

1 Deep Averaging Networks (DAN) 1.1.2 Results

1.1 a) DAN implementation

In this assignment, I implemented the **Deep Averaging Network (DAN)** for sentiment classification, as outlined in the lecture. The task involves classifying movie reviews as positive or negative based on the Rotten Tomatoes dataset.

1.1.1 Model Architecture

Embedding Layer: This layer transforms the input list of input indices into embeddings. By transforming the indices to the glove embeddings or a randomly initialized embedding layer.

Averaging: For each sentence, the embeddings of the words are averaged to produce a single fixed-size vector.

Feedforward Layers: The averaged embeddings are passed through multiple fully connected (dense) layers. I experimented with different dimensions, but ended up with the set of dimensions; [100, 500, 2].

Regularization: To prevent overfitting, I added dropout layers between the fully connected layers. I used a dropout rate of 50%.

Loss Function: I used cross entropy as my loss function.

For task **1a**, the goal was to implement the DAN and achieve at least 77% accuracy on the development set while keeping training time under five minutes. Using the glove embeddings I got the results presented in figure [1]. The training only took about 30 seconds on my computer, but it would be sufficient to stop after 5 epochs as the model started to overfit after this point. The performance was really good with minimal training, which indicates that the GloVe embeddings are good at catching the semantics of the words.

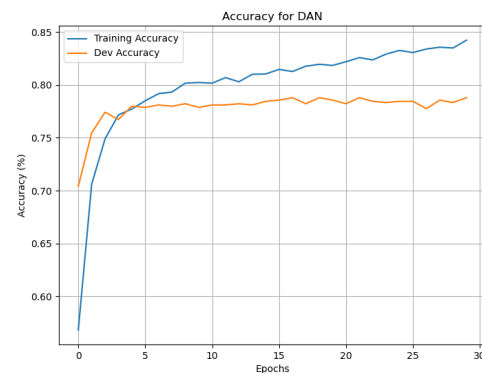


Figure 1: Train and dev accuracy during training using GloVe embeddings.

For task **1b**, I used a randomly initialized set of embeddings and trained them along with the rest of the network. The accuracy during training is visualized in figure [2]. With these embeddings, the training took 10 epochs instead of 5 epochs before it reached its peak on development accuracy. It also started to overfit significantly compared to the fixed GloVe embeddings, which indicates that the sentiment labels are getting memorized directly into the embeddings. It also never reached the same peak development accuracy as with the GloVe embeddings which is most

likely due to the embeddings being overfitted for the dataset. This could maybe be improved by implementing strict regularizations on the embeddings layer. It used almost twice the time to train compared to the GloVe model per epoch because of the increased number of parameters.

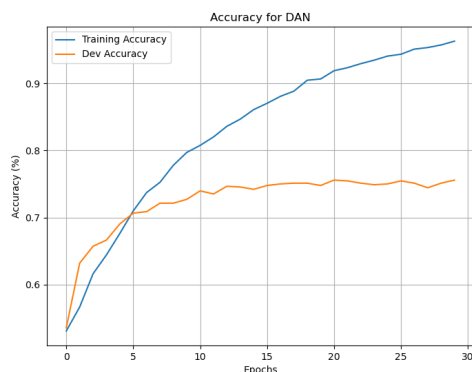


Figure 2: Train and dev accuracy during training using randomly initialized embeddings.

2 Byte-Pair Encoding (BPE)

I implemented BPE as suggested in the algorithm presented in [1]. The detailed implementation can be viewed in the code along with comments and explanations. I created a tokenizer in which used BPE on the training data to create a tokenizer function. This tokenizer were passed into the dataloaders so that the model had access to the tokenized version of the examples.

I experimented with different vocabulary sizes. Following is the results for the most interesting vocabulary sizes:

100 tokens: This is slightly more than the 65 unique ASCII characters in the dataset, meaning the BPE algorithm performed 35 merges. As a result, the vocabulary consists mainly of 1-2 character tokens. The DAN, therefore, averages over letter frequencies, which limits its ability to capture the complex

sentiments that humans typically understand by recognizing entire words or phrases. Despite this, it achieved an accuracy of 65%, which is decent. However, the model is likely memorizing word frequencies instead of learning the underlying sentiment patterns.

4096 tokens: A 4096-token BPE vocabulary balances word-level tokenization with subword representation. This larger size captures more complete words, which is crucial for sentiment analysis, while subword tokens improve generalization over rare words. This approach minimizes over-segmentation while retaining flexibility for less frequent terms, enhancing sentiment classification. It ended up converging at a 75% accuracy on the development set, which is very similar to what I acheived with the randomized embeddings.

12986 tokens: In this vocabulary, almost every word in the dataset was assigned its own token. This makes the tokenization quite similar to a randomized embedding model over the 14,923 unique words, with the difference that not all unique words are included and subwords are also represented. When compared to randomized embeddings, it performed slightly better, likely because the subword tokens allow the model to generalize better across unseen or rare words. This gives the model an advantage in handling words not fully present in the training data, while still benefiting from tokenizing frequent words efficiently. The results are visualized in figure [3].

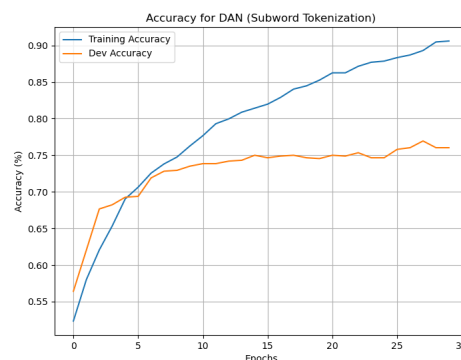


Figure 3: Train and dev accuracy during training using BPE tokenization and randomized embeddings.

3 Understanding Skip-Gram

From the problem statement, the context vectors are:

Q1 - 3a

Using the training data we get the following pairs of skipgrams:

- $(x = \text{the}, y = \text{dog})$
- $(x = \text{dog}, y = \text{the})$
- $(x = \text{the}, y = \text{cat})$
- $(x = \text{cat}, y = \text{the})$
- $(x = \text{a}, y = \text{dog})$
- $(x = \text{dog}, y = \text{a})$

To find the conditional probability that maximizes the data likelihood, we can use the approach suggested in the hint of choosing the model that matches the empirical distribution of the data. Given that $y = \text{the}$, we have the remaining skip-grams:

- $(x = \text{the}, y = \text{dog})$
- $(x = \text{the}, y = \text{cat})$

In which it's easy to see that the probability that maximizes this data is $P(y = \text{dog} | x = \text{the}) = 0.5$ and $P(y = \text{cat} | x = \text{the}) = 0.5$.

Q1 - 3b

For this task, we are tasked with finding a vector \mathbf{v}_{the} such that it gives nearly optimal probabilities for $P(\text{dog}|\text{the})$ and $P(\text{cat}|\text{the})$, both of which should be close to 0.5. The skip-gram equation is given by:

$$P(y|x) = \frac{\exp(\mathbf{v}_x \cdot \mathbf{c}_y)}{\sum_{y'} \exp(\mathbf{v}_x \cdot \mathbf{c}_{y'})}$$

Where:

\mathbf{v}_x is the vector for the center word x (here, "the"), and \mathbf{c}_y is the vector for the context word y (here, "dog", "cat", "the", "a").

$$\mathbf{c}_{\text{dog}} = \mathbf{c}_{\text{cat}} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad \mathbf{c}_{\text{the}} = \mathbf{c}_a = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

We are asked to find \mathbf{v}_{the} such that:

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

Because the embeddings for "dog" and "cat" are identical, we will have $P(\text{dog}|\text{the}) = P(\text{cat}|\text{the})$ for any \mathbf{v}_{the} . To ensure these probabilities are as close to 0.5 as possible, we need $P(\text{the}|\text{the})$ and $P(\text{a}|\text{the})$ to be as close to zero as possible.

Let us define \mathbf{v}_{the} as:

$$\mathbf{v}_{\text{the}} = \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{pmatrix}$$

We compute the dot products:

$$\begin{aligned} \mathbf{v}_{\text{the}} \cdot \mathbf{c}_{\text{dog}} &= \mathbf{v}_1 \cdot 0 + \mathbf{v}_2 \cdot 1 = \mathbf{v}_2 \\ \mathbf{v}_{\text{the}} \cdot \mathbf{c}_{\text{cat}} &= \mathbf{v}_1 \cdot 0 + \mathbf{v}_2 \cdot 1 = \mathbf{v}_2 \\ \mathbf{v}_{\text{the}} \cdot \mathbf{c}_{\text{the}} &= \mathbf{v}_1 \cdot 1 + \mathbf{v}_2 \cdot 0 = \mathbf{v}_1 \\ \mathbf{v}_{\text{the}} \cdot \mathbf{c}_a &= \mathbf{v}_1 \cdot 1 + \mathbf{v}_2 \cdot 0 = \mathbf{v}_1 \end{aligned}$$

The exponentials are:

$$\begin{aligned} \exp(\mathbf{v}_{\text{the}} \cdot \mathbf{c}_{\text{dog}}) &= \exp(\mathbf{v}_2) \\ \exp(\mathbf{v}_{\text{the}} \cdot \mathbf{c}_{\text{cat}}) &= \exp(\mathbf{v}_2) \\ \exp(\mathbf{v}_{\text{the}} \cdot \mathbf{c}_{\text{the}}) &= \exp(\mathbf{v}_1) \\ \exp(\mathbf{v}_{\text{the}} \cdot \mathbf{c}_a) &= \exp(\mathbf{v}_1) \end{aligned}$$

The normalization constant is:

$$\sum_{y'} \exp(\mathbf{v}_{\text{the}} \cdot \mathbf{c}_{y'}) = 2 \exp(\mathbf{v}_1) + 2 \exp(\mathbf{v}_2)$$

Therefore, the probabilities are:

$$\begin{aligned} P(\text{dog}|\text{the}) &= \frac{\exp(\mathbf{v}_2)}{2 \exp(\mathbf{v}_1) + 2 \exp(\mathbf{v}_2)} \\ P(\text{cat}|\text{the}) &= \frac{\exp(\mathbf{v}_2)}{2 \exp(\mathbf{v}_1) + 2 \exp(\mathbf{v}_2)} \\ P(\text{the}|\text{the}) &= \frac{\exp(\mathbf{v}_1)}{2 \exp(\mathbf{v}_1) + 2 \exp(\mathbf{v}_2)} \\ P(\text{a}|\text{the}) &= \frac{\exp(\mathbf{v}_1)}{2 \exp(\mathbf{v}_1) + 2 \exp(\mathbf{v}_2)} \end{aligned}$$

To achieve $P(\text{dog}|\text{the}) \approx 0.5$ and $P(\text{the}|\text{the}) \approx 0$, we need $\exp(\mathbf{v}_2) \gg \exp(\mathbf{v}_1)$. We can accomplish this by setting $\mathbf{v}_2 = 0$ and choosing a large negative value for \mathbf{v}_1 . Let's choose $\mathbf{v}_1 = -10$:

$$\mathbf{v}_{\text{the}} = \begin{pmatrix} -10 \\ 0 \end{pmatrix}$$

Now, compute the exponentials:

$$\begin{aligned} \exp(\mathbf{v}_1) &= \exp(-10) \approx 4.539 \times 10^{-5} \\ \exp(\mathbf{v}_2) &= \exp(0) = 1 \end{aligned}$$

The normalization constant becomes:

$$2 \exp(-10) + 2 \exp(0) \approx 0 + 2 = 2$$

Thus, the probabilities are:

$$\begin{aligned} P(\text{dog}|\text{the}) &\approx \frac{1}{2} \\ P(\text{cat}|\text{the}) &\approx \frac{1}{2} \\ P(\text{the}|\text{the}) &\approx 2.269 \times 10^{-5} \\ P(\text{a}|\text{the}) &\approx 2.269 \times 10^{-5} \end{aligned}$$

These probabilities satisfy our requirements of having a margin < 0.01 from the desired probability.

Q2 - 3c

With a window size $k = 1$, the skip-gram model considers the words immediately before and after the target word as context. For the given sentences, we

extract the (center word, context word) pairs as follows:

"the dog" - (the, dog), (dog, the)
"the cat" - (the, cat), (cat, the)
"a dog" - (a, dog), (dog, a)
"a cat" - (a, cat), (cat, a)

Which is all our pairs for our training data.

Q2 - 3d

We aim to find word vectors \mathbf{v}_x and context vectors \mathbf{c}_y such that the model's predicted probabilities $P(y|x)$ match the empirical probabilities from the training data.

Empirical Probabilities:

From the training examples, the empirical conditional probabilities are:

- For center words "the" and "a":

$$P(\text{dog}|x) = P(\text{cat}|x) = \frac{1}{2}, \quad x \in \{\text{the}, \text{a}\}$$

- For center words "dog" and "cat":

$$P(\text{the}|x) = P(\text{a}|x) = \frac{1}{2}, \quad x \in \{\text{dog}, \text{cat}\}$$

We can recognize that this is similar to problem **3b**, on which we've already found a solution to the word vector \mathbf{v}_{the} . Because our empirical data and context vectors are identical for \mathbf{v}_{the} and \mathbf{v}_a , and similarly with \mathbf{v}_{dog} and \mathbf{v}_{cat} , we can use the word vectors:

$$\mathbf{v}_{\text{the}} = \mathbf{v}_a = \begin{pmatrix} -10 \\ 0 \end{pmatrix}, \quad \mathbf{v}_{\text{dog}} = \mathbf{v}_{\text{cat}} = \begin{pmatrix} 0 \\ -10 \end{pmatrix}$$

Verification of Probabilities:

For center words $x \in \{\text{the}, \text{a}\}$ we have the word vector:

$$\mathbf{v}_x = \begin{pmatrix} -10 \\ 0 \end{pmatrix}$$

This will give us the following exponentials:

$$\mathbf{v}_x \cdot \mathbf{c}_y = \begin{cases} -10 \cdot 0 + 0 \cdot 1 = 0 & \text{if } y \in \{\text{dog}, \text{cat}\} \\ -10 \cdot 1 + 0 \cdot 0 = -10 & \text{if } y \in \{\text{the}, \text{a}\} \end{cases}$$

$$\exp(\mathbf{v}_x \cdot \mathbf{c}_y) = \begin{cases} \exp(0) = 1 & \text{if } y \in \{\text{dog, cat}\} \\ \exp(-10) \approx 0 & \text{if } y \in \{\text{the, a}\} \end{cases}$$

And the normalization constant will still be the same:

$$2 \times \exp(0) + 2 \times \exp(-10) \approx 2$$

This yields probabilities very close to the empirical data:

$$P(y|x) = \frac{\exp(\mathbf{v}_x \cdot \mathbf{c}_y)}{2} = \begin{cases} \frac{1}{2} & \text{if } y \in \{\text{dog, cat}\} \\ \approx 0 & \text{if } y \in \{\text{the, a}\} \end{cases}$$

For center words $x \in \{\text{dog, cat}\}$ we have the word vector:

$$\mathbf{v}_x = \begin{pmatrix} 0 \\ -10 \end{pmatrix}$$

In which we calculated the probabilities in **3b**) to be

$$P(y|x) = \frac{\exp(\mathbf{v}_x \cdot \mathbf{c}_y)}{2} = \begin{cases} \frac{1}{2} & \text{if } y \in \{\text{the, a}\} \\ \approx 0 & \text{if } y \in \{\text{dog, cat}\} \end{cases}$$

References

- [1] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In Katrin Erk and Noah A. Smith, editors, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics.