

Department of Artificial Intelligence and Machine Learning

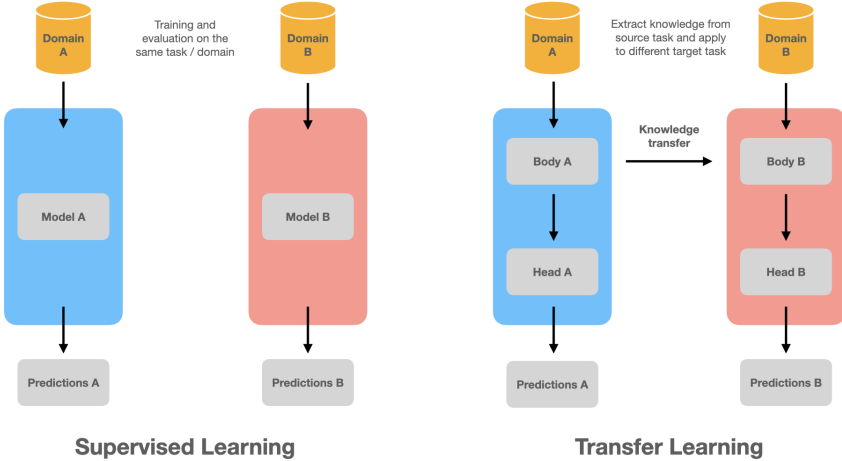
Natural Language Processing & Transformers (21AI53)

CIE 2

Scheme and Solution

SL. No	Questions	M
1	<p>a)</p> <ul style="list-style-type: none"> • Install NLTK • Import NLTK • Download Brown corpus --- nltk.download('brown') • Load the Categorized Corpus---from nltk.corpus import brown • Accessing Categories--- categories=brown.categoriesdocuments_in_news brown.fileids(categories='news') • Accessing Texts---text = brown.words(fileids=documents_in_news[0]) • Tokenization and Preprocessing---words = word_tokenize(text) 	5
	<p>b)</p> <p>Sequential backoff tagger is mainly used to increase the accuracy of the pos tagging process</p> <p>Backoff tagging is one of the core features of SequentialBackoffTagger. It allows you to chain taggers together so that if one tagger doesn't know how to tag a word, it can pass the word on to the next backoff tagger. If that one can't do it, it can pass the word on to the next backoff tagger, and so on until there are no backoff taggers left to check.</p> <pre> Defbackoff_tagger(train_sents,tagger_classes,backoff=None): for cls in tagger_classes: backoff = cls(train_sents, backoff=backoff) return backoff >>> from tag_util import backoff_tagger >>> backoff = DefaultTagger('NN') >>> tagger = backoff_tagger(train_sents, [UnigramTagger, BigramTagger, TrigramTagger], backoff=backoff) >>> tagger.evaluate(test_sents) 0.8806820634578028 </pre>	5
2	<p>a)</p> <ul style="list-style-type: none"> • Pos ---process of converting words into list of tuples (i.e., word,tag) • Tag signifies whether the word is noun, adjective, verb et cetera. • It is needed to extract phrases from sentence using tag patterns, grammar analysis and word sense disambiguation. • Most taggers are trainable and use list of tagged sentence as their training data <p>Any 2 approaches---default tagger and unigram tagger with example code</p>	5
	<p>b)</p> <pre> from nltk.tag import SequentialBackoffTagger from nltk.corpus import names class NamesTagger(SequentialBackoffTagger): def __init__(self, *args, **kwargs): SequentialBackoffTagger.__init__(self, *args, **kwargs) self.name_set = set([n.lower() for n in names.words()]) def choose_tag(self, tokens, index, history): word = tokens[index] </pre>	5

Department of Artificial Intelligence and Machine Learning

		<pre> if word.lower() in self.name_set: return 'NNP' else: return None >>> from taggers import NamesTagger >>> nt = NamesTagger() >>> nt.tag(['Jacob']) [('Jacob', 'NNP')] </pre>	
3	a)	 <p>Supervised Learning</p> <p>Transfer Learning</p> <p>It is common practice in NLP use <i>transfer learning</i> to train a convolutional neural network on one task and then adapt or <i>fine-tune</i> it on a new task, thus making use of the knowledge learned in the original task. Architecturally, this usually works by splitting the model in terms of a <i>body</i> and <i>head</i>, where the head is a task-specific network. During pretraining, the weights of the body learn broad features of the source domain, and it is these weights which are used to initialize the new model for the new task. Compared to traditional supervised learning, this approach typically produces high-quality models that can be trained much more efficiently on a variety of downstream tasks, and with much less labelled data.</p>	5
	b)	<p>Subword Tokenization is a Natural Language Processing technique(NLP)_in which a word is split into subwords and these subwords are known as tokens. This technique is used in any NLP task where a model needs to maintain a large vocabulary and complex word structures. The concept behind this, frequently occurring words should be in the vocabulary whereas rare words are split into frequent subwords. For example, the word “unwanted” might be split into “un”, “want”, and “ed”. The word “football” might be split into “foot”, and “ball”.</p> <pre> from transformers import AutoTokenizer model_ckpt = "distilbert-base-uncased" tokenizer = AutoTokenizer.from_pretrained(model_ckpt) text="I love machine learning!. Tokenization is awesome" encoded_text=tokenizer(text) print(encoded_text) tokens=tokenizer.convert_ids_to_tokens(encoded_text.input_ids) print(tokens) ['[CLS]', 'i', 'love', 'machine', 'learning', '!', '!', 'token', '##ization', 'is', 'awesome', '[SEP]'] </pre>	5

Department of Artificial Intelligence and Machine Learning

4	a)	<p><i>Self-attention is a sequence-to-sequence operation: a sequence of vectors goes in, and a sequence of vectors comes out. Let's call the input vectors x_1, x_2, \dots, x_t and the corresponding output vectors y_1, y_2, \dots, y_t. The vectors all have dimension k. To produce output vector y_i, the self attention operation simply takes a weighted average over all the input vectors, the simplest option is the dot product.</i></p> $\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$	5
4	b)	<p>Language barrier---zero shot training Data Biases Blackbox models Long documents</p>	5
5	a)	<p>Let the input text is : text = ""Dear Amazon, last week I ordered an Optimus Prime action figure \from your online store in Germany. Unfortunately, when I opened the package, \I discovered to my horror that I had been sent an action figure of Megatron \ instead! As a lifelong enemy of the Decepticons, I hope you can understand my \ dilemma. To resolve the issue, I demand an exchange of Megatron for the \ Optimus Prime figure I ordered. Enclosed are copies of my records concerning \ this purchase. I expect to hear from you soon. Sincerely, Bumblebee.""</p> <pre> from transformers import set_seed set_seed(42) generator = pipeline("text-generation") response = "Dear Bumblebee, I am sorry to hear that your order was mixed up." prompt = text + "\n\nCustomer service response:\n" + response outputs = generator(prompt, max_length=200) print(outputs[0]['generated_text']) </pre> <p>Setting `pad_token_id` to `eos_token_id`:50256 for open-end generation. Dear Amazon, last week I ordered an Optimus Prime action figure from your online > store in Germany. Unfortunately, when I opened the package, I discovered to > my horror that I had been sent an action figure of Megatron instead! As a > lifelong enemy of the Decepticons, I hope you can understand my dilemma. To > resolve the issue, I demand an exchange of Megatron for the Optimus Prime > figure I ordered. Enclosed are copies of my records concerning this purchase. > I expect to hear from you soon. Sincerely, Bumblebee.</p> <p>Customer service response: Dear Bumblebee, I am sorry to hear that your order was mixed up. The order was > completely mislabeled, which is very common in our online store, but I can > appreciate it because it was my understanding from this site and our customer > service of the previous day that your order was not made correct in our mind > and that we are in a process of resolving this matter. We can assure you that > your order</p>	7
5	b)	<p>Transformer models have a maximum input sequence length that is referred to as the maximum context size. Let dataset be df and let the feature name be "text". We can find the maximum</p>	3

Department of Artificial Intelligence and Machine Learning

possible length of words in sequence and plot a graph using the following code:

```
df["Words Per Tweet"] = df["text"].str.split().apply(len)
df.boxplot("Words Per Tweet", by='label_name', grid=False, showfliers=False,
color='black', )
plt.suptitle("")
plt.xlabel("");
```

