

CSU44061 Machine Learning | Weekly Assignment 9
Brendan McCann | 20332615

Part (i)

(a)

Briefly describing the following datasets:

- `input_childSpeech_trainingSet.txt`

This dataset contains a series of short sentences with poor grammar and vocabulary - in other words, by a child. The vocabulary size is 40. The sentence length ranges from 13 to 85 characters long.

- `input_childSpeech_testSet.txt`

This dataset contains a similar style as the previous one, except with different and unseen sentences. The vocabulary size is also 40, with a sentence length ranging between 13 and 83 characters long; very similar to the training set.

- `input_shakespeare.txt`

The style is much different. This text is from a play, so we have character names preceding their lines - for example: `MARCIUS:...`. The language is richer and more advanced. The vocabulary size is 25,671 and sentences range from lengths of 2 to 63 characters.

(b)

Seeing as the vocabulary of `input_childSpeech_trainingSet.txt` is very small and limited, I chose to reduce the number of word embeddings in the model from the default 384 to 40, which is the vocabulary size (unique characters). This reduced the model parameters to 127,320. My intuition was that it is unnecessary to have more word embeddings than the vocabulary size.

I kept the other hyper-parameters at the same values because the model parameter size dropped drastically after reducing the number of word embeddings.

TODO

Figure 1: train and evaluation loss in my downsized model.

```
No nap I jump high
I lid Where bads I want more juice please
Look moon No like it
```

Figure 2: Sample output generated from my downsized model. I was quite impressed: it captured the simple syntax and grammar of child-speak very well, and the sentences are understandable.

(c)

- `n_embd=80`

Parameter size: 488,840.

- `n_embd=160, n_layer = 3`

Parameter size: 972,840.

(d)

There is one bias term for each of Key, Value and Query vectors.

The key vector is used to measure how relevant each element in the input sequence is to the query. The inclusion of a bias term would either increase or decrease the relevance of other tokens.

See 3 for a sample generation of my model with only the key bias enabled. I believe this is because there is too much emphasis placed on the relations between the tokens as opposed to the value of the token itself. This could be addressed by modifying the bias of the value vector.

```
Where ball I clib
Dady reade read sad me de y
ve you I seang
```

Figure 3: The results were very poor and indecipherable as you can see.

The value vector holds the actual information or meaning of the token that will be passed to the next layer. The inclusion of a bias term would reduce the importance of common words such as "the, and, etc.". In our case, this is what we want in the context of child speech - the simpler the language, the better. Even if the model fails to capture the semantics or grammar (which may be a by-product of introducing bias for value), it will likely work ok because child speech is often not grammatically correct.

See 4 for a sample generation of my model with only the value bias enabled. We can see the relationships between words are stronger - Read is clearly related to **book**, Big to **truck**, and so on.

```
I jump Big hug More bufore me
Read book Big truck Saw big flufy dont
I wash hands I dove you
```

Figure 4: As a whole, the entire sentence is not very coherent. This is likely because the query bias was not set.

The query vector represents what a token 'looks for' in other tokens. The inclusion of a bias term will impact the magnitude of the search, either making it more or less scoped. In the context of child speech, my intuition is that we would prefer to focus its scope on tokens very close to each other, as opposed to ones that are far away.

See 5 for a sample generation of my model with only the query bias enabled.

I want mas Dadddy read me dinosaur book before bed Bath time
More bubblas Whatt is I run fasst I hide
Help me please Big hug

Figure 5: We can see there are some spelling issues, however the grammatical and syntactic structure is quite solid.

(e)

Without skip connections in our transformer model, the result is incomprehensible. See for sample output. This happens because without skip connections, when the model tries to fine-tune weights through back propagation, the model updates start to 'vanish' back through the model, so the early layers don't receive significant updates.

```
step 0: train loss 3.6862, val loss 3.6876
step 500: train loss 2.0982, val loss 2.0853
step 1000: train loss 2.0284, val loss 2.0159
step 1500: train loss 2.0136, val loss 2.0022
step 2000: train loss 2.0090, val loss 1.9967
step 2500: train loss 2.0062, val loss 1.9955
step 3000: train loss 2.0038, val loss 1.9914
step 3500: train loss 2.0027, val loss 1.9892
step 4000: train loss 2.0025, val loss 1.9897
step 4500: train loss 2.0018, val loss 1.9909
step 4999: train loss 2.0009, val loss 1.9897
```

Figure 6: train and evaluation loss values in a model without skip connections. We can see the loss values are stagnant, an indicator that the model is not learning.

Logg y Uh d I mbbumomp p w fake d fat Nont b
Re t dy tht bove Mo blpl
Coorpimieame meseshuirery re ont What

Figure 7: Sample output from a model trained without skip connections. It is incomprehensible because the model fails to learn effectively.

Part (ii)

(a)

(b)