

NLP ASSIGNMENT 1

PID: A69034019

QUESTION 1

LOGIC: Got word embeddings from glove using `get_word_embedding` and took mean for getting the `embed_vec`.

EXPERIMENTS:

1. No of layers: 2 layers have more accuracy than 3 at the end because of overfitting.
2. Hidden layer size: 100 is the below image output. default size. by increasing to 200:

3 LAYERS

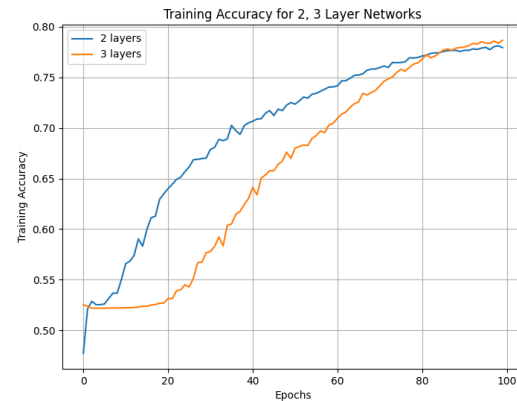
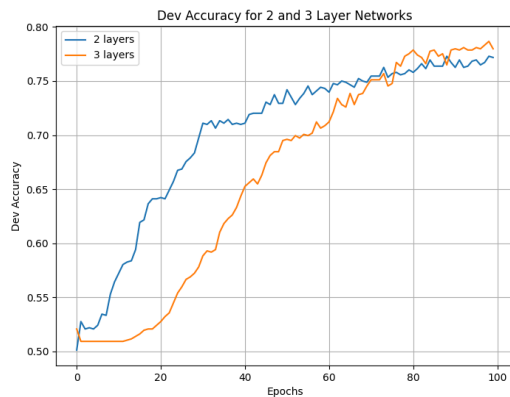
- a. Epoch #80: train accuracy 1.000, dev accuracy 0.784 train loss 0.001, test loss 1.511
- b. Epoch #90: train accuracy 1.000, dev accuracy 0.786 train loss 0.001, test loss 1.571
- c. Epoch #100: train accuracy 1.000, dev accuracy 0.790 train loss 0.001, test loss 1.679

2 LAYERS:

Epoch #80: train accuracy 0.886, dev accuracy 0.783 train loss 0.282, test loss 0.480
Epoch #90: train accuracy 0.899, dev accuracy 0.788 train loss 0.260, test loss 0.483
Epoch #100: train accuracy 0.911, dev accuracy 0.795 train loss 0.241, test loss 0.483
Increasing hidden layers increases accuracy.

3. source of embedding: 50 vs 300d - 300 has more accuracy than 50. 50 doesn't reach 77 percent accuracy for a simple 2 layer neural network but 300 reaches. Time taken for 77 percent accuracy < 2min.
4. Dropout layers: made less difference in DAN training compared to RANDAN. In DAN accuracy improved marginally.
5. Optimizer: SGD:
 - 2 layers:
 - i. Epoch #10: train accuracy 0.550, dev accuracy 0.564 train loss 0.683, test loss 0.683
 - ii. Epoch #20: train accuracy 0.635, dev accuracy 0.641 train loss 0.668, test loss 0.666
 - iii. Epoch #70: train accuracy 0.758, dev accuracy 0.749 train loss 0.529, test loss 0.523
 - iv. Epoch #100: train accuracy 0.779, dev accuracy 0.772 train loss 0.476, test loss 0.476
 - v. 3 layers:
 - vi. Epoch #10: train accuracy 0.522, dev accuracy 0.509 train loss 0.690, test loss 0.691
 - vii. Epoch #20: train accuracy 0.527, dev accuracy 0.524 train loss 0.686, test loss 0.687
 - viii. Epoch #70: train accuracy 0.737, dev accuracy 0.745 train loss 0.552, test loss 0.543
 - ix. Epoch #100: train accuracy 0.787, dev accuracy 0.780 train loss 0.453, test loss 0.465

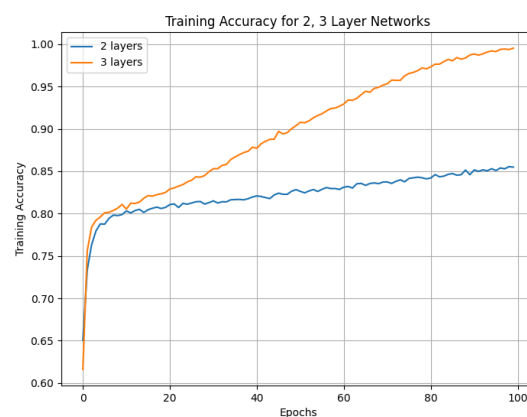
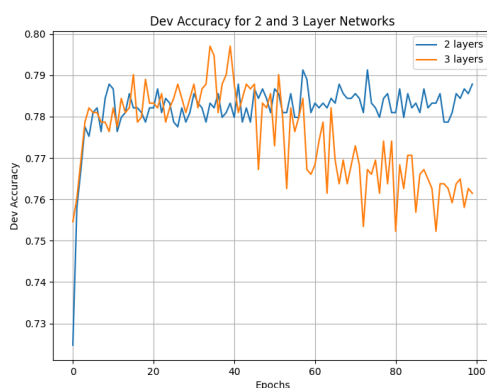
Adam vs SGD: More accuracy with SGD as it works well with large datasets, The way it converges is different from ADAM.



Adam:

```
2 layers:
Epoch #10: train accuracy 0.799, dev accuracy 0.788 train loss 0.431, test loss 0.454
Epoch #20: train accuracy 0.807, dev accuracy 0.782 train loss 0.414, test loss 0.459
Epoch #30: train accuracy 0.813, dev accuracy 0.781 train loss 0.402, test loss 0.458
Epoch #40: train accuracy 0.820, dev accuracy 0.783 train loss 0.392, test loss 0.456
Epoch #50: train accuracy 0.828, dev accuracy 0.781 train loss 0.382, test loss 0.456
Epoch #60: train accuracy 0.829, dev accuracy 0.781 train loss 0.371, test loss 0.457
Epoch #70: train accuracy 0.837, dev accuracy 0.784 train loss 0.361, test loss 0.469
Epoch #80: train accuracy 0.841, dev accuracy 0.781 train loss 0.351, test loss 0.460
Epoch #90: train accuracy 0.846, dev accuracy 0.783 train loss 0.341, test loss 0.466
Epoch #100: train accuracy 0.855, dev accuracy 0.788 train loss 0.330, test loss 0.466

3 layers:
Epoch #10: train accuracy 0.811, dev accuracy 0.776 train loss 0.413, test loss 0.447
Epoch #20: train accuracy 0.825, dev accuracy 0.783 train loss 0.383, test loss 0.455
Epoch #30: train accuracy 0.850, dev accuracy 0.784 train loss 0.342, test loss 0.464
Epoch #40: train accuracy 0.878, dev accuracy 0.797 train loss 0.294, test loss 0.481
Epoch #50: train accuracy 0.904, dev accuracy 0.786 train loss 0.246, test loss 0.531
Epoch #60: train accuracy 0.927, dev accuracy 0.766 train loss 0.193, test loss 0.589
Epoch #70: train accuracy 0.952, dev accuracy 0.768 train loss 0.143, test loss 0.651
Epoch #80: train accuracy 0.971, dev accuracy 0.774 train loss 0.100, test loss 0.770
Epoch #90: train accuracy 0.987, dev accuracy 0.763 train loss 0.061, test loss 0.904
Epoch #100: train accuracy 0.995, dev accuracy 0.761 train loss 0.032, test loss 1.086
```



1b) How does this compare to using GloVe? RANDAN has a bad accuracy in the beginning but increases as it learns.

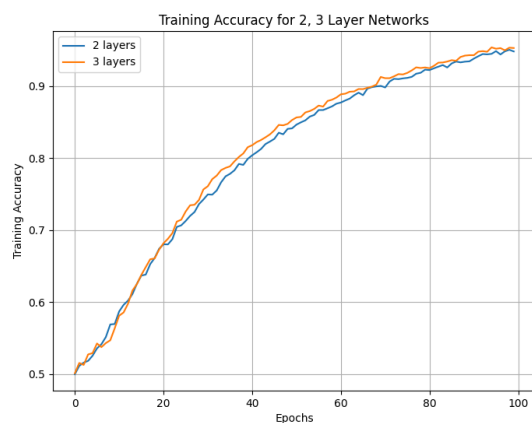
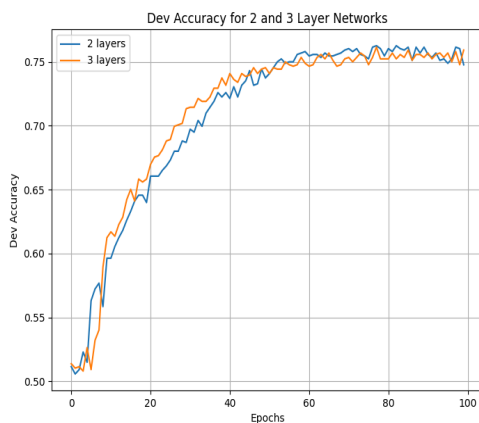
Discuss how training embeddings from scratch affect model convergence and performance:

RANDAN converges slower and less stable than DAN. Also, it has lesser accuracy(less generalization to unseen data) Have to give dropouts for it to train faster. else takes a lot of time.

RANDAN LOGIC: Got word index from the glove and initialised word embedding in Pytorch.

EXPERIMENTS

6. No of layers: 2 layers have more accuracy than 3 at the end because of overfitting.
7. Hidden layer size: increased to 200:
 - a. 2 layers
 - i. Epoch #70: train accuracy 0.919, dev accuracy 0.751 train loss 0.212, test loss 0.616
 - ii. Epoch #80: train accuracy 0.936, dev accuracy 0.776 train loss 0.172, test loss 0.651
 - b. 3 layers
 - i. Epoch #70: train accuracy 0.918, dev accuracy 0.758 train loss 0.200, test loss 0.717
 - ii. Epoch #80: train accuracy 0.938, dev accuracy 0.756 train loss 0.156, test loss 0.782
 - iii. Epoch #90: train accuracy 0.954, dev accuracy 0.758 train loss 0.125, test loss 0.893
 - c. greater size of hidden layer gives more accuracy.
8. Dropout layers: made less difference in DAN training compared to RANDAN. In DAN accuracy improved marginally. increasing the rate too much will underfit but too less results in overfitting.
 - a. dropout rate: 0.5 - images below show accuracy
 - b. dropout rate: 0.4: got better accuracy than 0.5
9. Optimizer: SGD gives bad accuracy since it can get stuck in local minima. Epoch #80: train accuracy 0.530, dev accuracy 0.518 train loss 0.690, test loss 0.690
10. Epoch #90: train accuracy 0.537, dev accuracy 0.521 train loss 0.689, test loss 0.690



QUESTION 2

For BPE 2layers give more accuracy than in DAN as compared to three layers since the model overfits for 3 layers.

VOCAB :2000

```

2 layers:
Read in 14923 vectors of size 300
Epoch #10: train accuracy 0.623, dev accuracy 0.673 train loss 0.653, test loss 0.639
Epoch #20: train accuracy 0.774, dev accuracy 0.737 train loss 0.488, test loss 0.533
Epoch #30: train accuracy 0.860, dev accuracy 0.774 train loss 0.333, test loss 0.501
Epoch #40: train accuracy 0.906, dev accuracy 0.787 train loss 0.236, test loss 0.537
Epoch #50: train accuracy 0.933, dev accuracy 0.764 train loss 0.171, test loss 0.644
Epoch #60: train accuracy 0.957, dev accuracy 0.779 train loss 0.117, test loss 0.710
Epoch #70: train accuracy 0.971, dev accuracy 0.780 train loss 0.083, test loss 0.829
Epoch #80: train accuracy 0.980, dev accuracy 0.782 train loss 0.062, test loss 0.942
Epoch #90: train accuracy 0.985, dev accuracy 0.771 train loss 0.047, test loss 1.026
Epoch #100: train accuracy 0.989, dev accuracy 0.786 train loss 0.034, test loss 1.136

```

```

3 layers:
Read in 14923 vectors of size 300
Epoch #10: train accuracy 0.735, dev accuracy 0.703 train loss 0.527, test loss 0.572
Epoch #20: train accuracy 0.866, dev accuracy 0.740 train loss 0.312, test loss 0.594
Epoch #30: train accuracy 0.927, dev accuracy 0.743 train loss 0.189, test loss 0.707
Epoch #40: train accuracy 0.957, dev accuracy 0.747 train loss 0.118, test loss 0.879
Epoch #50: train accuracy 0.973, dev accuracy 0.751 train loss 0.075, test loss 1.022
Epoch #60: train accuracy 0.984, dev accuracy 0.757 train loss 0.049, test loss 1.235
Epoch #70: train accuracy 0.991, dev accuracy 0.763 train loss 0.032, test loss 1.320
Epoch #80: train accuracy 0.991, dev accuracy 0.752 train loss 0.028, test loss 1.595
Epoch #90: train accuracy 0.995, dev accuracy 0.758 train loss 0.017, test loss 1.734
Epoch #100: train accuracy 0.992, dev accuracy 0.758 train loss 0.021, test loss 1.703

```

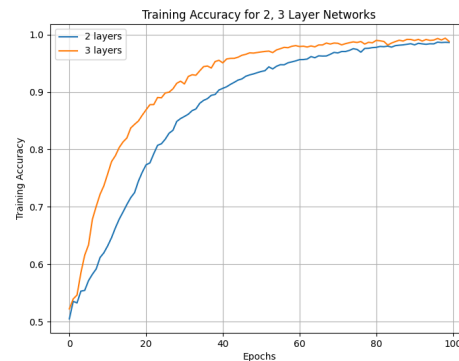
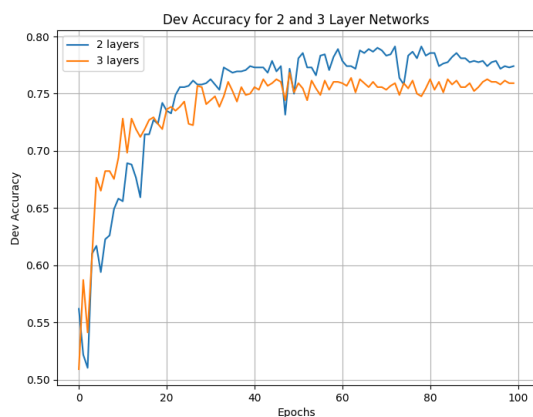
VOCAB : 1500

```

2 layers:
Read in 14923 vectors of size 300
Epoch #10: train accuracy 0.620, dev accuracy 0.658 train loss 0.654, test loss 0.643
Epoch #20: train accuracy 0.760, dev accuracy 0.742 train loss 0.513, test loss 0.544
Epoch #30: train accuracy 0.854, dev accuracy 0.759 train loss 0.353, test loss 0.509
Epoch #40: train accuracy 0.903, dev accuracy 0.774 train loss 0.248, test loss 0.527
Epoch #50: train accuracy 0.934, dev accuracy 0.750 train loss 0.170, test loss 0.642
Epoch #60: train accuracy 0.954, dev accuracy 0.789 train loss 0.123, test loss 0.666
Epoch #70: train accuracy 0.969, dev accuracy 0.788 train loss 0.091, test loss 0.774
Epoch #80: train accuracy 0.977, dev accuracy 0.783 train loss 0.065, test loss 0.887
Epoch #90: train accuracy 0.984, dev accuracy 0.778 train loss 0.050, test loss 0.975
Epoch #100: train accuracy 0.986, dev accuracy 0.774 train loss 0.042, test loss 1.064

3 layers:
Read in 14923 vectors of size 300
Epoch #10: train accuracy 0.737, dev accuracy 0.694 train loss 0.524, test loss 0.567
Epoch #20: train accuracy 0.860, dev accuracy 0.719 train loss 0.331, test loss 0.596
Epoch #30: train accuracy 0.919, dev accuracy 0.741 train loss 0.203, test loss 0.699
Epoch #40: train accuracy 0.955, dev accuracy 0.750 train loss 0.128, test loss 0.842
Epoch #50: train accuracy 0.969, dev accuracy 0.751 train loss 0.086, test loss 0.906
Epoch #60: train accuracy 0.981, dev accuracy 0.760 train loss 0.059, test loss 1.078
Epoch #70: train accuracy 0.985, dev accuracy 0.756 train loss 0.042, test loss 1.236
Epoch #80: train accuracy 0.986, dev accuracy 0.755 train loss 0.039, test loss 1.264
Epoch #90: train accuracy 0.991, dev accuracy 0.759 train loss 0.024, test loss 1.562
Epoch #100: train accuracy 0.988, dev accuracy 0.759 train loss 0.036, test loss 1.537

```



We can observe that reducing vocab size gives less accuracy. Reasons:

1. Small vocab size means less words: so the more frequent pairs maybe missed.
2. if a word is not present, BPE breaks into subwords which may not be correct semantic meaning.
3. smaller size => captures less meaning from the text
4. multiple words can have the same set of subwords after the split so it can not generalize well.

Compare the performance of your subword-based DAN model to the word-level DAN model you implemented earlier: BPE/SUBWORD-based model converges slower than DAN and gives higher accuracy than DAN. (expected trend)

QUESTION 3: Attached is a handwritten pdf after references.

INSTRUCTIONS TO EXECUTE CODE: Readme.md

REFERENCES:

1. <https://www.youtube.com/watch?v=zduSFxRajkE&t=20s>
2. Pytorch documentation
3. DAN: slides and sentiment.py understanding was enough to implement.
4. SKIPGRAM: using slides and online examples.
5. Stackoverflow, pytorch forums for errors: mostly debugging RANDAN.
 - a. `x = .xlong()` error
 - b. size issues: neural network
 - c. embedding errors for Randan
6. Chatgpt - for errors (though used checked forums and PyTorch subs for right error corrections, understand some parts of bpe, debugging RANDAN. to check some small mistakes like loop indices and syntax for IN bpe. (MY Loop was giving wrong answers for debugging. couldn't correct - debugged for a long time. then used this directly)

```
7.         w = w.replace(".join(p).p[0]+' '+p[1])
8.     final_txt.extend(w.split())
```

ASSIGN 1
A69034019

PART ③

Q①. (a):-

skip gram model:

$$P(q|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

\downarrow Context word \downarrow center word

u_w - context word vector
 v_w - center word vector.

$k=1$ $d=2$

$$u_w \text{ dog, cat} = (0, 1)$$

$$u_w \text{ a, the} = (1, 0)$$

$$\text{likelihood} = \sum_{(x,y)} \log(P(y|x))$$

WORD	CONTEXT WORD
the	dog, cat
dog	the, a
cat	the
a	dog

ANS: $P(\text{dog}|\text{the}) = 1/2$ $P(\text{cat}|\text{the}) = 1/2$

$$\Rightarrow P(y|\text{the}) = \begin{cases} 1/2 & y = \text{dog} \\ 1/2 & y = \text{cat} \\ 0 & \text{otherwise} \end{cases}$$

EXPLANATION:

the above is case because both dog & cat have equal likelihood of appearing after the. (they appear exactly once).

If we have $P(\text{dog}|\text{the}) = 0.8$ $P(\text{cat}|\text{the}) = 0.2$ then it means dog has more probability of appearing after the than cat. which is false.

Q. ① ②
lets say

PART ③

$u_{w, the} = (x, y)$
 $u_{w, dog, cat} = (0, 1)$ } take dot product $\Rightarrow y$

$$\Rightarrow \frac{P(dog | the)}{P(cat | the)} = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

$$\Rightarrow \frac{e^y}{e^x + e^x + e^y + e^y}$$

$$= \frac{e^y}{2e^x + 2e^y}$$

For it to be optimal, ① $P(dog | the) \approx P(cat | the)$
and

$$② P(dog | the) \approx 0.5$$

$$P(cat | the) \approx 0.5$$

since probability of them appearing is almost half using the given sentences.

Cases: (trying out values for (x, y)):-

$$x=0, y=0 \quad P(dog | the) = \frac{e^0}{2e^0 + 2e^0} = \frac{1}{4} = 0.25$$

$$x=1, y=0 \quad P(dog | the) = \frac{e^0}{2e^0 + 2e^1} = \frac{1}{2+2e} \approx 0.124$$

$$x=0, y=1 \quad P(dog | the) = \frac{e^1}{2e^0 + 2e^1} \approx 0.365$$

$$x=0, y=5 \quad P(dog | the) = \frac{e^5}{2e^0 + 2e^5} \approx 0.496$$

so x can't be large than y (since

denominator will be larger than Numerator so $P \downarrow$)
so keep least x ($x=0$)

$$x=0, y=10 \quad P(dog | the) = \frac{e^{10}}{2e^0 + 2e^{10}} \approx 0.5$$

So the ans: ($x=0$, $y=\text{any large no.}$) the larger the no. the better $p \rightarrow 0.5$.

$$\lim_{\substack{y \rightarrow \infty \\ \text{v small}}} \frac{e^y}{2e^x + 2e^y} \approx \frac{e^x}{2e^y} \approx 1/2.$$

Ans:- ($x=0$, $y=\text{large no.}$)

PART (3) - (3C)

$u = \text{context word}$ $v = \text{center word}$

Training Set:

- ($u = \text{the}$, $v = \text{dog}$)
- ($u = \text{dog}$, $v = \text{the}$)
- ($u = \text{the}$, $v = \text{cat}$)
- ($u = \text{a}$, $v = \text{dog}$)
- ($u = \text{cat}$, $v = \text{the}$)
- ($u = \text{dog}$, $v = \text{a}$)

PART (3) - (3D)

To maximize skip gram obj. $P(\text{dog}|\text{the})$
 $\approx P(\text{cat}|\text{the}) \approx 0.5$.

Taking c as $(0,1)$ (simple case)
and u as $(0,x)$ all

since from (3b) - 0, large no.

$$P(0|c) = \frac{\exp(u_0^T v_c)}{\sum_{w \in V} \exp(u_0^T v_w)}$$

gives $P(\text{cat}|\text{the}) \approx 0.5$ & $P(\text{dog}|\text{the})$

$v_{\text{the, dog, cat, a}} = (0,1)$

so x can be.

(if u 's are $(0,1)$ then u can't be $(x,0)$ since dot product 0)

$$P(0|c) \Rightarrow P(\text{dog}|\text{the}) \approx P(\text{cat}|\text{the}) \approx P(a|\text{the}) \approx 0.25$$

since its:
$$\frac{e^x}{e^x + e^x + e^x + e^x} \approx \frac{e^x}{4e^x} \approx \frac{1}{4}$$

But $P(\text{dog}|\text{the}) \approx P(\text{cat}|\text{the}) \approx 0.5$ 2. $P(a|\text{the}) \approx 0.1$ or less.

so u for: $(0, 10)$: can be 10 or large no.
 $\begin{matrix} \text{dog} \\ \text{cat} \end{matrix}$ should be 0. as $P(\text{the}|\text{the}) \approx 0$.
 $u_{\text{the}}, a \approx 0, \text{small} (0,1)$

$$\Rightarrow P(\text{dog}|\text{the}) = \frac{e^{10}}{e^1 + e^1 + 2e^{10}} \approx \frac{1}{2}$$

so greater the value of $(0, x)$
 more $P \rightarrow 0.5$

ANS: $u_{\text{dog}} = (0, \text{large no.})$
 $u_{\text{cat}} = (0, \text{large no.})$
 $u_{\text{the}} \approx (0, \text{small no.})$
 $u_a \approx (0, \text{small no.})$
 $\forall \text{ cat, dog, a, the} = (0, 1)$
 (context vectors)