

Cheat Sheet: Language Modeling with Transformers

Package/ Method	Description	Code Example
Dataset()	This code loads the IMDB data set and initializes iterators for both the training and validation sets. It then creates an iterator data_itr from the training iterator train_iter and retrieves the next data sample using the next() function.	<pre> 1 # Load the data set 2 train_iter, val_iter= IMDB() 3 data_itr=iter(train_iter) 4 next(data_itr) </pre>
Text Pipeline ()	You can utilize PyTorch's torchtext library to streamline the text processing pipeline for NLP tasks. Specifically, get_tokenizer("basic_english") from torchtext.data.utils is used for tokenizing the text data into a list of tokens. This tokenization process is essential for converting raw text into a format that your model can interpret. Furthermore, build_vocab_from_iterator, another utility from torchtext.vocab, is leveraged to construct the vocabulary from the tokenized text. This function iterates through the tokenized data, capturing the unique tokens and associating them with indices, including special symbols like <unk> for unknown tokens, <pad> for padding, and <eos> for end-of-sentence markers. By specifying specials and special_first=True, you ensure these special tokens are prioritized and properly indexed in your vocabulary, setting the groundwork for effective model training.	<pre> 1 # Define special symbols and indices 2 UNK_IDX, PAD_IDX, EOS_IDX = 0, 1, 2 3 # Make sure the tokens are in order of their indices to properly insert them in vocab 4 special_symbols = ['<unk>', '<pad>', '<eos>'] 5 tokenizer = get_tokenizer("basic_english") 6 def yield_tokens(data_iter): 7 for _, data_sample in data_iter: 8 yield tokenizer(data_sample) + ['<eos>'] 9 vocab = build_vocab_from_iterator(yield_tokens(train_iter), 10 specials=special_symbols, special_first=True) 11 vocab.set_default_index(UNK_IDX) 12 text_to_index = lambda text: [vocab[token] for token in 13 tokenizer(text)] + [EOS_IDX] 14 index_to_text = lambda seq_en: " ".join([vocab.get_itos()[index] for 15 index in seq_en if index != EOS_IDX]) </pre>
Creating data for next token prediction()	This code snippet demonstrates a critical step in preparing data for training a language model: generating input-target pairs, where each pair is used for the next token prediction. The function get_sample takes two parameters: block_size, which defines the maximum length of the text sample, and text, the input text from which the sample is generated. To create a diverse training set, torch.randint is used to select a random starting point within the text. This randomness ensures that the model encounters different segments of the text during training, which is vital for learning a robust representation of the language. The selected text segment (src_sequence) and its immediate next token (tgt_sequence) form a pair used to train the model to predict the next token given a sequence of tokens.	<pre> 1 def get_sample(block_size, text): 2 # Determine the length of the input text 3 sample_leg = len(text) 4 # Calculate the stopping point for randomly selecting a sample 5 # This ensures the selected sample doesn't exceed the text length 6 random_sample_stop = sample_leg - block_size 7 # Check if a random sample can be taken (if the text is longer than block_size) 8 if random_sample_stop >= 1: 9 # Randomly select a starting point for the sample 10 random_start = torch.randint(low=0, high=random_sample_stop, 11 size=(1,)).item() 12 # Define the endpoint of the sample 13 stop = random_start + block_size 14 # Create the input and target sequences 15 src_sequence = text[random_start:stop] 16 tgt_sequence= text[random_start + 1:stop + 1] 17 # Handle the case where the text length is exactly equal to or lesser than the block_size 18 else: 19 # Start from the beginning and use the entire text 20 random_start = 0 21 stop = sample_leg 22 src_sequence= text[random_start:stop] 23 tgt_sequence = text[random_start + 1:stop] 24 # Append an empty string to maintain sequence alignment 25 tgt_sequence.append('<endoftext>') 26 return src_sequence, tgt_sequence </pre>
(Src, Tgt) pairs ()	This code snippet generates a sample pair for training a language model, where block_size determines the maximum length of the sample and text represents the input text. It then prints the source sequence (src_sequences) and the corresponding target sequence (tgt_sequence). The function get_sample randomly selects a segment of the text of length block_size, ensuring proper alignment between the source and target sequences. If the text is shorter than block_size, it uses the entire text.	<pre> 1 block_size=10 2 src_sequences, tgt_sequence=get_sample(block_size, text) 3 print(src_sequences) 4 print(tgt_sequence) </pre>
		<pre> 1 BLOCK_SIZE=30 </pre>

Collate function()	<p>This code defines a function <code>collate_batch</code> to prepare batches of input-source and target sequences for training a language model. It iterates over a batch of data, generates source and target sequences using the <code>get_sample</code> function with a specified block size (<code>BLOCK_SIZE</code>), tokenizes the text using a tokenizer, and converts tokens to indices using a vocabulary. The sequences are then converted to PyTorch tensors and appended to the respective source and target batches. Finally, the function pads sequences within the batch to ensure uniform length and returns the processed batches.</p>	<pre> 2 def collate_batch(batch): 3 src_batch, tgt_batch = [], [] 4 DEVICE = torch.device("cuda" if torch.cuda.is_available() else "cpu") 5 for _, _textt in batch: 6 src_sequence, tgt_sequence = get_sample(BLOCK_SIZE, tokenizer(_textt)) 7 src_sequence = vocab(src_sequence) 8 tgt_sequence = vocab(tgt_sequence) 9 src_sequence = torch.tensor(src_sequence, dtype=torch.int64) 10 tgt_sequence = torch.tensor(tgt_sequence, dtype=torch.int64) 11 src_batch.append(src_sequence) 12 tgt_batch.append(tgt_sequence) 13 src_batch = pad_sequence(src_batch, padding_value=PAD_IDX, batch_first=False) 14 tgt_batch = pad_sequence(tgt_batch, padding_value=PAD_IDX, batch_first=False) 15 return src_batch.to(DEVICE), tgt_batch.to(DEVICE) 16 17 BATCH_SIZE=1 18 dataloader = DataLoader(train_iter, batch_size=BATCH_SIZE, shuffle=True, collate_fn=collate 19 val_dataloader= DataLoader(val_iter, batch_size=BATCH_SIZE, shuffle=True, collate_fn=collat </pre>
Masking future tokens: Causal mask()	<p>These functions create masks used in transformer-based models for self-attention:</p> <p><code>generate_square_subsequent_mask(sz, device=DEVICE)</code>: Generates a mask to prevent attending to subsequent positions in a sequence during self-attention.</p> <p><code>create_mask(src, device=DEVICE)</code>: Creates a mask for the source sequence to ensure that each position can only attend to previous positions during self-attention.</p>	<pre> 1 def generate_square_subsequent_mask(sz, device=DEVICE): 2 mask = (torch.triu(torch.ones((sz, sz), device=device)) == 1).transpose(0, 1) 3 mask = mask.float().masked_fill(mask == 0, float('-inf')).masked_fill(mask == 1, float(0.0)) 4 return mask 5 6 def create_mask(src, device=DEVICE): 7 src_seq_len = src.shape[0] 8 src_mask = 9 nn.Transformer.generate_square_subsequent_mask(None, src_seq_len).to(device) 10 return src_mask </pre>
Padding mask()	<p>The code generates a boolean mask, <code>src_padding_mask</code>, indicating the presence of padding tokens (True) in the source sequence <code>src</code>. Each True value corresponds to a padding token, while False represents non-padding tokens. The mask is structured to align with the sequence dimensions for proper masking.</p>	<pre> 1 src_padding_mask = (src == PAD_IDX).transpose(0, 1) 2 src = torch.tensor([[1, 1, 1, 1, 6, 169, 438, 709.. 3 padding_mask = tensor([[True, True, True, True, False, False, False, False, .. </pre>
Custom GPT model architecture()	<p>This forward pass function applies positional embeddings to input embeddings, incorporates source masks if provided, and passes the input through a transformer encoder. Finally, it passes the output through a linear layer (<code>lm_head</code>).</p>	<pre> 1 def forward(self, x, src_mask=None, key_padding_mask=None): 2 # Add positional embeddings to the input embeddings 3 x = self.embed(x) * math.sqrt(self.embed_size) 4 x = self.positional_encoding(x) 5 if src_mask is None: 6 src_mask, src_padding_mask = create_mask(x) 7 output = self.transformer_encoder(x, src_mask, key_padding_mask) 8 x = self.lm_head(x) 9 return x </pre>
Creating a small instance of the model	<p>This code snippet initializes a custom GPT (Generative Pre-trained Transformer) model with specified parameters such as embedding size, number of transformer encoder layers, number of attention heads, vocabulary size, and dropout probability. The model is then moved to the specified device (DEVICE).</p>	<pre> 1 ntokens = len(vocab) # size of vocabulary 2 emsize = 200 # embedding dimension 3 nlayers = 2 # number of ``nn.TransformerEncoderLayer`` in ``nn.TransformerEncoder`` 4 nhead = 2 # number of heads in ``nn.MultiheadAttention`` 5 dropout = 0.2 # dropout probability 6 model = CustomGPTModel(embed_size=emsize, num_heads=nhead, num_layers=nlayers, vocab_size = </pre>
Calculate the loss	<p>During loss calculation, the encoder model generates a source and a target. During prediction, the decoder model generates logits Class 1 and Class 2.</p>	<pre> 1 src= srctgt[0] 2 tgt= srctgt[1] 3 logits = model(src) </pre>
loss_fn	<p>In preparation for loss calculation, you can see the reshaping of logits, where each row corresponds to the prediction for a token, spanning across both the sequence and the batch dimensions. You can reshape the target tensor so that its elements correspond correctly to the logits. This process ensures that every row from the logits aligns with the appropriate target outcomes for accurate loss estimation.</p>	<pre> 1 logits_flat = logits.reshape(-1, logits.shape[-1]) 2 loss = loss_fn(logits_flat, tgt.reshape(-1)) </pre>
		<pre> 1 lr = 5 # learning rate 2 optimizer = torch.optim.SGD(model.parameters(), lr=lr) 3 scheduler = torch.optim.lr_scheduler.StepLR(optimizer, 10000, gamma=0.9) </pre>

Training ()	The training process is similar to other models, such as convolutional neural networks or CNNs, recurrent neural networks or RNNs, transformers, and generative models. It uses the modified loss shape and other functions, such as validation and checkpoint saving, that help in the optimization.	<pre> 4 def train(model: nn.Module, train_data) -> None: 5 model.train() # turn on train mode 6 total_loss = 0. 7 log_interval = 10000 8 start_time = time.time() 9 num_batches = len(list(train_data)) // block_size 10 for batch, srctgt in enumerate(train_data): 11 src= srctgt[0] 12 tgt= srctgt[1] 13 logits = model(src,src_mask=None, key_padding_mask=None) 14 logits_flat = logits.reshape(-1, logits.shape[-1]) 15 loss = loss_fn(logits_flat, tgt.reshape(-1)) 16 optimizer.zero_grad() 17 loss.backward() 18 torch.nn.utils.clip_grad_norm_(model.parameters(), 0.5) 19 optimizer.step() 20 total_loss += loss.item() 21 if (batch % log_interval == 0 and batch > 0) or batch==42060: 22 lr = scheduler.get_last_lr()[0] 23 ms_per_batch = (time.time() - start_time) * 1000 / log_interval 24 #cur_loss = total_loss / log_interval 25 cur_loss = total_loss / batch 26 ppl = math.exp(cur_loss) 27 print(f' epoch {epoch:3d} {batch//block_size:5d}/{num_batches:5d} 28 batches ' 29 f'lr {lr:02.4f} ms/batch {ms_per_batch:5.2f} ' 30 f'loss {cur_loss:5.2f} ppl {ppl:8.2f}') 31 #total_loss = 0 32 start_time = time.time() 33 scheduler.step() </pre>
Training: Validation function	The validation function is important to assess the model on a separate, invisible data set during training to gauge generalization.	<pre> 1 def validate(model, validation_loader, loss_fn): 2 model.eval() 3 total_loss = 0 4 with torch.no_grad(): 5 for src, tgt in validation_loader: 6 src, tgt = src.to(DEVICE), tgt.to(DEVICE) 7 logits = model(src) 8 loss = loss_fn(logits.reshape(-1, logits.shape[-1]), 9 tgt.reshape(-1)) 10 total_loss += loss.item() 11 return total_loss / len(validation_loader) </pre>
Checkpoint saving function	The checkpoint saving function is useful for saving the model's state after certain intervals or under specific conditions, like improved validation performance.	<pre> 1 def save_checkpoint(model, optimizer, 2 filename="my_checkpoint.pth"): 3 checkpoint = { 4 "model_state_dict": model.state_dict(), 5 "optimizer_state_dict": optimizer.state_dict(), 6 } 7 torch.save(checkpoint, filename) </pre>
Training: Evaluate function	The 'evaluate' function measures the performance of the model by computing its average loss in the validation data set. However, the trained model is useful to generate inferences.	<pre> 1 def evaluate(model: nn.Module, eval_data) -> float: 2 model.eval() # turn on evaluation mode 3 total_loss = 0. 4 with torch.no_grad(): 5 for src,tgt in eval_data: 6 tgt = tgt.to(DEVICE) 7 #seq_len = src.size(0) 8 logits = model(src,src_mask=None, key_padding_mask=None) 9 total_loss += loss_fn(logits.reshape(-1, logits.shape[-1]), 10 tgt.reshape(-1)).item() return total_loss / (len(list(eval_data)) -1) </pre>
Prompt()	Preparing an encoding prompt helps create a process for text generation. This process serves as a starting point for the model to generate subsequent tokens. Once this prompt is tokenized, the decoder model can process and generate the next tokens based on the input.	<pre> 1 def encode_prompt(prompt, block_size=BLOCK_SIZE): 2 # Handle None prompt 3 if prompt is None: 4 prompt = '<pad>' * block_size 5 else: 6 tokens = tokenizer(prompt) 7 number_of_tokens = len(tokens) 8 # Adjust prompt length to fit block_size 9 if number_of_tokens > block_size: 10 tokens = tokens[-block_size:] # Keep last block_size tokens 11 elif number_of_tokens < block_size: 12 padding = ['<pad>'] * (block_size - number_of_tokens) 13 tokens = padding + tokens # Prepend padding tokens 14 prompt_indices = vocab(tokens) 15 prompt_encoded = torch.tensor(prompt_indices, dtype=torch.int64).reshape(-1, 1) 16 return prompt_encoded </pre>

Prompt encoding	Tokenized decoded prompt	<pre> 1 prompt_encoded=encode_prompt("This is a prompt to get model 2 generate next words.") 3 prompt_encoded </pre>
Step 1: Generate function	The 'generate' function creates autoregressive text in the decoder model.	<pre> 1 #Autoregressive Language Model text generation 2 def generate(model, prompt=None, max_new_tokens=500, block_size=BLOCK_SIZE, vocab=vocab, to 3 model.to(DEVICE) 4 # Encode the input prompt 5 prompt_encoded = encode_prompt(prompt).to(DEVICE) 6 tokens = [] 7 # Generate new tokens up to max_new_tokens 8 for _ in range(max_new_tokens): 9 # Decode the encoded prompt using the model's decoder 10 logits = model(prompt_encoded) 11 # Bring the sequence length to the first dimension 12 logits = logits.transpose(0, 1) 13 # Select the logits of the last token in the sequence 14 logit_prediction = logits[:, -1] 15 # Choose the most probable next token from the logits(greedy decoding) 16 next_token_encoded = torch.argmax(logit_prediction, 17 dim=-1).reshape(-1, 1) 18 # If the next token is the end-of-sequence (EOS) token, stop generation 19 if next_token_encoded.item() == EOS_IDX: 20 break 21 # Append the next token to the prompt_encoded and keep only the last 'block_size' tokens 22 prompt_encoded = torch.cat((prompt_encoded, next_token_encoded), 23 dim=0)[-block_size:] 24 # Convert the next token index to a token string using the vocabulary 25 token_id = next_token_encoded.to('cpu').item() 26 tokens.append(vocab.get_itos()[token_id]) 27 # Join the generated tokens into a single string and return 28 return ' '.join(tokens) </pre>
Step 2: Generate a token	This function generates text using an autoregressive language model. It iteratively predicts the next token in the sequence based on the previous tokens, using greedy decoding. The generation stops when either the maximum number of new tokens is reached or an end-of-sequence token is predicted. Finally, it returns the generated text.	<pre> 1 #Autoregressive Language Model text generation 2 def generate(model, prompt=None, max_new_tokens=500, block_size=BLOCK_SIZE, vocab=vocab, to 3 model.to(DEVICE) 4 # Encode the input prompt 5 prompt_encoded = encode_prompt(prompt).to(DEVICE) 6 tokens = [] 7 # Generate new tokens up to max_new_tokens 8 for _ in range(max_new_tokens): 9 # Decode the encoded prompt using the model's decoder 10 logits = model.decoder(prompt_encoded,src_mask=None, key_padding_mask=None) 11 # Bring the sequence length to the first dimension 12 logits = logits.transpose(0, 1) 13 # Select the logits of the last token in the sequence 14 logit_prediction = logits[:, -1] 15 # Choose the most probable next token from the logits(greedy decoding) 16 next_token_encoded = torch.argmax(logit_prediction, dim=-1).reshape(-1, 1) 17 # If the next token is the end-of-sequence (EOS) token, stop generation 18 if next_token_encoded.item() == EOS_IDX: 19 break 20 # Append the next token to the prompt_encoded and keep only the last 'block_size' tokens 21 prompt_encoded = torch.cat((prompt_encoded, next_token_encoded), dim=0)[-block_size:] 22 # Convert the next token index to a token string using the vocabulary 23 token_id = next_token_encoded.to('cpu').item() 24 tokens.append(vocab.get_itos()[token_id]) 25 # Join the generated tokens into a single string and return 26 return ' '.join(tokens) </pre>
Step 3: Generate function	This function generates text using an autoregressive language model. It takes a pretrained model, an optional prompt, and parameters for controlling text generation. It iteratively predicts the next token in the sequence based on the previous tokens. The generation stops when either the maximum number of new tokens is reached or an end-of-sequence token is predicted. Finally, it returns the generated text.	<pre> 1 #Autoregressive Language Model text generation 2 def generate(model, prompt=None, max_new_tokens=500, block_size=BLOCK_SIZE, vocab=vocab, to 3 model.to(DEVICE) 4 # Encode the input prompt 5 prompt_encoded = encode_prompt(prompt).to(DEVICE) 6 tokens = [] 7 # Generate new tokens up to max_new_tokens 8 for _ in range(max_new_tokens): 9 # Decode the encoded prompt using the model's decoder 10 logits = model.decoder(prompt_encoded,src_mask=None, key_padding_mask=None) 11 # Bring the sequence length to the first dimension 12 logits = logits.transpose(0, 1) 13 # Select the logits of the last token in the sequence 14 logit_prediction = logits[:, -1] 15 # Choose the most probable next token from the logits(greedy decoding) 16 next_token_encoded = torch.argmax(logit_prediction, dim=-1).reshape(-1, 1) 17 # If the next token is the end-of-sequence (EOS) token, stop generation 18 if next_token_encoded.item() == EOS_IDX: 19 break </pre>

		<pre> 19 # Append the next token to the prompt_encoded and keep only the last 'block_size' tokens 20 prompt_encoded = torch.cat((prompt_encoded, next_token_encoded), dim=0)[-block_size:] 21 # Convert the next token index to a token string using the vocabulary 22 token_id = next_token_encoded.to('cpu').item() 23 tokens.append(vocab.get_itos()[token_id]) 24 # Join the generated tokens into a single string and return 25 return ' '.join(tokens) </pre>
Tokenization and vocabulary building	<p>This code sets up the necessary infrastructure for tokenizing text data and converting it into numerical indices, facilitating subsequent NLP tasks like model training and evaluation.</p>	<pre> 1 # Import the necessary libraries 2 tokenizer = get_tokenizer("basic_english") 3 # Define a function to yield tokenized samples 4 def yield_tokens(data_iter): 5 for label, data_sample in data_iter: 6 yield tokenizer(data_sample) 7 # Define special symbols and their indices 8 PAD_IDX, CLS_IDX, SEP_IDX, MASK_IDX, UNK_IDX = 0, 1, 2, 3, 4 9 special_symbols = ['[PAD]', '[CLS]', '[SEP]', '[MASK]', '[UNK]'] 10 # Split the data into training and testing sets using the IMDB data set 11 train_iter, test_iter = IMDB(split=('train', 'test')) 12 # Build the vocabulary from the training data 13 vocab = build_vocab_from_iterator(yield_tokens(train_iter), specials=special_symbols, speci 14 # Set the default index of the vocabulary to UNK_IDX 15 vocab.set_default_index(UNK_IDX) 16 # Get the size of the vocabulary 17 VOCAB_SIZE = len(vocab) 18 text_to_index=lambda text: [vocab(token) for token in tokenizer(text)] 19 index_to_en = lambda seq_en: " ".join([vocab.get_itos()[index] for index in seq_en]) </pre>
Text masking	<p>This code defines a function called Masking(token), which is responsible for applying masking to a token with a certain probability. If the mask decision is false (with an 80% chance), it returns the token with a '[PAD]' label. If masking is applied (with a 20% chance), it randomly selects between three cases.</p>	<pre> 1 def Masking(token): 2 # Decide whether to mask this token (20% chance) 3 mask = bernoulli_true_false(0.2) 4 # If mask is False, immediately return with '[PAD]' label 5 if not mask: 6 return token, '[PAD]' 7 # If mask is True, proceed with further operations 8 # Randomly decide on an operation (50% chance each) 9 random_opp = bernoulli_true_false(0.5) 10 random_swich = bernoulli_true_false(0.5) 11 # Case 1: If mask, random_opp, and random_swich are True 12 if mask and random_opp and random_swich: 13 # Replace the token with '[MASK]' and set label to a random token 14 mask_label = index_to_en(torch.randint(0, VOCAB_SIZE, (1,))) 15 token_ = '[MASK]' 16 # Case 2: If mask and random_opp are True, but random_swich is False 17 elif mask and random_opp and not random_swich: 18 # Leave the token unchanged and set label to the same token 19 token_ = token 20 mask_label = token 21 # Case 3: If mask is True, but random_opp is False 22 else: 23 # Replace the token with '[MASK]' and set label to the original token 24 token_ = '[MASK]' 25 mask_label = token 26 return token_, mask_label </pre>
MLM preparations	<p>This code defines a function Masking(token), which is responsible for applying masking to a token. It decides whether to mask the token with a 20% chance. If not, it returns the token with a '[PAD]' label. If masking is applied, it randomly chooses between three cases.</p> <p>Case 1: It replaces the token with '[MASK]' and assigns a random token as the label.</p> <p>Case 2: It retains the token unchanged and assigns the same token as the label.</p> <p>Case 3: It replaces the token with '[MASK]' and assigns the original token as the label.</p> <p>The choice between these cases is determined by two independent 50% chances (random_opp and random_swich). Finally, it returns the</p>	<pre> 1 def Masking(token): 2 # Decide whether to mask this token (20% chance) 3 mask = bernoulli_true_false(0.2) 4 # If mask is False, immediately return with '[PAD]' label 5 if not mask: 6 return token, '[PAD]' 7 # If mask is True, proceed with further operations 8 # Randomly decide on an operation (50% chance each) 9 random_opp = bernoulli_true_false(0.5) 10 random_swich = bernoulli_true_false(0.5) 11 # Case 1: If mask, random_opp, and random_swich are True 12 if mask and random_opp and random_swich: 13 # Replace the token with '[MASK]' and set label to a random token 14 mask_label = index_to_en(torch.randint(0, VOCAB_SIZE, (1,))) 15 token_ = '[MASK]' 16 # Case 2: If mask and random_opp are True, but random_swich is False 17 elif mask and random_opp and not random_swich: 18 # Leave the token unchanged and set label to the same token 19 token_ = token 20 mask_label = token 21 # Case 3: If mask is True, but random_opp is False 22 else: 23 # Replace the token with '[MASK]' and set label to the original token 24 token_ = '[MASK]' </pre>

	modified token and its corresponding label.	<div>25</div> <div>26</div> <div></div> <pre> mask_label = token return token_, mask_label </pre>
NSP preparations	<p>This code defines a function process_for_nsp that prepares inputs for training BERT for the next sentence prediction (NSP) task. It takes two inputs: input_sentences, a list of sentences, and input_masked_labels, a list of labels corresponding to masked tokens in the sentences.</p>	<div>1</div> <div>2</div> <div>3</div> <div>4</div> <div>5</div> <div>6</div> <div>7</div> <div>8</div> <div>9</div> <div>10</div> <div>11</div> <div>12</div> <div>13</div> <div>14</div> <div>15</div> <div>16</div> <div>17</div> <div>18</div> <div>19</div> <div>20</div> <div>21</div> <div>22</div> <div>23</div> <div>24</div> <div>25</div> <div>26</div> <div>27</div> <div>28</div> <div>29</div> <div>30</div> <div>31</div> <div>32</div> <div>33</div> <div>34</div> <div></div> <pre> def process_for_nsp(input_sentences, input_masked_labels): # Verify that both input lists are of the same length and have a sufficient number of sente if len(input_sentences) < 2: raise ValueError("Must have two same number of items.") if len(input_sentences) != len(input_masked_labels): raise ValueError("Both lists must have the same number of items.") bert_input = [] bert_label = [] is_next = [] available_indices = list(range(len(input_sentences))) while len(available_indices) >= 2: if random.random() < 0.5: # Choose two consecutive sentences to simulate the 'next sentence' scenario index = random.choice(available_indices[:-1]) # Exclude the last index # append list and add '[CLS]' and '[SEP]' tokens bert_input.append(['[CLS]']+input_sentences[index]+'[SEP]',input_sentences[index + 1]+'[bert_label.append(['[PAD]']+input_masked_labels[index]+'[PAD]', input_masked_labels[inde is_next.append(1) # Label 1 indicates these sentences are consecutive # Remove the used indices available_indices.remove(index) if index + 1 in available_indices: available_indices.remove(index + 1) else: # Choose two random distinct sentences to simulate the 'not next sentence' scenario indices = random.sample(available_indices, 2) bert_input.append(['[CLS]']+input_sentences[indices[0]]+'[SEP]',input_sentences[indices[bert_label.append(['[PAD]']+input_masked_labels[indices[0]]+'[PAD]', input_masked_labels is_next.append(0) # Label 0 indicates these sentences are not consecutive # Remove the used indices available_indices.remove(indices[0]) available_indices.remove(indices[1]) return bert_input, bert_label, is_next </pre>
Creating training-ready inputs for BERT	<p>This code defines a function prepare_bert_final_inputs, which prepares inputs for training a BERT model. It takes bert_inputs, bert_labels, and is_nexts as inputs. The function pads the inputs and labels with [PAD] tokens to ensure they are of equal length, creates segment labels for each pair of sentences, and converts the inputs, labels, and segment labels into tensors if to_tensor is set to True. Finally, it returns the processed inputs, labels, segment labels, and is_nexts.</p>	<div>1</div> <div>2</div> <div>3</div> <div>4</div> <div>5</div> <div>6</div> <div>7</div> <div>8</div> <div>9</div> <div>10</div> <div>11</div> <div>12</div> <div>13</div> <div>14</div> <div>15</div> <div>16</div> <div>17</div> <div>18</div> <div>19</div> <div>20</div> <div>21</div> <div>22</div> <div>23</div> <div>24</div> <div>25</div> <div>26</div> <div>27</div> <div>28</div> <div>29</div> <div>30</div> <div>31</div> <div>32</div> <div>33</div> <div>34</div> <div>35</div> <div></div> <pre> def prepare_bert_final_inputs(bert_inputs, bert_labels, is_nexts, to_tensor=True): def zero_pad_list_pair(pair_, pad='[PAD]'): pair = deepcopy(pair_) max_len = max(len(pair[0]), len(pair[1])) # Append [PAD] to each sentence in the pair until the maximum length is reached pair[0].extend([pad] * (max_len - len(pair[0]))) pair[1].extend([pad] * (max_len - len(pair[1]))) return pair[0], pair[1] # Flatten the tensor flatten = lambda l: [item for sublist in l for item in sublist] # Transform tokens to vocab indices tokens_to_index = lambda tokens: [vocab[token] for token in tokens] bert_inputs_final, bert_labels_final, segment_labels_final, is_nexts_final = [], [], [], [] for bert_input, bert_label, is_next in zip(bert_inputs, bert_labels, is_nexts): # Create segment labels for each pair of sentences segment_label = [[1] * len(bert_input[0]), [2] * len(bert_input[1])] # Zero-pad the bert_input, bert_label, and segment_label bert_input_padded = zero_pad_list_pair(bert_input) bert_label_padded = zero_pad_list_pair(bert_label) segment_label_padded = zero_pad_list_pair(segment_label, pad='') # Convert to tensors if to_tensor: # Flatten the padded inputs and labels, transform tokens to their corresponding vocab indic bert_inputs_final.append(torch.tensor(tokens_to_index(flatten(bert_input_padded))), dtype=to bert_labels_final.append(torch.tensor(tokens_to_index(flatten(bert_label_padded))), dtype=to segment_labels_final.append(torch.tensor(flatten(segment_label_padded), dtype=torch.int64)) is_nexts_final.append(is_next) else: # Flatten the padded inputs and labels bert_inputs_final.append(flatten(bert_input_padded)) bert_labels_final.append(flatten(bert_label_padded)) segment_labels_final.append(flatten(segment_label_padded)) is_nexts_final.append(is_next) return bert_inputs_final, bert_labels_final, segment_labels_final, is_nexts_final </pre>
		<div>1</div> <div>2</div> <div>3</div> <div>4</div> <div>5</div> <div></div> <pre> csv_file_path = 'train_bert_data_new.csv' # Open the CSV file for writing with open(csv_file_path, mode='w', newline='', encoding='utf-8') as file: csv_writer = csv.writer(file) </pre>

Creating training-ready CSV file from IMDB	This code writes the processed Internet Movie Database or IMDB data to a CSV file.	<pre> 6 # Write the header row 7 csv_writer.writerow(['Original Text', 'BERT Input', 'BERT Label', 'Segment Label', 'Is Next']) 8 # Wrap train_iter with tqdm for a progress bar 9 for n, (_, sample) in enumerate(tqdm(train_iter, desc="Processing samples")): 10 # Tokenize the sample input 11 tokens = tokenizer(sample) 12 # Create MLM inputs and labels 13 bert_input, bert_label = prepare_for_mlm(tokens, include_raw_tokens=False) 14 # Skip samples with insufficient input length 15 if len(bert_input) < 2: 16 continue 17 # Create NSP pairs, token labels, and is_next label 18 bert_inputs, bert_labels, is_nexts = process_for_nsp(bert_input, bert_label) 19 # Add zero-paddings, map tokens to vocab indices, and create segment labels 20 bert_inputs, bert_labels, segment_labels, is_nexts = prepare_bert_final_inputs(bert_inputs, 21 # Convert tensors to lists and then convert lists to JSON-formatted strings 22 for bert_input, bert_label, segment_label, is_next in 23 zip(bert_inputs, bert_labels, segment_labels, is_nexts): 24 bert_input_str = json.dumps(bert_input.tolist()) 25 bert_label_str = json.dumps(bert_label.tolist()) 26 segment_label_str = ','.join(map(str, segment_label.tolist())) 27 # Write the data to a CSV file row-by-row 28 csv_writer.writerow([sample, bert_input_str, bert_label_str, segment_label_str, is_next]) </pre>
<code>__init__</code>	Used to initialize objects of a class. It is also called a constructor.	<pre> 1 from pyspark.sql import SparkSession 2 spark = SparkSession.builder.appName("MyApp").getOrCreate() </pre>
<code>__len__</code>	Essentially used to implement the built-in <code>len()</code> function. Whenever you call <code>len()</code> , Python internally invokes the <code>__len__</code> magic method.	<pre> 1 def __len__(self): 2 return len(self.data) </pre>
<code>__getitem__</code>	Used to define the behavior of retrieving items from an object.	<pre> 1 def __getitem__(self, idx): 2 return self.data[idx] </pre>
<code>torch.tensor(...)</code>	Creates a PyTorch tensor from the Python object obtained from the JSON string. It converts the Python object into a PyTorch tensor.	<pre> 1 torch.tensor(json.loads(row['BERT Input'])) </pre>
<code>is_next</code>	A PyTorch tensor created from a value stored in a DataFrame row. Specifically, it's created from the value associated with the key 'Is Next'.	<pre> 1 is_next = torch.tensor(row['Is Next'], dtype=torch.long) </pre>
<code>collate_batch</code>	Responsible for collating individual samples into batches.	<pre> 1 def collate_batch(batch): 2 label_list, text_list, lengths = [], [], [] 3 for _label, _text in batch: 4 label_list.append(label_pipeline(_label)) 5 processed_text = torch.tensor(text_pipeline(_text), dtype=torch.int64) 6 text_list.append(processed_text) 7 lengths.append(processed_text.size(0)) 8 if CONFIG_USE_ROCM: 9 label_list = torch.tensor(label_list, device='cuda') 10 lengths = torch.tensor(lengths, device='cuda') 11 else: 12 label_list = torch.tensor(label_list) 13 lengths = torch.tensor(lengths) 14 padded_text_list = nn.utils.rnn.pad_sequence(text_list, batch_first=True) 15 padded_text_list.to('cuda') 16 #code.interact(local=locals()) 17 return padded_text_list, label_list, lengths </pre>
<code>forward</code>	Defines the forward pass computation, which includes applying the various embedding layers and dropout during training.	<pre> 1 def forward(self, bert_inputs, segment_labels=False): 2 my_embeddings = self.token_embedding(bert_inputs) 3 if self.train: 4 x = self.dropout(my_embeddings + self.positional_encoding(my_embeddings) + self.segment_embe 5 else: 6 x = my_embeddings + self.positional_encoding(my_embeddings) 7 return x </pre>
<code>torch.no_grad()</code>	Context manager provided by PyTorch that turns off gradients during validation or evaluation to save memory and computations.	<pre> 1 with torch.no_grad(): # Turning off gradients for validation saves memory and computations 2 for batch in dataloader: 3 bert_inputs, bert_labels, segment_labels, is_nexts = [b.to(device) for b in batch] </pre>
<code>evaluate</code>	Used for evaluating the BERT model's performance on the test data set. It calculates the average loss over all batches in the test data set and prints the average loss, average next	<pre> 1 def evaluate(dataloader=test_dataloader, model=model, loss_fn=loss_fn, device=device): 2 model.eval() # Turn off dropout and other training-specific behaviors 3 total_loss = 0 4 total_next_sentence_loss = 0 5 total_mask_loss = 0 </pre>

	the average loss, average next sentence loss, and average mask loss.	6 <code>total_batches = 0</code>
Adam	Initializes the Adam optimizer, which is a variant of stochastic gradient descent (SGD). It's commonly used for optimizing neural network models.	1 <code>optimizer = Adam(model.parameters(), lr=1e-4, weight_decay=0.01, betas=(0.9, 0.999))</code>
zero_grad()	Used to zero out the gradients of all parameters of the model. It's typically called before performing the backward pass to avoid accumulating gradients from previous iterations.	1 <code>import torch</code> 2 <code>from torch.autograd import Variable</code> 3 <code>import torch.optim as optim</code> 4 <code>def linear_model(x, W, b):</code> 5 <code> return torch.matmul(x, W) + b</code> 6 <code>data, targets = ...</code> 7 <code>a = Variable(torch.randn(4, 3), requires_grad=True)</code> 8 <code>b = Variable(torch.randn(3), requires_grad=True)</code> 9 <code>optimizer = optim.Adam([a, b])</code> 10 <code>for sample, target in zip(data, targets):</code> 11 <code> optimizer.zero_grad()</code> 12 <code> output = linear_model(sample, W, b)</code> 13 <code> loss = (output - target) ** 2</code> 14 <code> loss.backward()</code> 15 <code> optimizer.step()</code>
backward()	Computes gradients of the loss with respect to the model parameters.	1 <code>import torch</code> 2 <code>from torch.autograd import Variable</code> 3 <code>import torch.optim as optim</code> 4 <code>def linear_model(x, W, b):</code> 5 <code> return torch.matmul(x, W) + b</code> 6 <code>data, targets = ...</code> 7 <code>a = Variable(torch.randn(4, 3), requires_grad=True)</code> 8 <code>b = Variable(torch.randn(3), requires_grad=True)</code> 9 <code>optimizer = optim.Adam([a, b])</code> 10 <code>for sample, target in zip(data, targets):</code> 11 <code> optimizer.zero_grad()</code> 12 <code> output = linear_model(sample, W, b)</code> 13 <code> loss = (output - target) ** 2</code> 14 <code> loss.backward()</code> 15 <code> optimizer.step()</code>
<code>torch.nn.utils.clip_grad_norm_</code>	Used for gradient clipping, which is a technique to prevent the exploding gradient problem during training.	1 <code>torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)</code>
<code>step()</code>	Updates the parameters of the model using the gradients computed during backpropagation.	1 <code>optimizer.step()</code>
<code>torch.save</code>	Used to save the model's state dictionary to a file.	1 <code>torch.save(model.state_dict(), model_save_path)</code>
<code>plt.plot</code>	Used to plot data points on a graph. It takes the x-values, y-values, and optional arguments to customize the plot, such as line style, color, and label.	1 <code>plt.plot(range(1, num_epochs + 1), train_losses, label='Training Loss')</code>
<code>plt.xlabel</code>	Used to set the label for the x-axis of the plot.	1 <code>plt.xlabel('Epoch')</code>
<code>plt.ylabel</code>	Used to set the label for the y-axis of the plot.	1 <code>plt.ylabel('Loss')</code>
<code>plt.title</code>	Used to set the title of the plot. It specifies the text that will be displayed as the title above the plot.	1 <code>lt.title('Training and Evaluation Loss')</code>
<code>plt.legend</code>	Used to add a legend to the plot. It displays labels associated with each plot line.	1 <code>plt.legend()</code>
<code>plt.show</code>	Used to display the plot on the screen or in the output of the script.	1 <code>plt.show()</code>
<code>predict_nsp</code>	A function that takes two sentences, a BERT model, and a tokenizer as input. It tokenizes the input sentences using the tokenizer, then feeds the tokenized inputs to the BERT model to predict whether the second sentence follows the first one (Next Sentence Prediction task). The function returns a string indicating whether the second sentence follows the first one or not based on the model's prediction.	1 <code>sentence1 = "The cat is sitting on the chair."</code> 2 <code>sentence2 = "It is a rainy day"</code> 3 <code>print(predict_nsp(sentence1, sentence2, model, tokenizer))</code>

predict_mlm	Takes an input sentence, a BERT model, and a tokenizer as input. It tokenizes the input sentence using the tokenizer and converts it into token IDs. Then, it creates dummy segment labels filled with zeros and feeds the input tokens and segment labels to the BERT model. The function extracts the position of the [MASK] token and retrieves the predicted index for the [MASK] token from the model's predictions. Finally, it replaces the [MASK] token in the original sentence with the predicted token and returns the predicted sentence.	<pre> 1 def predict_mlm(sentence, model, tokenizer): 2 # Tokenize the input sentence and convert to token IDs, including special tokens 3 inputs = tokenizer(sentence, return_tensors="pt") 4 tokens_tensor = inputs.input_ids </pre>
generate_square_subsequent_mask	Generates a square subsequent mask for self-attention mechanisms in transformer-based models.	<pre> 1 def generate_square_subsequent_mask(sz, device=DEVICE): 2 mask = (torch.triu(torch.ones((sz, sz), device=device)) == 1).transpose(0, 1) 3 mask = mask.float().masked_fill(mask == 0, float('-inf')).masked_fill(mask == 1, float(0.0)) 4 return mask </pre>
create_mask	Creates masks for the source and target sequences, as well as padding masks for both sequences.	<pre> 1 def create_mask(src, tgt, device=DEVICE): 2 src_seq_len = src.shape[0] 3 tgt_seq_len = tgt.shape[0] 4 tgt_mask = generate_square_subsequent_mask(tgt_seq_len) 5 src_mask = torch.zeros((src_seq_len, src_seq_len), device=DEVICE).type(torch.bool) 6 src_padding_mask = (src == PAD_IDX).transpose(0, 1) 7 tgt_padding_mask = (tgt == PAD_IDX).transpose(0, 1) 8 return src_mask, tgt_mask, src_padding_mask, tgt_padding_mask </pre>
encode	Responsible for encoding the input source sequence into a fixed-dimensional representation that captures the contextual information of the input sequence.	<pre> 1 def encode(self, src: Tensor, src_mask: Tensor): 2 src_embedded = self.src_tok_emb(src) 3 src_pos_encoded = self.positional_encoding(src_embedded) 4 return self.transformer.encoder(src_pos_encoded, src_mask) </pre>
decode	Generates the output sequence based on the encoded source sequence and the target sequence.	<pre> 1 def decode(self, tgt: Tensor, memory: Tensor, tgt_mask: Tensor): 2 tgt_embedded = self.tgt_tok_emb(tgt) 3 tgt_pos_encoded = self.positional_encoding(tgt_embedded) 4 return self.transformer.decoder(tgt_pos_encoded, memory, tgt_mask) </pre>
train_epoch	Represents a training epoch in the training loop. It takes the model, optimizer, and training dataloader as input arguments and returns the average loss over the epoch.	<pre> 1 def train_epoch(model, optimizer, train_dataloader): 2 model.train() 3 losses = 0 4 for src, tgt in train_dataloader: 5 src = src.to(DEVICE) 6 tgt = tgt.to(DEVICE) 7 tgt_input = tgt[:-1, :] 8 src_mask, tgt_mask, src_padding_mask, tgt_padding_mask = create_mask(src, tgt_input) 9 src_mask = src_mask.to(DEVICE) 10 tgt_mask = tgt_mask.to(DEVICE) 11 src_padding_mask = src_padding_mask.to(DEVICE) 12 tgt_padding_mask = tgt_padding_mask.to(DEVICE) 13 logits = model(src, tgt_input, src_mask, tgt_mask, src_padding_mask, tgt_padding_mask, src_p 14 logits = logits.to(DEVICE) 15 optimizer.zero_grad() 16 tgt_out = tgt[1:, :] 17 loss = loss_fn(logits.reshape(-1, logits.shape[-1]), tgt_out. 18 Reshape(-1)) 19 loss.backward() 20 optimizer.step() 21 losses += loss.item() 22 return losses / len(list(train_dataloader)) </pre>
greedy_decode	Performs greedy decoding to generate an output sequence using the trained transformer model.	<pre> 1 def greedy_decode(model, src, src_mask, max_len, start_symbol): 2 src = src.to(DEVICE) 3 src_mask = src_mask.to(DEVICE) 4 memory = model.encode(src, src_mask) 5 ys = torch.ones(1, 1).fill_(start_symbol).type(torch.long).to(DEVICE) 6 for i in range(max_len-1): 7 memory = memory.to(DEVICE) 8 tgt_mask = (generate_square_subsequent_mask(ys.size(0))).type(torch.bool).to(DEVICE) 9 out = model.decode(ys, memory, tgt_mask) 10 out = out.transpose(0, 1) 11 prob = model.generator(out[:, -1]) 12 _, next_word = torch.max(prob, dim=1) 13 next_word = next_word.item() 14 ys = torch.cat([ys, torch.ones(1, 1).type_as(src.data).fill_(next_word)], dim=0) 15 if next_word == EOS_IDX: </pre>

		<pre> 16 break 17 return ys </pre>
<pre> translate(model: torch.nn.Module, src_sentence: str) </pre>	<p>Translates a given source sentence into the target language using the provided PyTorch model.</p>	<pre> 1 def translate(model: torch.nn.Module, src_sentence: str): 2 model.eval() 3 src = text_transform[SRC_LANGUAGE](src_sentence).view(-1, 1) 4 num_tokens = src.shape[0] 5 src_mask = (torch.zeros(num_tokens, num_tokens)).type(torch.bool) 6 tgt_tokens = greedy_decode(model, src, src_mask, max_len=num_tokens + 7 5, start_symbol=BOS_IDX).flatten() 8 return " ".join(vocab_transform[TGT_LANGUAGE].lookup_tokens(list(tgt_tokens.cpu().numpy()))) </pre>