## ⌄ Reason for Not Using Collate Function

- **Memory Efficiency**: When working with dense embeddings, memory efficiency becomes crucial. Utilizing offsets by concatenating indices allows for effective memory usage, particularly with PyTorch's `embedding_bag`, which minimizes computational load by summing embeddings based on these offsets.

- **Handling Variable Lengths**: Dense embeddings can represent sequences of varying lengths (like sentences or documents). Instead of padding these sequences, offsets are employed to indicate where each sequence begins, facilitating concatenation of indices. This approach avoids padding and enhances computational performance.

- **Sparse Embeddings**: Sparse embeddings, such as one-hot vectors or TF-IDF representations, are typically managed independently for each sample. Therefore, the use of concatenation and offset tracking is less critical since each document can be processed individually without needing special handling for sequence lengths.

```
1 import sys
2 import os
```

```
1 !pip install torchinfo
```

```
⯈  Collecting torchinfo
      Downloading torchinfo-1.8.0-py3-none-any.whl.metadata (21 kB)
    Downloading torchinfo-1.8.0-py3-none-any.whl (23 kB)
    Installing collected packages: torchinfo
    Successfully installed torchinfo-1.8.0
```

```
1 from google.colab import drive
2 drive.mount("/content/drive")
```

```
⯈  Mounted at /content/drive
```

```
 1 from datetime import datetime
 2 import matplotlib.pyplot as plt
 3 import seaborn as sns
 4 from pprint import pprint
 5 from torchinfo import summary
 6 import joblib
 7 from collections import Counter
 8 from functools import partial
 9 from pathlib import Path
10 import torch
```

```
11 import torch.nn as nn
12 import random
13 import numpy as np
14 from sklearn.metrics import confusion_matrix
```

```
1 # Load the dataset
2 data = joblib.load("/content/drive/MyDrive/NLP/df_multilabel_hw_cleaned.joblib
3 data.head()
```

|   | cleaned_text | Tags | Tag_Number |
|---|---|---|---|
| 0 | asp query stre dropdown webpage follow control... | c# asp.net | [0, 9] |
| 1 | run javascript code server java code want run ... | java javascript | [1, 3] |
| 2 | linq sql throw exception row find change hi li... | c# asp.net | [0, 9] |
| 3 | run python script php server run nginx web ser... | php python | [2, 7] |
| 4 | advice write function m try write function res... | javascript jquery | [3, 5] |

---

Next steps:   [ Generate code with  data ]   [ ⬤  View recommended plots ]   [ New interactive sheet ]

```
1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 47427 entries, 0 to 47426
Data columns (total 3 columns):
 #   Column        Non-Null Count  Dtype
---  ------        --------------  -----
 0   cleaned_text  47427 non-null  object
 1   Tags          47427 non-null  object
 2   Tag_Number    47427 non-null  object
dtypes: object(3)
memory usage: 1.1+ MB
```

```
1 import numpy as np
2 import ast
3
4 def process_data(data):
5     def safe_convert_tag(tag):
6         try:
7             return ast.literal_eval(tag)
8         except (ValueError, SyntaxError):
9             return None
10
11
12     y = [safe_convert_tag(tag) for tag in data['Tag_Number']]
13     x = np.array(data['cleaned_text'].astype(str)).reshape(-1, 1)
```

```
13       x = np.array(data['cleaned_text'].astype(str)).reshape(-1, 1)
14
15    return x, y
16
17
18 # Use the function to process your data
19 x, y = process_data(data)
```

```
1 from sklearn.preprocessing import MultiLabelBinarizer
2 mlb = MultiLabelBinarizer()
3
4 y = mlb.fit_transform(y)
5
6 print(type(y) , y.shape)
7 print(type(x) , x.shape)
```

```
<class 'numpy.ndarray'> (47427, 10)
<class 'numpy.ndarray'> (47427, 1)
```

```
1 from sklearn.model_selection import train_test_split
2
3 # Split the data into training and testing sets
4 X_train, X_test_valid, y_train, y_test_valid = train_test_split(x, y,test_size
5
6 # Further split the testing set into validation and testing sets
7 X_valid, X_test, y_valid, y_test = train_test_split(X_test_valid,y_test_valid,
```

```
1 if 'google.colab' in str(get_ipython()):
2     from google.colab import drive
3     drive.mount('/content/drive')
4
5     !pip install -U nltk -qq
6     !pip install -U spacy -qq
7     !python -m spacy download en_core_web_sm -qq
8
9     basepath = '/content/drive/MyDrive/NLP'
10    sys.path.append('/content/drive/MyDrive/NLP')
11 else:
12    basepath = '/Users/anxiousviking/Documents/course/Sem 3/NLP'
13    sys.path.append(
14    '/Users/anxiousviking/Documents/course/Sem 3/NLP/custom files')
```

```
Mounted at /content/drive
                                                      12.8/12.8 MB 80.2 MB/s eta 0:00
✔ Download and installation successful
You can now load the package via spacy.load('en_core_web_sm')
⚠ Restart to reload dependencies
If you are in a Jupyter or Colab notebook, you may need to restart Python in
order to load all the package's dependencies. You can do this by selecting the
'Restart kernel' or 'Restart runtime' option.
```

```
1 import CustomPreprocessorSpacy as cp
```

```
1 import spacy
2 nlp = spacy.load('en_core_web_sm')
3 cpp = cp.SpacyPreprocessor(model = 'en_core_web_sm', batch_size=1000)
```

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 vectorizer = TfidfVectorizer(analyzer='word', token_pattern=r"[\S]+",max_featu
3 vectorizer.fit(cpp.transform([str(x) for x in X_train.tolist()]))
```

```
 1 from scipy.sparse import csr_matrix
 2 class CustomDataset(torch.utils.data.Dataset):
 3     """
 4     Custom Dataset class for loading text and labels.
 5
 6     Attributes:
 7         X (numpy.ndarray): Feature data, an array of texts.
 8         y (list or array-like): Target labels.
 9         vectorizer (TfidfVectorizer): The TF-IDF vectorizer used to transform
10     """
11
12     def __init__(self, X, y, vectorizer,cpp):
13         """
14         Initialize the dataset with feature and target data.
15
16         Args:
17             X (list or array-like): The feature data (texts).
18             y (list or array-like): The target labels.
19             vectorizer (TfidfVectorizer): The TF-IDF vectorizer used to transf
20         """
21         X = [str(x) for x in X.tolist()]
22
23         # Storing feature data (texts)
24         self.X = cpp.transform(X)
25
26         # Storing the target labels
27         self.y = y
28
```

```
28
29          # Storing the TF-IDF vectorizer
30          self.vectorizer = vectorizer
31
32          # Transforming the texts to TF-IDF vectors
33          self.X_tfidf = self.vectorizer.transform(self.X)
34
35      def __len__(self):
36          """
37          Return the number of samples in the dataset.
38
39          Returns:
40              int: The total number of samples.
41          """
42          return len(self.X)
43
44      def __getitem__(self, idx):
45          """
46          Fetch and return a single sample from the dataset at the given index.
47
48          Args:
49              idx (int): Index of the sample to fetch.
50
51          Returns:
52              tuple: A tuple containing the label and the TF-IDF vector for the
53          """
54          # Retrieve the TF-IDF vector and corresponding label from the dataset
55          tfidf_vector = csr_matrix.toarray(self.X_tfidf[idx])
56          label = self.y[idx]
57
58          # Convert label to tensor of type float
59          label = torch.tensor(label, dtype=torch.float)
60
61          # Convert TF-IDF vector to tensor
62          tfidf_vector = torch.tensor(tfidf_vector, dtype=torch.float)
63
64          # Packing them into a tuple before returning
65
66          return tfidf_vector, label
```

```
1 # Create an instance of the CustomDataset class for the training set
2 trainset = CustomDataset(X_train, y_train,vectorizer,cpp)
3
4 # Create an instance of the CustomDataset class for the validation set
5 validset = CustomDataset(X_valid, y_valid,vectorizer,cpp)
6
7 # Create an instance of the CustomDataset class for the test set
8 testset = CustomDataset(X_test, y_test,vectorizer,cpp)
```

```
1 batch size = 2
```

```
 2 check_loader = torch.utils.data.DataLoader(dataset=trainset,
 3                                             batch_size=batch_size,
 4                                             shuffle=True,
 5                                             )
```

```
 1 class CustomModel(nn.Module):
 2     def __init__(self, input_dim, hidden_dim1, hidden_dim2, drop_prob1, drop_p
 3         super().__init__()
 4         self.hidden1 = nn.Linear(input_dim, hidden_dim1)
 5         self.relu1 = nn.ReLU()
 6         self.dropout1 = nn.Dropout(p=drop_prob1)
 7         self.batchnorm1 = nn.BatchNorm1d(num_features=hidden_dim1)
 8         self.hidden2 = nn.Linear(hidden_dim1, hidden_dim2)
 9         self.relu2 = nn.ReLU()
10         self.dropout2 = nn.Dropout(p=drop_prob2)
11         self.batchnorm2 = nn.BatchNorm1d(num_features=hidden_dim2)
12         self.output = nn.Linear(hidden_dim2, output_dim)
13
14     def forward(self, x):
15         x = self.hidden1(x)
16         x = self.relu1(x)
17         x = self.dropout1(x)
18         x = self.batchnorm1(x)
19         x = self.hidden2(x)
20         x = self.relu2(x)
21         x = self.dropout2(x)
22         x = self.batchnorm2(x)
23         x = self.output(x)
24         return x
```

```
 1 INPUT_DIM=5000
 2 HIDDEN_DIM1=200
 3 HIDDEN_DIM2=100
 4 DROP_PROB1=0.5
 5 DROP_PROB2=0.5
 6 NUM_OUTPUTS = 10
 7 EPOCHS=5
 8 BATCH_SIZE=128
 9 LEARNING_RATE=0.001
10 WEIGHT_DECAY=0.000
11 CLIP_VALUE = 10
12 PATIENCE = 5
13 dropout_p = 0.3
```

```
 1 import torch.optim as optim
 2 from torch.utils.data import DataLoader
 3 from tqdm import tqdm
 4
```

```
 5 model = CustomModel(input_dim=INPUT_DIM,
 6                     hidden_dim1=HIDDEN_DIM1,
 7                     hidden_dim2=HIDDEN_DIM2,
 8                     drop_prob1=0.5,
 9                     drop_prob2=0.5,
10                     output_dim=NUM_OUTPUTS)
```

```
 1 summary(model,(1, 5000))
```

```
===============================================================================
Layer (type:depth-idx)                     Output Shape              Param #
===============================================================================
CustomModel                                [1, 10]                   --
├─Linear: 1-1                              [1, 200]                  1,000,200
├─ReLU: 1-2                                [1, 200]                  --
├─Dropout: 1-3                             [1, 200]                  --
├─BatchNorm1d: 1-4                         [1, 200]                  400
├─Linear: 1-5                              [1, 100]                  20,100
├─ReLU: 1-6                                [1, 100]                  --
├─Dropout: 1-7                             [1, 100]                  --
├─BatchNorm1d: 1-8                         [1, 100]                  200
├─Linear: 1-9                              [1, 10]                   1,010
===============================================================================
Total params: 1,021,910
Trainable params: 1,021,910
Non-trainable params: 0
Total mult-adds (M): 1.02
===============================================================================
Input size (MB): 0.02
Forward/backward pass size (MB): 0.00
Params size (MB): 4.09
Estimated Total Size (MB): 4.11
===============================================================================
```

```
 1 # Define the device
 2 device = torch.device('cuda:0' if torch.cuda.is_available() else 'cpu')
 3
 4 # Move the model to the device
 5 model = model.to(device)
 6
 7 # Generate some dummy input data and offsets, and move them to the device
 8 data = trainset[0][0].squeeze(1).to(device)
```

```
 1 # Switch the model to evaluation mode
 2 model.eval()
 3
 4 # Perform inference
 5 output = model(data)
 6
 7 print(output)
```

```
    tensor([[ 0.0666,  0.0379,  0.0703, -0.0464,  0.0093,  0.0948, -0.0247, -0.06;
             -0.0508, -0.0895]], grad_fn=<AddmmBackward0>)
```

```python
 1 !pip install torchmetrics
 2 from torchmetrics import HammingDistance
 3
 4 def step(inputs, targets, model, device, loss_function=None, optimizer=None,cl
 5     """
 6     Performs a forward and backward pass for a given batch of inputs and targe
 7
 8     Parameters:
 9     - inputs (torch.Tensor): The input data for the model.
10     - targets (torch.Tensor): The true labels for the input data.
11     - model (torch.nn.Module): The neural network model.
12     - device (torch.device): The computing device (CPU or GPU).
13     - loss_function (torch.nn.Module, optional): The loss function to use.
14     - optimizer (torch.optim.Optimizer, optional): The optimizer to update mod
15
16     Returns:
17     - loss (float): The computed loss value (only if loss_function is not None
18     - outputs (torch.Tensor): The predictions from the model.
19     - train_hamming_distance (torchmetrics.HammingDistance): The Hamming dista
20     """
21     # Move the model and data to the device
22     model = model.to(device)
23     inputs, targets = inputs.squeeze(1).to(device), targets.to(device)
24
25     # Step 1: Forward pass to get the model's predictions
26     outputs = model(inputs)
27
28     # Step 2a: Compute the loss using the provided loss function
29     if loss_function:
30         loss = loss_function(outputs, targets)
31
32     # Step 2b: Update Hamming Distance metric
33     train_hamming_distance = HammingDistance(task="multilabel", num_labels=10)
34     y_pred = (outputs > 0.5).float()
35     train_hamming_distance.update(y_pred, targets)
36
37     # Step 3 and 4: Perform backward pass and update model parameters if an op
38     if optimizer:
39         optimizer.zero_grad()
40         loss.backward()
41         if clip_type == 'value':
42             torch.nn.utils.clip_grad_value_(model.parameters(), clip_value)
43         optimizer.step()
44
45     # Return relevant metrics
46     if loss_function:
47         return loss, outputs, train_hamming_distance
```

```
48      else:
49          return outputs, train_hamming_distance
```

```
Collecting torchmetrics
  Downloading torchmetrics-1.4.2-py3-none-any.whl.metadata (19 kB)
Requirement already satisfied: numpy>1.20.0 in /usr/local/lib/python3.10/dist-
Requirement already satisfied: packaging>17.1 in /usr/local/lib/python3.10/di:
Requirement already satisfied: torch>=1.10.0 in /usr/local/lib/python3.10/dis
Collecting lightning-utilities>=0.8.0 (from torchmetrics)
  Downloading lightning_utilities-0.11.7-py3-none-any.whl.metadata (5.2 kB)
Requirement already satisfied: setuptools in /usr/local/lib/python3.10/dist-p;
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.10,
Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-pacl
Requirement already satisfied: sympy in /usr/local/lib/python3.10/dist-package
Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-pacl
Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packag
Requirement already satisfied: fsspec in /usr/local/lib/python3.10/dist-packag
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/d:
Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.1(
Downloading torchmetrics-1.4.2-py3-none-any.whl (869 kB)
                                        869.2/869.2 kB 16.4 MB/s eta 0:00
Downloading lightning_utilities-0.11.7-py3-none-any.whl (26 kB)
Installing collected packages: lightning-utilities, torchmetrics
Successfully installed lightning-utilities-0.11.7 torchmetrics-1.4.2
```

```python
1 def train_epoch(train_loader, model, device, loss_function, optimizer):
2     """
3     Trains the model for one epoch using the provided data loader and updates
4
5     Parameters:
6     - train_loader (torch.utils.data.DataLoader): DataLoader object for the tr
7     - model (torch.nn.Module): The neural network model to be trained.
8     - device (torch.device): The computing device (CPU or GPU).
9     - loss_function (torch.nn.Module): The loss function to use for training.
10    - optimizer (torch.optim.Optimizer): The optimizer to update model paramet
11
12    Returns:
13    - train_loss (float): Average training loss for the epoch.
14    - epoch_hamming_distance (float): Hamming distance for the epoch.
15    """
16    # Set the model to training mode
17    model.train()
18
19    # Initialize variables to track running training loss and correct predicti
20    running_train_loss = 0.0
21    running_train_correct = 0
22
23    # Initialize Hamming Distance metric
24    hamming = HammingDistance(task="multilabel", num_labels=10).to(device)
25
26    # Iterate over all batches in the training data
27    for inputs, targets in train_loader:
```

```
27     for inputs, targets in train_loader:
28         # Move data to the appropriate device
29         inputs, targets = inputs.squeeze(1).to(device), targets.to(device)
30
31         # Perform a forward and backward pass, updating model parameters
32         loss, _, _ = step(inputs, targets, model, device, loss_function, optim
33
34         # Update running loss
35         running_train_loss += loss.item()
36
37         # Compute Hamming Distance for the epoch
38         y_pred = (model(inputs) > 0.5).float()
39         hamming.update(y_pred, targets)
40
41     # Compute average loss for the entire training set
42     train_loss = running_train_loss / len(train_loader)
43
44     # Compute Hamming Distance for the epoch
45     epoch_hamming_distance = hamming.compute()
46
47     return train_loss, epoch_hamming_distance
```

```
 1 from torchmetrics import HammingDistance
 2
 3 def val_epoch(valid_loader, model, device, loss_function):
 4     """
 5     Validates the model for one epoch using the provided data loader.
 6
 7     Parameters:
 8     - valid_loader (torch.utils.data.DataLoader): DataLoader object for the va
 9     - model (torch.nn.Module): The neural network model to be validated.
10     - device (torch.device): The computing device (CPU or GPU).
11     - loss_function (torch.nn.Module): The loss function to evaluate the model
12
13     Returns:
14     - val_loss (float): Average validation loss for the epoch.
15     - val_hamming_distance (float): Hamming distance for the epoch.
16     """
17     # Set the model to evaluation mode
18     model.eval()
19
20     # Initialize variables to track running validation loss and Hamming Distan
21     running_val_loss = 0.0
22     val_hamming_distance = HammingDistance(task="multilabel", num_labels=10).t
23
24     # Disable gradient computation
25     with torch.no_grad():
26         # Iterate over all batches in the validation data
27         for inputs, targets in valid_loader:
28             # Move data to the appropriate device
29             inputs, targets = inputs.squeeze(1).to(device), targets.to(device)
```

```python
29              inputs, targets = inputs.squeeze(1).to(device), targets.to(device)
30
31              # Perform a forward pass to get loss and number of correct predict
32              outputs = model(inputs)
33              loss = loss_function(outputs, targets)
34
35              # Update running loss
36              running_val_loss += loss.item()
37
38              # Update Hamming Distance metric
39              val_hamming_distance.update(torch.round(torch.sigmoid(outputs)), t
40
41      # Compute average loss and Hamming Distance for the entire validation set
42      val_loss = running_val_loss / len(valid_loader)
43      val_hamming_distance = val_hamming_distance.compute()
44
45      return val_loss, val_hamming_distance
46
```

```python
 1 def train(train_loader, valid_loader, model, optimizer, loss_function, epochs,
 2     """
 3     Trains and validates the model, and returns history of train and validatio
 4
 5     Parameters:
 6     - train_loader (torch.utils.data.DataLoader): DataLoader for the training
 7     - valid_loader (torch.utils.data.DataLoader): DataLoader for the validatio
 8     - model (torch.nn.Module): Neural network model to train.
 9     - optimizer (torch.optim.Optimizer): Optimizer algorithm.
10     - loss_function (torch.nn.Module): Loss function to evaluate the model.
11     - epochs (int): Number of epochs to train the model.
12     - device (torch.device): The computing device (CPU or GPU).
13     - patience (int): Number of epochs to wait for improvement before early st
14
15     Returns:
16     - train_loss_history (list): History of training loss for each epoch.
17     - train_hamm_history (list): History of training Hamming distance for each
18     - valid_loss_history (list): History of validation loss for each epoch.
19     - valid_hamm_history (list): History of validation Hamming distance for ea
20     """
21
22     # Initialize lists to store metrics for each epoch
23     train_loss_history = []
24     valid_loss_history = []
25     train_hamm_history = []
26     valid_hamm_history = []
27
28     # Initialize variables for early stopping
29     best_valid_loss = float('inf')
30     no_improvement = 0
31
```

```
32      # Loop over the number of specified epochs
33      for epoch in range(epochs):
34          # Train model on training data and capture metrics
35          train_loss, train_hamm = train_epoch(
36              train_loader, model, device, loss_function, optimizer)
37
38          # Validate model on validation data and capture metrics
39          valid_loss, valid_hamm = val_epoch(
40              valid_loader, model, device, loss_function)
41
42          # Store metrics for this epoch
43          train_loss_history.append(train_loss)
44          valid_loss_history.append(valid_loss)
45          train_hamm_history.append(train_hamm)
46          valid_hamm_history.append(valid_hamm)
47
48          # Output epoch-level summary
49          print(f"Epoch {epoch+1}/{epochs}")
50          print(f"Train Loss: {train_loss:.4f} | Train Hamming Distance: {train_
51          print(f"Valid Loss: {valid_loss:.4f} | Valid Hamming Distance: {valid_
52          print()
53
54          # Check for early stopping
55          if valid_loss < best_valid_loss:
56              best_valid_loss = valid_loss
57              no_improvement = 0
58          else:
59              no_improvement += 1
60              if no_improvement == patience:
61                  print(f"No improvement for {patience} epochs. Early stopping..
62                  break
63
64      return train_loss_history, train_hamm_history, valid_loss_history, valid_h
65
```

```
1 # training
2 EPOCHS=5
3 BATCH_SIZE=128
4 LEARNING_RATE=0.001
5 WEIGHT_DECAY=0.0
6 PATIENCE=10
```

```
1 # Fixing the seed value for reproducibility across runs
2 SEED = 2345
3 random.seed(SEED)
4 np.random.seed(SEED)
5 torch.manual_seed(SEED)
6 torch.cuda.manual_seed(SEED)
7 torch.backends.cudnn.deterministic = True
```

```
 8
 9
10 # Define the device for model training (use CUDA if available, else CPU)
11 device = torch.device('cuda:0' if torch.cuda.is_available() else 'cpu')
12
13 # Data Loaders for training, validation, and test sets
14 train_loader = torch.utils.data.DataLoader(trainset, batch_size = BATCH_SIZE,
15 valid_loader = torch.utils.data.DataLoader(validset, batch_size=BATCH_SIZE, sh
16 test_loader = torch.utils.data.DataLoader(testset, batch_size=BATCH_SIZE, shuf
17
18 # Define the loss function for the model, using cross-entropy loss
19 loss_function = nn.BCEWithLogitsLoss()
20
21 # Define the model with specified hyperparameters
22 model = CustomModel(input_dim=INPUT_DIM,
23                        hidden_dim1=HIDDEN_DIM1,
24                        hidden_dim2=HIDDEN_DIM2,
25                        drop_prob1=0.5,
26                        drop_prob2=0.5,
27                        output_dim=NUM_OUTPUTS)
28
29 model = model.to(device)
30
31 # Define the optimizer
32 optimizer = optim.AdamW(model.parameters(), lr=LEARNING_RATE, weight_decay=WEI
```

```
    /usr/local/lib/python3.10/dist-packages/torch/utils/data/dataloader.py:557: U
      warnings.warn(_create_warning_msg(
```

```
 1 for inputs, targets in train_loader:
 2   print(type(inputs))
 3   print(inputs.shape)
 4   print(targets.shape)
 5   break
```

```
    /usr/local/lib/python3.10/dist-packages/torch/utils/data/dataloader.py:557: U
      warnings.warn(_create_warning_msg(
    <class 'torch.Tensor'>
    torch.Size([128, 1, 5000])
    torch.Size([128, 10])
```

```
 1 for inputs, targets in train_loader:
 2     # Move inputs and targets to the CPU.
 3     inputs = inputs.squeeze(1).to(device)
 4     targets = targets.to(device)
 5     model = model.to(device)
 6     model.eval()
 7     # Forward pass
 8     with torch.no_grad():
 9         output = model(inputs)
```

```
10           loss = loss_function(output, targets)
11           print(f'Actual loss: {loss.item()}')
12       break
13
14 print(f'Expected Theoretical loss: {np.log(2)}')
```

```
    Actual loss: 0.6789126396179199
    Expected Theoretical loss: 0.6931471805599453
```

```
 1 CLIP_VALUE = 10
 2 # Call the train function to train the model
 3 train_losses, train_hamm, valid_losses, valid_hamm = train(
 4     train_loader, valid_loader, model, optimizer, loss_function, EPOCHS, devic
 5 )
```

```
    Epoch 1/5
    Train Loss: 0.3474 | Train Hamming Distance: 0.0855
    Valid Loss: 0.1411 | Valid Hamming Distance: 0.0491

    Epoch 2/5
    Train Loss: 0.1348 | Train Hamming Distance: 0.0474
    Valid Loss: 0.1161 | Valid Hamming Distance: 0.0437

    Epoch 3/5
    Train Loss: 0.1074 | Train Hamming Distance: 0.0391
    Valid Loss: 0.1085 | Valid Hamming Distance: 0.0392

    Epoch 4/5
    Train Loss: 0.0920 | Train Hamming Distance: 0.0330
    Valid Loss: 0.1041 | Valid Hamming Distance: 0.0372

    Epoch 5/5
    Train Loss: 0.0827 | Train Hamming Distance: 0.0303
    Valid Loss: 0.1045 | Valid Hamming Distance: 0.0372
```

```
 1 def plot_history(train_losses, train_metrics, val_losses=None, val_metrics=Non
 2     """
 3     Plot training and validation loss and metrics over epochs.
 4
 5     Args:
 6         train_losses (list): List of training losses for each epoch.
 7         train_metrics (list): List of training metrics (e.g., accuracy) for ea
 8         val_losses (list, optional): List of validation losses for each epoch.
 9         val_metrics (list, optional): List of validation metrics for each epoc
10
11     Returns:
12         None
13     """
14     # Determine the number of epochs based on the length of train_losses
15     epochs = range(1, len(train_losses) + 1)
16
```

```
16
17     # Plotting training and validation losses
18     plt.figure()
19     plt.plot(epochs, train_losses, label="Train")  # Plot training losses
20     if val_losses:  # Check if validation losses are provided
21         plt.plot(epochs, val_losses, label="Validation")
22     plt.xlabel("Epochs")
23     plt.ylabel("Loss")
24     plt.legend()
25     plt.show()
26
27     # Plotting training and validation metrics
28     if train_metrics[0] is not None:  # Check if training metrics are availabl
29         plt.figure()
30         plt.plot(epochs, train_metrics, label="Train")
31         if val_metrics:
32             plt.plot(epochs, val_metrics, label="Validation")
33         plt.xlabel("Epochs")
34         plt.ylabel("Metric")
35         plt.legend()
36         plt.show()
37
```

```
1 train_hamm
```

```
    [tensor(0.0855),
     tensor(0.0474),
     tensor(0.0391),
     tensor(0.0330),
     tensor(0.0303)]
```

```
1 import numpy as np
2 # Plot the training and validation losses and metrics
3 train_hamm_np = [ham.cpu().numpy() for ham in train_hamm]
4 valid_hamm_np = [ham.cpu().numpy() for ham in valid_hamm]
5 plot_history(train_losses, train_hamm_np, valid_losses, valid_hamm_np)
```

```
1 def get_acc_pred(data_loader, model, device):
2     """
3     Function to get predictions and accuracy for a given data using a trained
4     Input: data iterator, model, device
5     Output: predictions and accuracy for the given dataset
6     """
7     model = model.to(device)
8     # Set model to evaluation mode
9     model.eval()
10
11    # Create empty tensors to store predictions and actual labels
12    predictions = torch.Tensor().to(device)
13    y = torch.Tensor().to(device)
14
15    # Iterate over batches from data iterator
16    with torch.no_grad():
```

```
16        with torch.no_grad():
17            for inputs, targets in data_loader:
18                # Process the batch to get the loss, outputs, and correct predicti
19                outputs, _ = step(inputs, targets, model,
20                                  device, loss_function=None, optimizer=None)
21
22                # Choose the label with maximum probability
23                # Correct prediction using thresholding
24                y_pred = (outputs.data>0.5).float()
25
26                # Add the predicted labels and actual labels to their respective t
27                predictions = torch.cat((predictions, y_pred))
28                y = torch.cat((y, targets.to(device)))
29
30        # Calculate accuracy by comparing the predicted and actual labels
31        accuracy = (predictions == y).float().mean()
32
33        # Return tuple containing predictions and accuracy
34        return predictions, accuracy, y
```

```
1 # Get the prediction and accuracy
2 predictions_test, acc_test, y_test = get_acc_pred(test_loader, model, device)
3 predictions_train, acc_train, y_train = get_acc_pred(train_loader, model, devi
4 predictions_valid, acc_valid, y_valid = get_acc_pred(valid_loader, model, devi
```

```
1 # Print Test Accuracy
2 print('Valid accuracy', acc_valid * 100)
```

```
Valid accuracy tensor(96.1191)
```

```
1 from sklearn.metrics import multilabel_confusion_matrix
2
3 def plot_confusion_matrix(valid_labels, valid_preds, class_labels):
4     """
5     Plots a confusion matrix.
6
7     Args:
8         valid_labels (array-like): True labels of the validation data.
9         valid_preds (array-like): Predicted labels of the validation data.
10        class_labels (list): List of class names for the labels.
11    """
12    # Compute the confusion matrix
13    cm = multilabel_confusion_matrix(valid_labels, valid_preds)
14
15    # Plot the confusion matrix using Seaborn
16    fig, axs = plt.subplots(1, len(class_labels), figsize=(15, 5))
17    for i, (label, matrix) in enumerate(zip(class_labels, cm)):
18        sns.heatmap(matrix, annot=True, fmt="d", cmap="Reds", xticklabels=['0'
19        axs[i].set_title(f"Confusion Matrix for Class {label}")
```

```
20          axs[i].set_xlabel('Predicted Labels')
21          axs[i].set_ylabel('True Labels')
22
23      # Display the plot
24      plt.tight_layout()
25      plt.show()
```

```
1 plot_confusion_matrix(y_test.cpu().numpy(), predictions_test.cpu().numpy(), cl
```

```
1 test_hamming_distance = HammingDistance(task="multilabel", num_labels=10).to(dev
2 test_hamming_distance.update(y_test, predictions_test)
```

```
1 test_hamming_distance.compute()
```

```
    tensor(0.0379)
```

## ∨ Inferences

**Confusion Matrix**:

- **Negative Class**: True negatives are high at 7,215, indicating the model effectively predicts negative labels. However, there are 442 false negatives, showing instances where the model incorrectly predicts a sample as negative.
- **Positive Class**: True positives are also significant at 1,468, reflecting the model's strong performance in predicting positive labels. It has a low false positive rate of 65 and a moderate 199 false negatives, suggesting better performance on the positive class compared to the negative.

**Curves**:

- **Train Loss**: There is clear evidence of learning, with a steady decline in training loss. A sharp drop between epochs 1 and 2 indicates quick adaptation to the training data.
- **Validation Loss**: This loss decreases and plateaus around epoch 3, suggesting the model is nearing its generalization limit. The small gap between training and validation loss indicates good generalization without overfitting.
- **Stable Validation Loss**: The relatively flat curve after epoch 3 suggests that further training may not significantly enhance performance on the validation set.
- **Train Metric**: The model shows improved performance on the training data, correlating with the loss curve's significant drop between epochs 1 and 2.
- **Validation Metric**: This metric mirrors the training trend, decreasing sharply in early epochs and flattening around epoch 3, indicating strong performance on unseen data, although improvements may slow after this point.
- **Generalization**: As the validation metric flattens while the training metric continues to drop slightly, the model shows minimal overfitting, indicating potential benefits from early stopping after epochs 3 or 4.

1 Start coding or generate with AI.