	For this homework, we will work on (NLI) [https://nlp.stanford.edu/projects/snli/] The task is, give two sentences: a premise and a hypothesis, to classify the relation between them. We have three classes to describe this relationship. 1. Entailment: the hypothesis follows from the fact that the premise is true 2. Contradiction: the hypothesis contradicts the fact that the premise is true 3. Neutral: There is not relationship between premise and hypothesis See below for examples
[1]:	First let's load the Stanford NLI dataset from the huggingface datasets hub using the dataset package Explore the dataset! # Imports for most of the notebook import torch from transformers import BertModel from transformers import AutoTokenizer from typing import Dict, List from util import get_snli
[2]:	<pre>dataset = load_dataset("snli") print("Split sizes (num_samples, num_labels):\n", dataset.shape) print("\nExample:\n", dataset['train'][0]) Found cached dataset snli (/Users/xingyuchen/.cache/huggingface/datasets/snli/plain_text/1.0.0/1f60b67533b6 275561ff7828aad5ee4282d0e6f844fd148d05d3c6ea251b)</pre>
[3]:	 horse for a competition.', 'label': 1} Each example is a dictionary with the keys: (premise, hypothesis, label). Data Fields premise: a string used to determine the truthfulness of the hypothesis hypothesis: a string that may be true, false, or whose truth conditions may not be knowable when compared to the premise label: an integer whose value may be either 0, indicating that the hypothesis entails the premise, 1, indicating that the premise and hypothesis neither entail nor contradict each other, or 2, indicating that the hypothesis contradicts the premise. from util import get_snli train_dataset, validation_dataset, test_dataset = get_snli()
[4]:	Found cached dataset snli (/Users/xingyuchen/.cache/huggingface/datasets/snli/plain_text/1.0.0/1f60b67533b6275561ff7828aad5ee4282d0e6f844fd148d05d3c6ea251b) 0%
	<pre>print(Counter([t['label'] for t in validation_dataset])) print(Counter([t['label'] for t in test_dataset])) # We have a perfectly balanced dataset 9999 999 999 Counter({0: 3333, 1: 3333, 2: 3333}) Counter({0: 333, 1: 333, 2: 333}) Counter({0: 333, 1: 333, 2: 333}) We want a function to load samples from the huggingface dataset so that they can be batched and encoded for our model.</pre>
[5]:	Now let's reimplement our tokenizer using the huggingface tokenizer. Notice that our call method (the one called when we call an instance of our class) takes both premise batch and a hypothesis batch. The HuggingFace BERT tokenizer knows to join these with the special sentence seperator token between them. We let HuggingFace do most of the work here for making batches of tokenized and encoded sentences. # Nothing to do for this class! class BatchTokenizer:
	<pre>definit(self): """Initializes the tokenizer Args:</pre>
[5]: [6]:	<pre>"[CLS] this is also a premise [SEP] this is a second hypothesis [SEP]'] We can batch the train, validation, and test data, and then run it through the tokenizer def generate_pairwise_input(dataset: List[Dict]) -> (List[str], List[str], List[int]): """ TODO: group all premises and corresponding hypotheses and labels of the datapoints</pre>
[7]:	<pre>a datapoint as seen earlier is a dict of premis, hypothesis and label """ premises = [] hypothesis = [] label = [] for row in dataset: x, y, z = row.values() premises.append(x) hypothesis.append(y) label.append(z) return premises, hypothesis, label train_premises, train_hypotheses, train_labels = generate_pairwise_input(train_dataset) validation_premises, validation_hypotheses, validation_labels = generate_pairwise_input(validation_dataset) test premises, test hypotheses, test labels = generate_pairwise input(test dataset)</pre>
[8]:	<pre>def chunk(lst, n): """Yield successive n-sized chunks from lst.""" for i in range(0, len(lst), n): yield lst[i:i + n] def chunk_multi(lst1, lst2, n): for i in range(0, len(lst1), n): yield lst1[i: i + n], lst2[i: i + n] # Notice that since we use huggingface, we tokenize and # encode in all at once! tokenizer = BatchTokenizer() batch_size = 4</pre>
[9]:	<pre>train_input_batches = [b for b in chunk_multi(train_premises, train_hypotheses, batch_size)] # Tokenize + encode train_input_batches = [tokenizer(*batch) for batch in train_input_batches] Let's batch the labels, ensuring we get them in the same order as the inputs def encode_labels(labels: List[int]) -> torch.FloatTensor: """Turns the batch of labels into a tensor Args: labels (List[int]): List of all labels in the batch Returns:</pre>
10]:	<pre>torch.FloatTensor: Tensor of all labels in the batch """ return torch.LongTensor([int(1) for 1 in labels]) train_label_batches = [b for b in chunk(train_labels, batch_size)] train_label_batches = [encode_labels(batch) for batch in train_label_batches] Now we implement the model. Notice the TODO and the optional TODO (read why you may want to do this one.) class NLIClassifier(torch.nn.Module):</pre>
	<pre>definit(self, output_size: int, hidden_size: int): super()init() self.output_size = output_size self.hidden_size = hidden_size # Initialize BERT, which we use instead of a single embedding layer. self.bert = BertModel.from_pretrained("prajjwall/bert-small") # TODO (OPTIONAL]: Updating all BERT parameters can be slow and memory intensive. # Freeze them if training is too slow. Notice that the learning # rate should probably be smaller in this case. # Uncommenting out the below 2 lines means only our classification layer will be updated. for param in self.bert.parameters(): param_requires_grad = False self.bert hidden_dimension = self.bert.config.hidden_size # TODO: Add an extra hidden layer in the classifier, projecting # from the BERT hidden dimension to hidden size. # TODO: Add a relu nonlinearity to be used in the forward method # https://pytorch.org/docs/stable/generated/torch.nn.ReLU.html self.relu = torch.nn.ReLU() self.classifier = torch.nn.Linear(self.hidden_size, self.output_size) self.log_softmax = torch.nn.LogSoftmax(dim=2) def encode_text(self, symbols: Dict) -> torch.Tensor: """Encode the (batch of) sequence(s) of token symbols with an LSTM. Then, get the last (non-padded) hidden state for each symbol and return that. Args: symbols (Dict): The Dict of token specifications provided by the HuggingFace tokenizer</pre>
	<pre>Returns: torch.Tensor: The final hiddens tate of the LSTM, which represents an encoding of</pre>
	<pre>def forward(self, symbols: Dict,) -> torch.Tensor: """_summary_ Args: symbols (Dict): The Dict of token specifications provided by the HuggingFace tokenizer Returns: torch.Tensor: _description_ """ encoded_sents = self.encode_text(symbols)</pre>
	<pre>import numpy as np def precision(predicted_labels, true_labels, which_label=1): """ Precision is True Positives / All Positives Predictions """ pred_which = np.array([pred == which_label for pred in predicted_labels]) true_which = np.array([lab == which_label for lab in true_labels]) denominator = t_sum(pred_which) if denominator: return t_sum(logical_and(pred_which, true_which))/denominator else: return 0.</pre>
	<pre>def recall(predicted_labels, true_labels, which_label=1): """ Recall is True Positives / All Positive Labels """ pred_which = np.array((pred == which_label for pred in predicted_labels)) true_which = np.array((lab == which_label for lab in true_labels)) denominator = t_sum(true_which) if denominator: return t_sum(logical_and(pred_which, true_which))/denominator else: return 0. def f1_score(predicted_labels: List[int], true_labels: List[int], which_label: int): """ F1 score is the harmonic mean of precision and recall """ P = precision(predicted_labels, true_labels, which_label=which_label) R = recall(predicted_labels, true_labels, which_label=which_label) if P and R: return 2*P*R/(P+R) else: return 0. def macro_f1(predicted_labels: List[int], true_labels: List[int], possible_labels: List[int] possible_labels: List[int]</pre>
13]:	<pre>scores = [f1_score(predicted_labels, true_labels, l) for l in possible_labels] # Macro, so we take the uniform avg. return sum(scores) / len(scores) Training loop. from tqdm import tqdm_notebook as tqdm import random def training_loop(num_epochs, train_features, train_labels,</pre>
	<pre>train_labels, dev_sents, dev_labels, optimizer, model,): print("Training") loss_func = torch.nn.NLLLoss() batches = list(zip(train_features, train_labels)) random.shuffle(batches) for i in range(num_epochs): losses = [] for features, labels in tqdm(batches): # Empty the dynamic computation graph optimizer.zero_grad()</pre>
	<pre>preds = model(features).squeeze(1) loss = loss_func(preds, labels) # Backpropogate the loss through our model loss.backward() optimizer.step() losses.append(loss.item()) print(f"epoch {i}, loss: {sum(losses)/len(losses)}") # Estimate the f1 score for the development set print("Evaluating dev") all_preds = [] all_labels = [] for sents, labels in tqdm(zip(dev_sents, dev_labels), total=len(dev_sents)): pred = predict(model, sents) all_preds.extend(pred) all_labels_extend(list(labels_numpy()))</pre>
20]:	all_labels.extend(list(labels.numpy())) dev_f1 = macro_f1(all_preds, all_labels, list(set(all_labels))) print(f"Dev F1 (dev_f1)") # Return the trained model return model # You can increase epochs if need be epochs = 10 # TODO: Find a good learning rate LR = 0.0005 possible_labels = len(set(train_labels)) model = NLIClassifier(output_size=possible_labels, hidden_size = 64) optimizer = torch.optim.AdamW(model.parameters(), LR) validation_input_batches = [b for b in chunk_multi(validation_premises, validation_hypotheses, batch_size)] # Tokenize + encode validation_input_batches = [tokenizer(*batch) for batch in validation_input_batches] validation_batch_labels = [b for b in chunk(validation_labels, batch_size)] validation_batch_labels = [encode_labels(batch) for batch in validation_batch_labels] Some weights of the model checkpoint at prajjwall/bert-small were not used when initializing BertModel: ['credictions.transform.LayerNorm.ight', 'cls.seq relationship.bias', 'cls.predictions.decoder.weight', 'cls.seq relationship.weight', 'cls.predictions.
21]:	<pre>ions.transform.dense.weight', 'cls.predictions.bias', 'cls.predictions.transform.dense.bias'] - This IS expected if you are initializing BertModel from the checkpoint of a model trained on another task with another architecture (e.g. initializing a BertForSequenceClassification model from a BertForPreTraining del) This IS NOT expected if you are initializing BertModel from the checkpoint of a model that you expect to xactly identical (initializing a BertForSequenceClassification model from the checkpoint of a model that you expect to xactly identical (initializing a BertForSequenceClassification model from a BertForSequenceClassification model from a BertForSequenceClassification model from a BertForSequenceClassification model from the checkpoint of a model that you expect to xactly identical (initializing a BertForSequenceClassification model from the checkpoint of a model that you expect to xactly identical (initializing a BertForSequenceClassification model from the checkpoint of a model that you expect to xactly identical (initializing a BertForSequenceClassification model from the checkpoint of a model that you expect to xactly identical (initializing a BertForSequenceClassification model from the checkpoint of a model that you expect to xactly identical (initializing a BertForSequenceClassification model from the checkpoint of a model that you expect to xactly identical (initializing a BertForSequenceClass</pre>
	<pre>optimizer, model,) Training /var/folders/31/yzh9j02x7bxd463c11x0_21h0000gn/T/ipykernel_21469/408791926.py:18: TqdmDeprecationWarning: Tfunction will be removed in tqdm==5.0.0 Please use `tqdm.notebook.tqdm` instead of `tqdm.tqdm_notebook` for features, labels in tqdm(batches): 0% </pre>
	Please use `tqdm.notebook.tqdm` instead of `tqdm.tqdm_notebook`
	0% 0/2500 [00:00 , ?it/s] epoch 6, loss: 0.9836426736235618 Evaluating dev 0% 0/250 [00:00<?, ?it/s] Dev F1 0.48121076978999605 0% 0/2500 [00:00<?, ?it/s] epoch 7, loss: 0.9784472920656204 Evaluating dev 0% 0/250 [00:00<?, ?it/s] Dev F1 0.48489830961978936 0% 0/2500 [00:00<?, ?it/s] epoch 8, loss: 0.9726479669570923 Evaluating dev 0% 0/250 [00:00<?, ?it/s] Dev F1 0.48397059252223357</td
21]:	<pre>0% </pre>
	<pre>(encoder): BertEncoder((layer): ModuleList((0): BertLayer(</pre>
	<pre>(output): BertSelfOutput((dense): Linear(in_features=512, out_features=512, bias=True) (LayerNorm): LayerNorm((512,), eps=1e-12, elementwise_affine=True) (dropout): Dropout(p=0.1, inplace=False) } (intermediate): BertIntermediate((dense): Linear(in_features=512, out_features=2048, bias=True) (intermediate_act_fn): GELUActivation() } (output): BertOutput((dense): Linear(in_features=2048, out_features=512, bias=True) (LayerNorm): LayerNorm((512,), eps=1e-12, elementwise_affine=True) (dropout): Dropout(p=0.1, inplace=False) } (2): BertLayer((attention): BertAttention(</pre>
	<pre>(dense): Linear(in_reatures=512, out_reatures=2248, bias=True) (intermediate_act_fn): GELUActivation() } (output): BertOutput((dense): Linear(in_features=2048, out_features=512, bias=True) (LayerNorm): LayerNorm((512,), eps=le-12, elementwise_affine=True) (dropout): Dropout(p=0.1, inplace=False) }) (3): BertLayer((attention): BertAttention((self): BertSelfAttention((query): Linear(in_features=512, out_features=512, bias=True) (key): Linear(in_features=512, out_features=512, bias=True) (dropout): Dropout(p=0.1, inplace=False)) (output): BertSelfOutput((dense): Linear(in_features=512, out_features=512, bias=True) (dapenNorm): LayerNorm((512,), eps=le-12, elementwise_affine=True) (dropout): Dropout(p=0.1, inplace=False) } (intermediate): BertIntermediate((dense): Linear(in_features=512, out_features=2048, bias=True) (intermediate_act_fn): GELUActivation()) (output): BertOutput((dense): Linear(in_features=2048, out_features=512, bias=True) (LayerNorm): LayerNorm((512,), eps=le-12, elementwise_affine=True) (dapenNorm): LayerNorm((512,), eps=le-12, elementwise_affine=True) (dropout): Dropout(p=0.1, inplace=False) } </pre>
22]:	<pre> } } } } (pooler): BertPooler((dense): Linear(in_features=512, out_features=512, bias=True) (activation): Tanh() } (hidden_layer): Linear(in_features=512, out_features=64, bias=True) (relu): ReLU() (classifier): Linear(in_features=64, out_features=3, bias=True) (log_softmax): LogSoftmax(dim=2)) # TODO: Get a final macro F1 on the test set. # You should be able to mimic what we did with the validation set. # test_premises, test_hypotheses, test_labels test_input_batches = [b for b in chunk_multi(test_premises, test_hypotheses, batch_size)] # Tokenize + encode test_input_batches = [tokenizer(*batch) for batch in test_input_batches] test_batch_labels = [b for b in chunk(test_labels, batch_size)] test_batch_labels = [encode_labels(batch) for batch in test_batch_labels] all_preds = [] for sents, labels in tqdm(zip(test_input_batches, test_batch_labels), total=len(test_input_batches)): pred = predict(model, sents) all preds.extend(pred) all preds.extend(pred) all preds.extend(pred) all preds.extend(pred) all preds.extend(pred) all preds.extend(pred)</pre>
	<pre>all_preds.extend(pred) all_labels.extend(list(labels.numpy())) dev_f1 = macro_f1(all_preds, all_labels, list(set(all_labels))) print(f"Dev F1 {dev_f1}") /var/folders/31/yzh9j02x7bxd463cl1x0_2lh0000gn/T/ipykernel_21469/2551025514.py:16: TqdmDeprecationWarning: function will be removed in tqdm==5.0.0 Please use `tqdm.notebook.tqdm` instead of `tqdm.tqdm_notebook` for sents, labels in tqdm(zip(test_input_batches, test_batch_labels), total=len(test_input_batches)): 0% </pre>
	 Describe the task and what capability is required to solve it. Given two sentences: a premise and a hypothesis, classify the relationship between them. Three relationship: Entailment, contradiction, neutral. The dataset has these three values for training. We use HuggingFace BERT tokenizer to tokenize and encounted the sentences in batch. For NLI Classifier, we add layers to the model such as ReLU, Linear, LogSoftmax. For encode texrt we use pooler_output from the bertmodel output, which is the last layer hidden-state of the first token of the sequence (classification token) after further processing through the layers used for the auxiliary pretraining task. Then we just hyperparameter tuning for model to get relatively good f1-score. How does the method of encoding sequences of words in our model differ here compared.
	 How does the method of encoding sequences of words in our model differ here, compared to the word embeddings in HW 4. What is different? Why benefit does this method have? The differences are: 1. we dont need LSTM to encode context and packed the sequence from the LSTM output. The BERT model already encode context and gives us a single vector describing the sequence in the form of the token. 2. We need to preprocess the word embedding in HW4 to get a tensor that indexes the last non-PAD position in each tensor in the batch. However, in BERT we just reshape the dimension to the required output shape. The benefits are we dont need to preprocess the embeddings on our own method and no need to change the dimension of output from model just reshape them. Make the code clean and process quicker. Discuss your results. Did you do any hyperparameter tuning? Did the model solve the task?
	3. Discuss your results. Did you do any hyperparameter tuning? Did the model solve the task? I did hyperparameter tuning for hidden_size and learning rate, the result is showed below, the model solve the task: epoch = 10 batch size = 8 hidden size = 128 loss dev test Ir = 0.001 0.9846 0.5019 0.4429 Ir = 0.002 1.0997 0.1666 0.1666 Ir = 0.005 1.0469 0.4377 0.4612
	lr = 0.005 1.0469 0.4377 0.4612 $ lr = 0.01 1.1041 0.1666 0.1666 0.1666 0.3910 $ $ lr = 0.0005 0.9856 0.4998 0.4929 $ $ epoch = 10 learning rate = 0.0005 learning rate = 128 $ $ loss dev test $
	batch_size = 8
	batch_size = 4 loss dev test