

Scalable Data Engineering, Exercise 9

HENRY HAUSTEIN

Task 1

(a) Python

```
1 print(torch.sum(x, 0))
2 print(torch.sum(x, 1))
3 print(torch.sum(x, 2))
```

(b) Python

```
1 ones = torch.ones([1,2,4,5])
2 random = torch.randn(10,4)
3 random + ones.reshape([10, 4])
```

Task 2

Right answers are:

- Matrix Multiplication Backwards
- Backward Step of the linear layer

Task 3

Right answer is [0.27, 0.73]

Task 4

(a) Python

```
1 class NGramLanguageModeler(nn.Module):
2
3     def __init__(self, vocab_size, embedding_dim, context_size,
4                   hidden_size):
5         super(NGramLanguageModeler, self).__init__()
6         self.embeddings = nn.Embedding(vocab_size, embedding_dim)
7         self.linear1 = nn.Linear(context_size * embedding_dim,
8                                   hidden_size)
9         self.linear2 = nn.Linear(hidden_size, vocab_size)
```

```

9     def forward(self, inputs):
10         embeds = self.embeddings(inputs).view((1, -1))
11         out = F.relu(self.linear1(embeds))
12         out = self.linear2(out)
13         log_probs = F.log_softmax(out, dim=1)
14         return log_probs

```

(b) Python, only the loop

```

1  for epoch in tqdm.tqdm(range(EPOCHS), total=EPOCHS):
2      total_loss = 0
3      for context, target in trigrams:
4          # Step 1. Prepare the inputs to be passed to the model (i.
              e, turn the words into integer indices and wrap them
              in tensors)
5          context_idxs = torch.tensor([word_to_ix[w] for w in
              context], dtype=torch.long)
6
7          # Step 2. Recall that torch *accumulates* gradients.
              Before passing in a new instance, you need to zero out
              the gradients from the old instance
8          model.zero_grad()
9
10         # Step 3. Run the forward pass, getting log probabilities
              over next words
11         log_probs = model(context_idxs)
12
13         # Step 4. Compute your loss function. (Again, Torch wants
              the target word wrapped in a tensor)
14         loss = loss_function(log_probs, torch.tensor([word_to_ix[
              target]], dtype=torch.long))
15
16         # Step 5. Do the backward pass and update the gradient
17         loss.backward()
18         optimizer.step()
19
20         # Get the Python number from a 1-element Tensor by calling
              tensor.item()
21         total_loss += loss.item()
22
23     #print("\t", total_loss)
24     losses.append(total_loss)

```

(c) Python, only the second function

```

1  def most_similar(word_to_test, word_to_ix):
2      test_embedding = get_word_embedding_for_word(word_to_test,
              word_to_ix)
3
4      # get embeddings for all other possible words like aaa bbb
              ccc
5      cos = torch.nn.CosineSimilarity()

```

```

6     results = {}
7     for c in string.ascii_lowercase:
8         c_embedding = get_word_embedding_for_word(c+c+c,
9             word_to_ix)
10        cosine_similarity = cos(test_embedding, c_embedding)
11        results[c+c+c] = cosine_similarity.item()
12
13    sorted_results = dict(sorted(results.items(), key=lambda
14        item: -item[1]))
15    return sorted_results

```

Task 5

```

1  # Let's load the model
2  model = gensim.models.KeyedVectors.load_word2vec_format("
3      word2vec_embeddings.bin", binary=True)
4
5  # Print the length fo the whole vocabulary
6  print(len(model.wv.vocab))
7
8  # Print the embedding of the word good
9  print(model["good"])
10
11 # Get the 10 most similar words of "good"
12 print(model.most_similar("good", topn = 10))
13
14 # Check whether our embeddings are good at the analogy task, what
15   is the result to "?-man" when looking at "queen-king"?
16 print(model.most_similar(positive = ["woman", "king"], negative =
17     ["man"], topn = 1))

```