

6 LATE HOURS

## 1 SVD and PCA [35 Points]

Problem A [3 points]:

Solution A:

A. We have  $X$ , an  $N \times N$  matrix. We also have  $X = U \Sigma V^T$   
So,  $X^T = V \Sigma^T U^T$   
Thus, we have  $XX^T = U \Sigma V^T V \Sigma^T U^T$   
 $XX^T = U \Sigma^2 U^T$   
We recognize this as the definition of PCA.  
So by definition  $U$  are the principal components of  $X$  and  $\Sigma^2$  is the diagonal matrix whose entries are the eigenvalues  $XX^T$ .  
 $\Sigma$  contains the singular values of  $X$ . So, because  $\Sigma^2$  are the eigen values of  $XX^T$ , the singular values of  $X$  are the square root of the eigen values of  $XX^T$ .

**Problem B [4 points]:**

**Solution B:**

B. Intuitively, the eigenvalues of the PCA of  $XX^T$  are non-negative because the entries are variance of components and variance is always non-negative. Mathematically, we know the eigenvalues of  $XX^T$  are the square of the singular values of  $X$  and any number squared is non-negative.

Problem C [5 points]:

Solution C:

C. We want to show that the trace is invariant under cyclic permutations for any number of square matrices. We will show this for 2 matrices at first as the hint suggests and then generalize. Note that  $\text{tr}(AB) = \sum_{i=1}^N (AB)_{ii}$

$$\begin{aligned}\Rightarrow \text{tr}(AB) &= \sum_{i=1}^N \sum_{j=1}^N A_{ij} B_{ji} \\ &= \sum_{j=1}^N \sum_{i=1}^N B_{ji} A_{ij} \\ &= \sum_{j=1}^N (BA)_{jj} \\ \text{tr}(AB) &= \text{tr}(BA)\end{aligned}$$

Now we can take  $\text{tr}(ABC)$  and describe  $D=BC$ . So, we have  $\text{tr}(AD) = \text{tr}(DA)$  which extends to  $\text{tr}(ABC) = \text{tr}(BCA)$ . So we see we can generalize any number of square matrices to the 2 matrix case. So we have shown that trace is invariant under cyclic permutations of any number of square matrices.

**Problem D [3 points]:**

**Solution D:**

D. We need to store the truncated SVD with the  $k$  largest singular values of  $\Sigma$  and the corresponding columns of  $U$  and  $V$ . We need  $k$  values from  $\Sigma$ , the  $k$  largest which are in the diagonal. We need to store  $N \times k$  values from  $U$  and  $N \times k$  values from  $V$  to get the data from the corresponding columns. This requires storage of  $k(2N+1)$  values in total. For values of  $k \leq \frac{N^2}{2N+1}$  it is more efficient to store the truncated SVD.



Problem E [3 points]: .

Solution E:

E. Assuming  $D > N$  and  $X$  has rank  $N$ , we want to show  $U\Sigma = U'\Sigma'$ .

$U$  is  $D \times D$ ,  $\Sigma$  is  $D \times N$   
 $U'$  is  $D \times N$ ,  $\Sigma'$  is  $N \times N$

So we can see  $U\Sigma$  and  $U'\Sigma'$  are both  $D \times N$ .  
 So we have

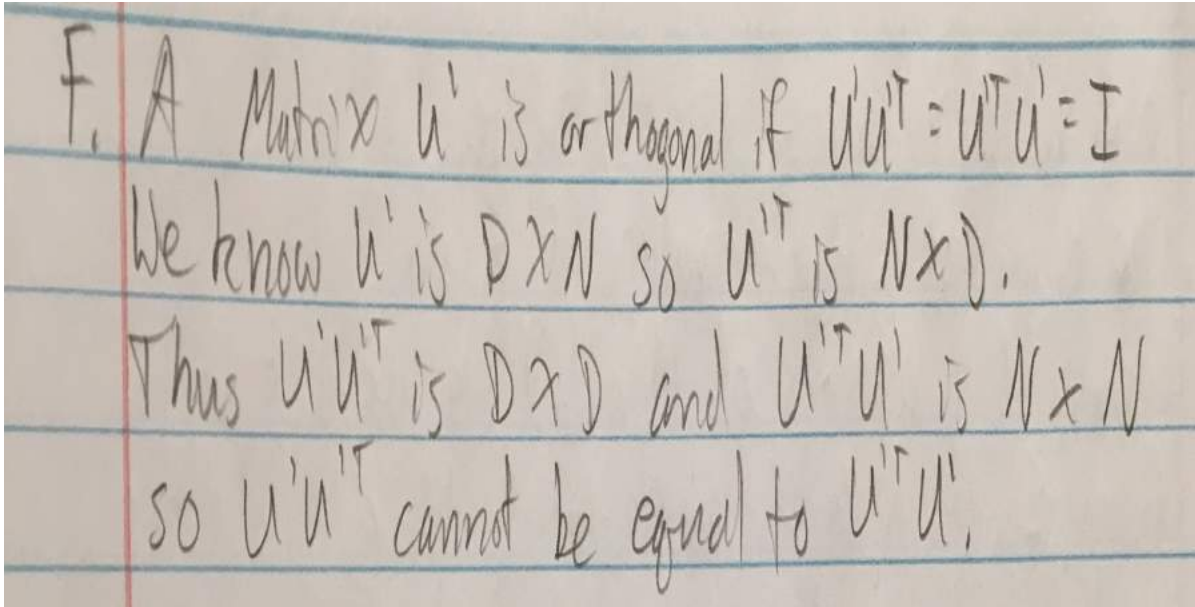
$$U = \begin{bmatrix} u_{00} & \dots & u_{0D} \\ \vdots & & \vdots \\ u_{D0} & \dots & u_{DD} \end{bmatrix} \quad \Sigma = \begin{bmatrix} s_{00} & 0 & 0 & \dots & 0 \\ 0 & s_{11} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & s_{NN} \\ 0 & \dots & 0 & \dots & 0 \end{bmatrix}$$

$$U' = \begin{bmatrix} u'_{00} & \dots & u'_{0N} \\ \vdots & & \vdots \\ u'_{D0} & \dots & u'_{DN} \end{bmatrix} \quad \Sigma' = \begin{bmatrix} s_{00} & s_{01} & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & s_{NN} \end{bmatrix}$$

Thus,  $U\Sigma = U'\Sigma'$

**Problem F [3 points]:**

**Solution F:**



F. A Matrix  $U'$  is orthogonal if  $U'U'^T = U'^TU' = I$   
We know  $U'$  is  $D \times N$  so  $U'^T$  is  $N \times D$ .  
Thus  $U'U'^T$  is  $D \times D$  and  $U'^TU'$  is  $N \times N$   
so  $U'U'^T$  cannot be equal to  $U'^TU'$ .

Problem G [4 points]:

Solution G:

G. We know the columns of  $U'$  are orthonormal and the rows of  $U'^T$  are orthonormal and when two corresponding orthonormal vectors are multiplied together the result is 1 and a orthonormal column with a non corresponding row results in 0. Thus, clearly,  $U'^T U' = I_{n \times n}$ .  $U' U'^T \neq I_{p \times p}$  because the columns of  $U'$ , not the rows, are orthonormal. Thus, we are not multiplying orthonormal vectors together so we won't get  $I_{p \times p}$ .

**Problem H [4 points]:**

**Solution H:**

H. Assuming that  $\Sigma$  is invertible we want to show  $X^+ = V\Sigma^+U^T$  is equivalent to  $V\Sigma^{-1}U^T$ . Note that  $\Sigma$  is a diagonal matrix with  $\sigma_i$  entries on the diagonal and 0 everywhere else. Note then that  $\Sigma^{-1}$  will have  $\frac{1}{\sigma_i}$  at every diagonal entry and 0 everywhere else. From lecture we know  $\sigma^+$  is when the entry is  $1/\sigma$  if  $\sigma > 0$  and 0 otherwise and thus, we see that  $\Sigma^{-1} = \Sigma^+$ . So,  $V\Sigma^+U^T = V\Sigma^{-1}U^T$  as desired.



**Problem I [4 points]:**

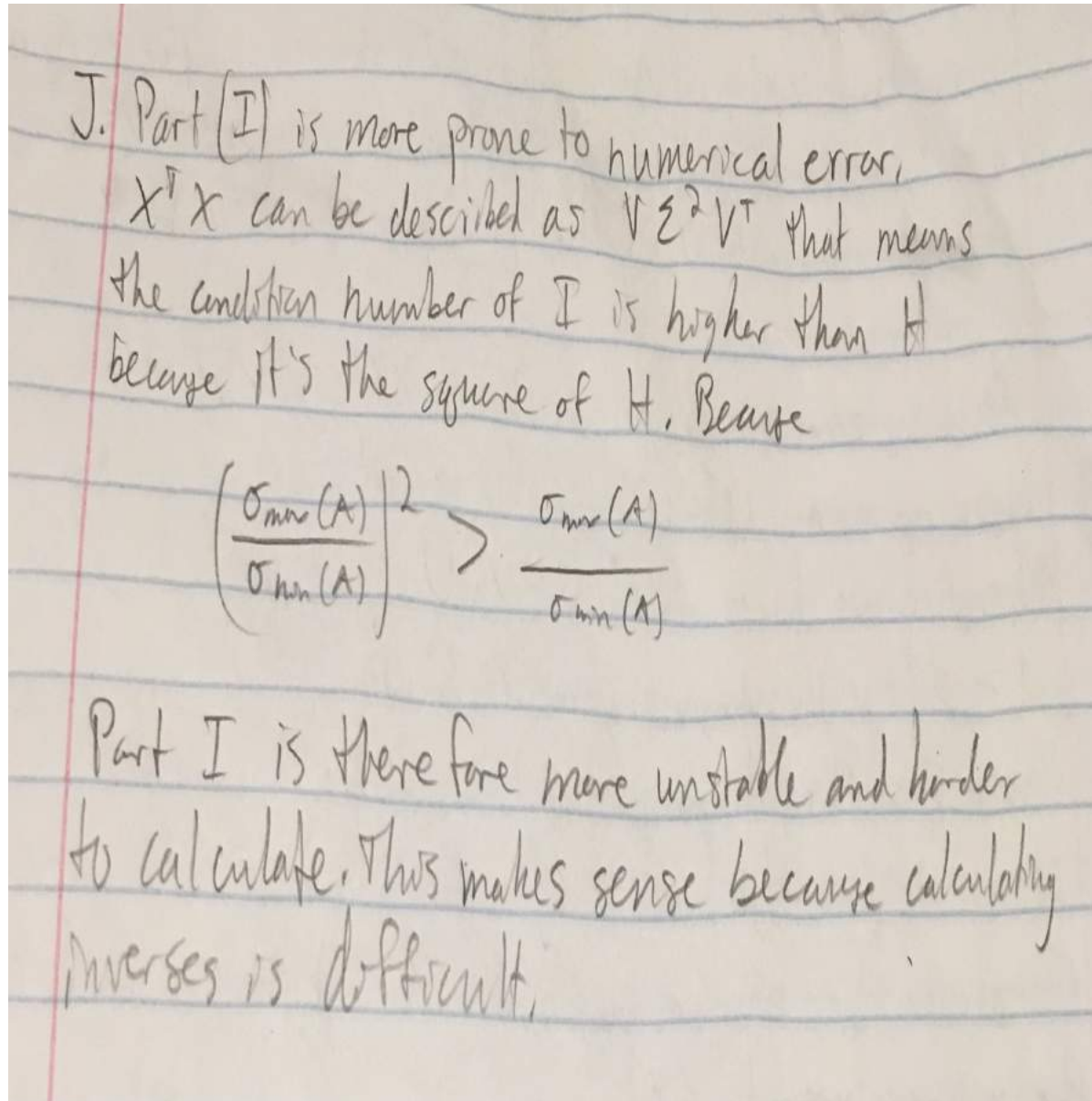
**Solution I:**

The image shows a handwritten solution on lined paper. It starts with the SVD decomposition of matrix  $X$  as  $X = U \Sigma V^T$ . Then, it derives the pseudoinverse  $(X^+ X)^{-1}$  by taking the inverse of both sides, resulting in  $(X^+ X)^{-1} = (V \Sigma^2 V^T)^{-1} = (V^T)^{-1} (\Sigma^2)^{-1} (V)^{-1}$ . Next, it multiplies this by  $X^T$  to get  $(X^+ X)^{-1} X^T = (V^T)^{-1} (\Sigma^2)^{-1} V^T \Sigma U^T$ . Finally, it simplifies this expression to  $(X^+ X)^{-1} X^T = V \Sigma^{-1} U^T$ , which is the desired result.

$$\begin{aligned} \text{I. } X &= U \Sigma V^T \\ (X^+ X)^{-1} &= (V \Sigma^2 V^T)^{-1} = (V^T)^{-1} (\Sigma^2)^{-1} (V)^{-1} \\ (X^+ X)^{-1} X^T &= (V^T)^{-1} (\Sigma^2)^{-1} V^T \Sigma U^T \\ (X^+ X)^{-1} X^T &= V \Sigma^{-1} U^T \text{ as desired} \end{aligned}$$

Problem J [2 points]:

Solution J:



## 2 Matrix Factorization [30 Points]

Problem A [5 points]:

Solution A:

B. To find the optimal  $u_i$  we set  $\partial u_i = 0$  and solve for  $u_i$ . So, we have

$$0 = \lambda u_i - \sum_j v_j (y_{ij} - u_i^T v_j)^T$$

$$\Rightarrow 0 = \lambda u_i - \sum_j v_j y_{ij}^T + \sum_j v_j v_j^T u_i$$

$$\sum_j v_j y_{ij} = u_i (\lambda I + \sum_j v_j v_j^T)$$

$$\Rightarrow u_i = (\lambda I + \sum_j v_j v_j^T)^{-1} (\sum_j y_{ij} v_j)$$

For optimal  $v_j$  we set  $\partial v_j = 0$

$$0 = \lambda v_j - \sum_i u_i (y_{ij} - u_i^T v_j)^T$$

$$0 = \lambda v_j - \sum_i u_i y_{ij}^T + \sum_i u_i u_i^T v_j$$

$$\sum_i u_i y_{ij} = v_j (\lambda I + \sum_i u_i u_i^T)$$

$$\Rightarrow v_j = (\lambda I + \sum_i u_i u_i^T)^{-1} (\sum_i u_i y_{ij})$$

Problem B [5 points]:

Solution B:

B. To find the optimal  $u_i$  we set  $\partial u_i = 0$  and solve for  $u_i$ . So, we have

$$0 = \lambda u_i - \sum_j v_j (y_{ij} - u_i^T v_j)^T$$
$$\Rightarrow 0 = \lambda u_i - \sum_j v_j y_{ij}^T + \sum_j v_j v_j^T u_i$$
$$\sum_j v_j y_{ij} = u_i (\lambda I + \sum_j v_j v_j^T)$$
$$\Rightarrow u_i = (\lambda I + \sum_j v_j v_j^T)^{-1} (\sum_j y_{ij} v_j)$$

For optimal  $v_j$  we set  $\partial v_j = 0$

$$0 = \lambda v_j - \sum_i u_i (y_{ij} - u_i^T v_j)^T$$
$$0 = \lambda v_j - \sum_i u_i y_{ij}^T + \sum_i u_i u_i^T v_j$$
$$\sum_i u_i y_{ij} = v_j (\lambda I + \sum_i u_i u_i^T)$$
$$\Rightarrow v_j = (\lambda I + \sum_i u_i u_i^T)^{-1} (\sum_i u_i y_{ij})$$



**Problem C [10 points]:**

**Solution C:** *See 2D.py and prob2utils.py for the solution code.*

**Problem D [5 points]:**

**Solution D:**

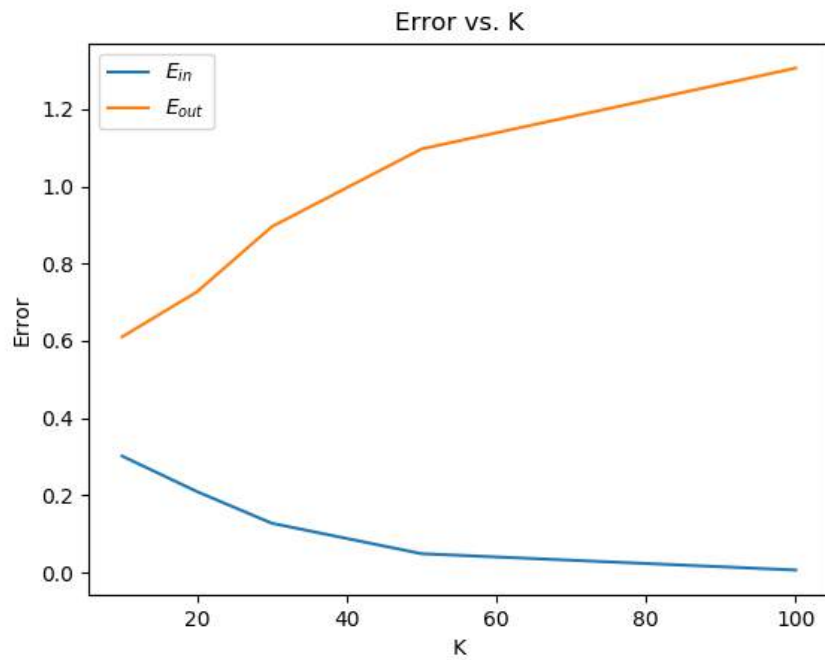


Figure 1: Unregularized factorization

*As  $K$  increases  $E_{out}$  goes up and  $E_{in}$  goes down. This makes sense because as the number of latent factors increases, the more the complex the model becomes. The model increases in complexity because  $K$  increases the dimensionality of the problem. This clearly leads to overfitting that we can see in the graph.*

**Problem E [5 points]:**

**Solution E:**

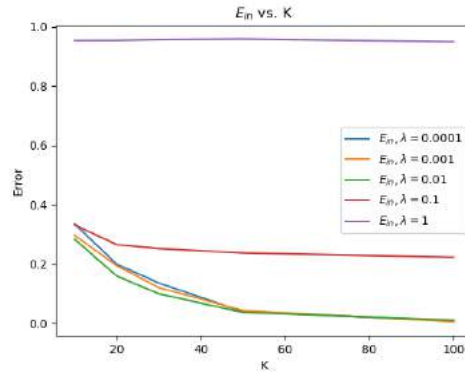


Figure 2:  $E_{in}$  vs  $k$  for different  $\lambda$

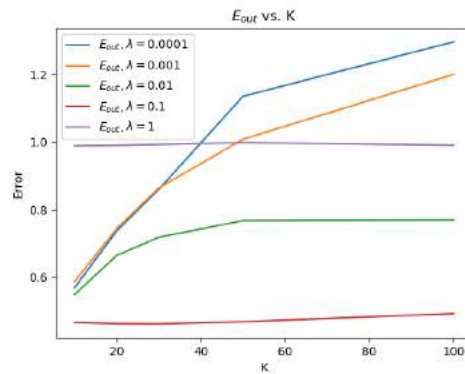


Figure 3:  $E_{out}$  vs  $k$  for different  $\lambda$

With the  $E_{in}$  graphs we see regularization of 1 causes underfitting and that the rest see a decrease in  $E_{in}$  as the complexity of the model increases. With  $E_{out}$  we see that regularization of 1 is super underfit and that regularization of .1 produces the best results. As regularization decreases we get a plot that looks more and more similar to the  $E_{out}$  from the previous part which was unregularized. For regularization under .1 we see that as the number of latent factors increases overfitting starts to occur.

### 3 Word2Vec Principles [35 Points]

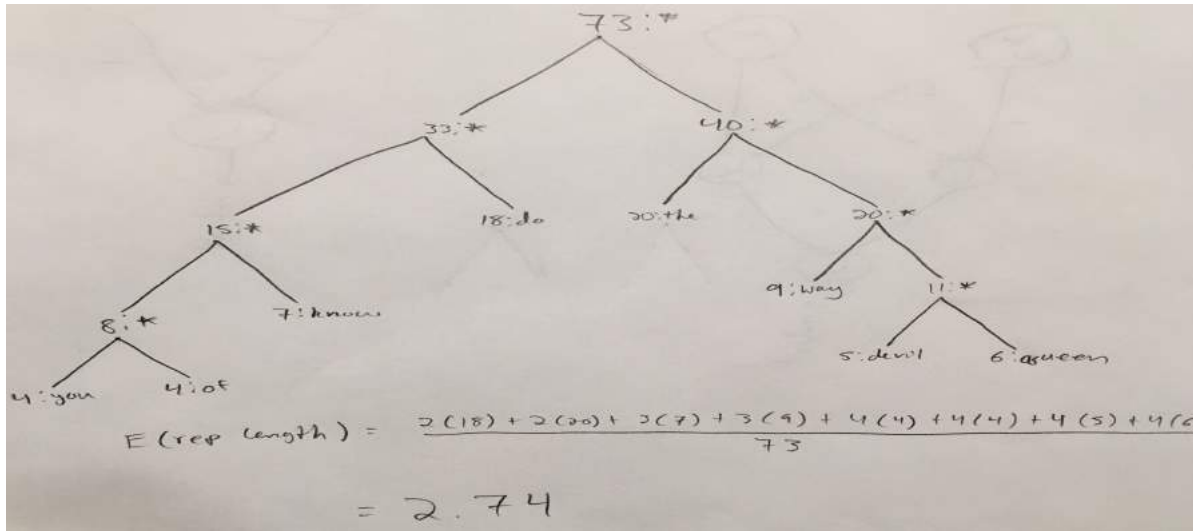
**Problem A** [5 points]:

**Solution A:** *Computing these gradients scales approximately linearly with  $W$  and  $D$ . Each new word will be an input word once and an output approximately  $2 \cdot D$  times. So  $2 \cdot D$  more pairs will get made. With an increase in  $D$  we see that it will create  $2 \cdot W$  more pairs because it increases the window by one on each side of the input word.*

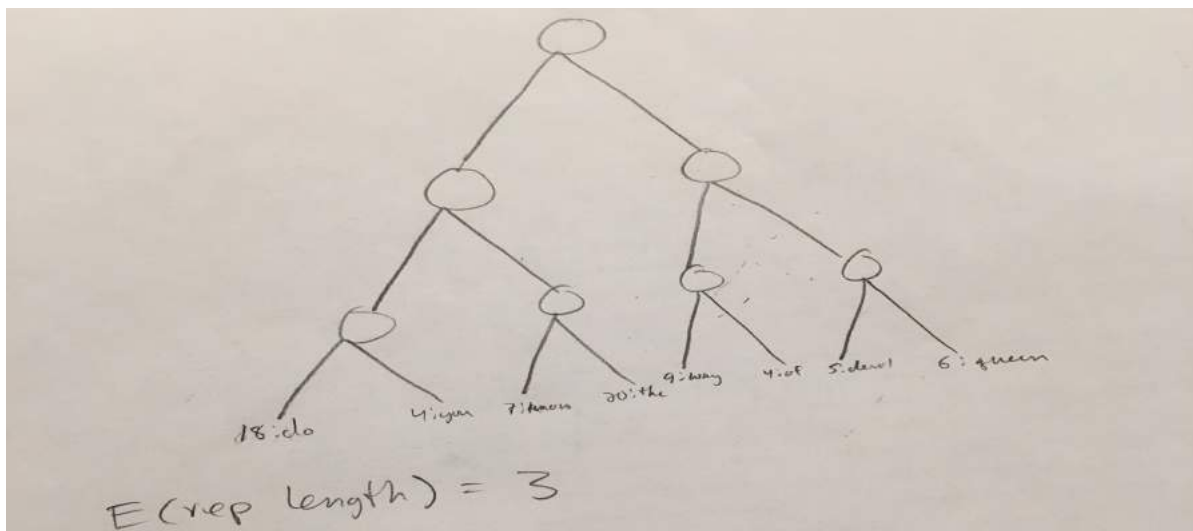


**Problem B [10 points]:**

**Solution B: Huffman tree**



**Binary tree**



**Problem C [3 points]:**

**Solution C:** *I would expect the value of the training objective to increase as  $D$  increases. Too large of a  $D$  would be computationally expensive and inefficient and could lead to overfitting which is undesirable.*

**Problem D [10 points]:**

**Solution D:** *See solution code in P3C.py*

**Problem E [2 points]:**

<b>Solution E:</b> $308 \times 10$
------------------------------------



**Problem F [2 points]:**

<b>Solution F:</b> $10 \times 308$
------------------------------------

**Problem G [1 points]:**

**Solution G:**

```
Pair(green, eggs), Similarity: 0.99271554
Pair(eggs, green), Similarity: 0.99271554
Pair(likes, drink), Similarity: 0.98298174
Pair(drink, likes), Similarity: 0.98298174
Pair(box, goat), Similarity: 0.9789746
Pair(goat, box), Similarity: 0.9789746
Pair(fox, goat), Similarity: 0.976961
Pair(boat, goat), Similarity: 0.9768292
Pair(some, slow), Similarity: 0.974435
Pair(slow, some), Similarity: 0.974435
Pair(there, here), Similarity: 0.97167295
Pair(here, there), Similarity: 0.97167295
Pair(tree, goat), Similarity: 0.9714184
Pair(samiam, anywhere), Similarity: 0.9709146
Pair(anywhere, samiam), Similarity: 0.9709146
Pair(brush, comb), Similarity: 0.9708585
Pair(comb, brush), Similarity: 0.9708585
Pair(wink, drink), Similarity: 0.9693825
Pair(not, box), Similarity: 0.96915597
Pair(five, eight), Similarity: 0.96764493
Pair(eight, five), Similarity: 0.96764493
Pair(mouse, fox), Similarity: 0.9672581
Pair(them, or), Similarity: 0.96560454
Pair(or, them), Similarity: 0.96560454
Pair(ham, green), Similarity: 0.9652538
Pair(blue, old), Similarity: 0.9648802
Pair(old, blue), Similarity: 0.9648802
Pair(heads, grows), Similarity: 0.9645252
Pair(grows, heads), Similarity: 0.9645252
Pair(in, dark), Similarity: 0.96441853
```

**Problem H [2 points]:**

**Solution H:** *One pattern is that if  $(x, y)$  shows up as similar then it is very likely that  $(y, x)$  also shows up as a pair right after it. Another thing to notice is that common and well know Dr. Seuss phrases like green eggs and anywhere samiam occur in this set. Some odd pairs like (blue, old) appear which could be a function of it being Dr. Seuss writing again.*