

## **1 SVD and PCA [35 Points]**

**Problem A [3 points]:**

### 1 SVD and PCA [35 Points]

**Problem A [3 points]:** Let  $X$  be a  $N \times N$  matrix. For the singular value decomposition (SVD)  $X = U\Sigma V^T$ , show that the columns of  $U$  are the principal components of  $X$ . What relationship exists between the singular values of  $X$  and the eigenvalues of  $XX^T$ ?

From the slides, we know PCA gives the following

$$XX^T = U\Sigma U^T \leftarrow \text{PCA}$$

From SVD decomposition, we have:

$$X = U\Sigma V^T$$

$$\Sigma^T = \Sigma \leftarrow \text{diag matrix}$$

$$\begin{aligned} X \cdot X^T &= U\Sigma V^T V \Sigma^T = U\Sigma V^T (\Upsilon \Sigma V^T)^T \\ &= U\Sigma V^T V \Sigma^T U^T \end{aligned}$$

Considering  $V$  is an orthogonal Matrix,  $V^T V = I$

$$\therefore X \cdot X^T = U\Sigma^2 U^T.$$

From the previous equation, we have the matrix  $U$  solves PCA with  $\Lambda = \Sigma^2$ .

Thus, each column of  $U$  is an eigenvector of  $XX^T$ , and we know  $XX^T$  is a principal component of  $X$ .

$\therefore$  The columns of  $U$  in the SVD equation are the principal components of  $X$ .

The relationship between the singular values of  $X$  and the eigenvalues of  $XX^T$  is that the latter are squared of the former.

**Problem B [4 points]:**

The mathematical justification is that they are squared of the singular values of X and thus, non negative. A more intuitive explanation is that eigenvalues represent the total amount of variance that can be explained by a given principal component. In theory they can be positive or negative, but as in practice they represent the variance they can not be negative, as variance is always positive.

**Problem C [5 points]:**

**Problem C [5 points]:** In calculating the Frobenius and trace matrix norms, we claimed that the trace is invariant under cyclic permutations (i.e.,  $\text{Tr}(ABC) = \text{Tr}(BCA) = \text{Tr}(CAB)$ ). Prove that this holds for any number of square matrices.

*Hint:* First prove that the identity holds for two matrices and then generalize. Recall that  $\text{Tr}(AB) = \sum_{i=1}^N (AB)_{ii}$ . Can you find a way to expand  $(AB)_{ii}$  in terms of another sum?

from the hint:

$$\text{Tr}(AB) = \sum_{i=1}^N (AB)_{ii} = \sum_{i=1}^N \sum_{j=1}^N A_{ij} \cdot B_{ji} \quad (1)$$

$$\text{tr}(\beta A) = \sum_{i=1}^N (\beta A)_{ii} = \sum_{i=1}^N \sum_{j=1}^N \beta_{ij} \cdot A_{ji} = \sum_{i=1}^N \sum_{j=1}^N A_{ji} \beta_{ij}$$

$$\text{tr}(\beta A) = \sum_{j=1}^N \sum_{i=1}^N A_{ji} \beta_{ij} \quad (2)$$

From (1) and (2), we have that  $\text{tr}(AB) = \text{tr}(\beta A)$  because both of them give the same sum.

For the general case, we can show that this relationship holds by grouping matrices in order to have the product of just 2 matrices.

$$\therefore \text{tr}(ABC) = \sum_{i=1}^N \sum_{j=1}^N A_{ij} \cdot (BC)_{ji}$$

$$\text{tr}(\beta CA) = \sum_{j=1}^N \sum_{i=1}^N (\beta CA)_{ji} \cdot A_{ij}$$

$\therefore \text{tr}(ABC) = \text{tr}(\beta CA)$  and similar for more matrices.

**Problem D [3 points]:**

For  $k$  singular values, we need to store  $k$  columns of both  $U$  and  $V$  (which are initially of length  $N$ ), so we will end with  $N*k*2$  values (from  $U$  and  $V$ ) plus  $k$  values from the matrix  $\Sigma$  (which has  $k$  term on the diagonal). Thus, we will end with  $2*N*k + k = k*(1+2*N)$  values stored in the truncated SVD. Storing the truncated SVD will be more efficient than storing the whole SVD matrix when  $k < N^2/(2*N+1)$

**Dimensions & Orthogonality**

**Problem E [3 points]:**

**Problem E [3 points]:** Assume that  $D > N$  and that  $X$  has rank  $N$ . Show that  $U\Sigma = U'\Sigma'$ , where  $\Sigma'$  is the  $N \times N$  matrix consisting of the first  $N$  rows of  $\Sigma$ , and  $U'$  is the  $D \times N$  matrix consisting of the first  $N$  columns of  $U$ . The representation  $U'\Sigma'V^T$  is called the "thin" SVD of  $X$ .

$$X = U\Sigma V^T$$

$$X = D \begin{bmatrix} U \\ \Sigma \\ V^T \end{bmatrix}$$

$D = N \times D$   
 $U = D \times N$   
 $\Sigma = N \times N$

Assume  $D \geq N$ .  $X$  has rank  $N$

We have to show that  $U\Sigma = U'\Sigma'$

$$(U\Sigma)_{ij} = \sum_{k=1}^D U_{ik} \Sigma_{kj} = \sum_{k=1}^N U_{ik} \cdot \Sigma_{kj} + \sum_{k=N+1}^D U_{ik} \cdot \underline{\Sigma_{kj} \quad 0}$$

Since  $\Sigma$  only has non-zero entries in the diagonal elements until the position  $\Sigma_{nn}$ , the summation highlighted in blue will be zero, because for  $k > N$  the element  $U_{ik}$  will be multiplied by terms  $\Sigma_{kj}$ , which are zero.

$$\text{Thus } (U\Sigma)_{ij} = \sum_{k=1}^N U_{ik} \cdot \Sigma_{kj} = (U'\Sigma')_{ij} \quad \square$$

for every  $i \in [1, D]$ ,  $j \in [1, N]$

**Problem F [3 points]:**

We know that  $U'$  is a  $D \times N$  matrix so  $U'^T$  is a  $N \times D$  matrix. Thus  $U'^*U'^T$  is a  $D \times D$  matrix, meanwhile  $U'^T*U'$  is a  $N \times N$  matrix. Hence  $U'^*U'^T$  and  $U'^T*U'$  are never the same matrices because they are of different dimensions and as a consequence the definition of orthogonality can not hold.

**Problem G [4 points]:**

**Problem G [4 points]:** Even though  $U'$  is not orthogonal, it still has similar properties. Show that  $U'^T U' = I_{N \times N}$ . Is it also true that  $U' U'^T = I_{D \times D}$ ? Why or why not? Note that the columns of  $U'$  are still orthonormal. Also note that orthonormality implies linear independence.

Since  $U'$  has orthonormal columns, we know that for  $U'^T U$  we will have 1 on diagonal elements ( $\langle U_i, U_j \rangle = 1 \forall i, j \in \mathbb{N}^n$ ) and 0 in the off-diagonal elements ( $\langle U_i, U_j \rangle = 0 \forall i \neq j, i, j \in \mathbb{N}^n$ ). Thus, the result is

$I_{N \times N}$   $\square$

For the case of  $U' U'^T$ , we know that  $U'$  is a  $D \times N$  matrix, and that  $D > N$ . So, considering  $U'$  has more rows than columns, they are not all linearly independent and hence, they can not be all orthonormal, and thus, the dot product between two different rows can be different of zero

$$\therefore U' U'^T \neq I_D$$

## Pseudoinverses

### Problem H [4 points]:

Considering the properties of diagonal matrices and inverses, we know that the inverse of a diagonal matrix is a diagonal matrix with each entry equal to the reciprocal of the corresponding entry in the original matrix.

Taking into account that  $\Sigma$  is invertible, and considering what was mentioned in lecture 10 regarding pseudo inverses, we know that  $\Sigma^+$  with entries  $\sigma^+ = 1/\sigma$  if  $\sigma > 0$  and 0 otherwise is equivalent to  $\Sigma^{-1}$ , which as mentioned above, has entries equal to the reciprocal of the corresponding entry in the original matrix.

$$A^+ = V\Sigma^+U^T$$

#### Pseudoinverse

$$\Sigma^+ = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \sigma_D \end{bmatrix} \quad \sigma^+ = \begin{cases} 1/\sigma & \text{if } \sigma > 0 \\ 0 & \text{otherwise} \end{cases}$$

**Problem I [4 points]:**

**Problem I [4 points]:** Another expression for the pseudoinverse is the least squares solution  $X^{+'} = (X^T X)^{-1} X^T$ . Show that (again assuming  $\Sigma$  invertible) this is equivalent to  $V \Sigma^{-1} U^T$ .

$$X^{+'} = (X^T X)^{-1} X^T$$

$$X^{+'} = X^{-1} (X^T)^{-1} X^T$$

$$X^{+'} = (U \Sigma V^T)^{-1} (V \Sigma V^T)^{-1} (V \Sigma U^T)$$

$$X^{+'} = (V^T)^{-1} \Sigma^{-1} U^{-1} U^T \Sigma^{-1} V^T V \Sigma U^T$$

$$X^{+'} = (V^T)^{-1} \Sigma^{-1} U^{-1} \quad \text{I}$$

$$X^{+'} = V \Sigma^{-1} U^T \quad \square$$


---

$$\begin{aligned} V &= V^T \\ U^T &= U^{-1} \end{aligned} \quad \left. \begin{array}{l} \text{orthogonal} \\ \text{ } \end{array} \right\}$$

$$V \cdot V^T = I$$

$$V^T \cdot V = I$$

**Problem J [2 points]:**

**Problem J [2 points]:** One of the two expressions in problems H and I for calculating the pseudoinverse is highly prone to numerical errors. Which one is it, and why? Justify your answer using condition numbers.

$$x^t = V \Sigma^t U^t \quad \leftarrow \text{Part II}$$

$$x^+ = (X^T X)^{-1} X^T \quad \leftarrow \text{Part I}$$

$$\text{We have that } X^T X = (U \Sigma V^T)^T (U \Sigma V^T)$$

$$\begin{aligned} &= V \Sigma U^T U \Sigma V^T \quad U^T U = I \text{ (orth)} \\ &= V \Sigma^2 V^T \end{aligned}$$

$\therefore X^T X$  has singular values equal to the entries of  $\Sigma^2$

Using condition number to analyze the stability

$$\text{For } (X^T X) \text{ (case Part I)} \rightarrow \frac{\sigma_{\max}}{\sigma_{\min}} = k_1$$

$$\text{For } \Sigma \text{ (case Part II)} \rightarrow \frac{\sigma_{\max}}{\sigma_{\min}} = k_2$$

According to condition number theory, a matrix is well-conditioned while  $k$  is not significantly larger than 1, at the contrary the matrix is said to be ill-conditioned.

$\therefore$  considering  $k_1 \geq k_2$ , computing the inverse of  $(X^T X)$  is more prone to error than calculating the inverse of  $\Sigma$  and thus, Part I is more prone to error.

## **2 Matrix Factorization [30 Points]**

**Problem A [5 points]:**

## 2 Matrix Factorization [30 Points]

In the setting of collaborative filtering, we derive the coefficients of the matrices  $U \in \mathbb{R}^{M \times K}$  and  $V \in \mathbb{R}^{N \times K}$  by minimizing the regularized squared error:

$$\arg \min_{U, V} \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2) + \frac{1}{2} \sum_{i,j} (y_{ij} - u_i^T v_j)^2$$

where  $u_i^T$  and  $v_j^T$  are the  $i^{\text{th}}$  and  $j^{\text{th}}$  rows of  $U$  and  $V$ , respectively, and  $\|\cdot\|_F$  represents the Frobenius norm. Then  $Y \in \mathbb{R}^{M \times N} \approx UV^T$ , and the  $ij$ -th element of  $Y$  is  $y_{ij} \approx u_i^T v_j$ .

**Problem A [5 points]:** Derive the gradients of the above regularized squared error with respect to  $u_i$  and  $v_j$ , denoted  $\partial_{u_i}$  and  $\partial_{v_j}$  respectively. We can use these to compute  $U$  and  $V$  by stochastic gradient descent using the usual update rule:

$$u_i = u_i - \eta \partial_{u_i}$$

$$v_j = v_j - \eta \partial_{v_j}$$

where  $\eta$  is the learning rate.



We have

$$\|V\|_F^2 = \sum_i \|v_i\|^2 \leftarrow \text{Frobenius Norm.}$$

$$\partial U_i = \frac{\lambda}{2} (2U_i) + \frac{1}{2} \sum_j (-v_j) \ell(y_{ij} - u_i^T v_j)$$

$$\underline{\partial U_i = \lambda U_i - \sum_j v_j (y_{ij} - u_i^T v_j)^T}$$

and similarly for  $\partial V_j$

$$\underline{\partial V_j = \lambda V_j - \sum_i U_i (y_{ij} - u_i^T v_j)^T}$$

**Problem B [5 points]:**

**Problem B [5 points]:** Another method to minimize the regularized squared error is alternating least squares (ALS). ALS solves the problem by first fixing  $U$  and solving for the optimal  $V$ , then fixing this new  $V$  and solving for the optimal  $U$ . This process is repeated until convergence.

Derive closed form expressions for the optimal  $u_i$  and  $v_j$ . That is, give an expression for the  $u_i$  that minimizes the above regularized square error given fixed  $V$ , and an expression for the  $v_j$  that minimizes it given fixed  $U$ .

## 1) Expression for $u_i$ given fixed $V$

For minimization, we must make the last derivative equal to zero.

$$f(u_i) = \frac{1}{2} \|v_j(y_{ij} - u_i^\top v_j)\|^2 = 0$$

$$\begin{aligned} \frac{\partial f}{\partial u_i} &= \sum_j v_j u_i^\top - \sum_j v_j v_j^\top u_i = \sum_j v_j u_i^\top - \left(\sum_j v_j v_j^\top\right) u_i \\ \left(\frac{1}{2} I + \sum_j v_j v_j^\top\right) u_i &= \sum_j v_j y_{ij} \end{aligned}$$

$$u_i = \left(\frac{1}{2} I + \sum_j v_j v_j^\top\right)^{-1} \left(\sum_j v_j y_{ij}\right)$$

## 2) Expression for $v_j$ given fixed $U$

$$f(v_j) = \frac{1}{2} \|u_i(y_{ij} - v_j^\top u_i)\|^2 = 0$$

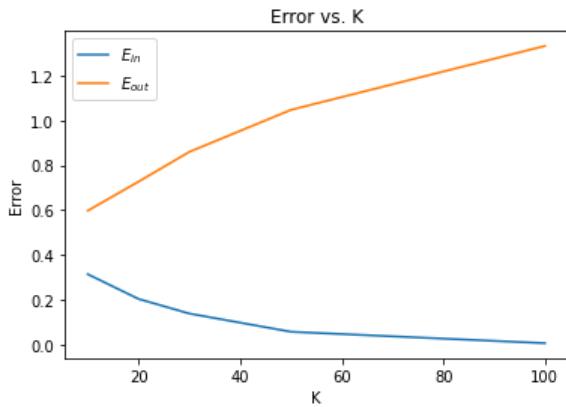
$$\frac{\partial f}{\partial v_j} = \sum_i u_i y_{ij}^\top - \sum_i u_i u_i^\top v_j = 0$$

$$\left(\frac{1}{2} I + \sum_i u_i u_i^\top\right) v_j = \sum_i u_i y_{ij}^\top$$

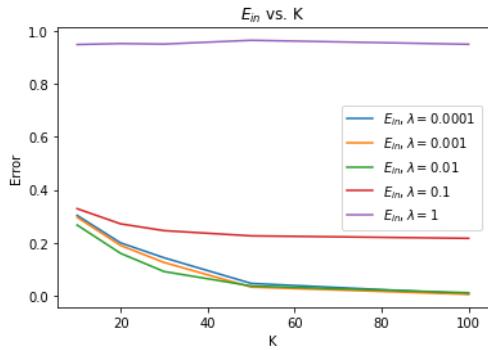
$$v_j = \left(\frac{1}{2} I + \sum_i u_i u_i^\top\right)^{-1} \left(\sum_i u_i y_{ij}^\top\right)$$

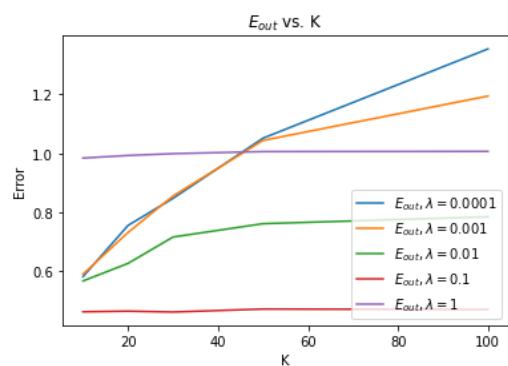
**Problem C [10 points]:**

<https://colab.research.google.com/drive/18PvFSwLD0ICcgxmkL3zq6wDTw3gaZzeR?usp=sharing>

**Problem D [5 points]:**

As we can see in the plot, as  $K$  goes up  $E_{in}$  decreases and  $E_{out}$  increases. This can be explained considering that  $K$  determines the dimensionality of matrices  $U$  and  $V$ . When  $K$  is lower, the input data is smaller and we are summarizing the original data, making the model not to learn the singularities of the training data and with this, prevent overfitting. On the contrary, when  $k$  increases and dimensionality of matrices  $U$  and  $V$  goes up, the model tends to capture these singularities from the training data when the test error goes up, which means that the model is overfitting and it decreases its capacity of generalizing for new data.

**Problem E [5 points]:**



The pattern is similar to part C. In general, as K goes up  $E_{in}$  tends to decrease meanwhile  $E_{out}$  tends to increase. For increasing regularization, such as for  $\lambda = 1$  or  $0.1$ , both curves  $E_{in}$  and  $E_{out}$  are flatten, but for smaller  $\lambda$ 's both of them follow the pattern from previous part that where  $\lambda = 0$ , and when K increases the models start overfitting.

### **3 Word2Vec Principles [35 Points]**

**Problem A [5 points]:**

### Problem #3

**Problem A [5 points]:** If we wanted to train this model with naive gradient descent, we'd need to compute all the gradients  $\nabla \log p(w_O|w_I)$  for each  $w_O, w_I$  pair. How does computing these gradients scale with  $W$ , the number of words in the vocabulary, and  $D$ , the dimension of the embedding space? To be specific, what is the time complexity of calculating  $\nabla \log p(w_O|w_I)$  for a single  $w_O, w_I$  pair?

$$\text{We have } \log p(w_O|w_I) = \log \left( \frac{\exp(U_{w_O}^T V_{w_I})}{\sum_{w=1}^W \exp(U_w^T V_{w_I})} \right)$$

$$\begin{aligned} f_{w_O}(\log(w_O|w_I)) &= \underbrace{f_{w_O}(\log(\exp(U_{w_O}^T V_{w_I})))}_{U_{w_O}^T V_{w_I}} - \log \left( \sum_{w=1}^W \exp(U_w^T V_{w_I}) \right) \\ &= V_{w_I} - f_{w_O} \left( \log \left( \sum_{w=1}^W \exp(U_w^T V_{w_I}) \right) \right) \end{aligned}$$

$$f_{w_O}(\log(w_O|w_I)) = V_{w_I} - \frac{V_{w_I}}{\sum_{w=1}^W \exp(U_w^T V_{w_I})} \cdot (\exp(U_{w_O}^T V_{w_I}))$$

Similarly for  $f_{w_I}$  we have:

$$f_{w_I}(\log(w_O|w_I)) = U_{w_O} - \frac{1}{\sum_{w=1}^W \exp(U_w^T V_{w_I})} \cdot \sum_{w=1}^W (U_w^T V_{w_I}) \cdot U_w$$

The time complexity for computing gradients of  $w_I$  and  $w_O$  is dominated by the numerator term in  $f_{w_I}$ . We can see that, the dot product of  $U_w$  and  $V_{w_I}$  ( $D$ ), and then multiplies it by  $U_w$  ( $D$ ). Considering this operation is made for each  $w$  in  $\{1, \dots, W\}$  ( $W$ ), the time complexity for calculating a gradient for each pair is  $O(W^*D)$ .

**Problem B [10 points]:**

**Problem B [10 points]:** When the number of words in the vocabulary  $W$  is large, computing the regular softmax can be computationally expensive (note the normalization constant on the bottom of Eq. 2). For reference, the standard fastText pre-trained word vectors encode approximately  $W \approx 218000$  words in  $D = 100$  latent dimensions. One trick to get around this is to instead represent the words in a binary tree format and compute the hierarchical softmax.

When the words have all the same frequency, then any balanced binary tree will minimize the average representation length and maximize computational efficiency of the hierarchical softmax. But in practice,

Table 1: Words and frequencies for Problem B

Word	Occurrences
do	18
you	4
know	7
the	20
way	9
of	4
devil	5
queen	6

words occur with very different frequencies — words like “a”, “the”, and “in” will occur many more times than words like “representation” or “normalization”.

The original paper (Mikolov et al. 2013) uses a Huffman tree instead of a balanced binary tree to leverage this fact. For the 8 words and their frequencies listed in the table below, build a Huffman tree using the algorithm found [here](#). Then, build a balanced binary tree of depth 3 to store these words. Make sure that each word is stored as a *leaf node* in the trees.

The representation length of a word is then the length of the path (the number of edges) from the root to the leaf node corresponding to the word. For each tree you constructed, compute the expected representation length (averaged over the actual frequencies of the words).

1)

frequency

word

4

you

4

of

5

devil

6

queen

7

know

9

way

18

do

20

the

1)

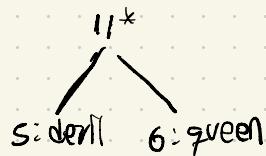
z)

8: \*

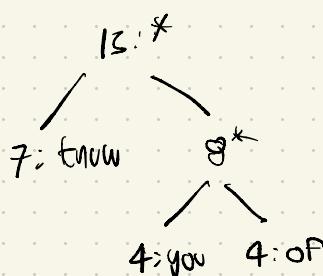
4: you

4: of

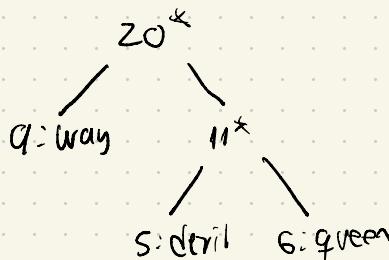
2) S devil  
6 green  
7 know  
8 \*  
9 way  
10 do  
11 the



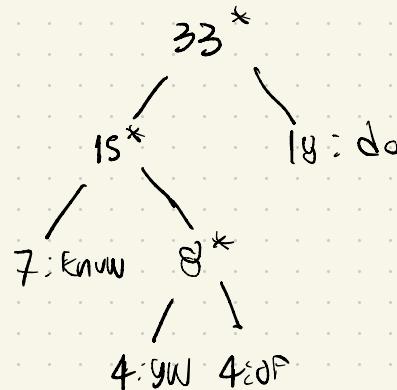
3) 7 know  
8 \*  
9 way  
11 \*  
10 do  
12 the



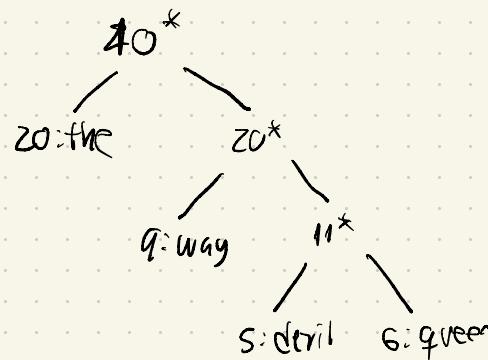
4) 9 way  
11 \*  
15 \*  
18 do  
20 the



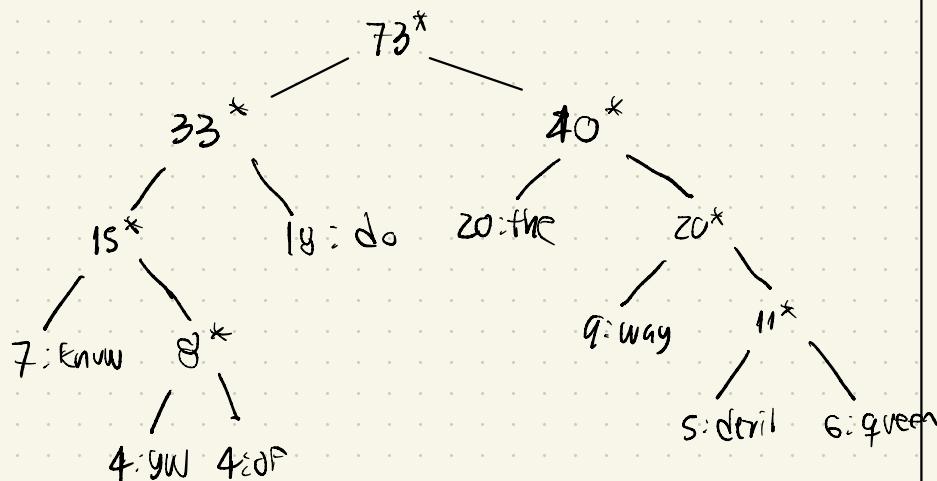
5) 15 \*  
18 do  
20 the  
21 \*



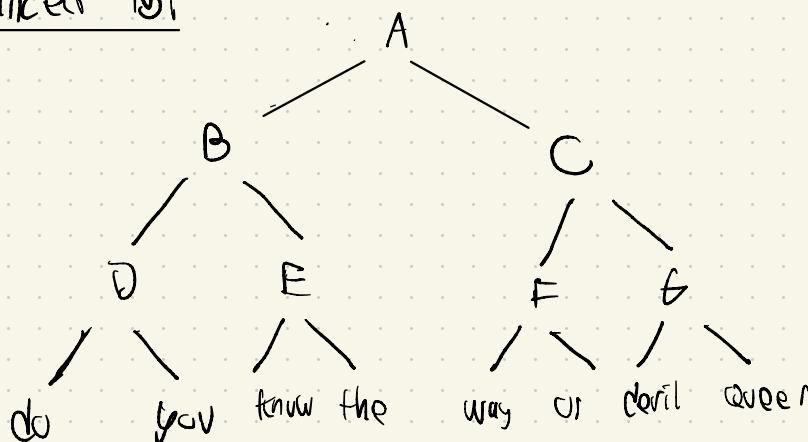
7) 20  
20 \*  
33 \*



9) 33 \*  
40 \*



Balanced BT



The expected representation length is 2.74 for the Huffman tree and 3 for balanced tree

**Problem C [3 points]:**

As we increase D, the value of the training objective will also increase, considering that  $p(w_0|w_i)$  involves multiplying  $v^T$  and  $v$ , where  $v \in \mathbb{R}^D$ . Thus, as D increases, it becomes more computationally expensive to make these calculations. Another aspect to consider when increasing D, is that the model will be more likely to overfit the training data, decreasing its capacity to generalize for new data.

**Problem D [10 points]:**

<https://colab.research.google.com/drive/1gRL03mP0kC-8vVBBYBBeweISSWbspfbl?usp=sharing>

**Problem E [2 points]:**

The dimension of the weigh matrix in the hidden layer is 308x10, where 308 is the number of unique words in the text file and 10 is the number of neurons in the hidden layer

**Problem F [2 points]:**

The dimension of the weight matrix of the output layer is 10x308

**Problem G [1 points]:**

```

Pair(finger, top), Similarity: 0.9980181
Pair(top, finger), Similarity: 0.9980181
Pair(foot, cold), Similarity: 0.9855466
Pair(cold, foot), Similarity: 0.9855466
Pair(shoe, off), Similarity: 0.98533446
Pair(off, shoe), Similarity: 0.98533446
Pair(eggs, ham), Similarity: 0.9851607
Pair(ham, eggs), Similarity: 0.9851607
Pair(green, ham), Similarity: 0.9828462
Pair(every, day), Similarity: 0.98120433
Pair(day, every), Similarity: 0.98120433
Pair(today, gone), Similarity: 0.9808445
Pair(gone, today), Similarity: 0.9808445
Pair(eight, nine), Similarity: 0.979169
Pair(nine, eight), Similarity: 0.979169
Pair(four, six), Similarity: 0.9791032
Pair(six, four), Similarity: 0.9791032
Pair(drink, thing), Similarity: 0.97856116
Pair(thing, drink), Similarity: 0.97856116
Pair(long, way), Similarity: 0.97736496
Pair(way, long), Similarity: 0.97736496
Pair(wire, goodbye), Similarity: 0.9767371
Pair(goodbye, wire), Similarity: 0.9767371
Pair(likes, drink), Similarity: 0.9729196
Pair(sad, glad), Similarity: 0.97162807
Pair(glad, sad), Similarity: 0.97162807
Pair(wink, pink), Similarity: 0.97104865
Pair(pink, wink), Similarity: 0.97104865
Pair(hot, sun), Similarity: 0.9701199
Pair(sun, hot), Similarity: 0.9701199

```

**Problem H [2 points]:**

The first thing that I noticed is that there appear many words that rhyme between them, such as sad and glad, or wink and pink. This makes sense considering the way the poem was written, where words that rhyme with others tends to be close to each other in the text. Another pattern is that many pairs appears twice in both ways ((a,b) and (b,a)), and usually they are immediately below the other. Another pattern is the appearance of words that have the same context, such as (hot, sun) or (eggs, ham).