

(2a) By definition of Frobenius norm from last slide 8.

$$\text{we know } \|U\|_F^2 = \sum_{ij} u_{ij}^2 \quad \& \quad \|V\|_F^2 = \sum_{ij} v_{ij}^2$$

$$\therefore \partial u_i \left(\frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2) + \frac{1}{2} \sum_j (y_{ij} - U^T v_i)^2 \right)$$

$$= \partial u_i \left(\frac{\lambda}{2} \left(\sum_j u_{ij}^2 + \sum_j v_{ij}^2 \right) + \frac{1}{2} \sum_j (y_{ij} - U^T v_i)^2 \right)$$

$$= \frac{\lambda}{2} \left(2u_i + 0 \right) + \frac{1}{2} \sum_j (2(y_{ij} - U^T v_i) (-v_i))$$

$$= \lambda u_i - \sum_j (u_{ij} - U^T v_i)^2 v_i$$

and

$$\therefore \partial v_i \left(\frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2) + \frac{1}{2} \sum_j (y_{ij} - U^T v_i)^2 \right)$$

$$= \partial v_i \left(\frac{\lambda}{2} \left(\sum_j u_{ij}^2 + \sum_j v_{ij}^2 \right) + \frac{1}{2} \sum_j (y_{ij} - U^T v_i)^2 \right)$$

$$= \frac{\lambda}{2} (0 + 2v_i) + \frac{1}{2} \sum_j (2(y_{ij} - U^T v_i)^2 (-u_i))$$

$$= \lambda v_i - \sum_j (u_{ij} - U^T v_i)^2 u_i$$

(2B) we'll find the closed form each partial derivative when each is equal to zero as that is our "optimal" v_i & u_i (the "critical point")

$$\therefore \partial u_i \left(\frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2) + \frac{1}{2} \sum_j (y_{ij} - U^T v_i)^2 \right) = 0$$

$$= \lambda u_i - \sum_j (u_{ij} - U^T v_i)^2 v_i = 0 \quad \therefore \lambda u_i = \sum_j (u_{ij} - U^T v_i)^2 v_i$$

$$\rightarrow \lambda u_i = \sum_j u_{ij} v_i - \sum_j u_i v_j v_i \rightarrow \lambda u_i + \sum_j u_i v_j v_i = \sum_j u_{ij} v_i$$

$$\rightarrow u_i (\lambda + \sum_j v_j v_i) = \sum_j u_{ij} v_i \quad \therefore u_i = \frac{\sum_j u_{ij} v_i}{(\lambda + \sum_j v_j v_i)}$$

and for v_i :

$$= \partial v_i \left(\frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2) + \frac{1}{2} \sum_j (y_{ij} - U^T v_i)^2 \right) = 0$$

$$= \lambda v_i - \sum_j (u_{ij} - U^T v_i)^2 u_i = 0 \quad \therefore \lambda v_i = \sum_j (u_{ij} - U^T v_i)^2 u_i$$

$$\rightarrow \lambda v_i = \sum_j u_{ij} u_i - \sum_j u_i u_j u_i \rightarrow \lambda v_i = \sum_j u_{ij} u_i - \sum_j u_i u_j u_i$$

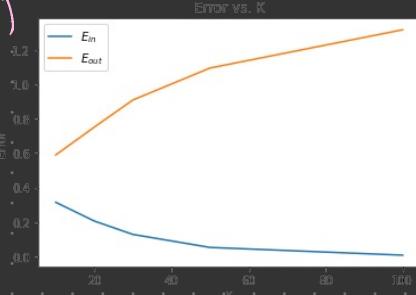
$$\rightarrow \lambda v_i + \sum_j u_i u_j u_i = \sum_j u_{ij} u_i \rightarrow v_i (\lambda + \sum_j u_i u_j) = \sum_j u_{ij} u_i$$

$$\therefore v_i = \frac{\sum_j u_{ij} u_i}{(\lambda + \sum_j u_i u_j)}$$

(2C)

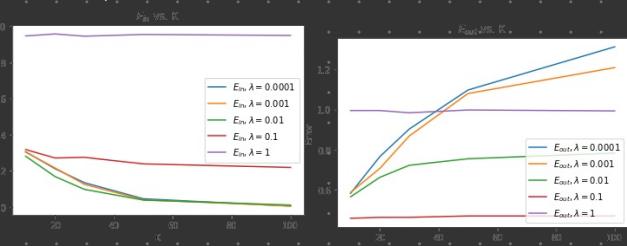
See Code For Prob2

(2d)



The E_{out} (out of sample error) is significantly higher than the E_{in} (in sample error) for all K values and as K values increase E_{in} decreases while E_{in} increases meaning we have overfitted with all models which makes sense as our regularization parameter is set to 0. This lack of regularization means our models simply memorize the training set hence fitting noise and making them bad generalizers and the higher the K value, the greater this overfitting is as we are increasing model complexity hence allowing more noise to be fit in our training set.

(2E)



The trend is the same as in the previous part for most Lambda values with the E_{out} being generally higher than E_{in} as a result of overfitting and this overfitting worsens for some of the lambda values (10^{***-4} , 10^{***-3} , 10^{***-2}) as K increases, making them weak regularizers. For Lambda equals 1, the regularization term has too much weight therefore all our models want to only minimize U and V norms, yielding relatively the same predictions for all inputs and poor performance for both training and testing error hence increasing K barely affects the models therefore consistently performing poorly for all K values as the both the E_{in} and E_{out} are consistently extremely high. Lambda equals 0.1 is our best regularizer because it yields the lowest E_{out} and as our values for K increase, E_{in} continually slightly decreases and E_{out} seems to slightly decrease before stagnating to stay consistently low hence avoiding increased overfitting as seen for all other values.

$$\begin{aligned}
 (3a) \quad \log P(w_0|w_I) &= \log \left(\frac{\exp(V_{w_0}^T V_{w_I})}{\sum_{w \in I} \exp(V_w^T V_{w_I})} \right) = \log \left(\frac{\exp(V_{w_0}^T V_{w_I})}{\sum_{w \in I} \exp(V_w^T V_{w_I})} \right) \\
 &= \log(C^{V_{w_0}^T V_{w_I}}) - \log \left(\sum_{w \in I} C^{V_w^T V_{w_I}} \right) \\
 &= V_{w_0}^T V_{w_I} - \log \left(\sum_{w \in I} C^{V_w^T V_{w_I}} \right)
 \end{aligned}$$

Now to find the gradient:

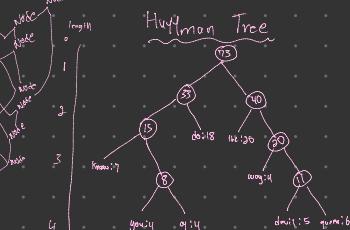
$$\nabla_{V_{w_0}} \log P(w_0|w_I) = V_{w_I} - \frac{V_{w_I} C^{V_{w_0}^T V_{w_I}}}{\sum_{w \in I} C^{V_w^T V_{w_I}}}$$

$$\nabla_{V_{w_I}} \log P(w_0|w_I) = V_{w_0} - \frac{V_{w_0} C^{V_w^T V_{w_I}}}{\sum_{w \in I} C^{V_w^T V_{w_I}}}$$

The time complexity of $V_{w_0}^T V_{w_I}$ is $O(D)$ as we are taking the dot product for all D dimensions w.r.t. the summation $\sum C^{V_w^T V_{w_I}}$. The complexity of $O(D^2)$ as we are calculating $V_w^T V_{w_I}$ many times. For a single w_0, w_I pair the time complexity is $O(D^2)$.

(3B)

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



The expected length of the binary tree is 3, the depth of our tree:

$$\begin{aligned}
 \text{The expected length of the Huffman tree is } 3, \text{ the depth of our tree:} \\
 \text{The expected length of the binary tree is } 3, \text{ the depth of our tree:} \\
 \text{Pair(fear, take), Similarity: 0.9954883} \\
 \text{Pair(take, fear), Similarity: 0.9954883} \\
 \text{Pair(feet, low), Similarity: 0.99465114} \\
 \text{Pair(low, feet), Similarity: 0.99465114} \\
 \text{Pair(called, sleep), Similarity: 0.98782295} \\
 \text{Pair(sleep, called), Similarity: 0.98782295} \\
 \text{Pair(head, wave), Similarity: 0.98776114} \\
 \text{Pair(wave, head), Similarity: 0.98776114} \\
 \text{Pair(made, only), Similarity: 0.9875726} \\
 \text{Pair(only, made), Similarity: 0.9875726} \\
 \text{Pair(little, took), Similarity: 0.9871121} \\
 \text{Pair(took, little), Similarity: 0.9871121} \\
 \text{Pair(teeth, boat), Similarity: 0.9870644} \\
 \text{Pair(boat, teeth), Similarity: 0.9870644} \\
 \text{Pair(ned, pink), Similarity: 0.9868588} \\
 \text{Pair(pink, ned), Similarity: 0.9868588} \\
 \text{Pair(but, grow), Similarity: 0.98576736} \\
 \text{Pair(grow, but), Similarity: 0.98576736} \\
 \text{Pair(story, low), Similarity: 0.9852524} \\
 \text{Pair(it, foot), Similarity: 0.9851571} \\
 \text{Pair(foot, it), Similarity: 0.9851571} \\
 \text{Pair(dad, thin), Similarity: 0.98467755} \\
 \text{Pair(thin, dad), Similarity: 0.98467755} \\
 \text{Pair(fat, ride), Similarity: 0.98459214} \\
 \text{Pair(ride, fat), Similarity: 0.98459214} \\
 \text{Pair(too, take), Similarity: 0.98443496} \\
 \text{Pair(cant, way), Similarity: 0.9842897} \\
 \text{Pair(way, cant), Similarity: 0.9842897} \\
 \text{Pair(nine, boat), Similarity: 0.9842129} \\
 \text{Pair(cans, yink), Similarity: 0.984207}
 \end{aligned}$$

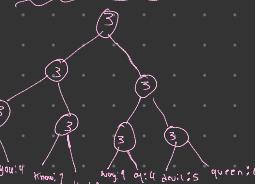
(3C)

As D increases, we expect the value of the training objectives to increase. As an increased D allows for better model selection, hence allowing us to fit the training data better. However, we might not want to use a very large D in order to avoid overly complex models that fit noise resulting in overfitting. This increased in dimensionality could also be very computationally expensive hence we deal with a very large D .

(3E) The dimensions are 308x10 showing the number of words and embedding dimensions.

(3F) The dimensions are 10x308.

Balanced Binary Tree:



$$\begin{aligned}
 \text{The expected length of the binary tree is } 3, \text{ the depth of our tree:} \\
 \text{The expected length of the binary tree is } 3, \text{ the depth of our tree:} \\
 \text{Pair(fear, take), Similarity: 0.9954883} \\
 \text{Pair(take, fear), Similarity: 0.9954883} \\
 \text{Pair(feet, low), Similarity: 0.99465114} \\
 \text{Pair(low, feet), Similarity: 0.99465114} \\
 \text{Pair(called, sleep), Similarity: 0.98782295} \\
 \text{Pair(sleep, called), Similarity: 0.98782295} \\
 \text{Pair(head, wave), Similarity: 0.98776114} \\
 \text{Pair(wave, head), Similarity: 0.98776114} \\
 \text{Pair(made, only), Similarity: 0.9875726} \\
 \text{Pair(only, made), Similarity: 0.9875726} \\
 \text{Pair(little, took), Similarity: 0.9871121} \\
 \text{Pair(took, little), Similarity: 0.9871121} \\
 \text{Pair(teeth, boat), Similarity: 0.9870644} \\
 \text{Pair(boat, teeth), Similarity: 0.9870644} \\
 \text{Pair(ned, pink), Similarity: 0.9868588} \\
 \text{Pair(pink, ned), Similarity: 0.9868588} \\
 \text{Pair(but, grow), Similarity: 0.98576736} \\
 \text{Pair(grow, but), Similarity: 0.98576736} \\
 \text{Pair(story, low), Similarity: 0.9852524} \\
 \text{Pair(it, foot), Similarity: 0.9851571} \\
 \text{Pair(foot, it), Similarity: 0.9851571} \\
 \text{Pair(dad, thin), Similarity: 0.98467755} \\
 \text{Pair(thin, dad), Similarity: 0.98467755} \\
 \text{Pair(fat, ride), Similarity: 0.98459214} \\
 \text{Pair(ride, fat), Similarity: 0.98459214} \\
 \text{Pair(too, take), Similarity: 0.98443496} \\
 \text{Pair(cant, way), Similarity: 0.9842897} \\
 \text{Pair(way, cant), Similarity: 0.9842897} \\
 \text{Pair(nine, boat), Similarity: 0.9842129} \\
 \text{Pair(cans, yink), Similarity: 0.984207}
 \end{aligned}$$

(3G)

(3H) We notice that most pairs occur twice s.t. for words x, y then we will have the pair (x, y) and (y, x) show up twice next to each other with the same similarity percentage as our model associates words with their neighbors in both directions.