

Neural Networks and Deep Learning HW #3

1. Given $x \in \mathbb{R}^n$, $W \in \mathbb{R}^{m \times n}$ where $m < n$.

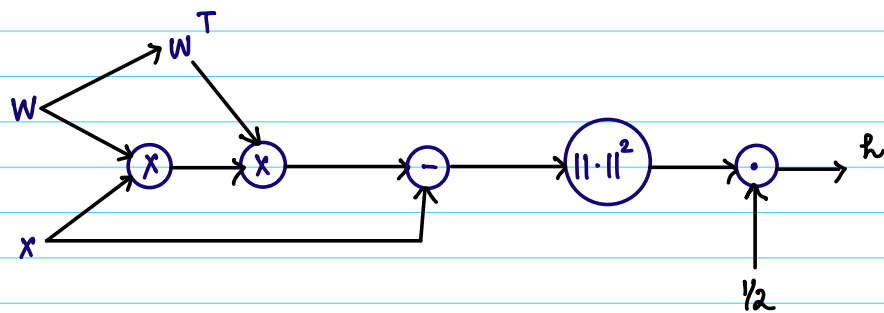
a) First let's consider the loss equation

$$L = \frac{1}{2} \|W^T W x - x\|^2$$

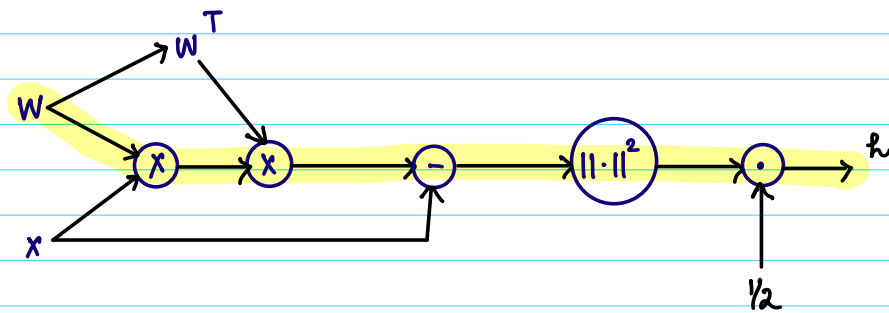
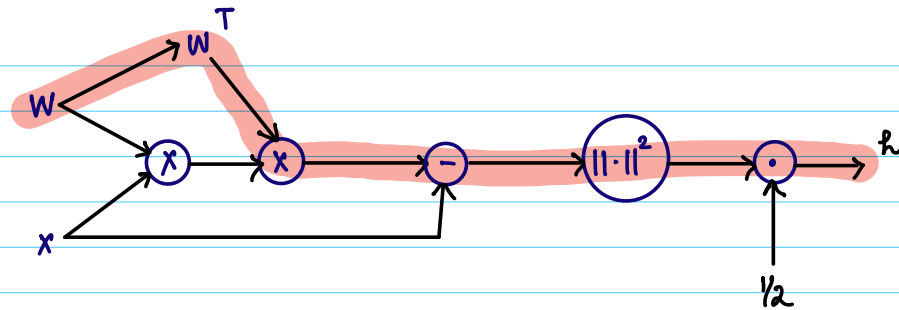
* We know that the loss L should be minimized, i.e. $W^T W x - x$ must be minimized.

* Wx is the hidden representation and for the loss to be minimized, Wx will have to preserve information about x .

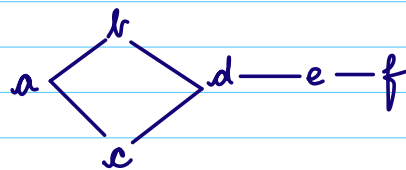
b) Computation graph:-



c) The two paths to W are highlighted below:



Since W contributes to two paths, the ∇_W^h will be expressed as a sum of derivatives along each path. let's consider this case:



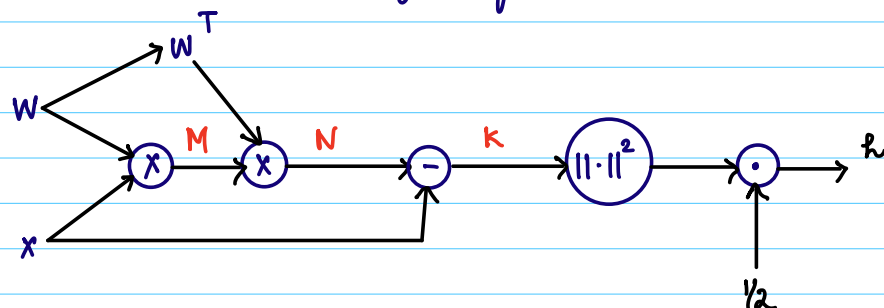
This graph represents the above scenario and consists of two paths to reach a .

$$a \rightarrow b \rightarrow d \rightarrow e \rightarrow f$$

$$a \rightarrow c \rightarrow d \rightarrow e \rightarrow f$$

$$\begin{aligned} \text{hence } \frac{\partial f}{\partial a} &= \frac{\partial f}{\partial e} \cdot \frac{\partial e}{\partial d} \cdot \frac{\partial d}{\partial b} \cdot \frac{\partial b}{\partial a} \\ &\quad + \frac{\partial f}{\partial e} \cdot \frac{\partial e}{\partial d} \cdot \frac{\partial d}{\partial c} \cdot \frac{\partial c}{\partial a} \end{aligned}$$

d) Let's consider the computation graph from 8):



* $K = W^T W x - x$
 where $W^T = n \times m$
 $W = m \times n$
 $x = n \times 1$

Hence $K = W^T W x - x \in \mathbb{R}^n$

* $N = W^T M$
 $N = W^T W x$
 where
 $W^T = n \times m$
 $W x = m \times 1$

Hence $N = W^T M \in \mathbb{R}^n$

* $M = W x$
 where $W = m \times n$
 $x = n \times 1$
 Hence $M = W x \in \mathbb{R}^m$

* $h = \frac{1}{2} \|K\|^2$, $\partial h / \partial K = K \rightarrow \textcircled{0}$

* $K = N - x$, $\partial h / \partial N = \partial K / \partial N \times \partial h / \partial K$ (chain rule)

We know that $\partial h / \partial K = K$

$\partial K / \partial N = \partial (N - x) / \partial N = I$

Hence $\partial h / \partial N = K \rightarrow \textcircled{1}$

We could also say that $\partial h / \partial K = \partial h / \partial N$ as $-$ operator distributes the gradient.

$$* \quad N = W^T M \quad \text{where } M = Wx \in \mathbb{R}^m$$

$$\partial h / \partial W^T = \partial h / \partial N \cdot \partial N / \partial W^T$$

$$\partial h / \partial N = k \in \mathbb{R}^n$$

$$\partial N / \partial W^T = \partial (W^T M) / \partial W^T = M^T = (Wx)^T \in \mathbb{R}^{1 \times m}$$

$$\text{hence } \partial h / \partial W^T = k (Wx)^T \longrightarrow \textcircled{2}$$

$$* \quad \partial h / \partial M = \partial h / \partial N \cdot \partial N / \partial M$$

$$\partial h / \partial N = k \in \mathbb{R}^n$$

$$\text{we know that } M = Wx \in \mathbb{R}^m, \text{ hence } \partial h / \partial M \in \mathbb{R}^m$$

$$\partial N / \partial M = \partial (W^T M) / \partial M = W \in \mathbb{R}^{m \times n}$$

$$\text{hence } \partial h / \partial M = Wk$$

$$* \quad \partial h / \partial W = \partial h / \partial M \times \partial M / \partial W$$

$$\partial h / \partial W \text{ must be } \mathbb{R}^{m \times n}$$

$$\partial h / \partial M = Wk \in \mathbb{R}^m, \quad \partial M / \partial W = \partial (Wx) / \partial W = x^T \in \mathbb{R}^{1 \times n}$$

$$\therefore \partial h / \partial W = WKX^T \in \mathbb{R}^{m \times n} \longrightarrow (3)$$

From (2) and (3), we know both contribute to ∇_W^h i.e

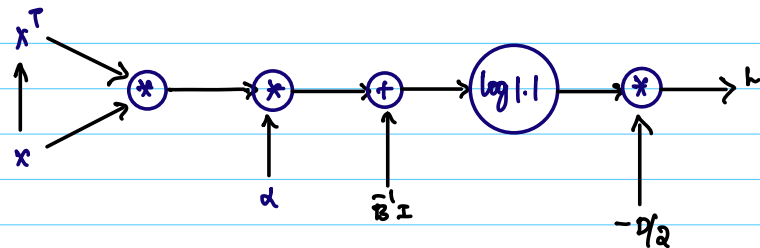
$$\nabla_W^h = \text{Backprop of } W + \text{Backprop of } W^T$$

Hence

$$\nabla_W^h = WKX^T + K(WX)^T = WKX^T + WXK^T$$

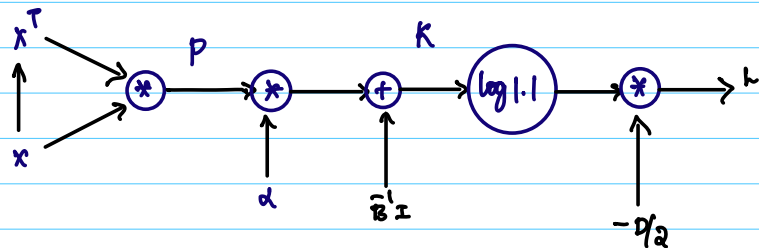
Question 2:

a) Computation graph for h_1 :



b) We need to compute $\partial h_1 / \partial x$

From the question, we know that $K = \alpha x x^T + \beta^{-1} I$



$$h_1 = -D/2 \log(|K|)$$

$$\partial h_1 / \partial K = -D/2 (K^{-1})^T = -D/2 (K^T)^{-1} \quad (\text{From the Matrix cookbook})$$

Now, we propagate backwards encountering a \oplus operation,

we know \oplus distributes the gradients and hence the gradient remains the same.

Now, we encounter a \otimes operator, hence the gradient @ P would be:

$$\partial h / \partial P = -\alpha D/2 \cdot (K^T)^{-1}$$

$$\partial h / \partial x x^T = -\alpha D/2 \cdot (K^T)^{-1}$$

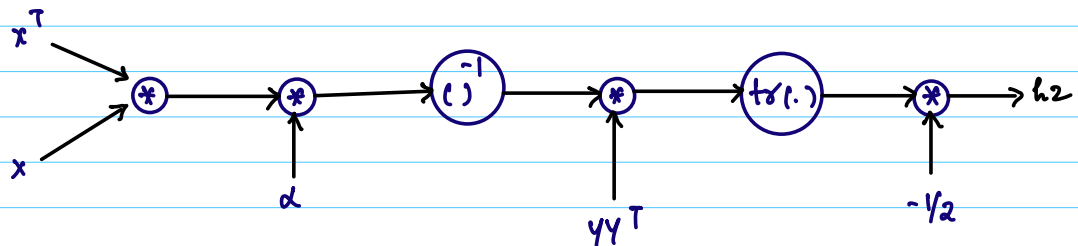
$$P = x x^T$$

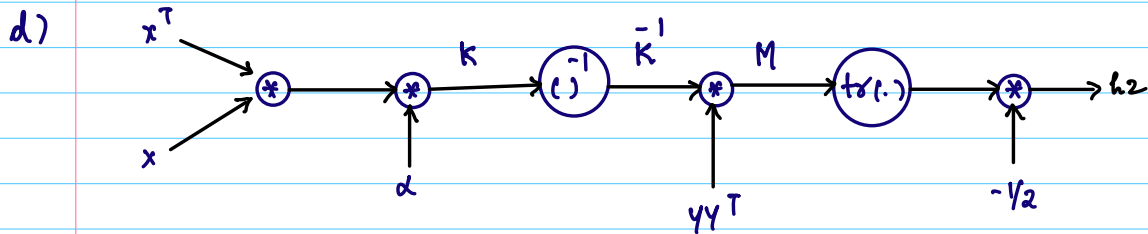
$$\partial h / \partial x = \partial P / \partial x \cdot \partial h / \partial P$$

$$\partial P / \partial x = x x^T = 2x$$

$$\therefore \partial h / \partial x = -\alpha D (K^T)^{-1} x$$

c) Computational graph for h_2 :





$$h_2 = -1/2 \operatorname{tr} \left((\alpha x x^T + \bar{K}^{-1} \mathbf{I})^{-1} y y^T \right)$$

$$h_2 = -1/2 \operatorname{tr} (M)$$

$$\partial h_2 / \partial M = -1/2 \mathbf{I} \quad (\text{From Matrix Cookbook})$$

$$\partial h_2 / \partial \bar{K}^{-1} = -1/2 \cdot y y^T$$

$$\partial h_2 / \partial K = 1/2 \bar{K}^{-1} y y^T \bar{K}^{-1}$$

Hence from the last part, we could say that:

$$\partial h_2 / \partial x = \alpha \bar{K}^{-1} y y^T \bar{K}^{-1} x$$

e) Hence $\partial h / \partial x = \partial h_1 / \partial x + \partial h_2 / \partial x$

$$\partial h / \partial x = \alpha \bar{K}^{-1} y y^T \bar{K}^{-1} x - \alpha D K^{-1} x$$