# CS771: Mini-Project 1

**Chaitanya Vishwas Bramhapurikar**
230305

**Om Ji Gupta**
230719

**Pranshu Agarwal**
230776

**Sunij Singh Gangwar**
231051

## 1   Task 1

Let us start with the first part of task 1, i.e. we shall train our model for the dataset involving emoticons.

### 1.1   Emoticon Dataset

We started off with unicode representation[1] of emojis and svm as our model (well, we have tried other models too but svm is giving us the best accuracy) for learning and predicting labels for the validation set. Given below is the accuracy plot for the same:
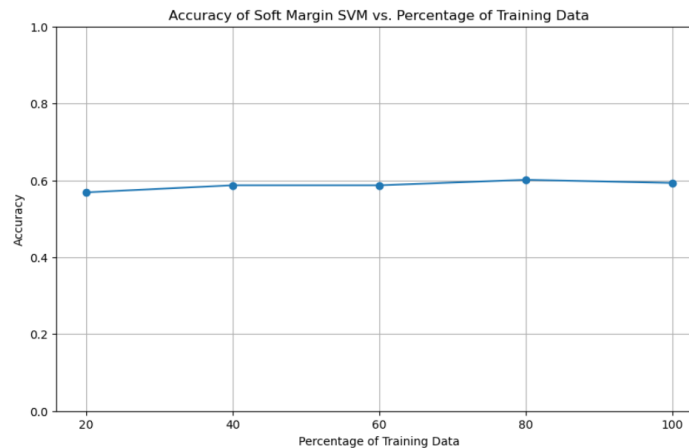


Figure 1: Accuracy plot for unicode representation

We then went on try one-hot encoding of the emoticons. The code snippet of the same is:

```
unique_emojis = set()
for i in train_emoticon_X:
    unique_emojis.update(list(i))
unique_emojis = list(unique_emojis)

#All the unique emojis collected from each string of emoticons. Next part
    is creating emoji-to-index dictionary
```

---

[1]the main point is that we had to convert the emojis into some viable representation for the training of our the model, so we started with unicode representation

```
7
8 emoji_to_index = {emoji: idx for idx, emoji in enumerate(unique_emojis)}
9
10 # One-hot encoding function for a single emoji
11 def one_hot(emoji):
12     encoding = np.zeros(len(unique_emojis), dtype=int)
13     if emoji in emoji_to_index:
14         index = emoji_to_index[emoji]
15         encoding[index] = 1
16     return encoding
17
18 # We finally need one_hot encoding function for a sequence of 13 emojis
19 def one_hot_encode_and_concatenate(emoji_sequence):
20     concatenated_vector = np.zeros(13 * len(unique_emojis), dtype=int)
21     # 13 one-hot vectors concatenated
22     for idx, emoji in enumerate(emoji_sequence[:13]):
23         one_hot_vector = one_hot(emoji)  # One-hot encode each emoji
24         concatenated_vector[idx * len(unique_emojis):(idx + 1) * len(
25     unique_emojis)] = one_hot_vector  # Concatenate
25     return concatenated_vector
26
27 train_X = input_emoticon_X.apply(one_hot_encode_and_concatenate)
28
29 # Convert the result into a 2D numpy array for training (batch_size x 13
        * num_unique_emojis)
30 train_X = np.stack(train_X)
31
32 valid_X = valid_emoticon_X.apply(one_hot_encode_and_concatenate)
33 valid_X = np.stack(valid_X)
```

We have calculated the number of unique emojis as 214 for the training dataset (as well as test dataset) so the size of individual input is $214 * 13 = 2782$ (this will be useful later when we calculate the number of trainable parameters in task 2). On this encoding, we tried different models like DTs, Logistic Regression, SVMs and it lead us to the following accuracy plot:
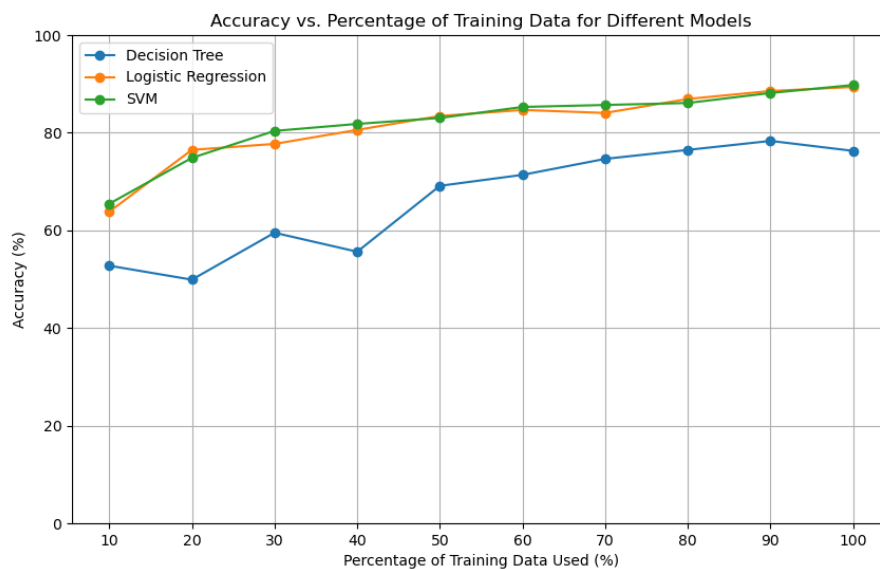


Figure 2: Plot of accuracy vs percentage of training data used

2

This was the maximum we could achieve using one-hot encoding. As it was allowed to use pre-trained deep learning model as feature extractor, we decided to learn[2] and use them directly for our feature transformation to improve our accuracy! That didn't go quite well (maybe we didn't use suitable neural-based encoding, who knows):

```
models_to_test = [
'all-MiniLM-L6-v2',
'distilbert-base-nli-stsb-mean-tokens',
'roberta-base-nli-stsb-mean-tokens',
'bert-base-nli-mean-tokens'
]
```

Encoding the dataset emoticons using the first one in the above models gave us the graph given below:
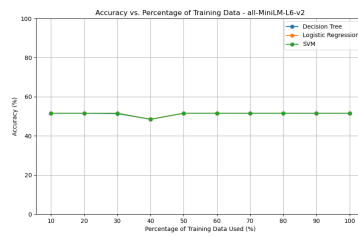


Figure 3: Using all-MiniLM-L6-v2 sentence transformer

We also tried KNN and LwP based on these transformers but their accuracy also wasn't good anyway. This wasn't nearly enough. We then tried DTs, Logistic Regression, SVM, etc for the distilbert-base-nli-stsb-mean-tokens neural-based transformer and other ones as well but they didn't yield any plots with decent accuracy. We realized later that the extracted features from transformers face lot of collisions i.e. for two different sequences of emoticons we get the same encoding, which results in bad accuracy. We also tried RNN and CNN for feature extraction but they also gave nearly 50% accuracy

## 1.2 Deep Features Dataset

For this dataset, we decided to use Averaging (which was mentioned in the question itself), Max Pooling, Concatenation, and a neural based feature extractor (CNN). For each of these encodings, we tried DTs, Logistic Regression, SVM and other models.

```
# 1. Averaging Embeddings
train_feat_X_avg = np.mean(train_feat_X, axis=1)
valid_feat_X_avg = np.mean(valid_feat_X, axis=1)
accuracies_avg = evaluate_decision_tree(train_feat_X_avg,
    valid_feat_X_avg)

# 2. Concatenation of Embeddings
train_feat_X_concat = train_feat_X.reshape(train_feat_X.shape[0], -1)
valid_feat_X_concat = valid_feat_X.reshape(valid_feat_X.shape[0], -1)
accuracies_concat = evaluate_decision_tree(train_feat_X_concat,
    valid_feat_X_concat)

# 3. MaxPooling of Embeddings
train_feat_X_maxpool = np.max(train_feat_X, axis=1)
valid_feat_X_maxpool = np.max(valid_feat_X, axis=1)
accuracies_maxpool = evaluate_decision_tree(train_feat_X_maxpool,
    valid_feat_X_maxpool)
# 4. Standard CNN-based Embedding Extraction
```

[2]Hugging face'sdocumentation was of big help

3

```
16  def build_cnn_model():
17      model = models.Sequential()
18      model.add(layers.Conv1D(filters=32, kernel_size=3, activation='relu',
         input_shape=(13, 768)))
19      model.add(layers.MaxPooling1D(pool_size=2))
20      model.add(layers.Conv1D(filters=64, kernel_size=3, activation='relu')
        )
21      model.add(layers.MaxPooling1D(pool_size=2))
22      model.add(layers.Flatten())
23      return model
24
25  cnn_model = build_cnn_model()
26  cnn_model.compile(optimizer='adam', loss='binary_crossentropy', metrics=[
        'accuracy'])
27
28  # Get CNN features
29  train_feat_cnn = cnn_model.predict(train_feat_X)
30  valid_feat_cnn = cnn_model.predict(valid_feat_X)
31  accuracies_cnn = evaluate_decision_tree(train_feat_cnn, valid_feat_cnn, "
        CNN")
```

As the number of trainable parameters in CNN exceed the limit, we have not included it in the final code.

### 1.2.1 Decision Tree

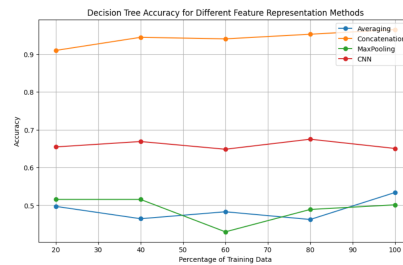We started off with the decision tree which gave us below plot:



Figure 4: Decision Tree accuracy based on different feature extraction methods.

### 1.2.2 Logistic Regression

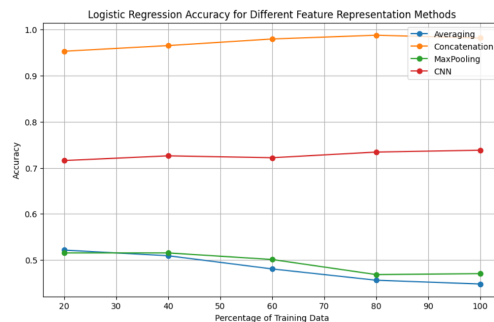We then went on to try Logistic Regression which gave us the following plot:



Figure 5: Logistic Regression accuracy based on different feature extraction methods.

4

### 1.2.3 Support Vector Machine

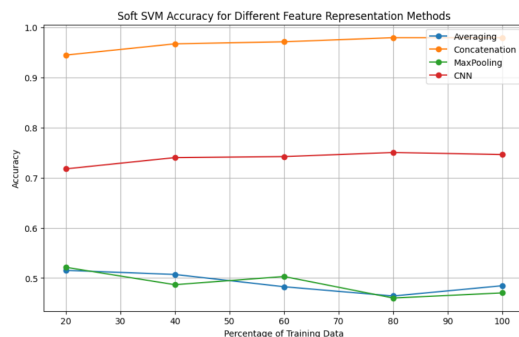The plot for support vector machine is given below:



Figure 6: Soft Margin SVM accuracy based on different feature extraction methods.

### 1.2.4 Final Code

We observe that concatenation is giving us maximum accuracy across all the feature extraction methods. We have combined the plots of all the models for concatenation below:
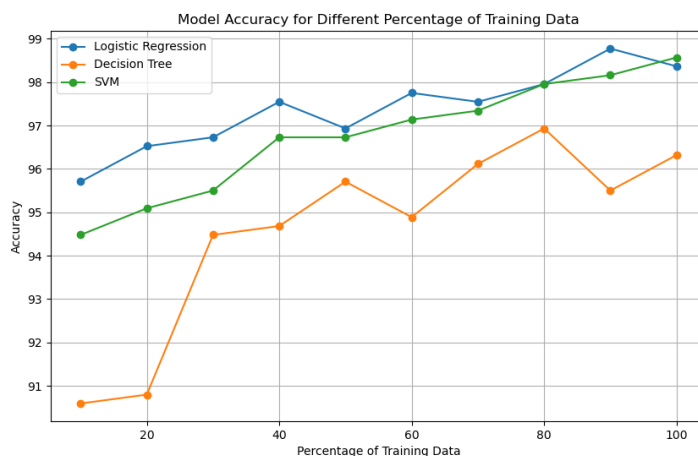


Figure 7: Accuracies of different models for concatenation

It can be seen from the above figure that the logistic regression is performing relatively better than the other models. The code snippet for logistic regression is given below:

```python
percentages = [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0]

def evaluate_LR(train_feat_X_method, valid_feat_X_method, method_name):
    accuracies = []
    for percentage in percentages:
        n_samples = int(len(train_feat_X_method) * percentage)
        x_train_subset = train_feat_X_method[:n_samples]
        y_train_subset = train_feat_Y[:n_samples]
        LR_classifier = LogisticRegression(random_state=42)
        LR_classifier.fit(x_train_subset, y_train_subset)

        # Predict on the validation set
        y_pred = LR_classifier.predict(valid_feat_X_method)
```

```
14
15          # Calculate accuracy
16          accuracy = accuracy_score(valid_feat_Y, y_pred)
17          accuracies.append(accuracy)
18
19      return accuracies
20  # Concatenation of Embeddings
21  train_feat_X_concat = train_feat_X.reshape(train_feat_X.shape[0], -1)
22  valid_feat_X_concat = valid_feat_X.reshape(valid_feat_X.shape[0], -1)
23  accuracies_concat = evaluate_LR(train_feat_X_concat, valid_feat_X_concat,
        "Concatenation")
```
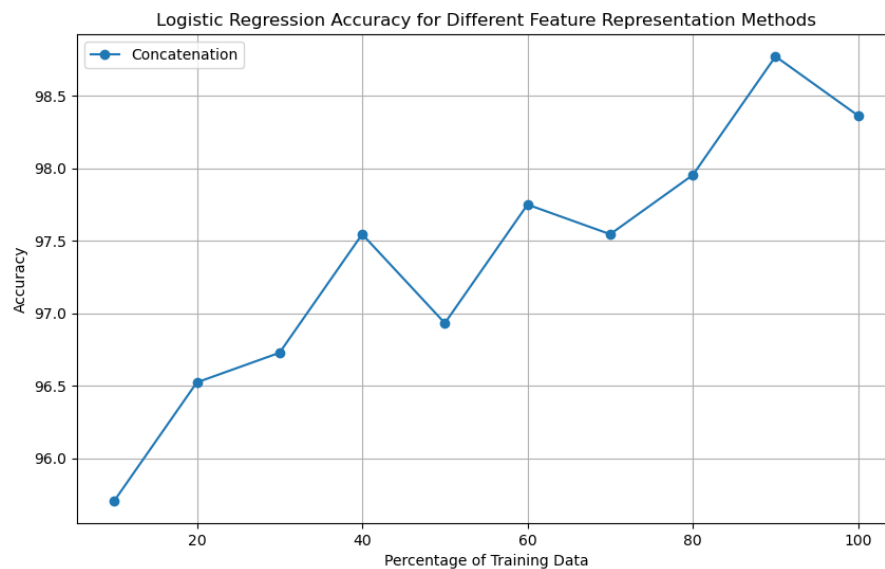
The corresponding plot is shown below:



Figure 8: Logistic Regression accuracy with concatenation feature extractor.

## 1.3 Text Sequence Dataset

We decided to first convert the string into a sequence of integers, rather we could have directly passed string as an input parameter in the model but these models generally work with numerical data inputs, so changing it into a sequence of integers is a better choice.

```
1  # Preprocessing: Convert input strings to sequences of integers
2  max_length = 50
3  train_seq_X = [[int(digit) for digit in seq] for seq in train_seq_X]
4  valid_seq_X = [[int(digit) for digit in seq] for seq in valid_seq_X]
5  # Preprocess validation data
6  train_seq_X = np.array(train_seq_X)
7  valid_seq_X = np.array(valid_seq_X)
8
9  #Convert it into numpy array
10 train_seq_Y = np.array(train_seq_Y)
11 valid_seq_Y = np.array(valid_seq_Y)
12
13 # Define CNN Model architecture with reduced parameters
14 def build_cnn_model(input_shape):
15     model = Sequential()
```

```
16      model.add(Embedding(input_dim=10, output_dim=16, input_length=
        input_shape[1]))   # Reduced embedding dimension
17      model.add(Conv1D(16, kernel_size=3, activation='relu'))  # Reduced
        number of filters
18      model.add(MaxPooling1D(pool_size=2))
19      model.add(Flatten())
20      model.add(Dense(16, activation='relu'))   # Reduced dense layer size
21      model.add(Dropout(0.3))
22      model.add(Dense(1, activation='sigmoid'))   # For binary
        classification
23      model.compile(optimizer='adam', loss='binary_crossentropy', metrics=[
        'accuracy'])
24      return model
25
26 # Function to train models and record validation accuracy
27 def train_model_on_subsets(model_func, model_type, X_train, y_train,
        X_val, y_val):
28      percentages = [0.2, 0.4, 0.6, 0.8, 1.0]
29      accuracies = []
30
31      for p in percentages:
32          subset_size = int(len(X_train) * p)
33          X_train_subset = X_train[:subset_size]
34          y_train_subset = y_train[:subset_size]
35
36          if model_type == "CNN":
37              # Build and train the CNN model
38              model = model_func(X_train_subset.shape)
39              model.fit(X_train_subset, y_train_subset, epochs=10,
        batch_size=64, verbose=0)
40              y_val_pred = (model.predict(X_val) > 0.5).astype("int32")
41              val_accuracy = accuracy_score(y_val, y_val_pred)
42          else:
43              # Train traditional machine learning models
44              model = model_func()
45              model.fit(X_train_subset, y_train_subset)
46              y_val_pred = model.predict(X_val)
47              val_accuracy = accuracy_score(y_val, y_val_pred)
48
49          accuracies.append(val_accuracy)
50
51      return percentages, accuracies
52
53 # Different models we have used
54 models = {
55      "Logistic Regression": lambda: LogisticRegression(max_iter=1000),
56      "kNN": lambda: KNeighborsClassifier(n_neighbors=5),
57      "Decision Tree": lambda: DecisionTreeClassifier(),
58      "Soft Margin SVM": lambda: SVC(kernel='linear', C=1),
59      "CNN": build_cnn_model
60 }
```

Initially, the number of trainable parameters exceeded the 10,000 trainable parameters constraint. We then reduced the number of filters and dense layer size to satisfy the constraint.

### 1.3.1 Further Exploration

We decided to use neural-based feature extractors with the hope of achieving better accuracy.

```
1 def build_embedding_model():
2      model = tf.keras.Sequential()
3      model.add(Embedding(input_dim=vocab_size, output_dim=embedding_dim,
        input_length=max_length))
4      model.add(GlobalAveragePooling1D())
```
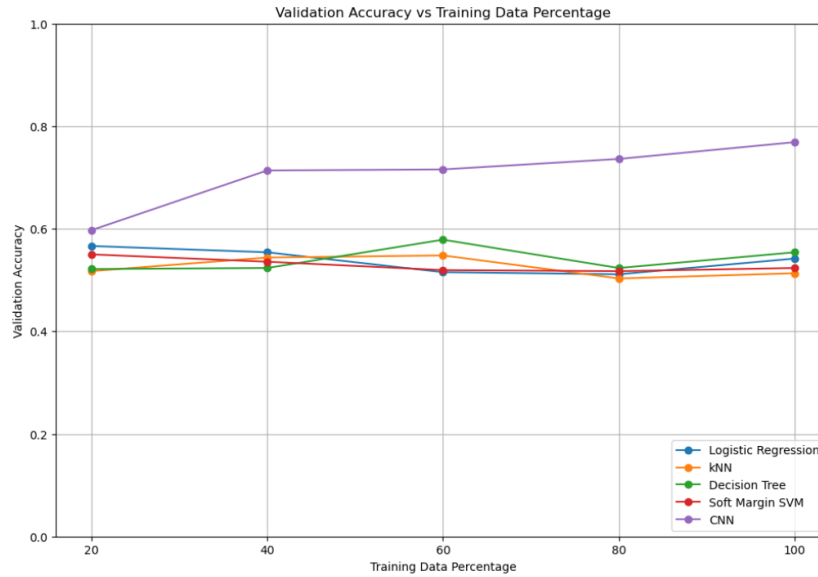
Figure 9: Final Plot of accuracies for different models.

```python
 5      return model
 6
 7  def build_cnn_model():
 8      model = tf.keras.Sequential()
 9      model.add(Embedding(input_dim=vocab_size, output_dim=embedding_dim,
        input_length=max_length))
10      model.add(Conv1D(filters=32, kernel_size=3, activation='relu'))
11      model.add(GlobalAveragePooling1D())
12      return model
13
14  def build_rnn_model():
15      model = tf.keras.Sequential()
16      model.add(Embedding(input_dim=vocab_size, output_dim=embedding_dim,
        input_length=max_length))
17      model.add(SimpleRNN(64))
18      return model
19
20  def build_lstm_model():
21      model = tf.keras.Sequential()
22      model.add(Embedding(input_dim=vocab_size, output_dim=embedding_dim,
        input_length=max_length))
23      model.add(LSTM(64))
24      return model
```

For the above feature extractors, we got the following accuracy plot for different models:

### 1.3.1.1. Decision Tree



Figure 10: Decision Tree accuracy for different neural network-based feature extractors.
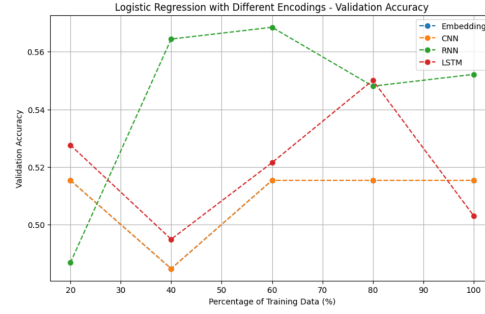
8

### 1.3.1.2. Logistic Regression



Figure 11: Logistic Regression accuracy for different neural network-based feature extractors.

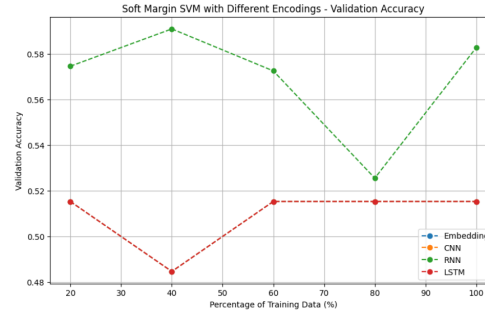### 1.3.1.3. Support Vector Machine



Figure 12: SVM accuracy for different neural network-based feature extractors.

## 2 Task 2

As we can see that in second part of task 1, we have achieved an accuracy of approximately $98.5\%$, so when we combine all the datasets from all the parts, we hope to get better accuracy than this. We tried to concatenate all the dataset inputs into one. The input size of first datset is 2782, second dataset is 9984, third dataset is 50 which in total exceeds our constraint of trainable parameters (1000). So we have used PCA to reduce the dimensionality of the second dataset to 7000.
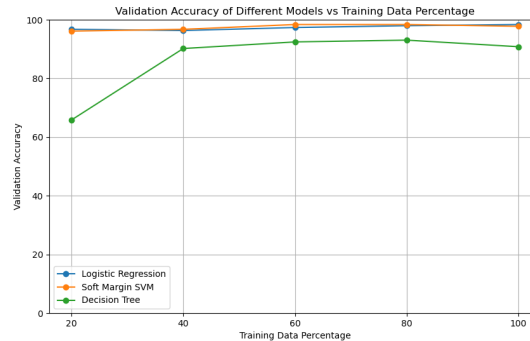


Figure 13: Accuracy plot for different models when the dimensionality of second dataset is 7000

The code snippet for the above plot is:

```python
# Prepare the third dataset: digit encoding
max_length = 50
train_seq_X_encoded = [[int(digit) for digit in seq] for seq in
    train_seq_X]
train_seq_X_encoded = np.array(train_seq_X_encoded)

valid_seq_X_encoded = [[int(digit) for digit in seq] for seq in
    valid_seq_X]
valid_seq_X_encoded = np.array(valid_seq_X_encoded)

# Feature Transformation
# One-Hot Encoding for Emoticons
encoder = OneHotEncoder(sparse_output=False, handle_unknown = 'ignore')
train_emoticon_X_encoded = encoder.fit_transform(np.array(
    train_emoticon_X).reshape(-1, 1))
valid_emoticon_X_encoded = encoder.transform(np.array(valid_emoticon_X).
    reshape(-1, 1))


# Flatten the feature matrix
train_feat_X_reshaped = train_feat_X.reshape(train_feat_X.shape[0], -1)

# Flatten the feature matrix
valid_feat_X_reshaped = valid_feat_X.reshape(valid_feat_X.shape[0], -1)

# Apply PCA to reduce deep features to 7,000 features
pca = PCA(n_components=7000)   # Adjust number of components as needed
train_feat_X_reduced = pca.fit_transform(train_feat_X_reshaped)
valid_feat_X_reduced = pca.transform(valid_feat_X_reshaped)

# Concatenate Deep Features (after PCA) with One-Hot Encoded Features
combined_train_features = np.concatenate((train_feat_X_reduced,
    train_emoticon_X_encoded), axis=1)
combined_valid_features = np.concatenate((valid_feat_X_reduced,
    valid_emoticon_X_encoded), axis=1)

# Concatenate the Combined Features with Digit Encoded Features
final_train_features = np.concatenate((combined_train_features,
    train_seq_X_encoded), axis=1)
final_valid_features = np.concatenate((combined_valid_features,
    valid_seq_X_encoded), axis=1)


final_train_labels = train_seq_Y
final_valid_labels = valid_seq_Y

# Define models
models = {
    "Logistic Regression": LogisticRegression(max_iter=1000),
    "Soft Margin SVM": SVC(kernel='linear', C=1),
    "Decision Tree": DecisionTreeClassifier()
}

# Store results for plotting
percentages = [0.2, 0.4, 0.6, 0.8, 1.0]   # Training percentages
results = {model_name: [] for model_name in models.keys()}   # Initialize
    results

# Train on subsets of the training data
for p in percentages:
    subset_size = int(len(final_train_features) * p)
    X_train_subset = final_train_features[:subset_size]
    y_train_subset = final_train_labels[:subset_size]
```

```
55
56     for model_name, model in models.items():
57         model.fit(X_train_subset, y_train_subset)
58
59         # Evaluate on the validation set
60         y_val_pred = model.predict(final_valid_features)  # Use the
       validation features
61         val_accuracy = accuracy_score(final_valid_labels, y_val_pred)
62
63         results[model_name].append(val_accuracy*100)
64
65 # Plotting the results
66 plt.figure(figsize=(10, 6))
67 for model_name, accuracies in results.items():
68     plt.plot([p * 100 for p in percentages], accuracies, marker='o',
       label=model_name)
69
70 plt.title("Validation Accuracy of Different Models vs Training Data
       Percentage")
71 plt.xlabel("Training Data Percentage")
72 plt.ylabel("Validation Accuracy")
73 plt.xticks([20, 40, 60, 80, 100])
74 plt.ylim(0, 100)  # Set y-axis limits for clarity
75 plt.grid()
76 plt.legend()
77 plt.show()
```

This PCA uses the entire dataset to extract features and reduce dimensionality of the deep features dataset, we have found the accuracy of logistic regression to be the highest across all the 3 models mentioned above, which is $98.16\%$ which was already not a marginal improvement over our original accuracy for the second dataset, so it can be inferred that on applying PCA on incremental fractions of data will definitely give even lesser accuracy.

We have tried above approaches but we coudn't get any marginal improvement over accuracy obtained in deep features dataset. The second dataset seems to be the most accurate representation of the distribution on which the all the datasets are based. Therefore, we will use the deep features dataset to make predictions on test data.

## References

- Used the sk-learn documentation for classification.
- Read about Sentence Transformers from Hugging Face.
- Learnt about tensor flow from its documentation.
- Took help of chat-gpt in debugging codes.