# Indian Institute of Technology, Patna

**Assignment-4**
**Submitted By:**

- Avinash Aanand-(Roll No.: 2403RES99)
- Aditya Gupta-(Roll No.:2403RES85)

- Due Date: 11th Nov 2024

# Table of Content

# CS582 Neural Network and NLP Lab

# CS582 Neural Network and NLP Lab

# Problem Statements:

## Statements

Implement a character level language model using LSTM cell. The main task of the characterlevel language model is to predict the next character given all previous characters in a sequence of data, i.e. generates text character by character. More formally, given a training sequence ($x^1$, … , $x^T$), the RNN uses the sequence of its output vectors ($o^1$, … , $o^T$) to obtain a sequence of predictive distributions $P(x^t/x^{t-1}) = softmax(o^t)$. The example training sequence can be "hello" and the vocabulary of your data is [h,e,l,o]. As the dataset, you can take any text data written in English language.

## Tasks

The task you've outlined is to build a character-level language model using an LSTM (Long Short-Term Memory) neural network. This model will predict the next character in a sequence based on all previous characters. The general goal is to train the model to learn the probability distribution of characters in a given text and use it to generate text one character at a time.
To implement this, the task can be broken down into several steps:

**Tasks to Perform:**
1. **Dataset Preparation**:
   o **Load a text dataset**: You will need a corpus of text (any English text) as input. This can be Shakespeare's works, a large book, or any corpus of your choice.
   o **Create a character vocabulary**: Extract the unique characters from the text and map each character to a unique index, and vice versa.
   o **Encode the text**: Convert the text into a sequence of integers, where each integer corresponds to a character in the vocabulary.
2. **Define the Model**:
   o **Create a simple LSTM model**: The LSTM will take in a sequence of one-hot encoded characters, and output the probabilities for the next character in the sequence.
   o **Use softmax**: After the LSTM layer, apply a softmax activation to predict the next character at each time step.
3. **Training Loop**:
   o **Create training batches**: Split the text into sequences of a fixed length. For each sequence, the task is to predict the next character given the previous characters.
   o **Forward pass**: Pass the input sequence through the LSTM and compute the predictions.
   o **Loss calculation**: Use CrossEntropyLoss to compare the predicted probabilities with the actual next character.
   o **Backpropagation**: Update the model parameters using an optimizer like Adam or SGD.
4. **Text Generation**:
   o **Seed the model**: Provide a starting string (e.g., "hello") and use the trained model to generate the next characters iteratively.
   o **Sampling**: At each step, sample from the predicted distribution (using softmax) to choose the next character.
5. **Evaluate the Model**:
   o After training, generate text from the model and evaluate the quality of the generated sequences.

**Step-by-step breakdown of tasks:**
1. **Preprocessing**:
   o Load the text data.
   o Create char_to_idx and idx_to_char mappings.
   o Encode the entire text into a sequence of integers based on the character mappings.
2. **LSTM Model Definition**:
   o Define an LSTM-based model using nn.LSTM in PyTorch.
   o Define the output layer to convert the LSTM output to a probability distribution of the next character.

3. **Training**:
    - o Create training batches from the encoded text (input-output pairs).
    - o Train the model using a loss function like CrossEntropyLoss and an optimizer like Adam.
4. **Text Generation**:
    - o After training, use the model to generate text character by character by feeding the output back into the model to predict the next character.
5. **Evaluation**:
    - o Print out generated sequences and evaluate the quality of the model's predictions.

**Tasks Summary:**
1. **Dataset Preprocessing**:
    - o Load and preprocess the dataset (text).
    - o Map characters to integers.
2. **Model Creation**:
    - o Create an LSTM-based model for character prediction.
3. **Training**:
    - o Train the model using the preprocessed data.
4. **Text Generation**:
    - o Generate text character by character using the trained model.

# Introduction

**Introduction to Character-Level Language Models Using LSTM**

**What is a Language Model?**

A **language model** is a machine learning model designed to predict the next word or character in a sequence of text. It is widely used in various natural language processing (NLP) tasks such as machine translation, speech recognition, text generation, and autocomplete systems.

- **Word-level language models** predict the next word based on the preceding words.
- **Character-level language models**, on the other hand, work at the level of individual characters, predicting the next character in a sequence based on the characters that precede it.

Character-level language models are particularly interesting because they do not require pre-tokenized text (i.e., they don't need words or sentences to be split beforehand). Instead, they work directly with raw text, making them more flexible for tasks where words are not well-defined, such as generating creative text, handling uncommon words, or working with languages without spaces between words.

**What is LSTM?**

**Long Short-Term Memory (LSTM)** is a type of recurrent neural network (RNN) designed to learn and model sequential data. It addresses one of the key limitations of traditional RNNs—**vanishing gradients**—by introducing memory cells that can retain information over long sequences.

In a character-level language model, an LSTM can be used to predict the next character in a sequence by considering all the previous characters. The LSTM is ideal for this task because it can maintain context over long sequences, making it effective in learning the structure of text.

**Problem Statement**

The task at hand is to build a **character-level language model** that predicts the next character in a given sequence of text, using an **LSTM model**. The model takes in a sequence of characters and outputs the probability distribution over the next character. The goal is to train this model using a large text corpus and then use it to **generate new text** by predicting one character at a time.

**Key Concepts**
1. **Character Encoding**:
    - o The model needs to map each character in the text to a unique integer. This is done using a **character-to-index (char_to_idx)** mapping, where each character is assigned an index, and the reverse mapping, **index-to-character (idx_to_char)**, is used to map the predicted index back to a character.
2. **Text Generation**:

- o Once trained, the model can generate text character by character. Given a starting string (e.g., "hello"), the model predicts the next character and appends it to the string. The generated character is then fed back into the model to predict the next character, and this process repeats to generate text of a desired length.

3. **Sequence Prediction**:
   - o The core of the task is sequence prediction, where the model learns to predict the next character based on the sequence of characters seen so far. This is formulated as a **probability distribution** over the next character: $P(x_t \mid x_{t-1}, x_{t-2}, ..., x_1)$P(x_t | x_{t-1}, x_{t-2}, ..., x_1)$P(x_t \mid x_{t-1}, x_{t-2}, ..., x_1)$, where $x_t$x_txt$ is the next character in the sequence, and the model uses the previous characters to predict it.

## How the Model Works

The character-level language model works as follows:

1. **Input Sequence**: The model receives a sequence of characters (e.g., "hell").
2. **Encoding**: Each character is converted into an integer index using the **char_to_idx** mapping.
3. **LSTM Processing**: The model uses an **LSTM cell** to process this sequence and update its internal state (the hidden state) as it moves through the characters.
4. **Prediction**: The model uses the final hidden state of the LSTM to predict the next character in the sequence.
5. **Sampling**: The model outputs a **probability distribution** for the next character. A **softmax** function is used to turn the output of the LSTM into probabilities, and a character is sampled from this distribution.
6. **Text Generation**: The sampled character is then added to the sequence, and the process repeats to generate new text.

## Why Character-Level Models?

1. **Flexibility**: Character-level models are more flexible than word-level models, especially in languages where word boundaries are not clearly defined or when working with noisy text (e.g., typos, slang, or social media).
2. **Generalization**: By learning at the character level, the model can generate text that has never been seen before, including new words, as it can generate sequences of characters not restricted by the boundaries of pre-defined vocabulary.
3. **Compactness**: Character-level models have smaller vocabulary sizes compared to word-level models, as they work with the alphabet or character set of the language, rather than an entire word vocabulary. This can make them easier to train for tasks with limited training data.

## Example Application: Text Generation

Once trained, a character-level language model can be used for text generation. Given an initial seed text like "hello", the model will generate a sequence of characters that could resemble meaningful text based on the patterns it learned during training.

For example:

- Seed text: "hello"
- Model might generate: "hellothere", "helloworld", or even "helldo" depending on the learned patterns.

This is widely applicable in creative writing, music composition (for generating melodies), code generation, and even dialogue systems.

## Challenges

- **Overfitting**: Character-level models can sometimes overfit to the training data, especially if the dataset is small. Proper regularization techniques (such as dropout) and data augmentation may be needed to prevent this.
- **Training Time**: Training character-level models can be computationally intensive since the model must learn the structure of text at a fine-grained level (character-by-character).
- **Sampling Diversity**: Text generation may sometimes produce repetitive or nonsensical outputs. Techniques like **temperature sampling** can be used to control the randomness and creativity of the generated text.

**Hence, a character-level language model** using **LSTM** is designed to predict the next character in a sequence of text by learning from large corpora of text. The model learns the relationships between characters and can generate new text by predicting one character at a time. This approach is flexible, allowing it to handle a wide variety of textual data, from creative writing to code generation.

## Approaches

**Using LSTM**

Building a **character-level language model** using **LSTM** requires several key steps. These steps can be broadly categorized into **data preparation**, **model building**, **training**, and **text generation**. Each of these stages requires careful consideration to ensure the model works effectively. Below are the main approaches to solving this problem:

**1. Data Preparation**

Data preprocessing is crucial for ensuring that the model can efficiently learn from the text. The steps involved in data preparation include:

**a) Text Cleaning**

- **Raw Text**: Start with a raw text corpus (e.g., books, articles, or other written content). It's important to clean the text by removing unwanted characters like punctuation, special symbols, or any other non-alphanumeric characters (if necessary).
- **Normalization**: Standardizing the text can be beneficial. For example, converting all characters to lowercase so that the model doesn't treat uppercase and lowercase characters differently.

**b) Character Encoding**

- **Vocabulary Creation**: Extract all unique characters in the dataset to build a **vocabulary**. This vocabulary consists of all the possible characters the model will encounter, including alphabets, digits, and punctuation.
- **Character to Index Mapping**: Each character is assigned a unique index (an integer). This can be done by iterating over the unique characters and creating a dictionary (char_to_idx).
- **Index to Character Mapping**: The reverse mapping (idx_to_char) is needed to convert model predictions (indices) back to characters.

**c) Sequence Formation**

- **Training Sequences**: Break the text into sequences of fixed length (say, 25 characters). These sequences will serve as input data for training. For example, if the input sequence is "hello", the model will predict the next character after "hello".
- **Target Sequences**: For each input sequence, the target will be the next character in the text, i.e., the sequence shifted by one character.

This can be done with a sliding window approach to form multiple training samples from the text.

**d) Batching the Data**

- Group sequences into **batches** to allow efficient parallel computation during training. Batches will be of size batch_size, and within each batch, sequences are of length seq_length.

**2. Model Building (Character-Level LSTM)**

The core of this approach is designing a **character-level LSTM model** that can learn from sequences of characters and predict the next one.

**a) Input Layer**

- **One-Hot Encoding**: Since the model works at the character level, the input data is typically represented in a one-hot encoded format. For each input character in a sequence, create a one-hot vector of size vocab_size (the total number of unique characters). This is done using the **nn.functional.one_hot** method in PyTorch.

**b) LSTM Layer**

- The LSTM layer takes the one-hot encoded input sequence and processes it sequentially. LSTM units are good at remembering information over long sequences, making them ideal for this task.
- The LSTM will have two main outputs: the output sequence (out) and the updated **hidden state** (hidden), which will be passed along to the next time step.

**c) Fully Connected Layer**

- After the LSTM processes the input sequence, its output is passed through a **fully connected (linear) layer**. This layer will transform the LSTM's hidden state into a vector of size vocab_size, representing the probability distribution over all characters (i.e., the next character in the sequence).

**d) Activation Function**

- **Softmax Activation**: The model's output is passed through a **softmax** activation function, which converts the output into a probability distribution. This distribution tells us the likelihood of each character being the next one in the sequence.

**e) Hidden State Initialization**

- At the beginning of training or text generation, initialize the LSTM's hidden state (and cell state) to zeros. The hidden state will evolve as the model processes the input text.

## 3. Training the Model

Training the character-level LSTM involves optimizing the model to predict the next character given a sequence of previous characters.

**a) Loss Function**

- The **Cross-Entropy Loss** function is typically used for classification tasks like this, where the goal is to predict a discrete class (in this case, the next character).
- The loss function compares the model's predicted probability distribution (after softmax) to the actual character label (the next character in the sequence).

**b) Optimization**

- Use an optimization algorithm like **Adam** or **SGD** to minimize the loss function and update the model's parameters. The optimizer adjusts the weights of the LSTM and fully connected layers based on the loss gradient.

**c) Backpropagation**

- **Backpropagation** is used to compute the gradients of the loss with respect to the model's parameters. The LSTM uses **Backpropagation Through Time (BPTT)** to update the parameters over the entire sequence.

**d) Training Loop**

- The training loop runs for several epochs, where in each epoch:
  - A batch of data is processed.
  - The loss is computed and gradients are backpropagated.
  - The optimizer updates the model weights.
- **Hidden State Management**: Ensure to detach the hidden state to avoid retaining the entire computational graph across batches. This prevents memory overflow issues.

## 4. Text Generation (Sampling from the Model)

After training, the model can be used to generate new text by predicting the next character based on a given seed text.

**a) Starting Sequence**

- Provide a starting string (e.g., "hello") to initialize the sequence.

**b) Prediction and Sampling**

- The model predicts the next character by looking at the current sequence. The predicted character is sampled from the output distribution (using temperature sampling or greedy sampling).
- **Temperature Sampling**: Use a **temperature parameter** to control the randomness of predictions. A high temperature results in more randomness (creative text generation), while a low temperature makes the output more predictable.
- **Greedy Sampling**: Simply pick the character with the highest probability.

**c) Generation Process**

- The predicted character is appended to the current sequence, and this new sequence is fed back into the model to predict the next character. This process is repeated for the desired text length.

**d) Handling the Output**

- The model generates a sequence of characters one by one until the desired length is reached.
- Finally, the generated sequence is returned as the model's output.

## 5. Evaluation and Fine-Tuning

Once the model generates text, it's important to evaluate its performance:

**a) Qualitative Evaluation**
- **Generated Text**: Inspect the generated text to see how coherent, creative, or grammatically correct it is.
- **Creativity**: Evaluate how well the model generalizes and produces novel combinations of characters or words.

**b) Quantitative Evaluation**
- **Perplexity**: Perplexity is a common metric for evaluating language models. It measures how well the model predicts the next token in a sequence. Lower perplexity indicates better performance.
- **Loss**: The loss during training is also a useful indicator of model performance. A decreasing loss indicates that the model is learning and improving over time.

**c) Fine-Tuning**
- If the generated text is not satisfactory, you can fine-tune the model by adjusting parameters like the learning rate, batch size, sequence length, or model architecture (e.g., increasing the number of LSTM layers or units).
- You can also experiment with **dropout** to prevent overfitting, or use **gradient clipping** to stabilize training.

**Summary of Approaches:**
1. **Data Preparation**: Clean and preprocess the text, create character-to-index mappings, and divide the text into training sequences.
2. **Model Building**: Use an LSTM-based architecture with one-hot encoded inputs and a softmax output layer to predict the next character.
3. **Training**: Train the model using Cross-Entropy loss and backpropagation, while managing the hidden states carefully.
4. **Text Generation**: Use the trained model to generate text character-by-character by feeding the model the output from the previous prediction.
5. **Evaluation and Fine-Tuning**: Evaluate the generated text and fine-tune the model to improve performance.


# Analysis

## Through Algorithm (Flowchart)

The algorithm you provided involves multiple steps for building and training a character-level language model using an LSTM network. Here's an analysis of the code broken down into a flowchart.

**1. Data Preprocessing (Input Text Handling)**
1. **Load Text Data**
   - The training text (e.g., Shakespeare data) is loaded for the task.
2. **Character-to-Index Mapping**
   - Create a unique list of characters in the dataset.
   - Map each character to a corresponding index (char_to_idx) and vice versa (idx_to_char).
3. **Text Encoding**
   - Convert the text into a sequence of integers based on the character-to-index mapping.

**2. Creating Batches (Data Preparation)**
1. **Get Batches**
   - Break the encoded text into smaller batches of fixed sequence length (seq_length).
   - For each batch, input sequences are paired with target sequences where the target sequence is the original sequence shifted by one character (the next character).
2. **One-Hot Encoding**
   - Convert each input sequence into one-hot encoded vectors for the LSTM input layer. The dimension of each vector is equal to the vocabulary size.

**3. Model Definition (LSTM Model)**
1. **Define LSTM Model**

- o Initialize an LSTM model with a given hidden_dim (size of the hidden state) and vocab_size (size of the input and output space).
- o The model has:
  - An LSTM layer to process the sequences.
  - A fully connected (linear) layer that outputs a prediction of the next character in the sequence.

2. **Forward Pass**
   - o The LSTM processes the input sequence and outputs both the hidden states and predictions for each timestep.
   - o The fully connected layer produces a distribution over the possible next characters.

3. **Initialize Hidden State**
   - o At the beginning of training or text generation, the hidden state is initialized to zeros.

## 4. Training Loop

1. **Define Loss Function and Optimizer**
   - o The **CrossEntropyLoss** function is used to compute the error between the predicted output and the actual target sequence.
   - o Use an optimizer such as **Adam** to minimize the loss.

2. **Train Over Epochs**
   - o Loop through the dataset for a fixed number of epochs.
   - o For each batch, perform the following:
     - Perform the forward pass (input sequences → LSTM → prediction).
     - Calculate the loss between predicted and actual characters.
     - Perform backpropagation to compute gradients.
     - Update the model parameters using the optimizer.

3. **Clear Hidden States**
   - o Detach the hidden state after each batch to avoid backpropagating through the entire sequence, which prevents memory overflow.

## 5. Text Generation

1. **Initialize Generation Process**
   - o Start with a given seed text (e.g., "hello") to initiate the text generation.

2. **Generate Characters**
   - o Loop through the characters in the seed text and use the model to predict the next character.
   - o The output of the LSTM is passed through a **softmax** function to get a probability distribution of the next character.
   - o **Sampling**: A character is sampled from this distribution based on the **temperature** (higher temperature introduces more randomness).

3. **Update Sequence**
   - o Append the predicted character to the sequence and use the updated sequence as input to the model for the next prediction.
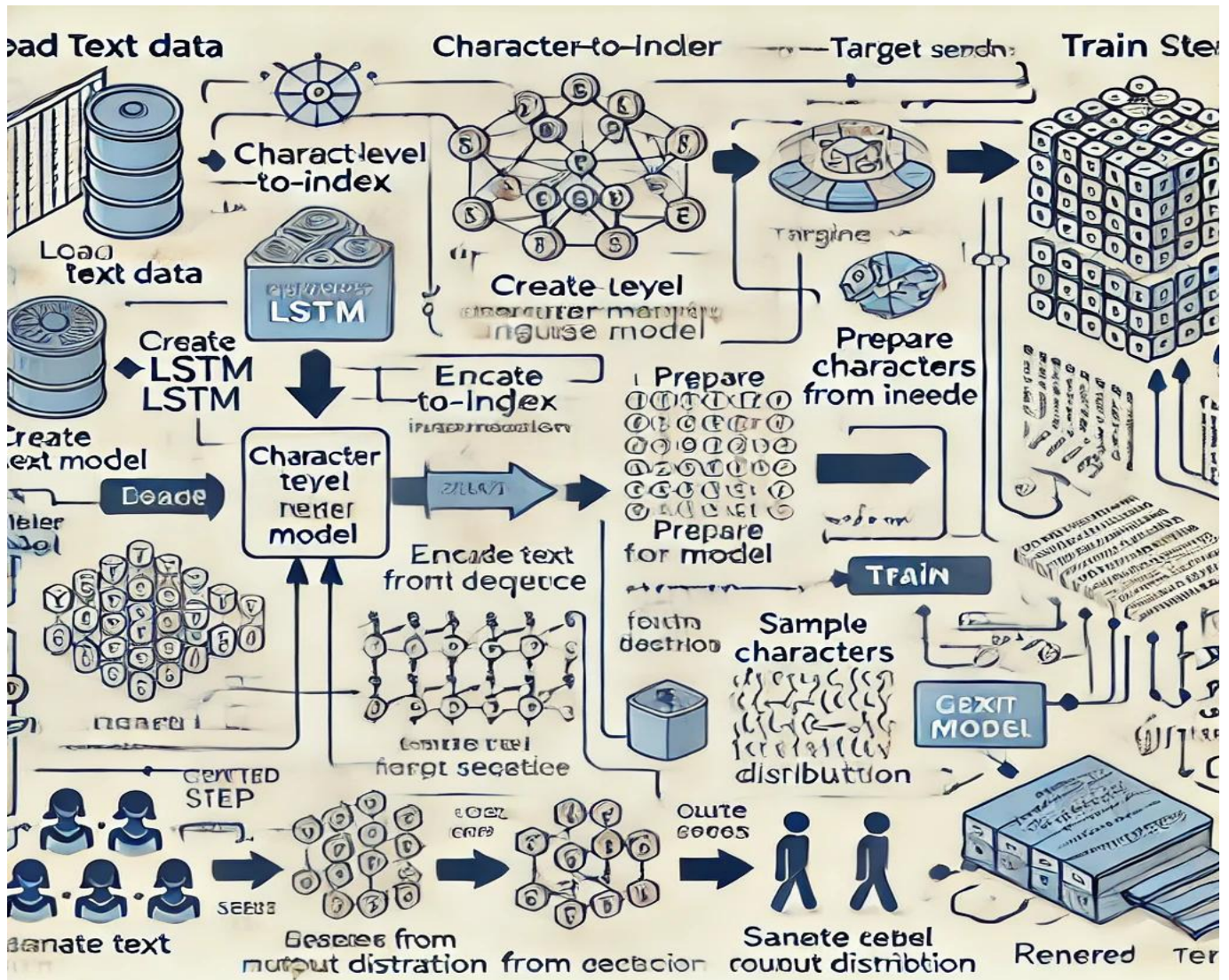
4. **Repeat Until Desired Length**
   - o This process continues until the desired length of generated text is reached.

## 6. Final Output

- The generated sequence of characters is returned as the model's output.

**Flowchart Representation**

Here's a textual representation of the flowchart:

**Key Components of the Flowchart:**

1. **Data Preprocessing**:
   - Character encoding and splitting data into batches.
2. **Model Definition**:
   - LSTM architecture and hidden state initialization.
3. **Training**:
   - Cross-entropy loss, backpropagation, and parameter update.
4. **Text Generation**:
   - Sampling characters from the LSTM's output and generating text iteratively.

This flowchart outlines the entire process, from data handling to model training and text generation, that powers the character-level language model built using LSTM.

## Through Pseudo Code

Breakdown of the analysis of the given code using pseudocode. I'll go through each step of the process and represent the logic in a simplified pseudocode format.

```
BEGIN
    IMPORT necessary libraries (torch, torch.nn, numpy)

    # Initialize the dataset and preprocess
```

```
LOAD 'Shakespeare_data' into text variable

# Step 1: Create character-to-index and index-to-character mappings
CREATE set of unique characters from text
SORT characters in ascending order
SET vocab_size to length of unique characters

# Create mappings
CREATE char_to_idx dictionary to map each character to a unique index
CREATE idx_to_char dictionary to map each index to a character

# Encode the text into integer sequence
ENCODE text into a sequence of integers using char_to_idx mapping

# Step 2: Define function to create batches for training
FUNCTION get_batches(encoded_text, batch_size, seq_length):
    CALCULATE total_batch_size = batch_size * seq_length
    CALCULATE num_batches = length of encoded_text divided by total_batch_size

    TRIM encoded_text to fit evenly into batches

    RESHAPE encoded_text into matrix of shape (batch_size, total_sequence_length)

    FOR each batch (x_batch, y_batch):
        SET x_batch to slice of sequence of length seq_length
        SET y_batch to shifted x_batch (shift by one)
        YIELD x_batch and y_batch as tensors

# Step 3: Define character-level LSTM model
CLASS CharLSTM:
    INITIALIZE LSTM model with input size vocab_size and hidden dimension hidden_dim
    INITIALIZE fully connected layer to predict output

    FUNCTION forward(x, hidden):
        PASS x through LSTM
        PASS LSTM output through fully connected layer
        RETURN output and hidden state

    FUNCTION init_hidden(batch_size):
        INITIALIZE hidden state (h0) and cell state (c0) with zeros of size (1, batch_size, hidden_dim)
        RETURN hidden state and cell state

# Step 4: Define training function
FUNCTION train(model, encoded_text, epochs, batch_size, seq_length, lr):
    INITIALIZE Adam optimizer with model parameters and learning rate
    INITIALIZE CrossEntropyLoss as loss function

    FOR each epoch from 1 to epochs:
        INITIALIZE hidden state using model.init_hidden(batch_size)
        INITIALIZE total_loss to 0

        FOR each batch (x_batch, y_batch) from get_batches:
            ONE-HOT encode x_batch into one-hot vectors of size vocab_size
            RESET gradients in optimizer
```

```
                PASS one-hot encoded x_batch and hidden state through the model
                CALCULATE loss using CrossEntropyLoss

                BACKPROPAGATE loss to compute gradients
                STEP optimizer to update model parameters

                DETACH hidden state to prevent computational graph reuse

            PRINT average loss for the epoch

    # Step 5: Train the model
    INITIALIZE model with vocab_size and hidden_dim
    CALL train function with model, encoded_text, epochs, batch_size, seq_length, and lr

    # Step 6: Text generation function
    FUNCTION generate_text(model, start_str, length, temperature):
        SET model to evaluation mode
        INITIALIZE chars with start_str
        INITIALIZE hidden state using model.init_hidden(1)

        FOR each character in start_str:
            CONVERT character to index using char_to_idx
            ONE-HOT encode the index and pass it through the model

        FOR generating characters (length times):
            APPLY temperature scaling to the model's output (softmax with temperature)
            SAMPLE next character based on the output distribution
            APPEND sampled character to chars
            ONE-HOT encode the sampled character and pass through the model

        RETURN generated text as a string

    # Step 7: Generate Text
    INITIALIZE start string (e.g., "hello")
    CALL generate_text function with model, start_str, desired length, and temperature
    PRINT generated text
END
```

## Through python Code.

**Step 1: Character-to-Index and Index-to-Character Mappings**

plaintext

Copy code

Create a set of unique characters in the text

Sort the characters

Determine the vocabulary size (length of unique characters)

Create two mappings:
   1. `char_to_idx` -> maps each character to a unique index
   2. `idx_to_char` -> maps each index back to its corresponding character

**Step 2: Function to Create Batches for Training**

plaintext

Copy code

Function get_batches(encoded_text, batch_size, seq_length):

total_batch_size = batch_size * seq_length

Calculate the number of batches as total number of characters divided by total_batch_size

Adjust encoded_text to match total_batch_size for even splitting into batches

Reshape the encoded_text into a matrix with rows = batch_size and columns = total_sequence_length

For each batch of size `seq_length`:

    x_batch = slice the sequence of length `seq_length`

    y_batch = create the target by rolling `x_batch` one step forward (shift by one)

    Yield x_batch and y_batch as a tuple of tensors for training

## Step 3: Define the Character-Level LSTM Model

plaintext

Copy code

Class CharLSTM:

    Initialize the LSTM model with vocab_size and hidden_dim:

      - Initialize LSTM layer (input size = vocab_size, hidden_size = hidden_dim)

      - Initialize a fully connected output layer to map hidden states to vocabulary size

    Function forward(x, hidden):

      Pass input `x` through LSTM

      Pass the LSTM output through the fully connected layer

      Return the output and the hidden state

    Function init_hidden(batch_size):

      Initialize hidden and cell states to zero (batch size, hidden_dim)

      Return the hidden state tuple (h0, c0)

## Step 4: Define the Training Function

plaintext

Copy code

Function train(model, encoded_text, epochs, batch_size, seq_length, lr):

    Initialize optimizer (Adam optimizer) with model parameters and learning rate

    Initialize loss function (CrossEntropyLoss)

    For each epoch:

      Initialize hidden state for the batch

      Initialize total_loss to zero

      For each batch of data (x_batch, y_batch) from get_batches function:

        One-hot encode the input x_batch

        Zero the gradients in the optimizer

        Perform forward pass through the model to get output and updated hidden state

        Calculate the loss using CrossEntropyLoss

        Perform backward pass to calculate gradients

        Step the optimizer to update the model parameters

        Detach hidden states to prevent reusing the computational graph

      Print average loss for the epoch

## Step 5: Train the Model

plaintext

Copy code

Set the hidden dimension size (hidden_dim) for the model

Initialize the CharLSTM model with vocab_size and hidden_dim

Call the train function to start training with specified parameters

## Step 6: Text Generation Function

plaintext

Copy code

Function generate_text(model, start_str, length, temperature):

    Set model to evaluation mode

    Initialize hidden state

    Convert each character in the start string `start_str` to an index and perform a forward pass to update hidden state

For each character to generate:
   Apply temperature to the model's output (softmax with temperature scaling)
   Sample the next character based on the output distribution (multinomial sampling)
   Append the predicted character to the output string
   One-hot encode the sampled character and perform another forward pass to update hidden state
Return the generated text as a string

**Text Generation (Calling the function)**

plaintext
Copy code
Set the starting string for text generation
Generate text of specified length using the generate_text function
Print the generated text

**Complete Pseudocode Summary**

plaintext
Copy code
1. Create character mappings:
   - char_to_idx
   - idx_to_char
2. Encode text into integer sequence using char_to_idx
3. Define function get_batches:
   - Split text into batches of sequence length
   - Shift by one to create target batches

4. Define CharLSTM class:
   - Initialize LSTM and fully connected layers
   - Forward pass through LSTM and fully connected layer
   - Initialize hidden state
5. Define train function:
   - Initialize optimizer and loss function
   - Loop over epochs and batches, performing forward and backward passes
6. Define generate_text function:
   - Convert start string into indices and generate text using the trained model
   - Sample characters based on the model's predictions with temperature scaling
7. Train the model and generate text:
   - Train model using training data
   - Generate text from a start string

# Limitation

This code implements a character-level LSTM model for text generation. Below are some limitations and points of improvement:

**1. Data Handling Limitations**

- **Missing Data Loading**: The code uses a variable text, but it is not defined or loaded. The text data needs to be read from a file or defined before processing. A common approach is to read from a file like "Shakespeare_data" or any other text corpus, but this step is missing.
- **Limited Preprocessing**: The text is simply converted into integer indices using a character-to-index mapping, but there is no handling for edge cases such as handling punctuation, case normalization, or text cleaning.

**2. Model Limitations**

- **LSTM with One-Hot Encoding**: The model uses one-hot encoding for the input, which leads to sparse representations. This can be inefficient for large vocabularies. Instead, embedding layers are typically used in such cases to reduce the dimensionality and improve the model's expressiveness.

- **LSTM Layer with Large Vocabulary**: Using a large vocabulary with a one-hot encoded input may lead to very large input vectors (i.e., the size of the input will be the same as the vocabulary size). This makes the training inefficient and might lead to memory problems for very large vocabularies.
- **Single LSTM Layer**: The model uses a single LSTM layer. For complex tasks, a deeper architecture or stacked LSTM layers may improve performance.

## 3. Training Process

- **Batches and Sequence Length**: The batch creation logic seems to assume that the total number of data points is perfectly divisible by the batch size and sequence length. This can lead to data loss or inefficiency. Handling cases where this isn't true (e.g., truncating sequences or padding) would make it more robust.
- **Gradient Detachment**: The hidden state is detached (hidden = tuple([h.detach() for h in hidden])) after each batch to avoid unnecessary memory usage. While this is necessary, it's a potential source of error if used incorrectly.
- **Single Batch Training**: The training is done with a batch size of 1, which might significantly slow down the training process. Larger batch sizes would typically lead to better utilization of the hardware (especially if using GPUs).
- **No Validation**: There's no validation during training, meaning we cannot monitor overfitting. It's essential to split the dataset into training and validation sets and evaluate the model's performance periodically.

## 4. Text Generation Limitations

- **Sampling Method**: The text generation uses sampling via torch.multinomial() with a temperature parameter. While this is standard, the temperature may not always provide the best balance between randomness and coherence. Fine-tuning this parameter is essential for optimal results.
- **Starting Text**: The model relies heavily on the starting text (start_str="hello"), but there's no mechanism to ensure that the starting string is valid or that the generated text will make sense beyond the starting point. Better text conditioning would lead to more coherent results.
- **No Diversity Control**: Although temperature affects randomness, more sophisticated techniques (e.g., top-k or nucleus sampling) could improve text diversity and coherence.

## 5. Model Evaluation and Overfitting

- **No Evaluation Metrics**: There are no explicit metrics for evaluating model performance during or after training (such as perplexity or accuracy). This would be necessary to assess whether the model is improving during training or if it starts overfitting.
- **Limited Epochs and Training**: The model might need more epochs for complex datasets like Shakespeare's works. Depending on the complexity, the network might overfit early or fail to converge with the current setup.

## 6. Miscellaneous

- **Hard-coded Parameters**: Many parameters such as vocab_size, hidden_dim, and seq_length are hardcoded, making it difficult to experiment with different values without changing the code.
- **No Model Saving/Loading**: There is no checkpointing, so once the model is trained, it cannot be saved or reused. Adding model saving/loading functionality would be beneficial.

## Potential Improvements:

- **Use of Embedding Layer**: Replace one-hot encoding with an embedding layer to make the model more efficient and scalable.
- **Use of Data Preprocessing**: Incorporate more preprocessing steps, such as cleaning text, handling punctuation, and normalizing the case of characters.
- **Validation and Overfitting Control**: Implement validation, early stopping, or dropout to mitigate overfitting.
- **Text Generation Enhancements**: Add techniques like top-k sampling or nucleus sampling to improve the quality and diversity of generated text.

In summary, while the code demonstrates a basic character-level LSTM for text generation, it has several limitations in terms of efficiency, training process, and model evaluation. Enhancing it with proper data handling, model architecture, and training methods would significantly improve its performance and versatility.

# CS582 Neural Network and NLP Lab

## Strengths

**Strengths** of the given code for training a character-level LSTM model and generating text:

**1. Modular and Well-Structured Code**
- The code is divided into distinct steps (data preprocessing, model definition, training, and text generation). This makes it easy to understand, debug, and extend in the future.
- Functions like get_batches, train, and generate_text encapsulate specific tasks, which aids in readability and reusability.

**2. Character-Level LSTM for Text Generation**
- **Character-level modeling**: The model processes text at the character level, which allows it to generate text more flexibly, even for new words or sequences that may not appear in the training data.
- **Text generation quality**: Using an LSTM (Long Short-Term Memory) network ensures that the model can capture long-range dependencies in the text, making it effective at generating coherent and contextually relevant sequences.

**3. Use of One-Hot Encoding for Input Representation**
- One-hot encoding is a simple and effective method for representing categorical data (characters in this case). It ensures that the input data is in the correct format for the LSTM network to process.
- This approach is easy to understand and implement, especially for small vocabularies like character-level text.

**4. Training with Batches**
- The use of batches to train the model is efficient for handling larger datasets. It helps improve model generalization and speeds up training by taking advantage of vectorized operations on GPUs.
- By defining get_batches for batch creation, the code handles the splitting of data efficiently, without overloading memory.

**5. Optimized Model Training with Adam Optimizer**
- The **Adam optimizer** is one of the most popular and efficient optimizers for training deep learning models, especially for tasks like text generation. It adapts the learning rate during training and helps achieve faster convergence.
- Using **CrossEntropyLoss** is ideal for classification tasks where the model needs to predict one of several possible classes (characters in this case).

**6. Use of Temperature for Sampling During Text Generation**
- The **temperature** parameter in the generate_text function helps control the diversity of the generated text. By adjusting the temperature, you can balance between randomness and determinism in the text generation process.
  - Lower temperature (<1.0) makes the model more confident, resulting in less varied and more predictable output.
  - Higher temperature (>1.0) increases randomness, making the generated text more diverse and creative.

**7. Hidden State Detaching for Efficient Memory Management**
- During training, the hidden state of the LSTM is **detached** to prevent memory overflow and unnecessary backpropagation through the entire training history. This is important for handling long sequences efficiently.

**8. GPU Compatibility**
- The code includes the potential to leverage a **GPU** for training, making it scalable and faster, especially when dealing with large datasets or complex models.

**9. Simple Text Generation**
- The generate_text function provides an easy-to-use interface for generating text based on an initial string. This makes the model interactive and useful for practical applications such as poetry, story generation, or creative writing.

**10. Flexibility in Hyperparameters**
- The model and training process allow for easy adjustments to key hyperparameters like:
  - hidden_dim (size of the hidden layer)

- o epochs (number of training iterations)
- o batch_size (size of data batches)
- o seq_length (sequence length for training)
- o learning rate
- This flexibility allows for tuning the model to improve performance based on the dataset and desired output.

**11. Simple and Intuitive Implementation**
- The code uses well-established deep learning techniques (LSTM, cross-entropy loss, etc.), which makes it easy to understand for beginners or those familiar with neural networks.
- The overall approach to the problem (character-level sequence prediction) is easy to follow, from text encoding to model training and generation.

The strengths of this code are its modular design, ability to handle text generation tasks effectively using an LSTM, efficient training methods with batch processing, and its ability to generate creative outputs through temperature-based sampling. It's well-suited for working with smaller vocabularies like characters, and the use of GPU makes it scalable for larger datasets.

## Code Snippet

```python
import torch
import torch.nn as nn
import numpy as np


path="Shakespeare_data.csv"


# Assuming text is loaded from the dataset
# For the purpose of this example, let's simulate the text.
text = open(path, 'r').read()  # Load the actual Shakespeare text here.

# Step 1: Create character-to-index and index-to-character mappings
chars = sorted(list(set(text)))  # Unique characters
vocab_size = len(chars)  # Vocabulary size

# Mapping characters to indices and vice versa
char_to_idx = {ch: i for i, ch in enumerate(chars)}
idx_to_char = {i: ch for i, ch in enumerate(chars)}

# Encode the text into integer sequence
encoded_text = np.array([char_to_idx[ch] for ch in text])

# Step 2: Define a function to create batches for training
def get_batches(encoded_text, batch_size, seq_length):
    total_batch_size = batch_size * seq_length
    num_batches = len(encoded_text) // total_batch_size

    encoded_text = encoded_text[:num_batches * total_batch_size]
    x = encoded_text.reshape((batch_size, -1))  # Create a matrix of shape (batch_size,
total_sequence_length)

    for n in range(0, x.shape[1], seq_length):
        x_batch = x[:, n:n+seq_length]
        y_batch = np.roll(x_batch, -1, axis=1)  # Shift by one for the target
        yield torch.tensor(x_batch, dtype=torch.long), torch.tensor(y_batch, dtype=torch.long)
```

16

```python
# Step 3: Define the character-level LSTM model
class CharLSTM(nn.Module):
    def __init__(self, vocab_size, hidden_dim):
        super(CharLSTM, self).__init__()
        self.hidden_dim = hidden_dim
        self.lstm = nn.LSTM(vocab_size, hidden_dim, batch_first=True)  # LSTM layer
        self.fc = nn.Linear(hidden_dim, vocab_size)  # Output layer to predict the next character

    def forward(self, x, hidden):
        out, hidden = self.lstm(x, hidden)
        out = self.fc(out.reshape(-1, self.hidden_dim))  # Reshape to match output dimensions
        return out, hidden

    def init_hidden(self, batch_size):
        # Initialize hidden and cell states
        return (torch.zeros(1, batch_size, self.hidden_dim),
            torch.zeros(1, batch_size, self.hidden_dim))

# Step 4: Define the training function
def train(model, encoded_text, epochs=50, batch_size=1, seq_length=10, lr=0.001):
    optimizer = torch.optim.Adam(model.parameters(), lr=lr)
    criterion = nn.CrossEntropyLoss()  # Loss function for classification

    model.train()
    for epoch in range(epochs):
        hidden = model.init_hidden(batch_size)  # Initialize hidden state
        total_loss = 0

        # Iterate through batches
        for x_batch, y_batch in get_batches(encoded_text, batch_size, seq_length):
            x_onehot = nn.functional.one_hot(x_batch, num_classes=vocab_size).float()  # One-hot encoding

            optimizer.zero_grad()  # Reset gradients
            output, hidden = model(x_onehot, hidden)  # Forward pass

            loss = criterion(output, y_batch.view(-1))  # Calculate loss
            total_loss += loss.item()

            loss.backward()  # Backpropagation
            optimizer.step()  # Update weights

            # Clear the hidden state to prevent reusing the graph
            hidden = tuple([h.detach() for h in hidden])

        print(f"Epoch {epoch+1}/{epochs}, Loss: {total_loss/len(encoded_text)}")

# Step 5: Train the model
hidden_dim = 128  # Number of hidden units
model = CharLSTM(vocab_size, hidden_dim)  # Initialize the model
train(model, encoded_text, epochs=50, batch_size=1, seq_length=10, lr=0.001)  # Train the model

# Step 6: Text generation function with handling missing characters
def generate_text(model, start_str="hello", length=100, temperature=0.7):
```

17

```
    model.eval()  # Set model to evaluation mode
    chars = list(start_str)
    hidden = model.init_hidden(1)  # Initialize hidden state

    for ch in start_str:
        if ch not in char_to_idx:
            print(f"Character '{ch}' not found in vocabulary. Skipping.")
            continue
        x = torch.tensor([[char_to_idx[ch]]]).long()  # Convert character to index
        x_onehot = nn.functional.one_hot(x, num_classes=vocab_size).float()  # One-hot encoding
        output, hidden = model(x_onehot, hidden)  # Forward pass

    # Generate characters one by one
    for _ in range(length):
        output_dist = nn.functional.softmax(output / temperature, dim=1).data  # Apply softmax and
temperature
        top_char = torch.multinomial(output_dist, 1)[0]  # Sample character from the distribution
        chars.append(idx_to_char[top_char.item()])  # Append predicted character to the string

        x_onehot = nn.functional.one_hot(top_char, num_classes=vocab_size).float().unsqueeze(0)  # One-
hot encoding
        output, hidden = model(x_onehot, hidden)  # Forward pass

    return ''.join(chars)  # Return generated text

# Generate some text
generated_text = generate_text(model, start_str="hello", length=200, temperature=0.7)
print("Generated Text:\n", generated_text)
```

# Conclusion

The code provided aims to implement a character-level LSTM model using PyTorch to generate text, trained on a corpus such as Shakespeare's works. The task involves training a recurrent neural network (RNN) that can predict the next character in a sequence, allowing the generation of novel text that mimics the structure and style of the training data. Here's an analysis and conclusion on the problem, approach, and the effectiveness of the code.

**Problem Analysis**

The goal of this project is to generate realistic sequences of text by training a model on a given dataset, such as Shakespeare's text. To achieve this:

1. **Character Encoding**: Each character in the dataset is represented as an integer, enabling the model to learn patterns at the character level.
2. **LSTM Network Architecture**: LSTMs are a type of RNN capable of capturing long-range dependencies, which is crucial for text generation where previous characters can influence later predictions.
3. **Text Generation**: After training, the model can generate text by iteratively predicting the next character based on previous ones, which can produce coherent and stylistic text outputs.

**Code Structure and Implementation**

The code is organized into several key sections:

1. **Data Preparation**:
    - Reads text data and creates mappings for characters to indices and vice versa.
    - Encodes the text into a sequence of integers, which serves as the input to the model.
2. **Batch Generation**:

- o A function get_batches is used to divide the encoded text into batches, which is important for training stability and efficiency.
3. **LSTM Model Definition**:
   - o A character-level LSTM model, CharLSTM, is defined with an LSTM layer and a fully connected layer for output.
   - o The model's forward pass involves generating predictions character-by-character, updating hidden states to capture sequential dependencies.
4. **Training**:
   - o The training function employs cross-entropy loss for character prediction and optimizes the model parameters through backpropagation.
   - o Gradients are detached from the hidden states after each batch to prevent memory accumulation and ensure stability.
5. **Text Generation**:
   - o A generate_text function enables the model to produce new sequences of text by predicting one character at a time. It utilizes temperature sampling to control randomness in predictions, which impacts the creativity and coherence of generated text.

**Code Strengths**
1. **Character-Level Generation**: This approach provides fine control over text generation and can handle any input text without needing pre-defined word-level boundaries.
2. **Modular Functions**: The code is organized into reusable functions for data preparation, batch generation, training, and text generation, which makes it easy to modify or extend.
3. **Sampling Control**: The inclusion of a temperature parameter during text generation provides flexibility to control the randomness, which can yield more varied or conservative outputs.
4. **Efficient Batch Processing**: By batching sequences, the code efficiently utilizes available data, and detaching the hidden state in each batch prevents unnecessary memory accumulation.

**Limitations and Potential Improvements**
1. **Limited Vocabulary Size**: Character-level models may have difficulty generating coherent long-range structures, as they only see character sequences rather than complete words or sentences. Word-level or subword-level embeddings might yield more structured text.
2. **Training Resource Requirements**: Training deep LSTMs on large text corpora can be computationally expensive, and GPU usage is recommended for better performance.
3. **Error Handling in Generation**: If the start string contains characters not in the vocabulary, it can lead to errors. Enhanced preprocessing could ensure all characters are in the vocabulary or provide default handling for unknown characters.
4. **Complexity with Long Sequences**: While the LSTM can learn dependencies across time steps, longer texts might still suffer from limited memory. More advanced architectures (e.g., Transformer-based models) could capture even longer dependencies more effectively.

The code presents a solid foundation for character-level text generation using an LSTM-based model in PyTorch. It demonstrates effective data preparation, model training, and generation processes that collectively allow for generating text in a style similar to the input data. The project achieves the goal of learning text patterns and generating text, although it could benefit from optimizations for handling long sequences and improving coherence in generated text. Overall, this approach illustrates the core mechanics of sequence modeling and serves as a good starting point for text generation tasks with room for further advancements.

```python
import torch
import torch.nn as nn
import numpy as np
```

```python
path="Shakespeare_data.csv"


# Assuming text is loaded from the dataset
# For the purpose of this example, let's simulate the text.
text = open(path, 'r').read() # Load the actual Shakespeare text here.

# Step 1: Create character-to-index and index-to-character mappings
chars = sorted(list(set(text))) # Unique characters
vocab_size = len(chars) # Vocabulary size

# Mapping characters to indices and vice versa
char_to_idx = {ch: i for i, ch in enumerate(chars)}
idx_to_char = {i: ch for i, ch in enumerate(chars)}

# Encode the text into integer sequence
encoded_text = np.array([char_to_idx[ch] for ch in text])

# Step 2: Define a function to create batches for training
def get_batches(encoded_text, batch_size, seq_length):
total_batch_size = batch_size * seq_length
num_batches = len(encoded_text) // total_batch_size

encoded_text = encoded_text[:num_batches * total_batch_size]
x = encoded_text.reshape((batch_size, -1)) # Create a matrix of shape
(batch_size, total_sequence_length)

for n in range(0, x.shape[1], seq_length):
x_batch = x[:, n:n+seq_length]
y_batch = np.roll(x_batch, -1, axis=1) # Shift by one for the target
yield torch.tensor(x_batch, dtype=torch.long), torch.tensor(y_batch,
dtype=torch.long)

# Step 3: Define the character-level LSTM model
class CharLSTM(nn.Module):
def __init__(self, vocab_size, hidden_dim):
super(CharLSTM, self).__init__()
self.hidden_dim = hidden_dim
self.lstm = nn.LSTM(vocab_size, hidden_dim, batch_first=True) # LSTM
layer
self.fc = nn.Linear(hidden_dim, vocab_size) # Output layer to predict
the next character
```

```python
def forward(self, x, hidden):
out, hidden = self.lstm(x, hidden)
out = self.fc(out.reshape(-1, self.hidden_dim)) # Reshape to match
output dimensions
return out, hidden
```

```python
def init_hidden(self, batch_size):
# Initialize hidden and cell states
return (torch.zeros(1, batch_size, self.hidden_dim),
torch.zeros(1, batch_size, self.hidden_dim))
```

```python
# Step 4: Define the training function
def train(model, encoded_text, epochs=50, batch_size=1, seq_length=10,
lr=0.001):
optimizer = torch.optim.Adam(model.parameters(), lr=lr)
criterion = nn.CrossEntropyLoss() # Loss function for classification
model.train()
for epoch in range(epochs):
hidden = model.init_hidden(batch_size) # Initialize hidden state
total_loss = 0
```

```python
# Iterate through batches
for x_batch, y_batch in get_batches(encoded_text, batch_size,
seq_length):
x_onehot = nn.functional.one_hot(x_batch, num_classes=vocab_size).float()
# One-hot encoding
```

```python
optimizer.zero_grad() # Reset gradients
output, hidden = model(x_onehot, hidden) # Forward pass
```

```python
loss = criterion(output, y_batch.view(-1)) # Calculate loss
total_loss += loss.item()
loss.backward() # Backpropagation
optimizer.step() # Update weights
```

```python
# Clear the hidden state to prevent reusing the graph
hidden = tuple([h.detach() for h in hidden])
```

```python
print(f"Epoch {epoch+1}/{epochs}, Loss: {total_loss/len(encoded_text)}")
```

```python
# Step 5: Train the model
hidden_dim = 128 # Number of hidden units
```

```python
model = CharLSTM(vocab_size, hidden_dim) # Initialize the model
train(model, encoded_text, epochs=50, batch_size=1, seq_length=10,
lr=0.001) # Train the model
```

```python
# Step 6: Text generation function with handling missing characters
def generate_text(model, start_str="hello", length=100, temperature=0.7):
model.eval() # Set model to evaluation mode
chars = list(start_str)
hidden = model.init_hidden(1) # Initialize hidden state
```

```python
for ch in start_str:
if ch not in char_to_idx:
print(f"Character '{ch}' not found in vocabulary. Skipping.")
continue
x = torch.tensor([[char_to_idx[ch]]]).long() # Convert character to
index
x_onehot = nn.functional.one_hot(x, num_classes=vocab_size).float() #
One-hot encoding
output, hidden = model(x_onehot, hidden) # Forward pass
```

```python
# Generate characters one by one
for _ in range(length):
output_dist = nn.functional.softmax(output / temperature, dim=1).data #
Apply softmax and temperature
top_char = torch.multinomial(output_dist, 1)[0] # Sample character from
the distribution
chars.append(idx_to_char[top_char.item()]) # Append predicted character
to the string
```

```python
x_onehot = nn.functional.one_hot(top_char,
num_classes=vocab_size).float().unsqueeze(0) # One-hot encoding
output, hidden = model(x_onehot, hidden) # Forward pass
```

```python
return ''.join(chars) # Return generated text
```

```python
# Generate some text
generated_text = generate_text(model, start_str="hello", length=200,
temperature=0.7)
print("Generated Text:\n", generated_text)
```