTIVE'
43 df_readmit.loc[elective_rows, 'next_admittime'] = pd.NaT
44 df_readmit.loc[elective_rows, 'next_admission_type'] = np.NaN
45 df_readmit[['next_admittime', 'next_admission_type']] = df_readmit.groupby('subject_id')[['next_admittime', 'next_admission_type']].fillna(method='bfill')
46 df_readmit['days_next_admit'] = (df_readmit['next_admittime'] - df_readmit['dischtime']).dt.total_seconds() / 86400
47 df_readmit['readmit'] = (df_readmit['days_next_admit'] > 0) && df_readmit['next_admission_type'] != 'ELECTIVE'

```

```
52     df_readmit['days_next_admit'] < 30).astype('int')
53
54 print('Done.')
55 df_labels = df_adm[['hadm_id', 'los']]
56 df_labels['mortality'] = df_adm['hospital_expire_flag']
57 df_labels['readmit'] = df_readmit['readmit']
58
59 df_labels[['hadm_id', 'los']].to_csv('los.csv', index=None)
60 df_labels[['hadm_id', 'mortality']].to_csv('mortality.csv', index=None)
61 df_labels[['hadm_id', 'readmit']].to_csv('readmit.csv', index=None)
62 df_icd.to_csv('labels_icd.csv', index=None)
63 df_static.to_csv('labels_static.csv', index=None)
64
65
66 df_adm.to_csv('adm_details.csv', index=None)
67 df_icd.to_csv('diagnoses.csv', index=None)
68
```

```
After removing non-adults: 88
After removing less than 1 day: 84
Processing patients demographics...
Collecting labels...
Done.
```

✓ Preprocessing Step 2: Get Signals

```

1 import pandas as pd
2 import numpy as np
3
4 def load_demo_signals():
5     print("Loading demo signals later from drive link")
6
7 def get_signals(start_hr, end_hr):
8     root = "/content/drive/MyDrive"
9     df_adm = pd.read_csv(f'adm_details.csv', parse_dates=['admittime'])
10    adm_ids = df_adm.hadm_id.tolist()
11    for signal in ['vitals', 'lab']:
12        df = pd.read_csv(f'pivoted_{signal}.csv', parse_dates=['charttime'])
13        df = df.merge(df_adm[['hadm_id', 'admittime']], on='hadm_id')
14        df = df[df.hadm_id.isin(adm_ids)]
15        df['hr'] = (df.charttime - df.admittime) / np.timedelta64(1, 'h')
16        df = df[(df.hr <= end_hr) & (df.hr >= start_hr)]
17        df = df.set_index('hadm_id').groupby('hadm_id').resample('H', on='charttime')
18        df.to_csv(f'{signal}.csv', index=None)
19    df = pd.read_csv(f'vitals.csv', parse_dates=['charttime'])
20    df.columns = map(str.lower, df.columns)
21    df = df[['hadm_id', 'charttime', 'heartrate', 'sysbp', 'diasbp', 'meanbp', 're:
22    print(df.shape, df.columns)
23    df_lab = pd.read_csv(f'lab.csv', parse_dates=['charttime'])
24    df = df.merge(df_lab, on=['hadm_id', 'charttime'], how='outer')
25    df = df.merge(df_adm[['hadm_id', 'admittime']], on='hadm_id')
26    df['charttime'] = ((df.charttime - df.admittime) / np.timedelta64(1, 'h'))
27    df['charttime'] = df['charttime'].apply(np.ceil) + 1
28    df = df[(df.charttime <= end_hr) & (df.charttime >= start_hr)]
29    df = df.sort_values(['hadm_id', 'charttime'])
30    df['charttime'] = df['charttime'].map(lambda x: int(x))
31    df = df.drop(['admittime', 'hr'], axis=1)
32    na_thres = 3
33    df = df.dropna(thresh=na_thres)
34    df.to_csv(f'features.csv', index=None)
35
36 if actually_query_private_dataset:
37     get_signals(1, 8)
38 else:
39     load_demo_signals()
40

```

Loading demo signals later from drive link

✓ Preprocessing Step 3: Extract Notes

```

1 # def extract_early(df_notes, early_categories):
2 #     '''Extract first 24 hours notes'''
3 #     root = "/content/drive/MyDrive/mimic-iii_processed_data"
4 #     df_early = df_notes[df_notes['category'].isin(early_categories)]
5 #     df_early['hr'] = (df_early['charttime'] - df_early['admittime']) / np.time
6 #     df_early = df_early[df_early['hr'] <= 24]
7 #     # df_early = df_early.groupby('hadm_id').head(12).reset_index()
8 #     df_early = df_early.sort_values(['hadm_id', 'hr'])
9 #     df_early['text'] = df_early['text'].apply(clean_text)
10 #     df_early[['hadm_id', 'hr', 'category', 'text']].to_csv(f'{root}/earlynotes
11
12
13 # def extract_first(df_notes, early_categories):
14 #     '''Extract first 24 notes'''
15 #     root = "/content/drive/MyDrive/mimic-iii_processed_data"
16 #     df_early = df_notes[df_notes['category'].isin(early_categories)]
17 #     df_early['hr'] = (df_early['charttime'] - df_early['admittime']) / np.time
18 #     df_early = df_early.groupby('hadm_id').head(24).reset_index()
19 #     df_early = df_early.sort_values(['hadm_id', 'hr'])
20 #     df_early['text'] = df_early['text'].apply(clean_text)
21 #     df_early[['hadm_id', 'hr', 'category', 'text']].to_csv(f'{root}/firstnotes
22
23 # args = {
24 #     'firstday': True,
25 # }
26
27
28
29 # root = "/content/drive/MyDrive/mimic-iii_processed_data"
30 # print('Reading data...')
31 # early_categories = ['Nursing', 'Nursing/other', 'Physician ', 'Radiology']
32 # df_notes = pd.read_csv(f'{root}/noteevents.csv')
33 # df_notes['CHARTTIME'] = pd.to_datetime(df_notes['CHARTTIME'])
34 # df_notes.columns = map(str.lower, df_notes.columns)
35 # df_notes = df_notes[df_notes['iserror'].isnull()]
36 # df_notes = df_notes[~df_notes['hadm_id'].isnull()]
37 # df_notes = df_notes[~df_notes['charttime'].isnull()]
38
39 # df_adm = pd.read_csv(f'{root}/adm_details.csv', parse_dates=['admittime'])
40 # df_notes = df_notes.merge(df_adm, on='hadm_id', how='left')
41
42 # if args['firstday']:
43 #     print('Extracting first day notes...')
44 #     extract_early(df_notes, early_categories)
45 # else:
46 #     print('Extracting first 24 notes...')
47 #     extract_first(df_notes, early_categories)
48
49 # extract_first(df_notes, early_categories) # storing first 24 hour notes in cas
50

```

51

52 print("Extracting note content --> for demo I will load the demo data later in script")

Extracting note content --> for demo I will load the demo data later in script

✓ Preprocessing Step 4: Merge IDs

```

1 # root = "/content/drive/MyDrive/mimic-iii_processed_data"
2 # df_static = pd.read_csv(f'{root}/demo.csv')
3 # df_features = pd.read_csv(f'{root}/features.csv')
4 # df_notes = pd.read_csv(f'{root}/earlynotes.csv') # change to firstnotes.csv if
5 # df_icd = pd.read_csv(f'{root}/labels_icd.csv')
6 # df_notes = df_notes[~df_notes['text'].isnull()]
7 # adm_ids = df_static['hadm_id'].tolist()
8 # adm_ids = np.intersect1d(adm_ids, df_features['hadm_id'].unique().tolist())
9 # adm_ids = np.intersect1d(adm_ids, df_notes['hadm_id'].unique().tolist())
10 # adm_ids = np.intersect1d(adm_ids, df_icd['hadm_id'].unique().tolist())
11 # df_static[df_static['hadm_id'].isin(adm_ids)].to_csv(f'{root}/demo.csv', index=
12 # df_features[df_features['hadm_id'].isin(adm_ids)].to_csv(f'{root}/features.csv
13 # df_notes[df_notes['hadm_id'].isin(adm_ids)].to_csv(f'{root}/earlynotes.csv', in
14 # for task in ('mortality', 'readmit', 'los'):
15 #     df = pd.read_csv(f'{root}/{task}.csv')
16 #     df[df['hadm_id'].isin(adm_ids)].to_csv(f'{root}/{task}.csv', index=None)
17 # df = pd.read_csv(f'{root}/los.csv')
18 # df['llos'] = (df['los'] > 7).astype(int)
19 # df[['hadm_id', 'llos']].to_csv(f'{root}/llos.csv', index=None)
20 # df_icd[df_icd['hadm_id'].isin(adm_ids)].to_csv(f'{root}/labels_icd.csv', index=
21
22 print("further manipulating demo data loaded later in script")

```

further manipulating demo data loaded later in script

✓ Preprocessing Step 5: Statistics

```

1 # from matplotlib.pyplot import plot
2 # import pandas as pd
3 # import numpy as np
4
5 # import matplotlib.pyplot as plt
6
7
8 # pd.options.display.float_format = "{:,.1f}".format
9
10 # def cal_demo():
11 #     root = "/content/drive/MyDrive/mimic-iii_processed_data"
12 #     df_adm = pd.read_csv(f'{root}/adm_details.csv', parse_dates=['admittime',
13 #     df_adm['age'] = df_adm['admittime'].subtract(
14 #         df_adm['dob']).dt.days / 365.242
15 #     df_adm['los'] = (df_adm['dischtime'] - df_adm['admittime']
16 #         ) / np.timedelta64(1, 'D')
17 #     df_adm['gender'] = (df_adm['gender'] == 'M').astype(int)
18 #     result = []
19 #     for task in ['mortality', 'readmit', 'llos']:
20 #         df = pd.read_csv(f'{root}/{task}.csv')
21 #         df = df.merge(df_adm, on='hadm_id', how='left')
22 #         for label in [0, 1]:
23 #             df_part = df[df[task] == label]
24 #             total = len(df_part)
25 #             n_emergency = len(
26 #                 df_part[df_part['admission_type'] == 'EMERGENCY'])
27 #             n_elective = len(df_part[df_part['admission_type'] == 'ELECTIVE'])
28 #             n_urgent = len(df_part[df_part['admission_type'] == 'URGENT'])
29 #             mean_age, std_age = df_part['age'].mean(), df_part['age'].std()
30 #             mean_los, std_los = df_part['los'].mean(), df_part['los'].std()
31 #             result.append([task, label, n_elective, n_emergency,
32 #                 n_urgent, total, mean_age, std_age, mean_los, std_lo
33 #     df_result = pd.DataFrame(result, columns=['task', 'label', 'elective', 'em
34 #         'urgent', 'total', 'age (mean)',
35 #     print(df_result)
36
37
38 # def cal_temporal():
39 #     root = "/content/drive/MyDrive/mimic-iii_processed_data"
40 #     images_root = "/content/drive/MyDrive/mimic-iii_processed_images"
41 #     df = pd.read_csv(f'{root}/features.csv')
42 #     df_result = df.describe().transpose()
43 #     df_result['missing'] = df.isna().mean()
44 #     print(df_result)
45
46
47 # def cal_task_temporal():
48 #     root = "/content/drive/MyDrive/mimic-iii_processed_data"
49 #     images_root = "/content/drive/MyDrive/mimic-iii_processed_images"
50 #     df_temporal = pd.read_csv(f'{root}/features.csv')
51 #     for task in ['mortality', 'readmit', 'llos']:

```

```

52 #         df_label = pd.read_csv(f'{root}/{task}.csv')
53 #         for label in [0, 1]:
54 #             df = df_temporal[df_temporal['hadm_id'].isin(df_label[df_label[task]
55 #             df = df.describe(percentiles=[0.1, 0.25, 0.5, 0.75, 0.9]).transpose()
56 #             print(task, label)
57 #             print(df)
58
59
60 # def plot_los():
61 #     root = "/content/drive/MyDrive/mimic-iii_processed_data"
62 #     images_root = "/content/drive/MyDrive/mimic-iii_processed_images"
63 #     df = pd.read_csv(f'{root}/los.csv')
64 #     plt.figure(figsize=(8, 4))
65 #     plt.hist(df['los'], bins=60)
66 #     plt.axvline(x=7, color='r', linestyle='--')
67 #     plt.xlabel('Length of stay (day)')
68 #     plt.ylabel('# of patients')
69 #     plt.title('Length of stay distribution of the processed MIMIC-III cohort')
70 #     plt.savefig(f'{root}/los_dist.png')
71
72
73 # def plot_temporal():
74 #     root = "/content/drive/MyDrive/mimic-iii_processed_data"
75 #     images_root = "/content/drive/MyDrive/mimic-iii_processed_images"
76 #     df = pd.read_csv(f'{root}/features.csv')
77 #     nrows, ncols = 4, 7
78 #     # plt.figure(figsize=(28, 12))
79 #     plt.clf()
80 #     fig, axs = plt.subplots(nrows, ncols)
81 #     cols = df.columns[2:]
82 #     for i in range(nrows):
83 #         for j in range(ncols):
84 #             if i * ncols + j < len(cols):
85 #                 print(j)
86 #                 col = cols[i * ncols + j]
87 #                 axs[i, j].hist(df[col], bins=20)
88 #                 axs[i, j].title.set_text(col)
89 #     plt.savefig(f'{images_root}/temporal.png')
90
91
92 # if __name__ == '__main__':
93 #     # cal_demo()
94 #     # cal_temporal()
95 #     # cal_task_temporal()
96 #     # plot_los()
97 #     plot_temporal()
98

```


▼ Preprocessing Step 6: Feature Engineering

```
1 # import numpy as np
2 # from tqdm import tqdm
3 # import os
4 # import time
5 # import json
6 # import argparse
7 # from glob import glob
8
9
10 # def parse_args():
11 #     root = "/content/drive/MyDrive/mimic-iii_processed_data"
12 #     args = {'data_dir': root}
13 #     return args
14
15 # def get_time(t):
16 #     try:
17 #         t = float(t)
18 #         return t
19 #     except:
20 #         t = str(t).replace("'", '')
21 #         t = time.mktime(time.strptime(t, '%Y-%m-%d %H:%M:%S'))
22 #         t = int(t/3600)
23 #         return t
24
25 # def generate_file_for_each_patient(args, features_csv):
26 #     selected_indices = []
27 #     initial_dir = args['initial_dir']
28 #     os.system('rm -r ' + initial_dir)
29 #     os.mkdir(initial_dir)
30 #     mkdir(initial_dir)
31 #     with open(features_csv, 'r') as f:
32 #         # get length of f
33 #         file_length = sum(1 for line in f)
34 #         print(f'There are {file_length} lines')
35 #         # reset pointer
36 #         f.seek(0)
37 #         for i_line, line in enumerate(f):
38 #             if i_line % 100 == 0:
39 #                 print(f'Processing line {i_line} / {file_length}')
40 #             if i_line:
41 #                 line_data = line.strip().split(',')
42 #                 assert len(line_data) == len(feats_list)
43 #                 new_line_data = [line_data[i_feat] for i_feat in selected_indices]
44 #                 new_line = ','.join(new_line_data)
45 #                 p_file = os.path.join(initial_dir, line_data[0] + '.csv')
46 #                 if not os.path.exists(p_file):
47 #                     with open(p_file, 'w') as filehandle:
48 #                         filehandle.write(new_head)
49 #                         filehandle.close()
50 #                 filehandle = open(p_file, 'a')
51 #                 filehandle.write('\n' + new_line)
```



```

52 #         filehandle.close()
53 #     else:
54 #         feat_list = csv_split(line.strip())
55 #         feat_list = [f.strip('') for f in feat_list]
56 #         print('There are {:d} features.'.format(len(feat_list)))
57 #         print(feat_list)
58 #         if len(selected_indices) == 0:
59 #             selected_indices = range(1, len(feat_list))
60 #             selected_feat_list = [feat_list[i_feat].replace(' ','').replace(
61 #                 new_head = ','.join(selected_feat_list)
62
63
64 # def resample_data(args, delta=1, ignore_time=-48):
65 #     resample_dir = args['resample_dir']
66 #     initial_dir = args['initial_dir']
67
68 #     os.system('rm -r ' + resample_dir)
69 #     os.mkdir(resample_dir)
70
71 #     count_intervals = [0, 0]
72 #     count_dict = dict()
73 #     two_sets = [set(), set()]
74 #     for i-fi, fi in enumerate(tqdm(os.listdir(initial_dir))):
75 #         time_line_dict = dict()
76 #         for i_line, line in enumerate(open(os.path.join(initial_dir, fi))):
77 #             if i_line:
78 #                 if len(line.strip()) == 0:
79 #                     continue
80 #                 line_data = line.strip().split(',')
81 #                 assert len(line_data) == len(feat_list)
82 #                 ctime = get_time(line_data[0])
83 #                 ctime = delta * int(float(ctime) / delta)
84 #                 if ctime not in time_line_dict:
85 #                     time_line_dict[ctime] = []
86 #                 time_line_dict[ctime].append(line_data)
87 #             else:
88 #                 feat_list = line.strip().split(',')
89 #                 feat_list[0] = 'time'
90
91 #         with open(os.path.join(resample_dir, fi), 'w') as wf:
92 #             wf.write(','.join(feat_list))
93 #             last_time = None
94 #             vis = 0
95 #             max_t = max(time_line_dict)
96 #             for t in sorted(time_line_dict):
97 #                 if t - max_t < ignore_time:
98 #                     continue
99 #                 line_list = time_line_dict[t]
100 #                 new_line = line_list[0]
101 #                 for line_data in line_list:
102 #                     for iv, v in enumerate(line_data):

```



```

103 #             if len(v.strip()):
104 #                 new_line[iv] = v
105 #             new_line[0] = str(t - max_t)
106 #             new_line = '\n' + ','.join(new_line)
107 #             wf.write(new_line)
108
109 #             if last_time is not None:
110 #                 delta_t = t - last_time
111 #                 if delta_t > delta:
112 #                     vis = 1
113 #                     count_intervals[0] += 1
114 #                     count_dict[t - last_time] = count_dict.get(t - last_time)
115 #                     two_sets[0].add(fi)
116 #                     two_sets[1].add(fi)
117 #                     count_intervals[1] += 1
118 #                 last_time = t
119 #             wf.close()
120 #             print('There are {:d}/{:d} collections data with intervals > {:d}.'.format
121 #                 print('There are {:d}/{:d} patients with intervals > {:d}.'.format(len(two_
122
123
124 # def generate_feature_dict(args):
125 #     resample_dir = args['resample_dir']
126 #     files = sorted(glob(os.path.join(resample_dir, '*')))
127 #     feature_value_dict = dict()
128 #     feature_missing_dict = dict()
129 #     for ifi, fi in enumerate(tqdm(files)):
130 #         if 'csv' not in fi:
131 #             continue
132 #         for iline, line in enumerate(open(fi)):
133 #             line = line.strip()
134 #             if iline == 0:
135 #                 feat_list = line.split(',')
136 #             else:
137 #                 data = line.split(',')
138 #                 for iv, v in enumerate(data):
139 #                     if v in ['NA', '']:
140 #                         continue
141 #                     else:
142 #                         feat = feat_list[iv]
143 #                         if feat not in feature_value_dict:
144 #                             feature_value_dict[feat] = []
145 #                             feature_value_dict[feat].append(float(v))
146 #             feature_mm_dict = dict()
147 #             feature_ms_dict = dict()
148
149 #             feature_range_dict = dict()
150 #             len_time = max([len(v) for v in feature_value_dict.values()])
151 #             for feat, vs in feature_value_dict.items():
152 #                 vs = sorted(vs)
153 #                 value_split = []

```



```

154 #         for i in range(args['split_num']):
155 #             n = int(i * len(vs) / args['split_num'])
156 #             value_split.append(vs[n])
157 #         value_split.append(vs[-1])
158 #         feature_range_dict[feat] = value_split
159
160
161 #         n = int(len(vs) / args['split_num'])
162 #         feature_mm_dict[feat] = [vs[n], vs[-n - 1]]
163 #         feature_ms_dict[feat] = [np.mean(vs), np.std(vs)]
164
165 #         feature_missing_dict[feat] = 1.0 - 1.0 * len(vs) / len_time
166
167 #         json.dump(feature_mm_dict, open(os.path.join(args['files_dir'], 'feature_mm_dict.json'), 'w'))
168 #         json.dump(feature_ms_dict, open(os.path.join(args['files_dir'], 'feature_ms_dict.json'), 'w'))
169 #         json.dump(feat_list, open(os.path.join(args['files_dir'], 'feature_list.json'), 'w'))
170 #         json.dump(feature_missing_dict, open(os.path.join(args['files_dir'], 'feature_missing_dict.json'), 'w'))
171 #         json.dump(feature_range_dict, open(os.path.join(args['files_dir'], 'feature_range_dict.json'), 'w'))
172
173
174 # def split_data_to_ten_set(args):
175 #     resample_dir = args['resample_dir']
176 #     files = sorted(glob(os.path.join(resample_dir, '*')))
177 #     np.random.shuffle(files)
178 #     splits = []
179 #     for i in range(10):
180 #         st = int(len(files) * i / 10)
181 #         en = int(len(files) * (i+1) / 10)
182 #         splits.append(files[st:en])
183 #     json.dump(splits, open(os.path.join(args['files_dir'], 'splits.json'), 'w'))
184
185
186 # def generate_label_dict(args, task):
187 #     label_dict = dict()
188 #     for i_line, line in enumerate(open(os.path.join(args['data_dir'], '%s.csv' % task), 'r')):
189 #         if i_line:
190 #             data = line.strip().split(',')
191 #             pid = data[0]
192 #             label = ''.join(data[1:])
193 #             pid = str(int(float(pid)))
194 #             label_dict[pid] = label
195 #     with open(os.path.join(args['files_dir'], '%s_dict.json' % task), 'w') as f:
196 #         json.dump(label_dict, f)
197
198
199
200 # def generate_demo_dict(args, demo_csv):
201 #     demo_dict = dict()
202 #     demo_index_dict = dict()
203 #     for i_line, line in enumerate(open(demo_csv, 'r')):
204 #         if i_line:

```



```

205 #         data = line.strip().split(',')
206 #         pid = str(int(float(data[0])))
207 #         demo_dict[pid] = []
208 #         for demo in data[1:]:
209 #             if demo not in demo_index_dict:
210 #                 demo_index_dict[demo] = len(demo_index_dict)
211 #                 demo_dict[pid].append(demo_index_dict[demo])
212 #         with open(os.path.join(args['files_dir'], 'demo_dict.json'), 'w') as json_
213 #             json.dump(demo_dict, json_file)
214 #         with open(os.path.join(args['files_dir'], 'demo_index_dict.json'), 'w') as
215 #             json.dump(demo_index_dict, json_file)
216
217
218 # def main():
219 #     args = parse_args()
220 #     args['files_dir'] = os.path.join(args['data_dir'], 'files')
221 #     args['initial_dir'] = os.path.join(args['data_dir'], 'initial_data')
222 #     args['resample_dir'] = os.path.join(args['data_dir'], 'resample_dir')
223 #     args['split_num'] = 4000
224 #     print(args.items())
225
226 #     for x in ['files', 'initial_data', 'resample_dir']:
227 #         if x not in os.listdir(args['data_dir']):
228 #             if x == 'files':
229 #                 os.mkdir(args['files_dir'])
230 #             elif x == 'initial_data':
231 #                 os.mkdir(args['initial_dir'])
232 #             elif x == 'resample_dir':
233 #                 os.mkdir(args['resample_dir'])
234
235 #     features_csv = os.path.join(args['data_dir'], 'features.csv')
236 #     demo_csv = os.path.join(args['data_dir'], 'demo.csv')
237 #     for task in ['mortality', 'readmit', 'llos']:
238 #         generate_label_dict(args, task)
239 #     generate_demo_dict(args, demo_csv)
240 #     generate_file_for_each_patient(args, features_csv)
241 #     resample_data(args)
242 #     generate_feature_dict(args)
243 #     split_data_to_ten_set(args)
244
245 # main()

```

✓ Preprocessing Step 7: Docs2Vec

```
1 import pandas as pd
2 import numpy as np
3 from tqdm import tqdm
4 from sklearn.utils import shuffle
5 from gensim.models.doc2vec import Doc2Vec, TaggedDocument
6 import json
7 import argparse
8
9 def parse_args(epochs=None, phase=None):
10     args = {}
11     if epochs is None:
12         args['epochs'] = 30
13     if phase is None:
14         args['phase'] = 'infer'
15     return args
16
17
18 processed_root = "/content/drive/MyDrive/mimic-iii_processed_data"
19 files_root = "/content/drive/MyDrive/mimic-iii_processed_data/files"
20 models_root = "/content/drive/MyDrive/models"
21 args = parse_args()
22 df = pd.read_csv(f'{processed_root}/earlynotes_demo.csv')
23 df['text'] = df['text'].astype(str).apply(text2words)
24 print(os.listdir(models_root))
```

✓ DATA PIPELINE

```
1 vector_dict = json.load(open(f'{args.data_dir}/files/vector_dict.json', 'r'))
2
3 def find_index(v, vs, i=0, j=-1):
4     if j == -1:
5         j = len(vs) - 1
6
7     if v > vs[j]:
8         return j + 1
9     elif v < vs[i]:
10        return i
11    elif j - i == 1:
12        return j
13
14    k = int((i + j)/2)
15    if v <= vs[k]:
16        return find_index(v, vs, i, k)
17    else:
18        return find_index(v, vs, k, j)
19
20
21 class DataBowl(Dataset):
22     def __init__(self, args, files, phase='train'):
23         assert (phase == 'train' or phase == 'valid' or phase == 'test')
24         self.args = args
25         self.phase = phase
26         self.files = files
27         self.feature_mm_dict = json.load(
28             open(os.path.join(args.files_dir, 'feature_mm_dict.json'), 'r'))
29         self.feature_value_dict = json.load(open(os.path.join(
30             args.files_dir, 'feature_value_dict_%d.json' % args.split_num), 'r'))
31         self.demo_dict = json.load(
32             open(os.path.join(args.files_dir, 'demo_dict.json'), 'r'))
33         self.label_dict = json.load(
34             open(os.path.join(args.files_dir, '%s_dict.json' % args.task), 'r'))
35
36         print('Use the last %d collections data' % args.n_visit)
37
38     def map_input(self, value, feat_list, feat_index):
39         index_start = (feat_index + 1) * (1 + self.args.split_num) + 1
40
41         if value in ['NA', '']:
42             return 0
43         else:
44             value = float(value)
45             vs = self.feature_value_dict[feat_list[feat_index]][1:-1]
46             v = find_index(value, vs) + index_start
47             return v
48
49     def map_output(self, value, feat_list, feat_index):
50         if value in ['NA', '']:
51             return 0
```



```

52     else:
53         value = float(value)
54         minv, maxv = self.feature_mm_dict[feat_list[feat_index]]
55         if maxv <= minv:
56             print(feat_list[feat_index], minv, maxv)
57         assert maxv > minv
58         v = (value - minv) / (maxv - minv)
59         v = max(0, min(v, 1))
60         return v
61
62     def get_mm_item(self, idx):
63         input_file = self.files[idx]
64         print(input_file)
65         pid = input_file.split('/')[-1].split('.')[0]
66
67         if input_file in args.resample_dir:
68             with open(input_file) as f:
69                 input_data = f.read().strip().split('\n')
70         else:
71             input_data = []
72
73         time_list, input_list = [], []
74
75         for iline in range(len(input_data)):
76             inp = input_data[iline].strip()
77             if iline == 0:
78                 feat_list = inp.split(',')
79             else:
80                 in_vs = inp.split(',')
81                 ctime = int(inp.split(',')[0])
82                 input = []
83                 for i, iv in enumerate(in_vs):
84                     if self.args.use_ve:
85                         input.append(self.map_input(iv, feat_list, i))
86                     else:
87                         input.append(self.map_output(iv, feat_list, i))
88                 input_list.append(input)
89                 time_list.append(- int(ctime))
90
91         if len(input_list) < self.args.n_visit:
92             for _ in range(self.args.n_visit - len(input_list)):
93                 # pad empty visit
94                 vs = [0 for _ in range(self.args.input_size + 1)]
95                 input_list = [vs] + input_list
96                 time_list = [time_list[0]] + time_list
97         else:
98             if self.use_first_records:
99                 input_list = input_list[: self.args.n_visit]
100                 time_list = time_list[: self.args.n_visit]
101             else:
102                 input_list = input_list[-self.args.n_visit:]

```

```

103         time_list = time_list[-self.args.n_visit:]
104
105     if self.args.value_embedding == 'no' or self.args.use_ve == 0:
106         input_list = np.array(input_list, dtype=np.float32)
107     else:
108         input_list = np.array(input_list, dtype=np.int64)
109     time_list = np.array(time_list, dtype=np.int64) + 1
110     assert time_list.min() >= 0
111     if self.args.value_embedding != 'no':
112         input_list = input_list[:, 1:]
113     else:
114         input_list = input_list.transpose()
115
116     label = np.array([int(l)
117                       for l in self.label_dict[pid]], dtype=np.float32)
118     # demo = np.array([self.demo_dict[pid] for _ in range(self.args.n_visit)]
119     demo = np.array(self.demo_dict.get(pid, 0), dtype=np.int64)
120
121     # content = self.unstructure_dict.get(pid, [])
122     # while len(content) < self.max_length:
123     #     content.append(0)
124     # content = content[: self.max_length]
125     # content = np.array(content, dtype=np.int64)
126     content = vector_dict[pid]
127     while len(content) < 12:
128         content.append([0] * 200)
129     content = content[:12]
130     content = np.array(content, dtype=np.float32)
131     # content = np.mean(content, axis=0)
132
133     return torch.from_numpy(input_list), torch.from_numpy(time_list), torch.
134
135     def __getitem__(self, idx):
136         return self.get_mm_item(idx)
137
138     def __len__(self):
139         return len(self.files)

```

✓ MODEL ARCHITECTURE

✓ LSTM

```

1
2 def value_embedding_data(d = 200, split = 200):
3     vec = np.array([np.arange(split) * i for i in range(int(d/2))], dtype=np.float32)
4     vec = vec / vec.max()
5     embedding = np.concatenate((np.sin(vec), np.cos(vec)), 1)
6     embedding[0, :d] = 0
7     embedding = torch.from_numpy(embedding)
8     return embedding
9
10
11 class LSTM(nn.Module):
12     def __init__(self, args):
13         super(LSTM, self).__init__()
14         self.args = args
15
16         # unstructure
17         if args.use_unstructure:
18             self.vocab_embedding = nn.Embedding (args.unstructure_size, args.embed_size)
19             self.vocab_lstm = nn.LSTM ( input_size=args.embed_size,
20                                         # hidden_size=args.hidden_size,
21                                         hidden_size=1,
22                                         num_layers=args.num_layers,
23                                         batch_first=True,
24                                         bidirectional=True)
25             self.vocab_mapping = nn.Sequential(
26                 nn.Linear(args.embed_size * 2, args.embed_size),
27                 nn.ReLU ( ),
28                 nn.Dropout ( 0.1),
29                 nn.Linear(args.embed_size, args.embed_size),
30             )
31             self.cat_output = nn.Sequential (
32                 nn.Linear (args.rnn_size * 3, args.rnn_size),
33                 nn.ReLU ( ),
34                 nn.Dropout ( 0.1),
35                 nn.Linear ( args.rnn_size, output_size),
36             )
37             self.cat_output = nn.Sequential (
38                 nn.ReLU ( ),
39                 nn.Dropout ( 0.1),
40                 nn.Linear (args.rnn_size * 3, output_size),
41             )
42
43         if args.value_embedding == 'no':
44             self.embedding = nn.Linear(args.input_size, args.embed_size)
45         else:
46             self.embedding = nn.Embedding (args.vocab_size, args.embed_size )
47         self.lstm1 = nn.LSTM (input_size=args.embed_size,
48                               hidden_size=args.hidden_size,
49                               num_layers=args.num_layers,
50                               batch_first=True,
51                               bidirectional=True)

```



```
52 self.lstm2 = nn.LSTM (input_size=args.embed_size,
53                       hidden_size=args.hidden_size,
54                       num_layers=args.num_layers,
55                       batch_first=True,
56                       bidirectional=True)
57 self.dd_embedding = nn.Embedding (args.n_ehr, args.embed_size )
58 self.value_embedding = nn.Embedding.from_pretrained(value_embedding_data
59 self.value_mapping = nn.Sequential(
60     nn.Linear ( args.embed_size * 2, args.embed_size),
61     nn.ReLU ( ),
62     nn.Dropout ( 0.1),
63 )
64 self.dd_mapping = nn.Sequential(
65     nn.Linear ( args.embed_size, args.embed_size),
66     nn.ReLU ( ),
67     nn.Dropout(0.1),
68     nn.Linear ( args.embed_size, args.embed_size),
69     nn.ReLU ( ),
70     nn.Dropout(0.1),
71 )
72 self.dx_mapping = nn.Sequential(
73     nn.Linear ( args.embed_size * 2, args.embed_size),
74     nn.ReLU ( ),
75     nn.Linear ( args.embed_size, args.embed_size),
76     nn.ReLU ( ),
77 )
78
79 self.tv_mapping = nn.Sequential (
80     nn.Linear ( args.embed_size * 2, args.embed_size),
81     nn.ReLU ( ),
82     nn.Linear ( args.embed_size, args.embed_size),
83     nn.ReLU ( ),
84     nn.Dropout ( 0.1),
85 )
86 self.relu = nn.ReLU ( )
87
88 lstm_size = args.rnn_size
89
90 lstm_size *= 2
91 self.output_mapping = nn.Sequential (
92     nn.Linear (lstm_size, args.rnn_size),
93     nn.ReLU ( ),
94     nn.Linear (args.rnn_size, args.rnn_size),
95     nn.ReLU ( )
96 )
97
98 self.output = nn.Sequential (
99     nn.Linear (args.rnn_size * 2, args.rnn_size),
100    nn.ReLU ( ),
101    nn.Dropout ( 0.1),
102    nn.Linear ( args.rnn_size, output_size),
```



```

103     )
104     self.pooling = nn.AdaptiveMaxPool1d(1)
105
106     self.one_output = nn.Sequential (
107         # nn.Linear (args.embed_size * 3, args.embed_size),
108         # nn.ReLU ( ),
109         nn.Dropout ( 0.1),
110         nn.Linear ( args.embed_size, output_size),
111     )
112
113
114     def visit_pooling(self, x):
115         output = x
116         size = output.size()
117         output = output.view(size[0] * size[1], size[2], output.size(3))    # (64, 30, 13)
118         output = torch.transpose(output, 1,2).contiguous()                  # (64, 30, 13)
119         output = self.pooling(output)                                         # (64, 30, 13)
120         output = output.view(size[0], size[1], size[3])                      # (64, 30, 13)
121         return output
122
123     def value_order_embedding(self, x):
124         size = list(x[0].size())                                              # (64, 30, 13)
125         index, value = x
126         xi = self.embedding(index.view(-1))                                  # (64*30*13, 200)
127         # xi = xi * (value.view(-1).float() + 1.0 / self.args.split_num)
128         xv = self.value_embedding(value.view(-1))                            # (64*30*13, 200)
129         x = torch.cat((xi, xv), 1)                                           # (64*30*13, 1024)
130         x = self.value_mapping(x)                                             # (64*30*13, 200)
131         size.append(-1)
132         x = x.view(size)                                                      # (64, 30, 13, 200)
133         return x
134
135
136     def forward(self, x, t, dd, content=None):
137
138         if 0 and content is not None:
139             content, _ = self.lstm1(content)
140             content = self.vocab_mapping(content)
141             content = torch.transpose(content, 1, 2).contiguous()
142             content = self.pooling(content)
143             content = content.view((content.size(0), -1))
144             return self.one_output(content)
145
146         # value embedding
147         x = self.value_order_embedding(x)
148         x = self.visit_pooling(x)
149
150         # demo embedding
151         dsize = list(dd.size()) + [-1]
152         d = self.dd_embedding(dd.view(-1)).view(dsize)
153         d = self.dd_mapping(d)

```

```

154     d = torch.transpose(d, 1,2).contiguous()           # (64*30, 200,
155     d = self.pooling(d)
156     d = d.view((d.size(0), -1))
157
158     # x = torch.cat((x, d), 2)
159     # x = self.dx_mapping(x)
160
161     # time embedding
162     # t = self.value_embedding(t)
163     # x = self.tv_mapping(torch.cat((x, t), 2))
164
165     # lstm
166     lstm_out, _ = self.lstm2( x )           # (64, 30, 1024)
167     output = self.output_mapping(lstm_out)
168     output = torch.transpose(output, 1,2).contiguous()   # (64,
169     # print('ouput.size', output.size())
170     output = self.pooling(output)           # (64,
171     output = output.view((output.size(0), -1))
172     out = self.output(torch.cat((output, d), 1))
173
174     # unstructure
175     if content is not None:
176         # print(content.size()) # [64, 1000]
177         content, _ = self.lstm1(content)
178         content = self.vocab_mapping(content)
179         content = torch.transpose(content, 1, 2).contiguous()
180         content = self.pooling(content)
181         content = content.view((content.size(0), -1))
182         out = self.cat_output(torch.cat((output, content, d), 1))
183
184
185     return out

```

✓ CNN

```
1 import os
2 import json
3 import torch
4 from torch import nn
5 import torch.nn.functional as F
6 from torch.autograd import *
7 import numpy as np
8 import sys
9 output_size = 1
10
11 def value_embedding_data(d = 200, split = 200):
12     vec = np.array([np.arange(split) * i for i in range(int(d/2))], dtype=np.float32)
13     vec = vec / vec.max()
14     embedding = np.concatenate((np.sin(vec), np.cos(vec)), 1)
15     embedding[0, :d] = 0
16     embedding = torch.from_numpy(embedding)
17     return embedding
18
19
20 def conv3(in_channels, out_channels, stride=1, kernel_size=3):
21     return nn.Conv1d(in_channels, out_channels, kernel_size=kernel_size,
22                      stride=stride, padding=1, bias=False)
23
24 class ResidualBlock(nn.Module):
25     def __init__(self, in_channels, out_channels, stride=1, downsample=None):
26         super(ResidualBlock, self).__init__()
27         self.conv1 = conv3(in_channels, out_channels, stride)
28         self.bn1 = nn.BatchNorm1d(out_channels)
29         self.relu = nn.ReLU(inplace=True)
30         self.conv2 = conv3(out_channels, out_channels)
31         self.bn2 = nn.BatchNorm1d(out_channels)
32         self.downsample = downsample
33
34     def forward(self, x):
35         residual = x
36         out = self.conv1(x)
37         out = self.bn1(out)
38         out = self.relu(out)
39         out = self.conv2(out)
40         out = self.bn2(out)
41         if self.downsample:
42             residual = self.downsample(x)
43         out += residual
44         out = self.relu(out)
```