

1. SVD and PCA:

Question A:

$$\begin{aligned}XX^T &= Y^T Y \\&= (U \Sigma V^T)^T (U \Sigma V^T) \\&= V \Sigma U^T U \Sigma V^T \\&= V \Sigma^2 V^T\end{aligned}$$

The Σ^2 in SVD is the Σ in PCA, so the columns of V are the PCA of X .

Question B:

$$\begin{aligned}(i) \quad \|A - B\|_F &= \text{trace}((A - B)^T (A - B)) \\&= \text{trace}(V^T (A - B)^T U U^T (A - B) V) \\&= \|U^T (A - B) V\|_F \\&= \|\Sigma - U^T B V\|_F \\&= \sum_{i,j} (\Sigma - U^T B V)_{ij}^2\end{aligned}$$

the ranks of U^T and V are longer than B , so $\text{rank}[U^T B V] = k$

$$\Sigma = \begin{pmatrix} s_1 & & & \\ & s_2 & & \\ & & \dots & \\ & & & s_k \\ & & & & 0 \end{pmatrix}$$

To minimize $\sum_{i,j} (\Sigma - U^T B V)_{ij}^2$, it is apparent that $U^T B V$ should be

$$\begin{pmatrix} s_1 & & & \\ & s_2 & & \\ & & \dots & \\ & & & s_k \\ & & & & 0 \end{pmatrix}$$

$$U^T B V = \sum_{i=1:k} I_k \Sigma I_k$$

$$B = (U I_k) \Sigma (I_k V^T)$$

$$= U_{1:k} \Sigma V_{1:k}^T$$

$$= A_k$$

$$\text{so, } \|A - A_k\|_F = \min_{\text{rank}(B) \leq k} \|A - B\|_F$$

(ii) We first prove $\|Y\|_* \geq \min_{Y=AB^T} \frac{1}{2} (\|A\|_F^2 + \|B\|_F^2)$, $Y = U \Sigma V^T$

$$\text{let } A = U \sqrt{\Sigma}, \quad B = V \sqrt{\Sigma}$$

$$\min_{Y=AB^T} \frac{1}{2} (\|A\|_F^2 + \|B\|_F^2)$$

$$\leq \frac{1}{2} (\|U \sqrt{\Sigma}\|_F^2 + \|V \sqrt{\Sigma}\|_F^2)$$

$$= \frac{1}{2} (\text{trace}(\sqrt{\Sigma} U^T U \sqrt{\Sigma}) + \text{trace}(\sqrt{\Sigma} V^T V \sqrt{\Sigma}))$$

$$= \frac{1}{2} (\text{trace}(\Sigma) + \text{trace}(\Sigma))$$

$$= \text{trace}(\Sigma)$$

$$= \|Y\|_*$$

Then we prove $\|Y\|_* \leq \min_{Y=AB^T} \frac{1}{2} (\|A\|_F^2 + \|B\|_F^2)$

$$\text{let } A' = U^T A, \quad B' = V^T B$$

$$A' B'^T = U^T A B V = U^T U \Sigma V^T V = \Sigma$$

$$\|A'\|_F^2 = \text{trace}(A'^T A') = \text{trace}(A^T U U^T A)$$

$$= \text{trace}(A^T A) = \|A\|_F^2$$

$$\|B'\|_F^2 = \|B\|_F^2$$

$$\text{So, } \min_{Y=AB^T} \frac{1}{2} (\|A\|_F^2 + \|B\|_F^2)$$

$$= \min_{\Sigma=A'B'^T} \frac{1}{2} (\|A'\|_F^2 + \|B'\|_F^2)$$

$$= \min_{\Sigma=A'B'^T} \frac{1}{2} (\text{trace}(A'A'^T) + \text{trace}(B'B'^T))$$

$$= \min_{\Sigma=A'B'^T} \frac{1}{2} \sum_i (\|A'_i\|^2 + \|B'_i\|^2)$$

$$\geq \min_{\Sigma=A'B'^T} \frac{1}{2} \sum_i (\|A'_i B'^T_i\|)$$

$$= \min_{\Sigma=A'B'^T} \sum_i (\sqrt{B'_i A'^T_i A'_i B'^T_i})$$

$$= \min_{\Sigma=A'B'^T} \sum_i (\sqrt{(A'B'^T)_{ii} (A'B'^T)_{ii}})$$

$$= \min_{\Sigma=A'B'^T} \sum_i (\sqrt{s_i^2})$$

$$= \min_{\Sigma=A'B'^T} \sum_i s_i$$

$$= \text{trace}(\Sigma)$$

$$= \|Y\|_*$$

$$\text{In Conclusion: } \|Y\|_* = \min_{Y=AB^T} \frac{1}{2} (\|A\|_F^2 + \|B\|_F^2)$$

2. Matrix Factorization:

Question A:

$$\frac{\partial}{\partial u_i} = \frac{\lambda}{N} u_i - (y_{ij} - u_i^T v_j) v_j$$

$$\frac{\partial}{\partial v_j} = \frac{\lambda}{N} v_j - (y_{ij} - u_i^T v_j) u_i$$

Question B:

$$\frac{\partial}{\partial u_i} \left(\frac{\lambda}{2} (\|u\|_F^2 + \|v\|_F^2) + \sum_{i,j} (y_{ij} - u_i^T v_j)^2 \right) = 0$$

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

$$u_i = (\lambda I_k + \sum_j v_j v_j^T)^{-1} \left(\sum_j v_j y_{ij} \right)$$

$$\frac{\partial}{\partial v_j} \left(\frac{\lambda}{2} (\|u\|_F^2 + \|v\|_F^2) + \sum_{i,j} (y_{ij} - u_i^T v_j)^2 \right) = 0$$

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

$$v_j = (\lambda I_k + \sum_i u_i u_i^T)^{-1} \left(\sum_i u_i y_{ij} \right)$$

Question C:

Question D:

~~with different latent factor k, E~~

Question E:

in "problem-2-plot.pdf"

3. Word2Vec:

Question A:

$$p(w_o | w_I) = \frac{\exp(v_{w_o}'^T v_{w_I})}{\sum_{w=1}^W \exp(v_{w_o}'^T v_{w_I})}$$

computing $\nabla \log p(w_o | w_I)$, the cost is proportional to W , which is always large ($10^5 \sim 10^7$ terms).

Question B:

when D increases to we using large D , our learning process may overfitting. And what happens to training objective?

- (i) The loss of input training objective will decrease;
- (ii) The loss of output training objective will increase due to overfitting.

Question C:

Yes. When using hierarchical softmax, the cost of computing $\log p(w_o | w_I)$ and $\nabla \log p(w_o | w_I)$ is proportional to $L(w_o)$, which on average is no greater than $\log W$. (the length of the path)
compared with Question A, it does improve.

Question D:

- (i) if W words are stored in binary tree, which are all leaves, There are $W-1$ internal nodes ^{does} the tree have.
(prove): let n be the depth of the tree

$$2^n = W$$

$$\begin{aligned} \text{then internal nodes: } 1 + 2^1 + 2^2 + \dots + 2^{n-1} &= 2^n - 1 \\ &= W - 1 \end{aligned}$$

- (ii) The hierarchical softmax formation has one representation v_w for each word w and one representation v_n' for every inner node n of the binary tree.

As we come up with the outcome of (i),

we have W words and $W-1$ internal nodes,

The hierarchical softmax model represent W v_w and $(W-1)$ v_n' .

In Softmax model where it assigns ~~two~~ two representations v_w and v_w' to each words. If we have W words, we will have W v_w and $W \times v_w'$.

Thus, the hierarchical softmax model does not require representing prohibitively many more vectors than the basic softmax model.

Question E:

- (i) we want to verify $\sum_{w=1}^W p(w|w_I) = 1$

$$\sum_{w=1}^W p(w|w_I) =$$

$$\sum_{w=1}^W \prod_{j=1}^{L(w)-1} \sigma([n(w, j+1) = ch(n(w, j))]) \cdot v_{n(w, j)}'^T v_{w_I}$$

$$\text{where as we have } [X] = \begin{cases} 1 & X \text{ is true} \\ -1 & X \text{ is not true} \end{cases}$$

we can divide whether $\underbrace{n(w, j+1) = ch(n(w, j))}_A$ is true or not, then $[A]$

$$= \begin{cases} 1 & A \text{ is true} \\ -1 & A \text{ is not true.} \end{cases}$$

then $\sigma([A]) \cdot v_{n(w, j)}'^T v_{w_I}$ can be divided to 2 conditions.

when A is true:

$$\sigma(v_{n(w,j)}'^T v_{w_I})$$

when A is not true:

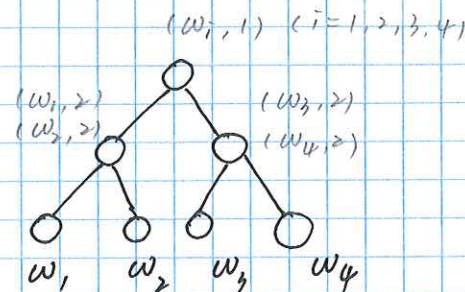
$$\sigma(-v_{n(w,j)}'^T v_{w_I})$$

we already know $\sigma(x) + \sigma(-x) = 1$

$$\begin{aligned} \sum_{w=1}^W p(w|w_I) &= \sum_{w=1}^W \prod_{j=1}^{L(w)-1} \sigma([n(w,j+1) = \text{ch}(n(w,j))]) \cdot v_{n(w,j)}'^T v_{w_I} \\ &= \sum_{w=1}^W [\sigma(v_{n(w_I,1)}'^T v_{w_I}) + \sigma(-v_{n(w_I,1)}'^T v_{w_I})] \\ &= \sigma(v_{n(w_I,1)}'^T v_{w_I}) + \sigma(-v_{n(w_I,1)}'^T v_{w_I}) \\ &= 1 \end{aligned}$$

Therefore, the hierarchical softmax defines a distribution.

(ii)



Set $\text{ch}(\cdot) = \text{left child}$, we have:

$$\begin{aligned} p(w_3|w_1) &= \sigma(\frac{1}{2}) \cdot \sigma(-\frac{1}{3}) \\ &= 0.2199 \end{aligned}$$

Question F: if the feature vectors cannot be initialized to all zeros.
we can't use the gradient descent methods as the
"gradient" of initial state are all the same and cannot
change to descending trend.

Question G:

we have $P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$

we want to prove if we know $f(w_i) > f(w_j)$,

$$\mathbb{E}[f(w_i)] \geq \mathbb{E}[f(w_j)]$$

we know:

$$\mathbb{E}[x] = \int x P(x) dx$$

therefore:

$$\begin{aligned}\mathbb{E}[f(w_i)] &= \int f(w_i) P(w_i) dw_i \\ &= \int f(w_i) \left(1 - \sqrt{\frac{t}{f(w_i)}}\right) dw_i\end{aligned}$$

we know, $f(w_i) > f(w_j)$

then,

$$1 - \sqrt{\frac{t}{f(w_i)}} > 1 - \sqrt{\frac{t}{f(w_j)}}$$

then,

$$f(w_i) \left(1 - \sqrt{\frac{t}{f(w_i)}}\right) > f(w_j) \left(1 - \sqrt{\frac{t}{f(w_j)}}\right)$$

therefore, we have

$\mathbb{E}(f(w_i)) > \mathbb{E}(f(w_j))$ ~~maintains~~ maintains, the rank of frequencies is preserved by this subsampling policy.