

Deep Learning Homework 5 Report

colab.research.google.com/drive/1FTVCpTF7wHnU4VlrkxTUgtpKyZ10-gj

Online C Compiler - ... What is an HTML Se...

152120191023 Selim Can Yazar Hw5.ipynb

Doğa Düzene Göster Ekle Çalma zamanı Açılır Yürüt Son yükleme tarihi: 12-12

+ Kod + Main

```
# Import necessary libraries
import pandas as pd
import numpy as np
import torch
from torch.utils.data import DataLoader, Dataset
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix, f1_score, accuracy_score
from sklearn.feature_extraction.text import CountVectorizer
import seaborn as sns
import matplotlib.pyplot as plt

# Step 1: Load the dataset
email_data = pd.read_excel('content/email_data.xlsx') # Path to the uploaded file

# Step 2: Preprocess the data
text = email_data['text'].values
labels = email_data['spam'].values

# Convert text to numeric using CountVectorizer
vectorizer = CountVectorizer(max_features=5000) # Limit to 5000 features
x = vectorizer.fit_transform(text).toarray()

# Train-test split
x_train, x_test, y_train, y_test = train_test_split(x, labels, test_size=0.2, random_state=42)

# Convert to PyTorch tensors
x_train_tensor = torch.tensor(x_train, dtype=torch.float32)
y_train_tensor = torch.tensor(y_train, dtype=torch.long)
x_test_tensor = torch.tensor(x_test, dtype=torch.float32)
y_test_tensor = torch.tensor(y_test, dtype=torch.long)

# Step 3: Create a custom dataset
class EmailDataset(Dataset):
    def __init__(self, x, y):
        self.x = x
        self.y = y

    def __len__(self):
        return len(self.y)

    def __getitem__(self, idx):
        return self.x[idx], self.y[idx]

train_dataset = EmailDataset(x_train_tensor, y_train_tensor)
test_dataset = EmailDataset(x_test_tensor, y_test_tensor) # Optimized batch size
train_loader = DataLoader(train_dataset, batch_size=64, shuffle=True)
test_loader = DataLoader(test_dataset, batch_size=64, shuffle=False)
```

colab.research.google.com/drive/1FTVCpTF7wHnU4VlrkxTUgtpKyZ10-gj

Online C Compiler - ... What is an HTML Se...

152120191023 Selim Can Yazar Hw5.ipynb

Doğa Düzene Göster Ekle Çalma zamanı Açılır Yürüt Son yükleme tarihi: 12-12

+ Kod + Main

```
# Step 4: Define the LSTM model with optimizations
def _init_(self, input_size, hidden_size, output_size, num_layers):
    super(_init_, self).__init__()
    self.lstm = nn.LSTM(input_size, hidden_size, num_layers=num_layers, batch_first=True)
    self.fc = nn.Linear(hidden_size, output_size)
    self.relu = nn.ReLU() # ReLU activation function

def forward(self, x):
    x, _ = self.lstm(x)
    x = self.relu(x) # Apply ReLU activation to the last LSTM output
    x = self.fc(x)
    return x

input_size = x_train.shape[1]
hidden_size = 256 # Increased hidden size
output_size = 2
num_layers = 2 # Multi-layer LSTM

model = _init_(input_size, hidden_size, output_size, num_layers)
optimizer = optim.Adam(model.parameters()) # Optimized learning rate

# Step 5: Train the model
num_epochs = 20 # Increased number of epochs
for epoch in range(num_epochs):
    model.train()
    for x_batch, y_batch in train_loader:
        outputs = model(x_batch.unsqueeze(1)) # Add sequence dimension
        loss = criterion(outputs, y_batch)
        optimizer.zero_grad_()
        loss.backward()
        optimizer.step()
    print(f"Epoch {epoch+1}/{num_epochs}, Loss: {loss.item()}")

# Step 6: Save the model
torch.save(model.state_dict(), 'content/optimized_lstm_model.pth')

# Step 7: Load the model and evaluate
model_load_state_dict(torch.load('content/optimized_lstm_model.pth'))
model.eval()

# Evaluate on test data
y_true = []
y_pred = []

with torch.no_grad():
    for x_batch, y_batch in test_loader:
        outputs = model(x_batch.unsqueeze(1))
        _, predicted = torch.max(outputs, 1)
        y_true.extend(y_batch.numpy())
        y_pred.extend(predicted.numpy())
```

colab.research.google.com/drive/1FTVCpTF7wHnU4VlrkxTUgtpKyZ10-gj

Online C Compiler - ... What is an HTML Se...

152120191023 Selim Can Yazar Hw5.ipynb

Doğa Düzene Göster Ekle Çalma zamanı Açılır Yürüt Son yükleme tarihi: 12-12

+ Kod + Main

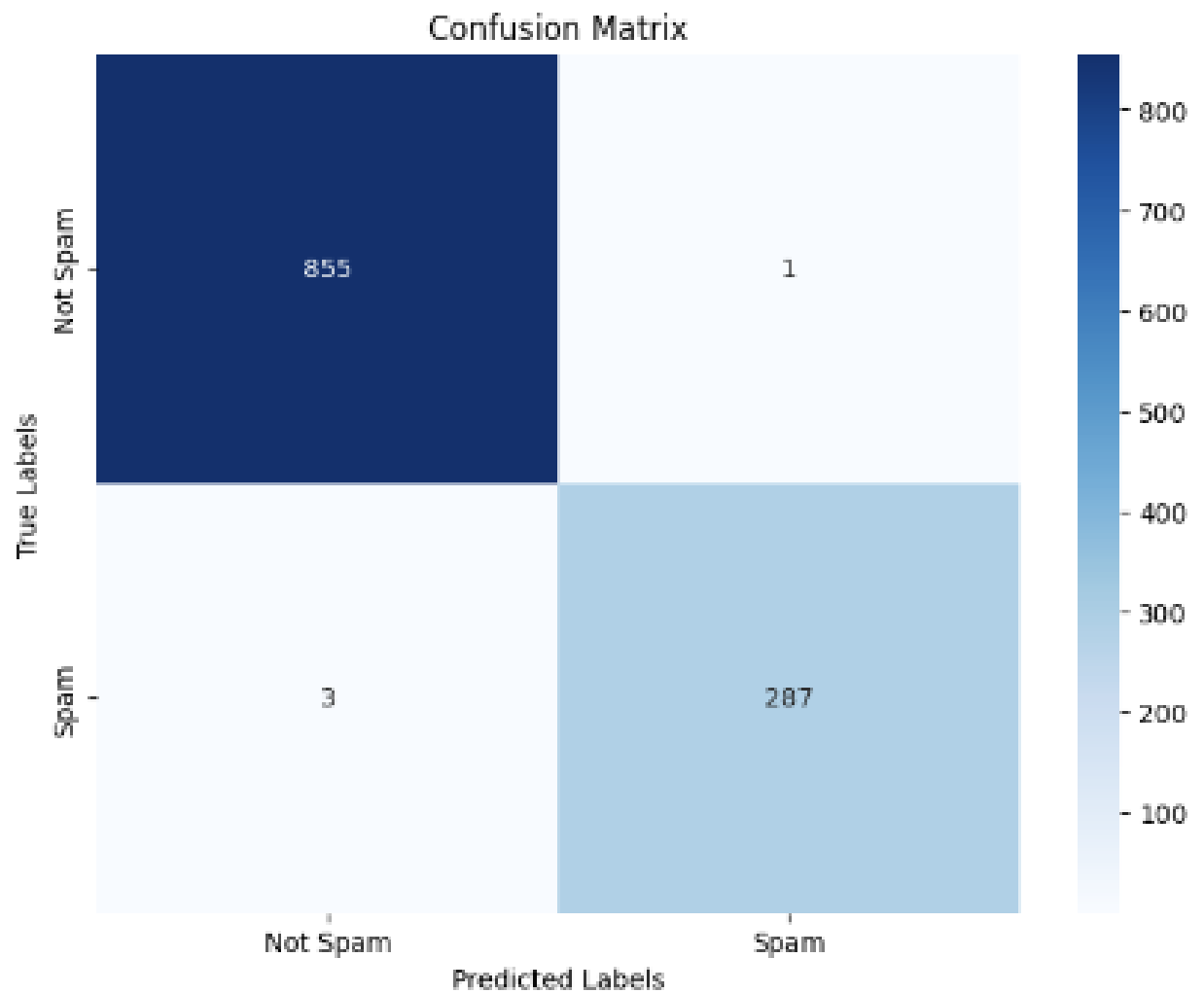
```
# Step 8: Calculate metrics
conf_matrix = confusion_matrix(y_true, y_pred)
f1 = f1_score(y_true, y_pred)
accuracy = accuracy_score(y_true, y_pred)

print("Confusion Matrix:")
print(conf_matrix)
print("F1 Score: (f1)")
print("Accuracy: (accuracy)")

# Step 9: Visualize the confusion matrix
plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', xticklabels=['Not Spam', 'Spam'], yticklabels=['Not Spam', 'Spam'])
plt.xlabel('Predicted Labels')
plt.ylabel('True Labels')
plt.title('Confusion Matrix')
plt.show()

Epoch 1/20, Loss: 0.8080059507601245
Epoch 2/20, Loss: 0.8001210817680884
Epoch 3/20, Loss: 0.80017351393178973
Epoch 4/20, Loss: 0.8001710013881077
Epoch 5/20, Loss: 0.80005202208970173
Epoch 6/20, Loss: 0.8001468017017176
Epoch 7/20, Loss: 0.80146786972205105
Epoch 8/20, Loss: 0.8021652589730122
Epoch 9/20, Loss: 0.80003464388285444
Epoch 10/20, Loss: 0.80077551296387894
Epoch 11/20, Loss: 2.18666752406748e-05
Epoch 12/20, Loss: 0.80019028178745646
Epoch 13/20, Loss: 0.80011013052370426
Epoch 14/20, Loss: 6.2264087862324e-05
Epoch 15/20, Loss: 1.50848834700484e-05
Epoch 16/20, Loss: 3.302486870755e-05
Epoch 17/20, Loss: 3.115984778111355e-06
Epoch 18/20, Loss: 6.7798152385255e-06
Epoch 19/20, Loss: 7.58427738115109e-06
Epoch 20/20, Loss: 4.6761087745238e-06
```

```
Confusion Matrix:  
[[855  1]  
 [ 3 287]]  
F1 Score: 0.9938795847750865  
Accuracy: 0.9965095986038395
```



Importing Libraries and Loading the Dataset

The code first imports the libraries needed for spam classification. Pandas and NumPy are used for data manipulation, while PyTorch is used for deep learning operations. Scikit-learn is included for training-test splitting, metric calculation, and text vectorization operations. Seaborn and Matplotlib are used as visualization tools. The dataset is loaded as an Excel file containing email texts and labels. This file is passed to the `emails_data` variable.

Processing and Digitizing the Data

In the dataset, the `text` column contains the email content, and the `spam` column contains the labels indicating whether these emails are spam or not. `CountVectorizer` is used to convert the email texts into a numeric form. With this method, the texts are converted into a numeric matrix according to word frequencies and limited to a maximum of 5000 features. Then, this digitized data is divided into training and test sets by 80% and 20%.

Conversion to PyTorch Tensors

The obtained training and test data are converted to tensors so that they can be processed by PyTorch. This process allows the data to be put into a form that can be processed by the model. The input data is converted to `float32` and the labels to `long` data type.

Creating Custom Dataset and Data Loaders

The `EmailDataset` class defines a customized dataset to be used in PyTorch's data loaders. This class allows the data to be separated into two parts: input (X) and label (y). Then, separate `DataLoader` objects are created for the training and test data sets. These loaders increase efficiency during model training by loading data in mini-batches.

Defining the LSTM Model

The code defines an LSTM-based classifier model. The model consists of an LSTM layer, a ReLU activation function, and a fully connected layer. LSTM handles the sequential nature of the input data while ReLU provides a nonlinear transformation. The output of the model provides an accuracy that can predict between two classes (spam and non-spam).

Training the Model

The model is trained using the CrossEntropyLoss loss function and the Adam optimization algorithm. During training, the loss between the model's predictions and the true values is calculated sequentially on mini-batches. This loss is then used to update the model weights using backpropagation. This process is repeated for 20 epochs.

Saving and Loading the Model

After the training is complete, the weights of the model are saved in a `.pth` file. The model is then loaded from this file and evaluated during the testing phase.

Evaluating the Results

The model is run on the test data and the predictions are obtained. The obtained predictions are compared with the actual labels to calculate a confusion matrix, F1 score and accuracy.

Visualization of Results

The confusion matrix is visualized to make the model's performance more understandable. The results show that the model classifies spam and non-spam emails with very high accuracy. The F1 score is 0.993 and the accuracy is as high as 0.996. Below, I have interpreted the table you provided as a result of the visualization:

- Confusion Matrix:

- There are 855 correctly classified non-spam emails and 287 correctly classified spam.

- Incorrect classifications: 1 non-spam email is labeled as "spam" and 3 spam are labeled as "non-spam". These results show that the model works very well and only makes a few incorrect classifications.

