

Report on Data Engineering and Model Development

Your Name

November 14, 2024

0.1 Data Engineering

The preprocessing pipeline involves the integration and transformation of large datasets. Below is the code used for data preprocessing and an explanation of each step:

```
1 import numpy as np
2 import tensorflow as tf
3
4 # Load embedding data
5 X = np.load('embeddings_1.npy')
6 X_2 = np.load('embeddings_2.npy')
7 X = np.concatenate((X, X_2), axis=0)
8
9 # Load label data
10 label_data = []
11 file_names = ['icd_codes_1.txt', 'icd_codes_2.txt']
12 for file_name in file_names:
13     with open(file_name, 'r') as file:
14         label_data.extend(line.strip() for line in file if
15                             line.strip())
16
17 # Create unique vocabulary for StringLookup
18 unique_codes = set()
19 for labels in label_data:
20     unique_codes.update(labels.split(";"))
21 unique_codes = sorted(unique_codes)
22
23 # Initialize StringLookup
24 lookup_layer = tf.keras.layers.StringLookup(vocabulary=unique_codes,
25     output_mode="multi_hot", mask_token=None, num_oov_indices=0)
26
27 # Process labels with tf.data.Dataset
28 label_data_ds = tf.data.Dataset.from_tensor_slices(label_data)
29 def encode_labels(labels):
30     return lookup_layer(tf.strings.split(labels, sep=";"))
31 multi_hot_labels_ds = label_data_ds.map(encode_labels,
32     num_parallel_calls=tf.data.AUTOTUNE).batch(1000)
33
34 # Concatenate results
35 y = tf.concat(list(multi_hot_labels_ds), axis=0)
36 print("Shape of y:", y.shape)
```

0.1.1 Explanation of Steps

- **Loading Embedding Data:** The embedding vectors are loaded from two files, `embeddings_1.npy` and `embeddings_2.npy`, and concatenated along the first axis to create a unified dataset. The dataset consists of 1024 features for each sample.
- **Reading Label Data:** Label data are extracted from text files `icd_codes_1.txt` and `icd_codes_2.txt`, where each line contains semicolon-separated labels.
- **Building Vocabulary:** A set of unique ICD-10 codes is created by iterating through all label sets and splitting them by the semicolon delimiter. This set is then sorted to ensure consistent encoding.

- **StringLookup Initialization:** A `StringLookup` layer is initialized with the vocabulary to map each label to a multi-hot encoded format.
- **Efficient Processing:** The `tf.data.Dataset` API is used to preprocess the label data efficiently. Each label set is split and encoded using the `StringLookup` layer, and batching is applied to optimize memory usage.
- **Concatenating Results:** The batched encoded labels are concatenated to form the final multi-hot encoded tensor `y`, which has dimensions (200000, 1400), representing the number of samples and unique labels.

0.2 Feature Engineering

A `StringLookup` layer was initialized using TensorFlow to map the ICD-10 codes to a multi-hot encoded format. This layer utilized the sorted vocabulary, ensuring consistent encoding of label sets. To optimize memory usage, the label data were processed using a `tf.data.Dataset` pipeline, leveraging batching and parallel processing to handle the large-scale dataset efficiently. The final encoded labels, stored in a tensor `y`, had dimensions (200000, 1400), indicating the number of samples and the total unique labels. Each embedding vector consists of 1024 features, representing the high-dimensional nature of the data.

0.3 Sampling

Train test split was used to split the training data and the validation data which is randomly done.

```

1 from sklearn.model_selection import train_test_split
2 X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2,
3         random_state=42)
4
5 print("X_train shape:", X_train.shape)
6 print("y_train shape:", y_train.shape)
7 print("X_test shape:", X_val.shape)
8 print("y_test shape:", y_val.shape)

```

Other than this, no other sampling was done because it was not required. As we used gpu to run out neural network model, the code worked sufficiently fast and so, no sampling was used. It was done to capture every information we can.

0.4 Exploratory Data Analysis (EDA)

0.4.1 Visualization

Below is the code snippet used to generate a box plot visualization of the first 10 embedding features:

```

1 import matplotlib.pyplot as plt
2
3 # Box plot for first 10 embedding features
4 plt.figure(figsize=(10, 6))

```

```

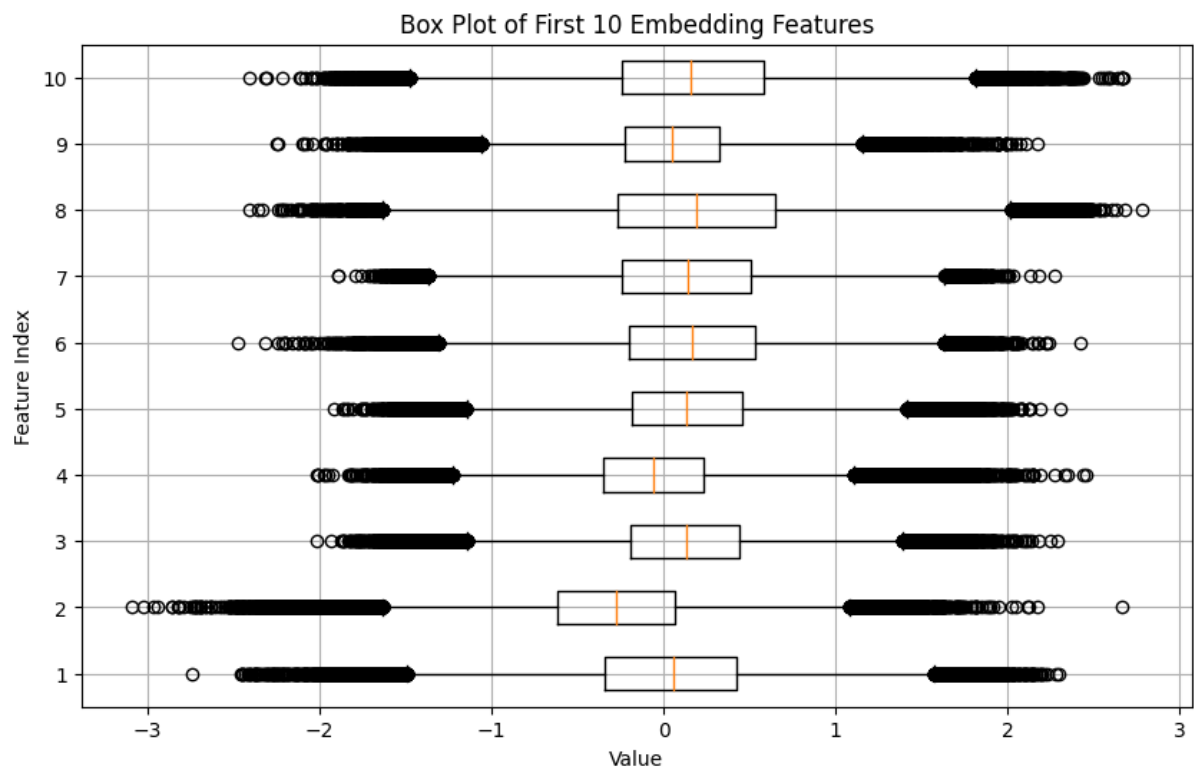
5 plt.boxplot(X[:, :10], vert=False)
6 plt.title("Box Plot of First 10 Embedding Features")
7 plt.xlabel("Value")
8 plt.ylabel("Feature Index")
9 plt.grid(True)
10 plt.show()

```

0.4.2 Explanation of the Visualization

The box plot displays the distribution of values for the first 10 embedding features out of a total of 1024 features. Key takeaways:

- The horizontal orientation makes it easier to compare feature distributions.
- The median, quartiles, and potential outliers for each feature are clearly visible.
- Variability in feature scales provides insights into normalization needs or model expectations.



Additionally, the following visualizations provide insights into label data:

```

1 # Distribution of label counts per sample
2 plt.figure(figsize=(10, 6))
3 plt.hist(label_counts, bins=30, edgecolor='k', alpha=0.7)
4 plt.title("Distribution of Label Counts per Sample")
5 plt.xlabel("Number of Labels")
6 plt.ylabel("Frequency")
7 plt.grid(True)
8 plt.show()
9
10 # Frequency of top 20 most common labels
11 plt.figure(figsize=(10, 6))

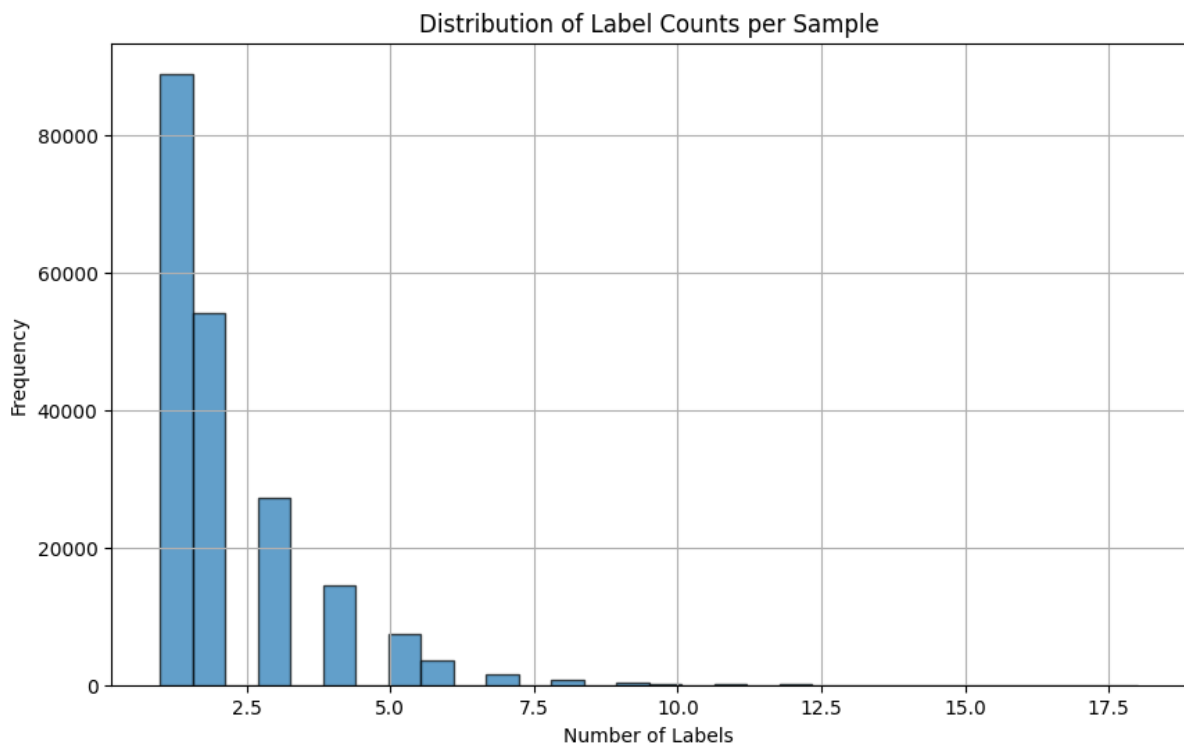
```

```

12 plt.barh(range(len(top_label_counts)), top_label_counts,
13          color='skyblue')
14 plt.yticks(range(len(top_label_counts)), [f"Label {i}" for i in
15          top_label_indices])
16 plt.xlabel("Frequency")
17 plt.ylabel("Labels")
18 plt.title("Top 20 Most Frequent Labels")
19 plt.grid(True, axis='x')
20 plt.show()

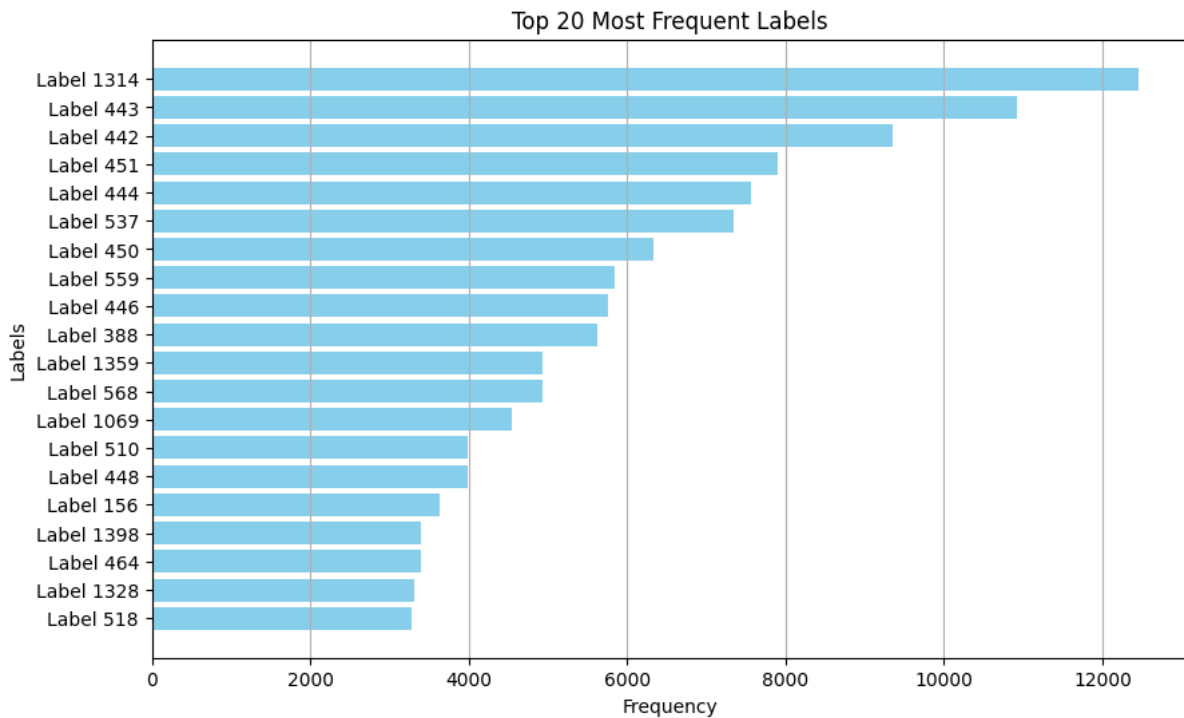
```

0.4.3 Distribution of Label Counts



The histogram above illustrates the distribution of label counts per sample. Most samples have a small number of associated labels, while a few have a larger number. This skewed distribution suggests that most instances are sparsely labeled.

0.4.4 Top 20 Most Frequent Labels



This visualization highlights the imbalance in label occurrences, with certain labels appearing significantly more frequently than others.

0.5 Models used

A number of models were considered for this data challenge. The first approach is to consider the multilabel problem as a multiple binary label classification problem. This process involves training around 1400 models (one for each label) and predict the test data on these 1400 models meaning that each model will predict 0 or 1 for its corresponding label.

The binary classifiers tried for this approach were logistic regression and lightbgm model.

Listing 1: Code for creating separate label files for multiple-binary classification

```
1 import pandas as pd
2 import os
3 import tensorflow as tf
4
5 # Define output folder for CSV files
6 output_folder = "label_files"
7 os.makedirs(output_folder, exist_ok=True)
8
9 # Save each column as a separate CSV file
10 for label_index in range(y_train.shape[1]):
11     # Extract each column (label) and save as a binary CSV
12     label_data = pd.DataFrame(y_train[:, label_index],
13                               columns=[f'label_{label_index}'])
14     label_data.to_csv(f"{output_folder}/label_{label_index}.csv",
15                       index=False, header=False)
```

```
15 print("CSV files created in the 'label_files' folder.")
```

the above code will create 1400 csv files of binary output for each label. Each csv will have 0s and 1s indicating the presence of that particular label in each train datapoint.

Listing 2: Code for implenting Logreg models for each label

```
1 from sklearn.linear_model import LogisticRegression
2 from joblib import Parallel, delayed
3 import joblib
4
5 # Folder for saving models
6 model_folder = "models"
7 os.makedirs(model_folder, exist_ok=True)
8
9 # Function to train a model for each label
10 def train_model_for_label(label_index):
11     # Load the binary target data for this label
12     y_train = pd.read_csv(f"{output_folder}/label_{label_index}.csv",
13                           header=None).values.ravel()
14
15     # Initialize and train the model
16     model = LogisticRegression(max_iter=100)
17     model.fit(X_train, y_train)
18
19     # Save the model
20     model_path = f"{model_folder}/model_{label_index}.pkl"
21     joblib.dump(model, model_path)
22
23     return model
24
25 # Train each model in parallel
26 Parallel(n_jobs=8)(delayed(train_model_for_label)(i) for i in
27                    range(y_train.shape[1]))
28
29 print("Training complete. Models saved in the 'models' folder.")
```

The above code trains the training dataset on logistic regression and saves 1400 models as output in a seperate folder.

We got 0.427 score (score refers to average micro f2 score) using this model. Nest, the main model for binary classification has been changed from logreg to lightgbm as we searched and found it as a better model in this scenaario. But, later we found that lightgbm model overpredicts (giving abnormally more number of labels for the test data). In fact, we got 0.000 score from using lightgbm.

From here, we switched our concentration to neural network models. Here is a general framework of the neural network model that we used:

Listing 3: sample neural network code

```
1 import tensorflow as tf
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import Dense, Dropout,
4   BatchNormalization, Input, LayerNormalization
5 from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint,
6   ReduceLROnPlateau
7 from tensorflow.keras.optimizers import Adam
8 import numpy as np
9
10 # Enhanced model architecture with larger layers and increased dropout
```

```

9  def create_model(input_shape):
10     model = Sequential([
11         Input(shape=(input_shape,)),
12         Dense(2048, activation='relu'),
13         LayerNormalization(),
14         Dropout(0.5),
15         Dense(1024, activation='relu'),
16         BatchNormalization(),
17         Dropout(0.5),
18         Dense(512, activation='relu'),
19         BatchNormalization(),
20         Dropout(0.4),
21         Dense(256, activation='relu'),
22         BatchNormalization(),
23         Dropout(0.4),
24         Dense(1400, activation='sigmoid')
25     ])
26     return model
27
28 # Initialize the model
29 model = create_model(X_train.shape[1])
30
31 # Compile model with Adam optimizer and custom focal loss
32 optimizer = Adam(learning_rate=0.0003)
33 model.compile(optimizer=optimizer, loss='binary_crossentropy',
34               metrics=[MicroF2Score()])
35
36 # Callbacks for early stopping and learning rate reduction
37 early_stopping = EarlyStopping(monitor='val_loss', patience=8,
38                                restore_best_weights=True, mode='min')
39 model_checkpoint = ModelCheckpoint('best_model.keras',
40                                   save_best_only=True, monitor='val_loss', mode='min')
41 lr_scheduler = ReduceLROnPlateau(monitor='val_loss', factor=0.5,
42                                  patience=4, mode='min')
43
44 # Train the model
45 history = model.fit(X_train, y_train,
46                     batch_size=32, # Optimized batch size
47                     epochs=100,
48                     validation_split=0.1,
49                     callbacks=[early_stopping, model_checkpoint,
50                                lr_scheduler],
51                     verbose=1)

```

Here are the important parts we tuned on the neural network model:

1. **Number of Layers:** Used no layer, 1 layer, 2 layer, 3 layers, 4 layers and 8 layers. It is found that no layer works best among all these (simply a sigmoid function between input and output gave the best results :))
2. **Activation function:** Tried relu, leaky relu, elu and gelu model. Found that leaky relu worked best among all.
3. **Dropout:** Changed dropout between 0.3, 0.4 and 0.5 and no dropout for each layer.
4. **Optimizer:** Tried adam, nadam, RMSprop, Adamax, Adagrad. Different optimizers worked well under different conditions (dependent on other parameters)

5. **Batch size:** Tried batch sizes 32, 64 and 128. Found 32 is ideal for our problem.
6. **Learning rate:** learning rate of 0.0003 was taken for all the models which was found to be ideal after doing trial and error.

Number of epochs is set to 100. Although, This will not be reached every time as we included early stopping criterion and the model is trained within 100 epochs usually.

Permutations and combinations were made with all these parameters and the optimal parameters were found which will be discussed with their corresponding scores in the further sections.

0.5.1 Additional things used in the Neural Network

- **Early stopping :** Early stopping criterion was used as we found out after about 50-60 epoch, validation test score was not improving and started declining (model started over fitting). This will automatically stops epoch once it reaches a steady state.

Listing 4: Early stopping

```
1 early_stopping = EarlyStopping(monitor='val_loss',
    patience=8, restore_best_weights=True, mode='min')
```

- **Dynamic learning rate change:** The ReduceLROnPlateau learning rate scheduler is a helpful tool for adjusting the learning rate dynamically during training. This setup helps to prevent overfitting and allows the model to converge more effectively by slowing down the learning rate when progress plateaus.

Listing 5: lr scheduler

```
1 lr_scheduler = ReduceLROnPlateau(monitor='val_loss',
    factor=0.5, patience=4, mode='min')
```

0.5.2 Metric Used

A custom micro average f2 score was written from scratch and used as a metric in training our neural network model. Following code snippet shows the function we have used.

Listing 6: Code for custom f2 metric

```
1 import tensorflow as tf
2
3 @tf.keras.saving.register_keras_serializable()
4 class MicroF2Score(tf.keras.metrics.Metric):
5     def __init__(self, name='micro_f2_score', beta=2, **kwargs):
6         super(MicroF2Score, self).__init__(name=name, **kwargs)
7         self.beta = beta
8         self.tp = self.add_weight(name='tp', initializer='zeros')
9         self.fp = self.add_weight(name='fp', initializer='zeros')
10        self.fn = self.add_weight(name='fn', initializer='zeros')
11
12    def update_state(self, y_true, y_pred, sample_weight=None):
13        # Threshold y_pred to get binary predictions
```

```

14     y_pred = tf.cast(y_pred > 0.5, tf.float32)
15
16     # Cast y_true to float32 to ensure compatibility
17     y_true = tf.cast(y_true, tf.float32)
18
19     # Calculate true positives, false positives, and false
20         negatives
21     true_positive = tf.reduce_sum(y_true * y_pred)
22     false_positive = tf.reduce_sum(y_pred * (1 - y_true))
23     false_negative = tf.reduce_sum((1 - y_pred) * y_true)
24
25     # Update the corresponding weights
26     self.tp.assign_add(true_positive)
27     self.fp.assign_add(false_positive)
28     self.fn.assign_add(false_negative)
29
30     def result(self):
31         precision = self.tp / (self.tp + self.fp +
32             tf.keras.backend.epsilon())
33         recall = self.tp / (self.tp + self.fn +
34             tf.keras.backend.epsilon())
35         f_beta = (1 + self.beta**2) * (precision * recall) /
36             (self.beta**2 * precision + recall +
37             tf.keras.backend.epsilon())
38         return f_beta
39
40     def reset_states(self):
41         self.tp.assign(0)
42         self.fp.assign(0)
43         self.fn.assign(0)

```

0.5.3 Prediction and Outputting as csv files

The following code generates the output csv file as required in the submission.

Listing 7: Code for prediction and outputting required csv file in logreg model

```

1  import tensorflow as tf
2  import numpy as np
3  import pandas as pd
4  import joblib
5  import os
6
7  num_labels = len(unique_codes) # Total number of labels
8
9  # Initialize the StringLookup layer for reverse lookup
10 lookup_layer = tf.keras.layers.StringLookup(
11     vocabulary=unique_codes, invert=True, output_mode="int",
12     mask_token=None, num_oov_indices=0
13 )
14
15 # Folder where models are saved
16 model_folder = "models"
17
18 # Array to store binary predictions for each label
19 y_pred = np.zeros((X_test.shape[0], num_labels), dtype=int)

```

```

20 # Predict for each label using the corresponding model
21 for label_index in range(num_labels):
22     # Load the saved model for this label
23     model_path = f"{model_folder}/model_{label_index}.pkl"
24     model = joblib.load(model_path)
25
26     # Predict probabilities for the test set and convert to binary
27     # labels
28     y_pred_probs = model.predict_proba(X_test)[: , 1] # Take
29     # probabilities for the positive class
30     y_pred[: , label_index] = (y_pred_probs > 0.5).astype(int)
31
32 # Convert binary predictions to ICD-10 codes
33 predicted_indices = [np.where(pred_row == 1)[0] for pred_row in y_pred]
34 predicted_codes = [lookup_layer(indices).numpy() for indices in
35 predicted_indices]
36 predicted_codes = [[code.decode('utf-8') for code in row] for row in
37 predicted_codes]
38
39 # Join ICD-10 codes with semicolons for each test instance
40 predicted_labels = [';'.join(row) for row in predicted_codes]
41
42 # Create the final submission DataFrame
43 submission_df = pd.DataFrame({
44     'id': range(1, len(predicted_labels) + 1),
45     'labels': predicted_labels
46 })
47
48 # Save to CSV file
49 submission_df.to_csv('submission.csv', index=False)
50 print("Predictions saved to 'submission.csv'.")

```

Listing 8: output for Neural network model

```

1 import tensorflow as tf
2 import numpy as np
3 import pandas as pd
4
5 # Load test data
6 X_test = np.load('test_data.npy')
7
8 # Predict probabilities and convert to binary
9 y_pred_probs = model.predict(X_test)
10 y_pred = (y_pred_probs > 0.5).astype(int)
11
12 # Prepare reverse lookup
13 lookup_layer = tf.keras.layers.StringLookup(
14     vocabulary=unique_codes, invert=True, output_mode="int",
15     mask_token=None, num_oov_indices=0
16 )
17
18 # Convert predictions to ICD10 codes
19 predicted_indices = [np.where(pred_row == 1)[0] for pred_row in y_pred]
20 predicted_codes = [lookup_layer(indices).numpy() for indices in
21 predicted_indices]
22 predicted_codes = [[code.decode('utf-8') for code in row] for row in
23 predicted_codes]
24 predicted_labels = [';'.join(row) for row in predicted_codes]

```

```

22
23 # Create and save submission DataFrame
24 submission_df = pd.DataFrame({
25     'id': range(1, len(predicted_labels) + 1),
26     'labels': predicted_labels
27 })
28 submission_df.to_csv('submission.csv', index=False)

```

0.5.4 Ensemble

Listing 9: Ensemble methods used

```

1 # Ensemble methods
2 def majority_voting(preds):
3     binary_preds = (preds > 0.5).astype(int)
4     majority_vote = np.mean(binary_preds, axis=0) >= 0.5
5     return majority_vote.astype(int)
6
7 def weighted_average(preds, weights):
8     weighted_preds = np.average(preds, axis=0, weights=weights)
9     return (weighted_preds > 0.5).astype(int)
10
11 def max_pooling(preds):
12     max_preds = np.max(preds, axis=0)
13     return (max_preds > 0.5).astype(int)
14
15 def threshold_adjustment(preds, threshold=0.4):
16     avg_preds = np.mean(preds, axis=0)
17     return (avg_preds > threshold).astype(int)

```

An ensemble of best models was done which turned out to be the best prediction.

0.6 Hacks and Workarounds

- Developed an efficient batching system using TensorFlow `tf.data.Dataset` to pre-process large-scale label data.
- Leveraged `StringLookup` with a custom vocabulary to ensure consistent multi-hot encoding for label sets.
- Designed memory-efficient pipelines to handle datasets with high-dimensional embeddings and large label vocabularies.
- Multiple ensemble models of neural networks were used, each configured with different layers and types of optimizers. Aggregation of predictions was performed using a voting mechanism. Final predictions for the test dataset were fine-tuned with threshold tuning based on the validation set.

0.7 Performance

The final model is an ensemble of various simple neural network models which produced a score of around 0.475 and the logreg model.

The final score we got is 0.490 with a public score of 0.489 placing in the sixth position overall.

As our best model is an ensemble of other models (we just used ensemble techniques on the test outputs), there is no particular training and validation score as no model is separately done.

However, for the best individual model trained (we got a score of 0.478 with this), the model is a no layer, no dropout adam optimizer neural network model. The training score (custom f2 score used) we obtained from training the model is 0.876 and the validation score is 0.837 for this model.

0.8 Conclusion

It is found that the models we train is very particular for our dataset. We have tried a lot of combinations but we are shocked to find that the more basic our codes became (reduced the complexity of the models), the better our scores were.

Also, domain knowledge also plays an important role in building our model. The idea of using simpler model occurred to us when we were searching about the origin of the dataset, how it is produced and understanding more about the features and labels.