# PTBxl data

November 23, 2021

```
[1]: import wfdb
     import numpy as np
     import os
     import pandas as pd
     import math
     import h5py
     import matplotlib.pyplot as plt
     import re
```

```
[2]: def print_object_attributes(obj): #https://stackoverflow.com/questions/192109/
     →is-there-a-built-in-function-to-print-all-the-current-properties-and-values-of-a
         for attr in dir(obj):
             print("obj.%s = %r" % (attr, getattr(obj, attr)))
```

### 0.0.1 Read MIT format .dat ecg data files and .hea headers

```
[3]: #Download first at: https://physionet.org/content/ptbdb/1.0.0/
     #BASE_DIR = '/media/julian/Volume/data/ECG/mit-bih-arrhythmia-database-1.0.0/'␣
     →#Arrhythmia
     BASE_DIR = '/media/julian/Volume/data/ECG/
     →ptb-xl-a-large-publicly-available-electrocardiography-dataset-1.0.1/'
     def get_file_list(BASE_DIR, relative=True, filter_function=None):
         record_files = []
         #file_endings = ['.dat', '.hea', '.xyz']
         with open(os.path.join(BASE_DIR, 'RECORDS')) as recs:
             record_files = recs.read().splitlines()
         if filter_function:
             record_files = list(filter(filter_function, record_files))
         if not relative:
             record_files = [os.path.join(BASE_DIR, f) for f in record_files]
         return record_files


     record_files = get_file_list(BASE_DIR)
     print(len(record_files), 'files found')
```

```
43674 files found
```

### 0.0.2 Extract signal from *.dat files & Read annotations & Read comments

```python
[4]: def read_comment_map_PTB(record_path):
         #print(record_path)
         record = wfdb.rdrecord(record_path)
         comment_map = {}
         for c in record.comments:
             e = c.split(':')
             comment_map[e[0]] = e[1].strip()
         return comment_map
```

```python
[5]: def filter_comment(comment, key):
         c = comment
         if key == 'Reason for admission':
             if 'Cardiomyopathy' in c or 'Heart failure' in c:
                 return 'Cardiomyopathy'
             elif 'n/a' in k or 'Palpitation' in k:
                 return 'Miscellaneous'
             elif 'angina' in k:
                 new_comments['Angina'] = comments[k]
             else:
                 new_comments[k] = comments[k]
```

```python
[6]: def read_comment(record_path):
         record = wfdb.rdrecord(record_path)
         return record.comments
```

```python
[7]: def read_header(record_path):
         record = wfdb.rdheader(record_path, rd_segments=True)
         return record.comments
```

```python
[8]: def read_signal(record_path, physical=True):
         #print(record_path)
         record = wfdb.rdrecord(record_path, physical=physical)
         #print_object_attributes(record)
         if physical:
             data = record.p_signal
         else:
             data = record.d_signal
         return data
```

```python
[9]: def read_annotation(record_path, physical=True):
         try:
             annotation = wfdb.rdann(record_path, 'hea',␣
         ↪return_label_elements=['symbol', 'label_store', 'description'])
             #print(record_path)
```

```python
        #print('sample:', annotation.sample, 'symbol', annotation.symbol,
 →'contained labels', annotation.description)
        return (annotation.sample, annotation.symbol, annotation.label_store,
 →annotation.description)
    except ValueError as ve:
        print(record_path, ' annotation read failed:', ve)
        return None
```

**Save all signals and attributes in file_data (also note how many had functioning annotations)**

```python
[10]: def read_ptbxl_database():
    csvfile = os.path.join(BASE_DIR, 'ptbxl_database.csv')
    dataframe = pd.read_csv(csvfile)
    return dataframe

def read_ptbxl_scp_statements():
    csvfile = os.path.join(BASE_DIR, 'scp_statements.csv')
    dataframe = pd.read_csv(csvfile)
    return dataframe

def train_test_split(record_files_relative):
    df = read_ptbxl_csv()
    train, val, test = [], [], []
    for rf in record_files_relative:
        temp = rf.replace('resampled', '') #resampled file contains 'resampled'
 →but csv not
        row = df.loc[(df['filename_hr'].str.contains(temp)) |
 →(df['filename_lr'].str.contains(temp))]
        if len(row) < 1:
            print('no row found containing file', rf)
        fold = row['strat_fold'].values[0]
        if fold <= 8: #https://physionet.org/content/ptb-xl/1.0.1/
 →#Cross-validation Folds
            train.append(rf)
        elif fold == 9:
            val.append(rf)
        elif fold == 10:
            test.append(rf)
        else:
            print('found unknown strat fold number', fold)
        print("final split: train %d; validation %d; test %d" % (len(train),
 →len(val), len(test)))
    return train, val, test

def read_label(ptbxl_database_dataframe, spc_codes_dataframe, rf,
 →likelihood_threshold=0.0):
```

```python
    df = ptbxl_database_dataframe
    spc_df = spc_codes_dataframe
    temp = rf.replace('resampled', '') #resampled file contains 'resampled' but␣
 →csv not
    row = df.loc[(df['filename_hr'].str.contains(temp)) | (df['filename_lr'].
 →str.contains(temp))]
    if len(row) < 1:
        print(filename, 'not found in dataframe')
        return
    code = row['scp_codes'].values[0]
    labels = [(re.sub(r'\W+', '', c.split(':')[0]), c.split(':')[1].
 →replace('}', '').strip()) for c in code.split(',')]
    diagnostic_classes = []
    for l, p in labels:
        if float(p) > likelihood_threshold:
            scp_row = spc_df.loc[(spc_df.iloc[:, 0] == l) |␣
 →(spc_df['diagnostic_subclass'] == l)]
            if len(scp_row) < 1:
                print(l, 'not found in scp_statements')
                break
            diagnostic_classes.append((scp_row['diagnostic_class'].values[0],␣
 →p))
    return sorted(diagnostic_classes, key=lambda x: x[1], reverse=True)

record_files500 = get_file_list(BASE_DIR, filter_function=lambda x:␣
 →'records500' in x)

db_df = read_ptbxl_database()
scp_df = read_ptbxl_scp_statements()
print(scp_df.iloc[:, 0])
read_label(db_df, scp_df, 'records500/00000/00001_hr')
```

```
0        NDT
1        NST_
2        DIG
3        LNGQT
4        NORM
         …
66       BIGU
67       AFLT
68       SVTAC
69       PSVT
70       TRIGU
Name: Unnamed: 0, Length: 71, dtype: object
```

[10]: [('NORM', '100.0')]

```
[75]: import ast
      from collections import defaultdict
      dbdf = read_ptbxl_database()
      all_labels = defaultdict(set)

      #Collect all possible labels
      codes = dbdf['scp_codes']
      for row in codes:
          lbl_dict = ast.literal_eval(row)
          for k,v in lbl_dict.items():
              all_labels[k].update({v})

      #Filter out labels that have only 0 probablity
      all_labels = {k: v for k, v in all_labels.items() if max(v) > 0.0}

      #Build a dict with filename as key and scp code as values
      filenames = dbdf[['filename_hr', 'scp_codes']]
      file_codes = defaultdict(dict)
      for i, (f, c) in filenames.iterrows():
          lbl_dict = ast.literal_eval(c)
          for k, v in lbl_dict.items():
              if k in all_labels and v > 0.0: #First check not necessary
                  file_codes[f][k] = v/100.0

      code_indices = dict(zip(all_labels.keys(), range(len(all_labels.keys()))))
      file_codes_onehot = dict()
      for k, v in file_codes.items():
          hot_prob = np.zeros(len(code_indices))
          for ck, cv in v.items():
              hot_prob[code_indices[ck]] = cv
          file_codes_onehot[k] = hot_prob
      print(len(filenames), len(file_codes_onehot))
```

21837 21837

```
[ ]: [a[0] for a in sorted(code_indices.items(), key=lambda x: x[1])]
```

```
[11]: import glob
      p = '/media/julian/Volume/data/ECG/
      ↪ptb-xl-a-large-publicly-available-electrocardiography-dataset-1.0.1/
      ↪generated/1000'
      for f in glob.glob(os.path.join(p, '*_hr.*')):
          fn = os.path.basename(f)
          fp = os.path.dirname(f)
          os.rename(f, os.path.join(fp, os.path.splitext(fn)[0]+'resampled'+os.path.
      ↪splitext(fn)[1]))
```
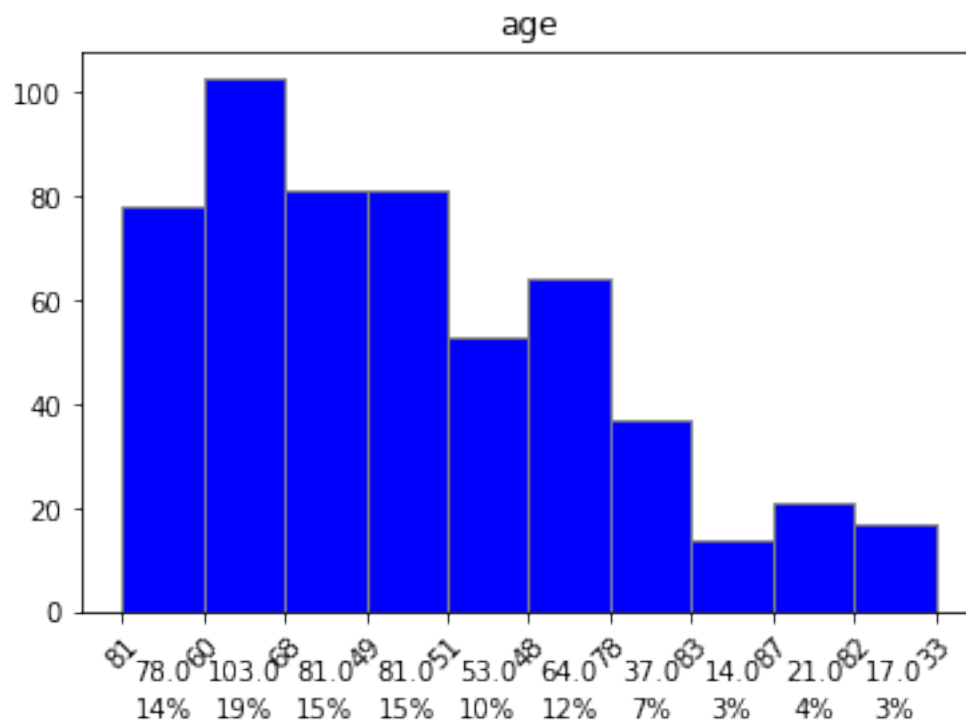
## 0.1 Playing around with PTB Comments
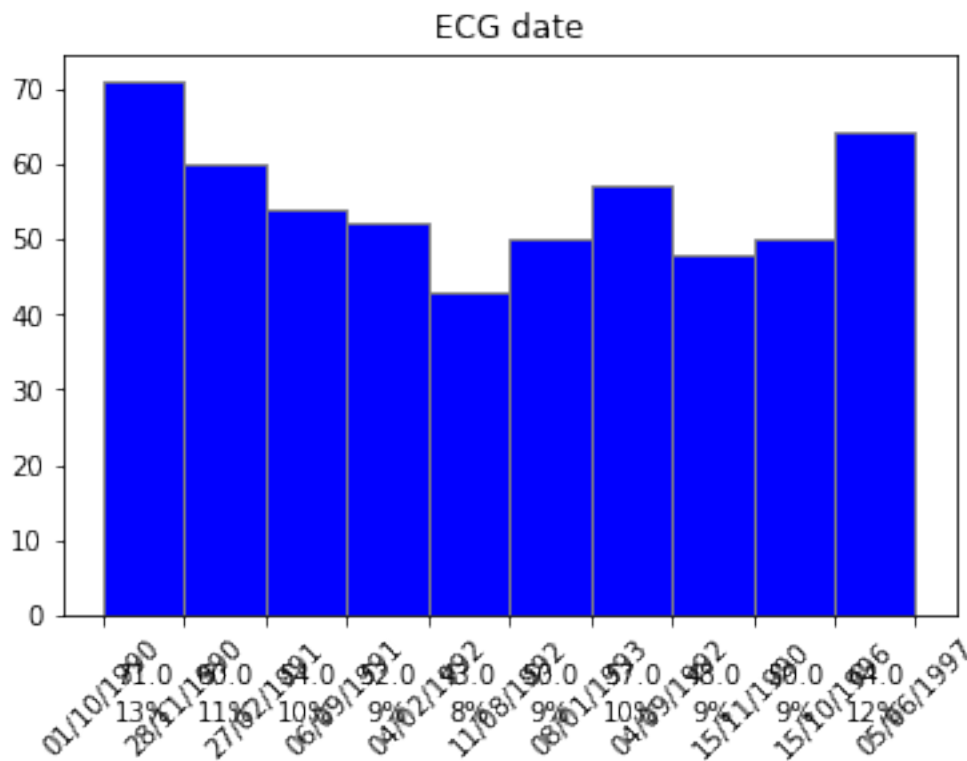
```
[20]: comment_data = []
      success = 0
      record_files = get_file_list(BASE_DIR, filter_function=lambda x: 'records500'
       ↪in x)
      for f in record_files:
          p = os.path.join(BASE_DIR, f)
          d = read_comment_map_PTB(p)
          comment_data.append(d)
      print(comment_data[0])
```

```
{}
```

```
[ ]: comment_data
```

```
[34]: for k in comment_data[0].keys():
          plot([c[k] for c in comment_data], k)
```

## sex



female 147.0 27% | female 0.0 0% | male 0.0 0% | male 391.0 71% | male 0.0 0% | n/a 0.0 0% | n/a 10.0 2% | n/a 0.0 0% | n/a 0.0 0% | 1.0 0%

## ECG date



01/10/1990 71.0 13% | 28/11/1990 60.0 11% | 27/02/1991 54.0 10% | 06/06/1991 52.0 9% | 04/02/1992 43.0 8% | 11/08/1992 50.0 9% | 08/01/1993 57.0 10% | 04/09/1992 48.0 9% | 15/11/1990 50.0 9% | 15/10/1996 64.0 12% | 05/06/1997

Diagnose

| | | | | | 549.0 | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 549.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 0% | 0% | 0% | 0% | 0% | 100% | 0% | 0% | 0% | 0% |

## Reason for admission



Myocardial infarction: 448.0, 82%
Healthy control: 8.0, 1%
Dysrhythmia: 18.0, 3%
Heart failure (NYHA 3): 27.0, 5%
Heart failure (NYHA 4): 17.0, 3%
Palpitation
Stable angina: 8.0, 2%
Hypertrophy: 17.0, 3%
Unstable angina: 5.0, 1%
Myocarditis

Acute infarction (localization)

## Former infarction (localization)



| | no | antero-septal | inferior | anterior | antero-lateral | anterior (2) | anterior (1), anterior (2) | anterior (1), inferior (2) | infero-posterolateral | anterior ? | inferior (1+2) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 426.0 | 28.0 | 31.0 | 5.0 | 7.0 | 6.0 | 14.0 | 28.0 | 2.0 | 3.0 | |
| | 78% | 5% | 6% | 1% | 1% | 1% | 3% | 5% | 0% | 1% | |

## Additional diagnoses

## Smoker



| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| no | no | yes | yes | yes | unknown | unknown | unknown | n/a | n/a |
| 218.0 | 0.0 | 0.0 | 190.0 | 0.0 | 0.0 | 114.0 | 0.0 | 0.0 | 27.0 |
| 40% | 0% | 0% | 35% | 0% | 0% | 21% | 0% | 0% | 5% |

## Number of coronary vessels involved



| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 2 | 3 | 3 | 3 | unknown | 0 | 0 | n/a |
| 149.0 | 0.0 | 87.0 | 0.0 | 64.0 | 0.0 | 190.0 | 0.0 | 32.0 | 27.0 |
| 27% | 0% | 16% | 0% | 12% | 0% | 35% | 0% | 6% | 5% |

## Infarction date (acute)



## Previous infarction (1) date

## Previous infarction (2) date



| n/a 543.0 99% | n/a 0.0 0% | 2-Jan-87 2.0 0% | 2-Jan-87 0.0 0% | 2-Jan-71 0.0 0% | 2-Jan-71 1.0 0% | 2-Jan-71 0.0 0% | 2-Jan-81 1.0 0% | 2-Jan-81 0.0 0% | 2-Jan-80 2.0 0% | 2-Jan-80 |

## Hemodynamics



| 0.0 0% | 0.0 0% | 0.0 0% | 0.0 0% | 0.0 0% | 549.0 100% | 0.0 0% | 0.0 0% | 0.0 0% | 0.0 0% |

Catheterization date

Ventriculography

Chest X-ray

| Label | Value | Percent |
|---|---|---|
| Heart size upper limit of norm | 21.0 | 4% |
| Enlarged right heart | 219.0 | 40% |
| Enlarged right heart, pulmonary venous congestion | 16.0 | 3% |
| Heart size upper limit of norm, pulmonary emphysema | 6.0 | 1% |
| Enlarged left heart, pulmonary venous congestion | 7.0 | 1% |
| Left ventricular enlargement | 14.0 | 3% |
| Normal sized heart, pulmonary venous congestion | 8.0 | 1% |
| | 9.0 | 2% |
| | 9.0 | 2% |
| Moderately enlarged heart | 240.0 | 44% |
| Globally enlarged heart, small pleural effusion (right-sided), Pulmonary venous congestion, postero-basal pneumonic infiltration | | |
| Enlarged left heart, severe pulmonary venous congestion, Marked pleural effusions on either side | | |

Peripheral blood Pressure (syst/diast)

| Label | Value | Percent |
|---|---|---|
| 140/80 mmHg | 315.0 | 57% |
| 110/70 mmHg | 45.0 | 8% |
| 130/80 mmHg | 48.0 | 9% |
| 100/60 mmHg | 53.0 | 10% |
| 130/100 mmHg | 16.0 | 3% |
| 120/90 mmHg | 27.0 | 5% |
| 100/75 mmHg | 10.0 | 2% |
| 100/50 mmHg | 4.0 | 1% |
| 140/100 mmHg | 17.0 | 3% |
| 160/95 mmHg | 10.0 | 2% |
| 180/100 mmHg | | |

## Pulmonary artery pressure (at rest) (syst/diast)



n/a 383.0 70%
23/6 cmH2O 7%
23/9 cmH2O 4%
12/3 cmH2O 4%
31/19 cmH2O 4%
43/21 cmH2O 3%
15/8 cmH2O 2%
23/4 cmH2O 1%
17/3 cmH2O 2%
30/7 cmH2O 4%
18/8 cmH2O

## Pulmonary artery pressure (at rest) (mean)



n/a 375.0 68%
33.0 cmH2O 6%
31.0 cmH2O 6%
44.0 cmH2O 8%
10.0 cmH2O 2%
6.0 cmH2O 1%
10.0 cmH2O 2%
14.0 cmH2O 3%
18.0 cmH2O 3%
6.0 cmH2O 1%

Pulmonary capillary wedge pressure (at rest)

Cardiac output (at rest)

# Cardiac index (at rest)



| Label | n/a | 3,2 l/min/sqmBSA | 2,46 l/min/sqmBSA | 3,66 l/min/sqmBSA | 2,76 l/min/sqmBSA | 3,51 l/min/sqmBSA | 2,5 l/min/sqmBSA | 2,04 l/min/sqmBSA | 2,07 l/min/sqmBSA | 3,8 l/min/sqmBSA | 3,89 l/min/sqrmBSA |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Value | 379.0 | 24 | 20 | 16 | 23 | 22 | 21 | 7 | 13 | 24 | |
| Percent | 69% | 4% | 4% | 3% | 4% | 4% | 4% | 1% | 2% | 4% | |

## Stroke volume index (at rest)



## Pulmonary artery pressure (laod) (syst/diast)

Pulmonary artery pressure (laod) (mean)

## Pulmonary capillary wedge pressure (load)



## Cardiac output (load)

Cardiac index (load)

## Stroke volume index (load)



| n/a | 387.0 | 36,3 ml/beat | 70 ml/beat | 39,2 ml/beat | 53,6 ml/beat | 38,7 ml/beat | 72,7 ml/beat | 31,2 ml/beat | 37,6 ml/beat | 65 ml/beat | 44,4 ml/beat |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 70% | 4% | 4% | 4% | 3% | 4% | 2% | 1% | 3% | 4% | |

## Aorta (at rest) (syst/diast)



| 160/64 cmH2O | 90/54 cmH2O | 106/69 cmH2O | 90/58 cmH2O | 170/80 cmH2O | 100/57 cmH2O | 132/62 cmH2O | 160/78 cmH2O | 120/64 cmH2O | 100/56 cmH2O | 116/68 cmH2O |
|---|---|---|---|---|---|---|---|---|---|---|
| 89% | 1% | 1% | 1% | 1% | 1% | 2% | 2% | 1% | 1% | |

28

Aorta (at rest) mean

# Left ventricular enddiastolic pressure

Left coronary artery stenoses (RIVA)

Left coronary artery stenoses (RCX)

No stenoses 182.0 33%
Ramus postero-lateralis of RCX 60% 16.0 3%
RCX distal to Ramus marginalis sinister_1 two serial stenoses of 60% and 95%, Ramus marginalis sinister_1 100% 12.0 2%
37.0 7%
not available 14.0 3%
RCX 239.0 44%
Ramus marginalis sinister_1 proximal 90% 9.0 2%
RCX 100% close to origin of ramus marginalis_1, Ramus marginalis_1 proximal 50% : Ramus marginalis_2 100% 15.0 3%
Ramus marginalis_1 proximal 70% 20.0 4%
RCX proximal 80%, PTCA 5.0 1%

Right coronary artery stenoses (RCA)

## Echocardiography



| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| n/a | | | | Echo | | | | | | |
| 496.0 | 11.0 | 5.0 | 19.0 | | 2.0 | 4.0 | 3.0 | 2.0 | 1.0 | |
| 90% | 2% | 1% | 3% | | 0% | 1% | 1% | | 1% | |

Labels:
- LV dilatation, Hypokinesia/Akinesia posterior, inferior and lateral wall, Hypokinesia of the septum, free anterior wall and apex
- Hypokinesia/Akinesia of the distal apikal septum, apex, adjacent parts of the anterior and the entire inferioen wall
- Normal LV dimensions. Akinesia of the inferior wall, the distal interventricular septums, apex and distal anterior wall
- LV dimensions within upper limits, slightly decreased contractility in the presence of LBBB
- Akinesia septal anterior, anterolateral. Reduced contrctility
- LV hypertrophy (septum), good contractility, slightly enlarged left atrium, mitral reguritation (grade 1)
- LV hypertrophy, slightly diminuished contractility, Hypokinesia od the interventricular septum (12mm). Left atrium slightly enlarged (50mm), normal valves
- Significantly decreased global LV contractility. Mitral regurgitation (grade 1-2), mild tricuspid regurgitation
- Significantly enlarged LV (81mm) with decreased contractility (FS 12%). Mitral regurgitation (grade 2-3), moderately enlarged left atrium (49mm). Pericardial effusion (ca. 4mm)

## Therapy



|  | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 549.0 | 0.0 | 0.0 | 0.0 | 0.0 |
|---|---|---|---|---|---|---|---|---|---|---|
|  | 0% | 0% | 0% | 0% | 0% | 100% | 0% | 0% | 0% | 0% |

## Infarction date



| 29-Sep-90 | 09-Nov-90 | 02-Jan-91 | 26-Feb-91 | 03-Jun-91 | 30-Aug-91 | 18-Feb-92 | 26-Oct-92 | n/a | 07-Jun-96 | 24-May-97 |
|---|---|---|---|---|---|---|---|---|---|---|
| 35.0 | 47.0 | 46.0 | 39.0 | 14.0 | 50.0 | 42.0 | 44.0 | 216.0 | 16.0 | |
| 6% | 9% | 8% | 7% | 3% | 9% | 8% | 8% | 39% | 3% | |

35

Admission date

29-Sep-90 26.0 5%
05-Nov-90 37.0 7%
17-Dec-90 42.0 8%
29-Jan-91 34.0 6%
03-Apr-91 29.0 5%
05-Jun-91 12.0 2%
24-Aug-91 35.0 6%
16-Jan-92 36.0 7%
24-Mar-92 33.0 6%
18-Nov-92 265.0 48%
n/a

## Medication pre admission



| Label | Value | Percent |
|---|---|---|
| L-Thyroxin | 144.0 | 26% |
| Isosorbit-Dinitrate Digoxin Glibenclamide | 15.0 | 3% |
| | 22.0 | 4% |
| Amiloride+Chlorothiazide | 15.0 | 3% |
| Metoprolol Metamizole | 265.0 | 48% |
| | 25.0 | 5% |
| Ca-antagonist | 5.0 | 1% |
| ASA Propafenon Naftidrofurylhydrogenoxalate Theophyllin | 15.0 | 3% |
| Nitrate-Spray Ca-antagonist | 19.0 | 3% |
| Metoclopramide | 24.0 | 4% |
| Ca-antagonist Captopril Isosorbit-Mononitrate Digoxin ASA Tilidin Nitrate-Spray Pindolol | | |

Start lysis therapy (hh.mm)

**Lytic agent**

| | Gamma-tPA | Streptokinase | Urokinase | n/a | Trihexylphenidyl-HCl | Tetrabenazin | Enalapril | Molsidomin | Benazepril |
|---|---|---|---|---|---|---|---|---|---|
| | 230.0 | 63.0 | 241.0 | 5.0 | 2.0 | 2.0 | 1.0 | 2.0 | 1.0 | 2.0 |
| | 42% | 11% | 44% | 1% | 0% | 0% | 0% | 0% | 0% | 0% |

Ca-antagonist Theophyllin 5-Fluoro-uracil Mitomycin Clonidin Doxazosin

Minoxidil Isosorbit-Dinitrate Atenolol Molsidomin Clonidin Doxazosin

Captopril Digoxin Furosemide Isosorbit-Dinitrate Lovastatin ASA Glibenclamide

39

Dosage (lytic agent)

Additional medication

In hospital medication

Medication after discharge

```python
from collections import Counter
k = "Reason for admission"
class_counts = Counter([c[k] for c in comment_data])
print(class_counts)
```

```
Counter({'Myocardial infarction': 368, 'Healthy control': 80, 'n/a': 27,
'Cardiomyopathy': 17, 'Bundle branch block': 17, 'Dysrhythmia': 16,
'Hypertrophy': 7, 'Valvular heart disease': 6, 'Myocarditis': 4, 'Stable
angina': 2, 'Heart failure (NYHA 2)': 1, 'Heart failure (NYHA 3)': 1, 'Heart
failure (NYHA 4)': 1, 'Palpitation': 1, 'Unstable angina': 1})
```

```python
def plot(data, title): #https://stackoverflow.com/questions/6352740/
    matplotlib-label-each-bin
    import matplotlib.pyplot as plt
    import numpy as np
```

```python
from matplotlib.ticker import FormatStrFormatter

fig, ax = plt.subplots()
counts, bins, patches = ax.hist(data, facecolor='blue', edgecolor='gray')
ax.set_title(title)
# Set the ticks to be at the edges of the bins.
ax.set_xticks(bins)
# Set the xaxis's tick labels to be formatted with 1 decimal place...
for tick in ax.get_xticklabels():
    tick.set_rotation(45)
# Label the raw counts and the percentages below the x-axis...
bin_centers = 0.5 * np.diff(bins) + bins[:-1]
for count, x in zip(counts, bin_centers):
    # Label the raw counts
    ax.annotate(str(count), xy=(x, 0), xycoords=('data', 'axes fraction'),
        xytext=(0, -18), textcoords='offset points', va='top', ha='center')

    # Label the percentages
    percent = '%0.0f%%' % (100 * float(count) / counts.sum())
    ax.annotate(percent, xy=(x, 0), xycoords=('data', 'axes fraction'),
        xytext=(0, -32), textcoords='offset points', va='top', ha='center')


# Give ourselves some more room at the bottom of the plot
plt.subplots_adjust(bottom=0.15)
plt.show()
```

**Plotting (does not seem to work for annotations)**

```
[1]: #wfdb.plot_wfdb(record=file_data[0], plot_sym=True, time_units='samples',
     ↪title='Test', figsize=(10,4), ecg_grids='all')
```

**Method for partitioning data into windows with a specific overlap** Warning: overlaps other than 0.5 have not been tested. Might also behave unexpected when 'shift' is not an Integer

```
[6]: def partition_data(data, window_size=3000, overlap=0.5, store_in_array=True,
     ↪align_right=True, verbose=True): #maybe allow non float overlap too
        samples, channels = data.shape
        if samples < window_size:
            print('too few samples (%d) to support window size of %d' % (samples,
     ↪window_size))
            return None
        if verbose: print('Input data has shape:', data.shape)
        shift = window_size*overlap
        offset = int(samples % shift)
        if align_right:
            used_data = data[offset:]
```

```
        else:
            used_data = data[:-offset]
        samples, _ = used_data.shape
        if verbose: print('The window of size %d will be shifted by %f. The total␣
 ↪data used is %d' % (window_size, shift, samples))
        partitioned = np.empty((int(samples/shift)-1, window_size, channels))
        if verbose: print('The partitioned data now has shape:', partitioned.shape)
        for i in range(len(partitioned)):
            index = int(i*shift)
            partitioned[i, :, :] = used_data[index:index+window_size, :]
        return partitioned
```

```
[ ]: partition_data(d) #nur zum testen
```

```
[8]: #testing
     print(partition_data(np.repeat(np.arange(10)[np.newaxis, :].T, 5, axis=1),␣
      ↪window_size=4))
     print(partition_data(np.repeat(np.arange(11)[np.newaxis, :].T, 5, axis=1),␣
      ↪window_size=4))
     print(partition_data(np.repeat(np.arange(9)[np.newaxis, :].T, 5, axis=1),␣
      ↪window_size=4))
     print(partition_data(np.repeat(np.arange(10)[np.newaxis, :].T, 5, axis=1),␣
      ↪window_size=10))
```

```
Input data has shape: (10, 5)
The window of size 4 will be shifted by 2.000000. The total data used is 10
The partitioned data now has shape: (4, 4, 5)
[[[0. 0. 0. 0. 0.]
  [1. 1. 1. 1. 1.]
  [2. 2. 2. 2. 2.]
  [3. 3. 3. 3. 3.]]

 [[2. 2. 2. 2. 2.]
  [3. 3. 3. 3. 3.]
  [4. 4. 4. 4. 4.]
  [5. 5. 5. 5. 5.]]

 [[4. 4. 4. 4. 4.]
  [5. 5. 5. 5. 5.]
  [6. 6. 6. 6. 6.]
  [7. 7. 7. 7. 7.]]

 [[6. 6. 6. 6. 6.]
  [7. 7. 7. 7. 7.]
  [8. 8. 8. 8. 8.]
  [9. 9. 9. 9. 9.]]]
Input data has shape: (11, 5)
```

```
The window of size 4 will be shifted by 2.000000. The total data used is 10
The partitioned data now has shape: (4, 4, 5)
[[[ 1.   1.   1.   1.   1.]
  [ 2.   2.   2.   2.   2.]
  [ 3.   3.   3.   3.   3.]
  [ 4.   4.   4.   4.   4.]]

 [[ 3.   3.   3.   3.   3.]
  [ 4.   4.   4.   4.   4.]
  [ 5.   5.   5.   5.   5.]
  [ 6.   6.   6.   6.   6.]]

 [[ 5.   5.   5.   5.   5.]
  [ 6.   6.   6.   6.   6.]
  [ 7.   7.   7.   7.   7.]
  [ 8.   8.   8.   8.   8.]]

 [[ 7.   7.   7.   7.   7.]
  [ 8.   8.   8.   8.   8.]
  [ 9.   9.   9.   9.   9.]
  [10. 10. 10. 10. 10.]]]
Input data has shape: (9, 5)
The window of size 4 will be shifted by 2.000000. The total data used is 8
The partitioned data now has shape: (3, 4, 5)
[[[1. 1. 1. 1. 1.]
  [2. 2. 2. 2. 2.]
  [3. 3. 3. 3. 3.]
  [4. 4. 4. 4. 4.]]

 [[3. 3. 3. 3. 3.]
  [4. 4. 4. 4. 4.]
  [5. 5. 5. 5. 5.]
  [6. 6. 6. 6. 6.]]

 [[5. 5. 5. 5. 5.]
  [6. 6. 6. 6. 6.]
  [7. 7. 7. 7. 7.]
  [8. 8. 8. 8. 8.]]]
Input data has shape: (10, 5)
The window of size 10 will be shifted by 5.000000. The total data used is 10
The partitioned data now has shape: (1, 10, 5)
[[[0. 0. 0. 0. 0.]
  [1. 1. 1. 1. 1.]
  [2. 2. 2. 2. 2.]
  [3. 3. 3. 3. 3.]
  [4. 4. 4. 4. 4.]
  [5. 5. 5. 5. 5.]
  [6. 6. 6. 6. 6.]
```

```
[7. 7. 7. 7. 7.]
[8. 8. 8. 8. 8.]
[9. 9. 9. 9. 9.]]]
```

**Generate a (hopefully correct) context task for any given partitioned data** *N* : The number of samples generated for each current window. If this number is too high (No more windows to randomly sample from) numpy will throw an error.

*observations* : The number of windows you have to predict the future.

*predictions*: The number of windows that will be predicted.

Total amount of data output will be (*Number of windows - observations - predictions*) * *N*

```python
[7]: def generate_context_task(partitioned_data, N, observations=5, predictions=3,
     →verbose=True, shuffle_all=False): #maybe use shuffle?
         #generate N-1 negative samples and 1 positive sample:
         windows, samples, channels = partitioned_data.shape
         positive_samples_x = []
         positive_samples_y = []
         negative_samples_x = []
         negative_samples_y = []
         for i in range(0, windows - observations - predictions):
             i_cur = i+observations
             positive_samples_x += [partitioned_data[i:i_cur, :, :]]
             positive_samples_y += [partitioned_data[i_cur:i_cur+predictions]]
             possible_choices = list(range(i))+list(range(i_cur+predictions,
     →windows))get_file_list(BASE_DIR)
             for _ in range(N-1):
                 choices = np.random.choice(possible_choices, predictions,
     →replace=False) #advanced use p param for different probabilities
                 negative_samples_x += [partitioned_data[i:i_cur, :, :]]
                 negative_samples_y += [partitioned_data[choices, :, :]]

         return positive_samples_x, positive_samples_y, negative_samples_x,
     →negative_samples_y
```

```python
[10]: #Testing
      px, py, nx, ny = generate_context_task(partition_data(d), 10)
```

```
Input data has shape: (120012, 15)
The window of size 3000 will be shifted by 1500.000000. The total data used is
120000
The partitioned data now has shape: (79, 3000, 15)
```

```python
[121]: (len(px), len(py)), (len(nx), len(ny))
```

```
[121]: ((71, 71), (639, 639))
```

```
[12]: px[0].shape, py[0].shape
```

```
[12]: ((5, 3000, 15), (3, 3000, 15))
```

```
[18]: def convert_dat_to_h5(storage_path, dat_file_paths, window_size=3000, overlap=0.
      ↪5, align_right=True, verbose=True):
          if not os.path.exists(storage_path):
              os.makedirs(storage_path)
          for f in dat_file_paths:
              data = read_signal(os.path.join(BASE_DIR, f))
              print(f)
              partitioned = partition_data(data, window_size, overlap, True,␣
      ↪align_right, verbose)
              target = os.path.join(storage_path, f.replace('/', '-') +'.h5')
              with h5py.File(target, 'w') as wf:
                  wf['windows'] = partitioned
                  wf.flush()
                  if verbose: print(target, 'file created and written. %d windows␣
      ↪saved.' % (len(partitioned)))
```

```
[ ]: record_files = get_file_list(BASE_DIR)
     convert_dat_to_h5('test', record_files)
```

```
[17]: BASE_DIR = '/media/julian/Volume/data/ECG/ptb-diagnostic-ecg-database-1.0.0/'
      record_files = get_file_list(BASE_DIR, relative=False)

      file_data = []
      success = 0
      for f in record_files:
          header_record = read_header(f)
          print(header_record.comments)
          break
```

['age: 81', 'sex: female', 'ECG date: 01/10/1990', 'Diagnose:', 'Reason for
admission: Myocardial infarction', 'Acute infarction (localization): infero-
latera', 'Former infarction (localization): no', 'Additional diagnoses: Diabetes
mellitus', 'Smoker: no', 'Number of coronary vessels involved: 1', 'Infarction
date (acute): 29-Sep-90', 'Previous infarction (1) date: n/a', 'Previous
infarction (2) date: n/a', 'Hemodynamics:', 'Catheterization date: 16-Oct-90',
'Ventriculography: Akinesia inferior wall', 'Chest X-ray: Heart size upper limit
of norm', 'Peripheral blood Pressure (syst/diast):  140/80 mmHg', 'Pulmonary
artery pressure (at rest) (syst/diast): n/a', 'Pulmonary artery pressure (at
rest) (mean): n/a', 'Pulmonary capillary wedge pressure (at rest): n/a',
'Cardiac output (at rest): n/a', 'Cardiac index (at rest): n/a', 'Stroke volume
index (at rest): n/a', 'Pulmonary artery pressure (laod) (syst/diast): n/a',
'Pulmonary artery pressure (laod) (mean): n/a', 'Pulmonary capillary wedge
pressure (load): n/a', 'Cardiac output (load): n/a', 'Cardiac index (load):

n/a', 'Stroke volume index (load): n/a', 'Aorta (at rest) (syst/diast): 160/64 cmH2O', 'Aorta (at rest) mean: 106 cmH2O', 'Left ventricular enddiastolic pressure: 11 cmH2O', 'Left coronary artery stenoses (RIVA): RIVA 70% proximal to ramus diagonalis_2', 'Left coronary artery stenoses (RCX): No stenoses', 'Right coronary artery stenoses (RCA): No stenoses', 'Echocardiography: n/a', 'Therapy:', 'Infarction date: 29-Sep-90', 'Catheterization date: 16-Oct-90', 'Admission date: 29-Sep-90', 'Medication pre admission: Isosorbit-Dinitrate Digoxin Glibenclamide', 'Start lysis therapy (hh.mm): 19:45', 'Lytic agent: Gamma-TPA', 'Dosage (lytic agent): 30 mg', 'Additional medication: Heparin Isosorbit-Mononitrate ASA Diazepam', 'In hospital medication: ASA Isosorbit-Mononitrate Ca-antagonist Amiloride+Chlorothiazide Glibenclamide Insulin', 'Medication after discharge: ASA Isosorbit-Mononitrate Amiloride+Chlorothiazide Glibenclamide']

```python
[44]: import datetime
      import time
      str(datetime.datetime.now().strftime("%d_%m_%y-%H"))
```

```
[44]: '19_11_20-15'
```

```python
[81]: def parse_comment_dict(wfdb_comment):
          comment_map = {}
          for c in wfdb_comment:
              e = c.lower().split(':')
              comment_map[e[0]] = e[1].strip()
          return comment_map
      label_mappings = {}
      def onehot_comment(wfdb_comment, terms: list, key_function_dict:dict = None):
          onehots = []
          if len(terms) > 0:
              terms_encoded = np.zeros(len(terms), dtype=bool)
              for i, term in enumerate(terms):
                  for c in wfdb_comment:
                      if term in c:
                          terms_encoded[i] = True
              onehots.append(terms_encoded)
          comment_dict = parse_comment_dict(wfdb_comment)
          for key, (func, n) in key_function_dict.items():
              if not key in label_mappings:
                  label_mappings[key] = {}
              value = func(comment_dict[key])
              if not value in label_mappings[key]:
                  label_mappings[key][value] = len(label_mappings[key])
              encoded_key = np.zeros(n, dtype=bool)
              encoded_key[label_mappings[key][value]] = value
              onehots.append(encoded_key)
          return np.concatenate(onehots)
```

```python
def filter_comment(key, comment_string):
    c = comment_string
    if key == 'reason for admission':
        if 'cardiomyopathy' in c or 'heart failure' in c:
            return 'cardiomyopathy'
        elif 'n/a' in c or 'palpitation' in c:
            return 'miscellaneous'
        elif 'angina' in c:
            return 'angina'
    return comment_string
```

```python
[82]: from collections import defaultdict
      default_key_function_dict = defaultdict(lambda y: (lambda x: True, 2))
      default_key_function_dict['age'] = lambda x: int(x)>50 if x.isnumeric() else
       ↪False, 2 #age
      default_key_function_dict['smoker'] = lambda x: x == 'yes', 2
      default_key_function_dict['reason for admission'] = lambda x:
       ↪filter_comment('reason for admission', x), 10
```

```python
[ ]: BASE_DIR = '/media/julian/Volume/data/ECG/ptb-diagnostic-ecg-database-1.0.0/'
     record_files = get_file_list(BASE_DIR, relative=False)

     for f in record_files:
         wfdb_comment = read_header(f)
         print(onehot_comment(wfdb_comment, [], default_key_function_dict))
```

```python
[85]: label_mappings
```

```python
[85]: {'age': {True: 0, False: 1},
       'smoker': {False: 0, True: 1},
       'reason for admission': {'myocardial infarction': 0,
        'healthy control': 1,
        'valvular heart disease': 2,
        'dysrhythmia': 3,
        'cardiomyopathy': 4,
        'miscellaneous': 5,
        'angina': 6,
        'hypertrophy': 7,
        'bundle branch block': 8,
        'myocarditis': 9}}
```

```python
[20]: import numpy as np
      def generate_sin_data(n, hz_low, hz_high):
          m_phase=0

          signal = np.empty(n)
```

```python
    f = hz_low
    fs= 44100

    phaseInc = 2*np.pi*f/fs


    for i in range(int(n/2)):
        signal[i] = np.sin(m_phase)
        m_phase = m_phase + phaseInc

    m_phase = m_phase % 2*np.pi


    f=hz_high
    phaseInc = 2*np.pi*f/fs


    for i in range(int(n/2), n):
        signal[i] = np.sin(m_phase)
        m_phase = m_phase + phaseInc

    m_phase = m_phase % 2*np.pi
    return signal
```
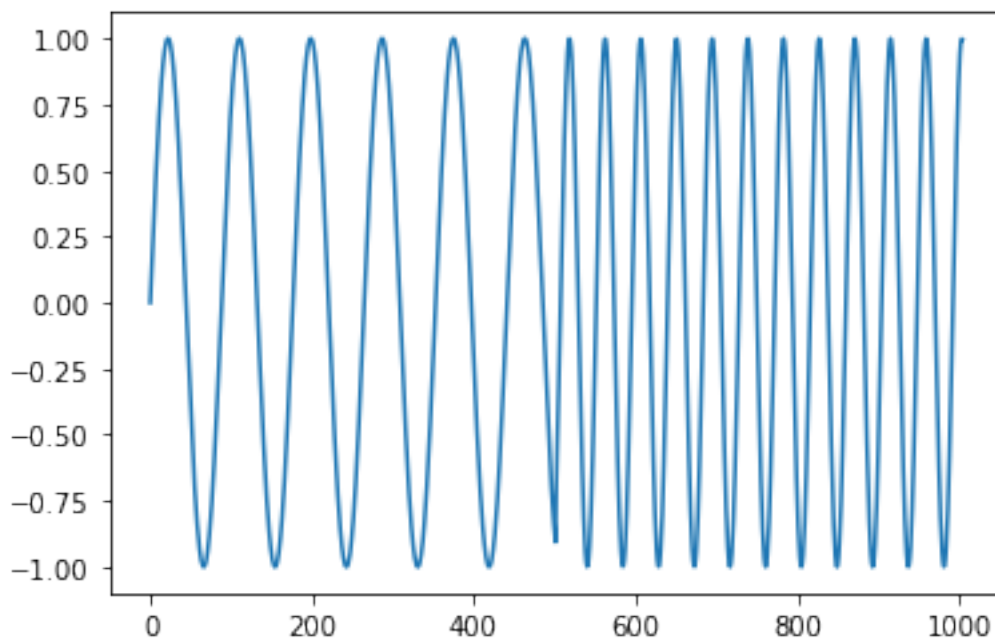
```python
[23]: import matplotlib.pyplot as plt

plt.plot(generate_sin_data(1005, 500, 1000))
plt.show()
```

```
[19]: record_files = get_file_list(BASE_DIR, filter_function=lambda x: 'records500'␣
      ↪in x)
      print(len(record_files), 'files found')
```

```
21837 files found
```

```
[6]: import torch as t
```

```
[25]: tar = t.zeros((4,8))
      tar[0,1] = 1.
      tar
      y = tar
```

```
[26]: pred = torch.rand_like(tar)
      pred[:, 1:5] = 0.
      pred
      logits = pred
```

```
[24]: t.sum((pred != 0.0) | (tar != 0.0))
```

```
[24]: tensor(17)
```

```
[37]: t.sum(t.abs(y-pred) <= 0.000005)/(4*8)
```

```
[37]: tensor(0.4688)
```

```python
[4]: import pandas as pd
     import numpy as np
     import wfdb
     import ast

     def load_raw_data(df, sampling_rate, path):
         if sampling_rate == 100:
             data = [wfdb.rdsamp(path+f) for f in df.filename_lr]
         else:
             data = [wfdb.rdsamp(path+f) for f in df.filename_hr]
         data = np.array([signal for signal, meta in data])
         return data

     path = '/media/julian/Volume/data/ECG/
      ↪ptb-xl-a-large-publicly-available-electrocardiography-dataset-1.0.1/'
     sampling_rate=500

     # load and convert annotation data
     Y = pd.read_csv(path+'ptbxl_database.csv', index_col='ecg_id')
     Y.scp_codes = Y.scp_codes.apply(lambda x: ast.literal_eval(x))

     # Load raw signal data
     #X = load_raw_data(Y, sampling_rate, path)

     # Load scp_statements.csv for diagnostic aggregation
     agg_df = pd.read_csv(path+'scp_statements.csv', index_col=0)
     agg_df = agg_df[agg_df.diagnostic == 1]

     def aggregate_diagnostic(y_dic):
         tmp = []
         for key in y_dic.keys():
             if key in agg_df.index:
                 tmp.append(agg_df.loc[key].diagnostic_class)
         return list(set(tmp))

     # Apply diagnostic superclass
     Y['diagnostic_superclass'] = Y.scp_codes.apply(aggregate_diagnostic)

     # Split data into train and test
     test_fold = 10
     # Train
     #X_train = X[np.where(Y.strat_fold != test_fold)]
     y_train = Y[(Y.strat_fold != test_fold)].diagnostic_superclass
     # Test
     #X_test = X[np.where(Y.strat_fold == test_fold)]
     y_test = Y[Y.strat_fold == test_fold].diagnostic_superclass
```

```
[5]: y_train
```

```
[5]: ecg_id
     1          [NORM]
     2          [NORM]
     3          [NORM]
     4          [NORM]
     5          [NORM]
                  …
     21833      [STTC]
     21834      [NORM]
     21835      [STTC]
     21836      [NORM]
     21837      [NORM]
     Name: diagnostic_superclass, Length: 19634, dtype: object
```