

Cheat Sheet: AI Models for NLP

Package/Method	Description	Code example
PyTorch/Embedding and EmbeddingBag	Embedding is a class that represents an embedding layer. It accepts token indices and produces embedding vectors. EmbeddingBag is a class that aggregates embeddings using mean or sum operations. Embedding and EmbeddingBag are part of the torch.nn module. The code example shows how you can use Embedding and EmbeddingBag in PyTorch.	<pre> 1 # Defining a data set 2 dataset = [3 "I like cats", 4 "I hate dogs", 5 "I'm impartial to hippos" 6] 7 #Initializing the tokenizer, iterator from the data set, and vocabulary 8 tokenizer = get_tokenizer('spacy', language='en_core_web_sm') 9 def yield_tokens(data_iter): 10 for data_sample in data_iter: 11 yield tokenizer(data_sample) 12 data_iter = iter(dataset) 13 vocab = build_vocab_from_iterator(yield_tokens(data_iter)) 14 #Tokenizing and generating indices 15 input_ids=lambda x:[torch.tensor(vocab(tokenizer(data_sample))) for data_sample in dataset] 16 index=input_ids(dataset) 17 print(index) 18 #Initiating the embedding layer, specifying the dimension size for the embeddings, 19 #determining the count of unique tokens present in the vocabulary, and creating the embedding layer 20 embedding_dim = 3 21 n_embedding = len(vocab) 22 n_embedding:9 23 embeds = nn.Embedding(n_embedding, embedding_dim) 24 #Applying the embedding object 25 i_like_cats=embeds(index[0]) 26 i_like_cats 27 impartial_to_hippos=embeds(index[-1]) 28 impartial_to_hippos 29 #Initializing the embedding bag layer 30 embedding_dim = 3 31 n_embedding = len(vocab) 32 n_embedding:9 33 embedding_bag = nn.EmbeddingBag(n_embedding, embedding_dim) 34 # Output the embedding bag 35 dataset = ["I like cats","I hate dogs","I'm impartial to hippos"] 36 index:[tensor([0, 7, 2]), tensor([0, 4, 3]), tensor([0, 1, 6, 8, 5])] 37 i_like_cats=embedding_bag(index[0],offsets=torch.tensor([0])) 38 i_like_cats </pre>
Batch function	Defines the number of samples that will be propagated through the network.	<pre> 1 def collate_batch(batch): 2 target_list, context_list, offsets = [], [], [0] 3 for _context, _target in batch: 4 target_list.append(vocab[_target]) 5 processed_context = torch.tensor(text_pipeline(_context), dtype=torch.int64) 6 context_list.append(processed_context) 7 offsets.append(processed_context.size(0)) 8 target_list = torch.tensor(target_list, dtype=torch.int64) 9 offsets = torch.tensor(offsets[:-1]).cumsum(dim=0) 10 context_list = torch.cat(context_list) 11 return target_list.to(device), context_list.to(device), offsets.to(device) 12 BATCH_SIZE = 64 # batch size for training 13 dataloader_cbow = DataLoader(cobw_data, batch_size=BATCH_SIZE, shuffle=True, collate_fn=collate_batch) </pre>
Forward pass	Refers to the computation and storage of intermediate variables (including outputs) for a neural network in order from the input to the output layer.	<pre> 1 def forward(self, text): </pre>
Stanford's pre-trained GloVe	Leverages large-scale data for word embeddings. It can be integrated into PyTorch for improved NLP tasks such as classification.	<pre> 1 from torchtext.vocab import GloVe,vocab 2 # Creating an instance of the 6B version of GloVe() model 3 glove_vectors_6B = GloVe(name='6B') # you can specify the model with the following format: GloVe(name='840B', dim=300) 4 # Build vocab from glove_vectors 5 vocab = vocab(glove_vectors_6B.stoi, 0, specials=('<unk>', '<pad>')) 6 vocab.set_default_index(vocab["<unk>"]) </pre>
	The vocab	<pre> 1 # Takes an iterator as input and extracts the next tokenized sentence. Creates a list of token indices using the vocab dicti </pre>

vocab	object is part of the PyTorch torchtext library. It maps tokens to indices. The code example shows how you can apply the vocab object to tokens directly.	<pre> 2 def get_tokenized_sentence_and_indices(iterator): 3 tokenized_sentence = next(iterator) 4 token_indices = [vocab[token] for token in tokenized_sentence] 5 return tokenized_sentence, token_indices 6 7 # Returns the tokenized sentences and the corresponding token indices. Repeats the process. 8 tokenized_sentence, token_indices = get_tokenized_sentence_and_indices(my_iterator) 9 next(my_iterator) 10 # Prints the tokenized sentence and its corresponding token indices. 11 print("Tokenized Sentence:", tokenized_sentence) 12 print("Token Indices:", token_indices) </pre>
Special tokens in PyTorch: <eos> and <bos>	Tokens introduced to input sequences to convey specific information or serve a particular purpose during training. The code example shows the use of <bos> and <eos> during tokenization. The <bos> token denotes the beginning of the input sequence, and the <eos> token denotes the end.	<pre> 1 # Appends <bos> at the beginning and <eos> at the end of the tokenized sentences 2 # using a loop that iterates over the sentences in the input data 3 tokenizer_en = get_tokenizer('spacy', language='en_core_web_sm') 4 tokens = [] 5 max_length = 0 6 for line in lines: 7 tokenized_line = tokenizer_en(line) 8 tokenized_line = ['<bos>'] + tokenized_line + ['<eos>'] 9 tokens.append(tokenized_line) 10 max_length = max(max_length, len(tokenized_line)) </pre>
Special tokens in PyTorch: <pad>	Tokens introduced to input sequences to convey specific information or serve a particular purpose during training. The code example shows the use of <pad> token to ensure all sentences have the same length.	<pre> 1 # Pads the tokenized lines 2 for i in range(len(tokens)): 3 tokens[i] = tokens[i] + ['<pad>'] * (max_length - len(tokens[i])) </pre>
Cross entropy loss	A metric used in machine learning (ML) to evaluate the performance of a classification model. The loss is measured as the probability value between 0 (perfect model) and 1. Typically, the aim is to bring the model as close to 0 as possible.	<pre> 1 from torch.nn import CrossEntropyLoss 2 model = TextClassificationModel(vocab_size, emsize, num_class) 3 loss_fn = CrossEntropyLoss() 4 predicted_label = model(text, offsets) 5 loss = criterion(predicted_label, label) </pre>
Optimization	Method to reduce losses in a model.	<pre> 1 # Creates an iterator object 2 optimizer = torch.optim.SGD(model.parameters(), lr=0.1) 3 scheduler = torch.optim.lr_scheduler.StepLR(optimizer, 1.0, gamma=0.1) 4 optimizer.zero_grad() 5 predicted_label = model(text, offsets) 6 loss = criterion(predicted_label, label) 7 loss.backward() 8 torch.nn.utils.clip_grad_norm_(model.parameters(), 0.1) 9 optimizer.step() </pre>
	NLTK (or Natural Language Toolkit) provides this function to	<pre> 1 from nltk.translate.bleu_score import sentence_bleu 2 def calculate_bleu_score(generated_translation, reference_translations): </pre>

sentence_bleu()	evaluate a hypothesis sentence against one or more reference sentences. The reference sentences must be presented as a list of sentences where each reference is a list of tokens.	<pre> 3 # Convert the generated translations and reference translations into the expected format for sentence_bleu 4 references = [reference.split() for reference in reference_translations] 5 hypothesis = generated_translation.split() 6 # Calculate the BLEU score 7 bleu_score = sentence_bleu(references, hypothesis) 8 return bleu_score 9 reference_translations = ["Asian man sweeping the walkway .", "An asian man sweeping the walkway .", "An Asian man sweeps the 10 bleu_score = calculate_bleu_score(generated_translation, reference_translations) </pre>
Encoder RNN model	The encoder-decoder seq2seq model works together to transform an input sequence into an output sequence. Encoder is a series of RNNs that process the input sequence individually, passing their hidden states to their next RNN.	<pre> 1 class Encoder(nn.Module): 2 def __init__(self, vocab_len, emb_dim, hid_dim, n_layers, dropout_prob): 3 super().__init__() 4 self.hid_dim = hid_dim 5 self.n_layers = n_layers 6 self.embedding = nn.Embedding(vocab_len, emb_dim) 7 self.lstm = nn.LSTM(emb_dim, hid_dim, n_layers, dropout = dropout_prob) 8 self.dropout = nn.Dropout(dropout_prob) 9 def forward(self, input_batch): 10 embed = self.dropout(self.embedding(input_batch)) 11 embed = embed.to(device) 12 outputs, (hidden, cell) = self.lstm(embed) 13 return hidden, cell </pre>
Decoder RNN model	The encoder-decoder seq2seq model works together to transform an input sequence into an output sequence. The decoder module is a series of RNNs that autoregressively generates the translation as one token at a time. Each generated token goes back into the next RNN along with the hidden state to generate the next token of the output sequence until the end token is generated.	<pre> 1 class Decoder(nn.Module): 2 def __init__(self, output_dim, emb_dim, hid_dim, n_layers, dropout): 3 super().__init__() 4 self.output_dim = output_dim 5 self.hid_dim = hid_dim 6 self.n_layers = n_layers 7 self.embedding = nn.Embedding(output_dim, emb_dim) 8 self.lstm = nn.LSTM(emb_dim, hid_dim, n_layers, dropout = dropout) 9 self.fc_out = nn.Linear(hid_dim, output_dim) 10 self.softmax = nn.LogSoftmax(dim=1) 11 self.dropout = nn.Dropout(dropout) 12 def forward(self, input, hidden, cell): 13 input = input.unsqueeze(0) 14 embedded = self.dropout(self.embedding(input)) 15 output, (hidden, cell) = self.lstm(embedded, (hidden, cell)) 16 prediction_logit = self.fc_out(output.squeeze(0)) 17 prediction = self.softmax(prediction_logit) 18 return prediction, hidden, cell </pre>
Skip-gram model	Predicts surrounding context words from a specific target word. It predicts one context word at a time from a target word.	<pre> 1 class SkipGram_Model(nn.Module): 2 def __init__(self, vocab_size, embed_dim): 3 super(SkipGram_Model, self).__init__() 4 # Define the embeddings layer 5 self.embeddings = nn.Embedding(num_embeddings=vocab_size, embedding_dim=embed_dim) 6 # Define the fully connected layer 7 self.fc = nn.Linear(in_features=embed_dim, out_features=vocab_size) 8 # Perform the forward pass 9 def forward(self, text): 10 # Pass the input text through the embeddings layer 11 out = self.embeddings(text) 12 # Pass the output of the embeddings layer through the fully connected layer 13 # Apply the ReLU activation function 14 out = torch.relu(out) 15 out = self.fc(out) 16 return out 17 model_sg = SkipGram_Model(vocab_size, emsize).to(device) 18 # Sequence generation function 19 CONTEXT_SIZE = 2 20 skip_data = [] 21 for i in range(CONTEXT_SIZE, len(tokenized_toy_data) - CONTEXT_SIZE): 22 context = (23 [tokenized_toy_data[i - j - 1] for j in range(CONTEXT_SIZE)] # Preceding words 24 + [tokenized_toy_data[i + j + 1] for j in range(CONTEXT_SIZE)] # Succeeding words) 25 target = tokenized_toy_data[i] 26 skip_data.append((target, context)) </pre>

		<pre> 27 skip_data=[('i', ['wish', 'i', 'was', 'little']), ('was', ['i', 'wish', 'little', 'bit']),... </pre>
collate_fn	Processes the list of samples to form a batch. The <code>batch</code> argument is a list of all your samples.	<pre> 1 def collate_fn(batch): 2 target_list, context_list = [], [] 3 for _context, _target in batch: 4 target_list.append(vocab[_target]) 5 context_list.append(vocab[_context]) 6 target_list = torch.tensor(target_list, dtype=torch.int64) 7 context_list = torch.tensor(context_list, dtype=torch.int64) 8 return target_list.to(device), context_list.to(device) </pre>
Training function	Trains the model for a specified number of epochs. It also includes a condition to check whether the input is for skip-gram or CBOW. The output of this function includes the trained model and a list of average losses for each epoch.	<pre> 1 def train_model(model, dataloader, criterion, optimizer, num_epochs=1000): 2 # List to store running loss for each epoch 3 epoch_losses = [] 4 for epoch in tqdm(range(num_epochs)): 5 # Storing running loss values for the current epoch 6 running_loss = 0.0 7 # Using tqdm for a progress bar 8 for idx, samples in enumerate(dataloader): 9 optimizer.zero_grad() 10 # Check for EmbeddingBag layer in the model CBOW 11 if any(isinstance(module, nn.EmbeddingBag) for _, module in model.named_modules()): 12 target, context, offsets = samples 13 predicted = model(context, offsets) 14 # Check for Embedding layer in the model skip gram 15 elif any(isinstance(module, nn.Embedding) for _, module in model.named_modules()): 16 target, context = samples 17 predicted = model(context) 18 loss = criterion(predicted, target) 19 loss.backward() 20 torch.nn.utils.clip_grad_norm_(model.parameters(), 0.1) 21 optimizer.step() 22 running_loss += loss.item() 23 # Append average loss for the epoch 24 epoch_losses.append(running_loss / len(dataloader)) 25 return model, epoch_losses </pre>
CBOW model	Utilizes context words to predict a target word and generate its embedding.	<pre> 1 class CBOW(nn.Module): 2 # Initialize the CBOW model 3 def __init__(self, vocab_size, embed_dim, num_class): 4 super(CBOW, self).__init__() 5 # Define the embedding layer using nn.EmbeddingBag 6 self.embedding = nn.EmbeddingBag(vocab_size, embed_dim, sparse=False) 7 # Define the fully connected layer 8 self.fc = nn.Linear(embed_dim, vocab_size) 9 def forward(self, text, offsets): 10 # Pass the input text and offsets through the embedding layer 11 out = self.embedding(text, offsets) 12 # Apply the ReLU activation function to the output of the first linear layer 13 out = torch.relu(out) 14 # Pass the output of the ReLU activation through the fully connected layer 15 return self.fc(out) 16 vocab_size = len(vocab) 17 emsize = 24 18 model_cbow = CBOW(vocab_size, emsize, vocab_size).to(device) </pre>
Training loop	Enumerates data from the DataLoader and, on each pass of the loop, gets a batch of training data from the DataLoader, zeros the optimizer's gradients, and performs an inference (gets predictions from the model for an input batch).	<pre> 1 for epoch in tqdm(range(1, EPOCHS + 1)): 2 model.train() 3 cum_loss=0 4 for idx, (label, text, offsets) in enumerate(train_dataloader): 5 optimizer.zero_grad() 6 predicted_label = model(text, offsets) 7 loss = criterion(predicted_label, label) 8 loss.backward() 9 torch.nn.utils.clip_grad_norm_(model.parameters(), 0.1) 10 optimizer.step() 11 cum_loss+=loss.item() 12 cum_loss_list.append(cum_loss) 13 accu_val = evaluate(valid_dataloader) 14 acc_epoch.append(accu_val) 15 if accu_val > acc_old: 16 acc_old= accu_val 17 torch.save(model.state_dict(), 'my_model.pth') </pre>