

```
In [1]: 1 import cupy as cp
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import matplotlib.path_effects as PathEffects
6 import matplotlib
7 from sklearn.utils import resample
8 from sklearn.preprocessing import MultiLabelBinarizer
9 #import seaborn as sns
10 from cuml.neighbors import KNeighborsClassifier # Importing MLKNN from RAPIDS
11 from cuml.decomposition import PCA
12 from keras.preprocessing.sequence import pad_sequences
13 from sklearn.model_selection import train_test_split
14 import warnings
15 from keras.layers import StringLookup
16 warnings.filterwarnings("ignore")
```

2024-11-05 07:53:37.343838: I tensorflow/core/util/port.cc:153] oneDNN custom operations are on. You may see slightly different numerical results due to floating-point round-off errors from different computation orders. To turn them off, set the environment variable `TF_ENABLE_ONEDNN_OPTS=0`.

2024-11-05 07:53:37.377153: E external/local_xla/xla/stream_executor/cuda/cuda_fft.cc:485] Unable to register cuFFT factory: Attempting to register factory for plugin cuFFT when one has already been registered

2024-11-05 07:53:37.435312: E external/local_xla/xla/stream_executor/cuda/cuda_dnn.cc:8454] Unable to register cuDNN factory: Attempting to register factory for plugin cuDNN when one has already been registered

2024-11-05 07:53:37.461823: E external/local_xla/xla/stream_executor/cuda/cuda_blas.cc:1452] Unable to register cuBLAS factory: Attempting to register factory for plugin cuBLAS when one has already been registered

2024-11-05 07:53:37.501641: I tensorflow/core/platform/cpu_feature_guard.cc:210] This TensorFlow binary is optimized to use available CPU instructions in performance-critical operations.

To enable the following instructions: SSE4.1 SSE4.2 AVX AVX2 AVX_VNNI FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.

```
In [2]: 1 # Load embeddings and codes
2 embeddings1 = np.load("embeddings_1.npy")
3 embeddings2 = np.load("embeddings_2.npy")
4
5 with open("icd_codes_1.txt", 'r') as file:
6     codes1 = [line.strip('\n').split(';') for line in file.readlines()]
7 with open("icd_codes_2.txt", 'r') as file:
8     codes2 = [line.strip('\n').split(';') for line in file.readlines()]
```

```
In [3]: 1 Code1 = pd.Series([label for sublist in codes1 for label in sublist ])
2 Code2 = pd.Series([label for sublist in codes2 for label in sublist ])
```

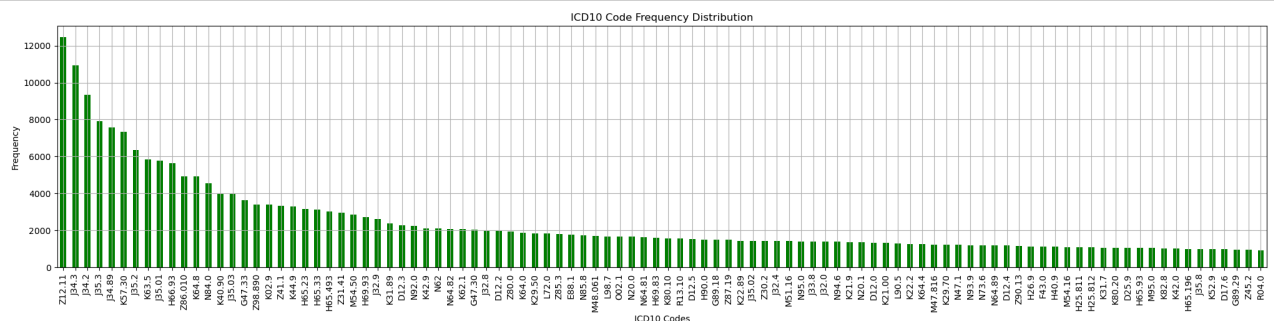
```
In [4]: 1 print("No of Unique codes in ICD_Codes_1:",len(Code1.unique()))
2 print("No of Unique codes in ICD_Codes_2:",len(Code2.unique()))
```

No of Unique codes in ICD_Codes_1: 1373
No of Unique codes in ICD_Codes_2: 1294

```
In [5]: 1 # Concatinating both the codes
2 CODE = pd.concat([Code1, Code2])
3 print("No of Unique ICD10 Codes:",len(CODE.unique()))
```

No of Unique ICD10 Codes: 1400

```
In [6]: 1 CODE.value_counts()[:100].plot(kind = 'bar', figsize= (25,5), color='green')
2 plt.title('ICD10 Code Frequency Distribution') # Top 100 ICD10 codes which is frequently repeated as Labels to embedding
3 plt.xlabel('ICD10 Codes')
4 plt.ylabel('Frequency')
5 plt.grid(True)
6 plt.xticks(rotation=90)
7 plt.show()
```



```
In [7]: 1 threshold = 100
2 underrepresented_class = CODE.value_counts()[CODE.value_counts() < threshold].index
```

```
In [8]: 1 # Extract samples (indices) that contain underrepresented codes
2 samples_to_oversample = [i for i, code_list in enumerate(codes1) if any(code in underrepresented_class for code in code_
3 len(samples_to_oversample)]
```

Out[8]: 15816

```
In [9]: 1 X_resampled, Y_resampled = resample(embeddings1[samples_to_oversample], [codes1[i] for i in samples_to_oversample], repl
```

```
In [10]: 1 codes1.extend(Y_resampled)
```

```
In [11]: 1 # Concatenate original and oversampled data
2 Embeddings1 = np.vstack((embeddings1, X_resampled))
3 Codes1 = codes1
```

```
In [12]: 1 X_resampled1, Y_resampled1 = resample(embeddings2[samples_to_oversample], [codes2[i] for i in samples_to_oversample], re
```

```
In [13]: 1 codes2.extend(Y_resampled1)
```

```
In [14]: 1 # Concatenate original and oversampled data
2 Embeddings2 = np.vstack((embeddings2, X_resampled1))
3 Codes2 = codes2
```

```
In [15]: 1 len(Embeddings1), len(Embeddings2), len(Codes1), len(Codes2)
```

Out[15]: (126491, 126491, 126491, 126491)

```
In [16]: 1 vocab = list(CODE.unique())
```

```
In [17]: 1 # Create the StringLookup layer with multi_hot output mode
2 lookup = StringLookup(vocabulary=vocab, output_mode='multi_hot')
```

WARNING: All log messages before absl::InitializeLog() is called are written to STDERR
I0000 00:00:1730773420.423799 3850 cuda_executor.cc:1001] could not open file to read NUMA node: /sys/bus/pci/devices/00:01:00.0/numa_node
Your kernel may have been built without NUMA support.
I0000 00:00:1730773420.424438 3850 cuda_executor.cc:1001] could not open file to read NUMA node: /sys/bus/pci/devices/00:01:00.0/numa_node
Your kernel may have been built without NUMA support.
I0000 00:00:1730773420.424460 3850 cuda_executor.cc:1001] could not open file to read NUMA node: /sys/bus/pci/devices/00:01:00.0/numa_node
Your kernel may have been built without NUMA support.
I0000 00:00:1730773420.428765 3850 cuda_executor.cc:1001] could not open file to read NUMA node: /sys/bus/pci/devices/00:01:00.0/numa_node
Your kernel may have been built without NUMA support.
I0000 00:00:1730773420.428789 3850 cuda_executor.cc:1001] could not open file to read NUMA node: /sys/bus/pci/devices/00:01:00.0/numa_node
Your kernel may have been built without NUMA support.
I0000 00:00:1730773420.428802 3850 cuda_executor.cc:1001] could not open file to read NUMA node: /sys/bus/pci/devices/00:01:00.0/numa_node
Your kernel may have been built without NUMA support.
I0000 00:00:1730773420.598736 3850 cuda_executor.cc:1001] could not open file to read NUMA node: /sys/bus/pci/devices/00:01:00.0/numa_node
Your kernel may have been built without NUMA support.
I0000 00:00:1730773420.598780 3850 cuda_executor.cc:1001] could not open file to read NUMA node: /sys/bus/pci/devices/00:01:00.0/numa_node
Your kernel may have been built without NUMA support.
2024-11-05 07:53:40.598787: I tensorflow/core/common_runtime/gpu/gpu_device.cc:2112] Could not identify NUMA node of platform GPU id 0, defaulting to 0. Your kernel may not have been built with NUMA support.
I0000 00:00:1730773420.598820 3850 cuda_executor.cc:1001] could not open file to read NUMA node: /sys/bus/pci/devices/00:01:00.0/numa_node
Your kernel may have been built without NUMA support.
2024-11-05 07:53:40.598835: I tensorflow/core/common_runtime/gpu/gpu_device.cc:2021] Created device /job:localhost/replica:0/task:0/device:GPU:0 with 9558 MB memory: -> device: 0, name: NVIDIA GeForce RTX 4070, pci bus id: 0000:01:00.0, compute capability: 8.9

```
In [18]: 1 import copy
2 C = copy.deepcopy(Codes1)
```

```
In [19]: 1 C.extend(Codes2) # contains codes1 followed by codes2 labels
```

```
In [20]: 1 len(C), len(codes2)+len(codes1) # to show both are of same length
```

Out[20]: (252982, 252982)

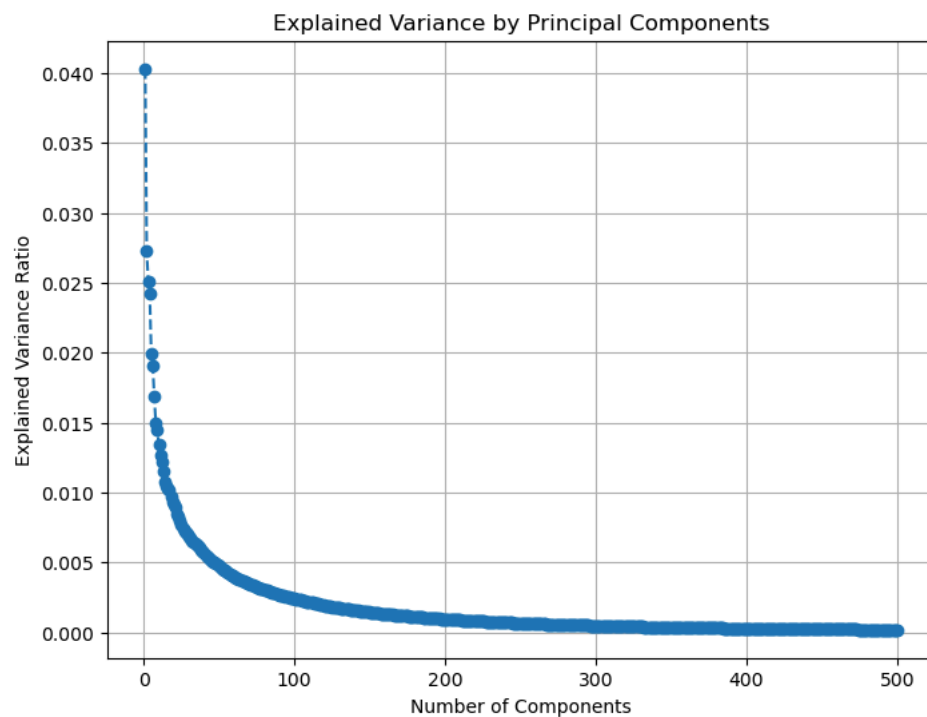
```
In [21]: 1 Embeddings = np.vstack((Embeddings1, Embeddings2))
```

```
In [22]: 1 len(Embeddings), Embeddings.shape
```

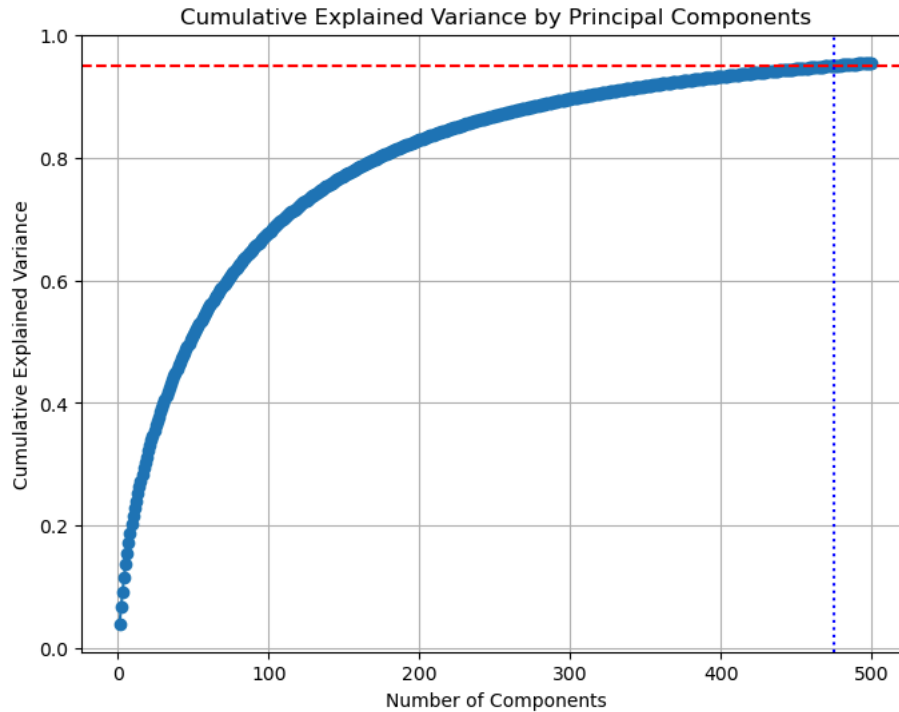
```
Out[22]: (252982, (252982, 1024))
```

```
In [23]: 1 p = PCA(n_components=500).fit(Embeddings)
```

```
In [24]: 1 explained_variance = p.explained_variance_ratio_  
2  
3 # Plot explained variance  
4 plt.figure(figsize=(8, 6))  
5 plt.plot(range(1, len(explained_variance) + 1), explained_variance, marker='o', linestyle='--')  
6 plt.title('Explained Variance by Principal Components')  
7 plt.xlabel('Number of Components')  
8 plt.ylabel('Explained Variance Ratio')  
9 plt.grid()  
10 plt.show()  
11  
12
```



```
In [25]: 1 cumulative_variance = np.cumsum(explained_variance)
2
3 plt.figure(figsize=(8, 6))
4 plt.plot(range(1, len(cumulative_variance) + 1), cumulative_variance, marker='o', linestyle='--')
5 plt.title('Cumulative Explained Variance by Principal Components')
6 plt.xlabel('Number of Components')
7 plt.ylabel('Cumulative Explained Variance')
8 plt.axhline(y=0.95, color='r', linestyle='--')
9 plt.axvline(x = 475, color = 'b', linestyle='dotted')
10 plt.grid()
11 plt.show()
```



```
In [26]: 1 Embeddings_transformed = PCA(n_components=500).fit_transform(Embeddings)
2
```

```
In [27]: 1 max_length = max(len(sublist) for sublist in C)
2 Codes_padded = pad_sequences(C, maxlen=max_length, padding='post', value='', dtype=object)
```

```
In [28]: 1 output = lookup(Codes_padded)
```

W0000 00:00:1730773430.271894 4024 gpu_backend_lib.cc:593] Can't find libdevice directory \${CUDA_DIR}/nvvm/libdevice. This may result in compilation or runtime failures, if the program we try to run uses routines from libdevice. Searched for CUDA in the following directories:

```
./cuda_sdk_lib
/home/conda/feedstock_root/build_artifacts/tensorflow-split_1729076995378/_build_env/targets/x86_64-linux
/usr/local/cuda
/home/student/miniforge3/envs/rapids-24.10/lib/python3.12/site-packages/tensorflow/python/platform/../../../../nvidia/cuda_nvcc
/home/student/miniforge3/envs/rapids-24.10/lib/python3.12/site-packages/tensorflow/python/platform/../../../../nvidia/cuda_nvcc
```

You can choose the search directory by setting `xla_gpu_cuda_data_dir` in `HloModule`'s `DebugOptions`. For most apps, setting the environment variable `XLA_FLAGS=-xla_gpu_cuda_data_dir=/path/to/cuda` will work.

```
In [29]: 1 output.numpy().shape, Embeddings_transformed.shape
```

```
Out[29]: ((252982, 1401), (252982, 500))
```

```
In [30]: 1 X_train, X_test, Y_train, Y_test = train_test_split(Embeddings_transformed, output.numpy(), test_size=0.2, random_state=
```

```
In [31]: 1 x_train, x_val, y_train, y_val = train_test_split(X_train, Y_train, test_size=0.2, random_state=42)
```

```
In [32]: 1 from cuml.multiclass import OneVsRestClassifier
2 from cuml.linear_model import LogisticRegression
```

```
In [33]: 1 #K = KNeighborsClassifier(weights='uniform', n_neighbors=5)
2 L = LogisticRegression(C=0.1, max_iter=10000)
3 K = OneVsRestClassifier(L)
```

```
In [34]: 1 K
```

```
Out[34]: OneVsRestClassifier()
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
In [35]: 1 K.fit(Embeddings_transformed, output.numpy())
2
3 from sklearn.metrics import f1_score
4 Y_pred_test = K.predict(X_test)
5 f1_macro = f1_score(Y_test, Y_pred_test, average='macro', zero_division=0)
6
7 print(f"F1 Macro Score for test data: {f1_macro}")
```

[W] [07:57:59.602338] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:06.581847] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:08.262423] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:12.564228] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:17.315985] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:22.031415] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:23.633626] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:25.972379] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:27.084904] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:45.678872] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:52.154148] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:54.467388] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:55.436228] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:56.690781] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:58:58.615645] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:59:03.059612] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:59:03.863460] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:59:07.932284] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:59:15.307046] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)
[W] [07:59:16.325874] L-BFGS stopped, because the line search failed to advance (step delta = 0.000000)

```
In [36]: 1 print("Training!!!")
```

Training!!!

```
In [ ]: 1
```

```
In [37]: 1 embeddings_test = np.load("test_data.npy")
2 pca = PCA(n_components=500).fit(Embeddings) # Fit PCA on training data
3 Embeddings_test_transformed = pca.transform(embeddings_test) # Apply same transformation to test data
4
5
6 Test_pred = K.predict(Embeddings_test_transformed)
7
```

```
In [38]: 1 V = lookup.get_vocabulary()
2 def codes_prediction(pred):
3     return [V[i] for i in np.where(pred == 1)[0]]
```

```
In [39]: 1
2 Test_decode = [codes_prediction(i) for i in Test_pred]
3
4 with open('decoded_labels_500.txt', 'w') as file:
5     for sublist in Test_decode:
6         file.write(','.join(sublist) + '\n')
```

```
In [40]: 1 def clean_predictions(pred):
2     if all(code == '[UNK]' for code in pred):
3         return ''
4     code_pred = [code for code in pred if code != '[UNK]']
5     code_pred.sort()
6     return ';'.join(code_pred)
```

```
In [41]: 1 Test_decode_cleaned = [clean_predictions(i) for i in Test_decode]
```

```
In [42]: 1 sub = pd.DataFrame({'id': range(1, len(Test_decode_cleaned)+1),
2                       'labels': Test_decode_cleaned})
```

```
In [43]: 1 sub.to_csv('submission15.csv', index=False)
```

In [44]: 1 `print("Done!")`

Done!

In []: 1