

ANEEK DAS - CSC2515 - HW4.

Q1. Multilayer Perceptron

i/p's $\rightarrow (x_1, x_2)$; x_1, x_2 are scalars.

o/p's $\rightarrow (y_1, y_2)$; $y_1 = \min(x_1, x_2)$
 $y_2 = \max(x_1, x_2)$

All hidden units \rightarrow ReLU activation ; $\text{ReLU}(x) = \max(0, x)$.

All o/p units \rightarrow Linear / no activation.

$$y_1 = \min(x_1, x_2)$$

$$\begin{aligned}\text{Now, } \min(x_1, x_2) &= -\max(0, x_1 - x_2) + x_1 \\ &= -\text{ReLU}(x_1 - x_2) + x_1.\end{aligned}$$

Now, we have to represent x_1 in terms of ReLU.

- If x is 've' we can represent it as $\max(0, x) = \text{ReLU}(x)$
- If x is '-ve' we can represent it as $-\max(0, -x) = -\text{ReLU}(-x)$

So, $y_1 \rightarrow \boxed{\min(x_1, x_2) = -\text{ReLU}(x_1 - x_2) + \text{ReLU}(x_1) - \text{ReLU}(-x_1)}$

$$\text{Now, } y_2 = \max(x_1, x_2)$$

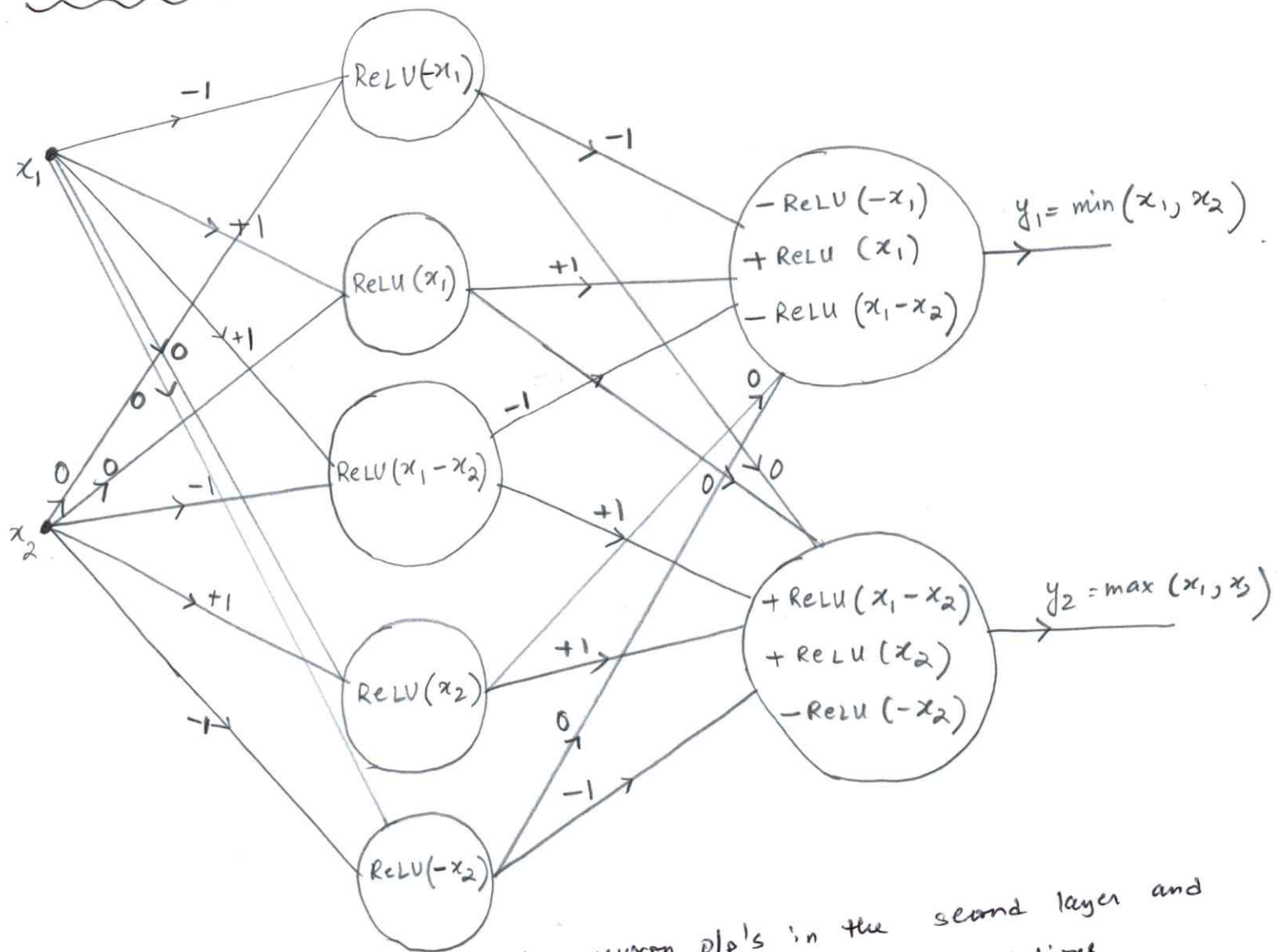
$$\max(x_1, x_2) = \max(0, x_1 - x_2) + x_2$$

$$\text{Now, similarly, } x_2 = \text{ReLU}(x_2) - \text{ReLU}(-x_2)$$

$$\text{So, } \max(x_1, x_2) = \max(0, x_1 - x_2) + \max(0, x_2) - \max(0, -x_2)$$

$y_2 \rightarrow \boxed{\max(x_1, x_2) = \text{ReLU}(x_1 - x_2) + \text{ReLU}(x_2) - \text{ReLU}(-x_2)}$

Network diagram -



The weights of each of the neuron o/p's in the second layer and the weights of the i/p's are written along the connections.

All biases are set to 0. { Ans.

So, the wt. vector of $x_1 \rightarrow [-1, +1, +1, 0, 0]$
 the wt. vector of $x_2 \rightarrow [0, 0, -1, +1, -1]$.

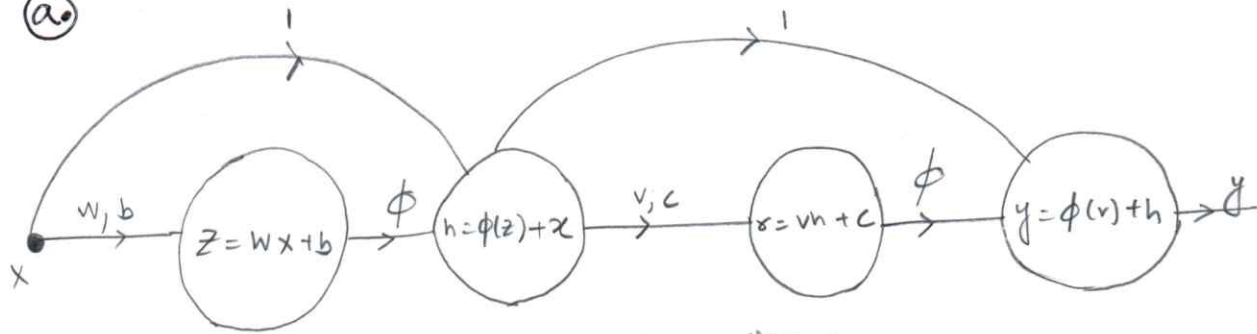
The learned quantity is written inside the neuron.

wts. of $\text{ReLU}(-x_1) \rightarrow [-1, 0]$
 wts. of $\text{ReLU}(x_1) \rightarrow [+1, 0]$
 wts. of $\text{ReLU}(x_1 - x_2) \rightarrow [-1, 1]$
 wts. of $\text{ReLU}(x_2) \rightarrow [0, +1]$
 wts. of $\text{ReLU}(-x_2) \rightarrow [0, -1]$

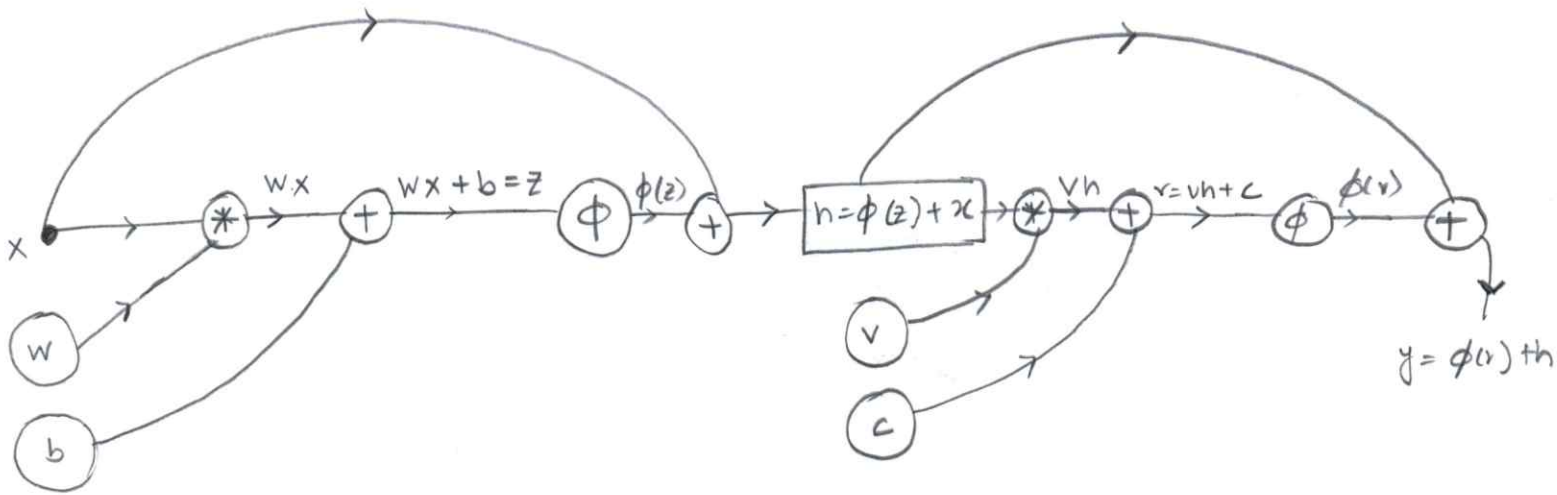
} Ans.

Q2. Back propagation through Residual Block.

(a)



So, the computational graph looks like :



(b) Determine the backprop rules for gradient w.r.t. w, b, v & c .

Let's say the o/p (actual o/p) is \bar{y} , and the loss l is defined as $l(y, \bar{y})$.

So, we have to find $\frac{\partial l}{\partial v}, \frac{\partial l}{\partial c}, \frac{\partial l}{\partial w}, \frac{\partial l}{\partial b}$.

first, $\frac{\partial l}{\partial v}$

$$\frac{\partial l}{\partial v} = \frac{\partial l}{\partial y} \cdot \frac{\partial y}{\partial r} \cdot \frac{\partial r}{\partial v}$$

putting the values,

$$= \frac{\partial l}{\partial y} \cdot \frac{\partial (\phi(r) + h)}{\partial r} \cdot \frac{\partial (vh + c)}{\partial v}$$

$$= \frac{\partial l}{\partial y} \cdot \left[\phi'(r) + \frac{\partial y}{\partial r} \right] \cdot h = \boxed{\frac{\partial l}{\partial y} \cdot \phi'(r) \cdot h = \frac{\partial l}{\partial v}}$$

$$\text{Now, } \frac{dl}{dc} = \frac{dl}{dy} \cdot \frac{dy}{dv} \cdot \frac{dv}{dc}$$

$$= \frac{dl}{dy} \cdot \frac{d}{dv} (\phi(v) + h) \cdot \frac{d}{dc} (vh + c) = \boxed{\frac{dl}{dy} \cdot \phi'(v) = \frac{dl}{dc}}$$

Now, $\frac{dl}{dw} =$ gradient through ① wt. layer (v, c) + gradient through skip connection. ②

$$\text{② } \frac{dl}{dw} = \frac{dl}{dy} \cdot \frac{dy}{dh} \cdot \frac{dh}{dz} \cdot \frac{dz}{dw}$$

$$= \frac{dl}{dy} \cdot \frac{d}{dh} (\phi(v) + h) \cdot \frac{d}{dz} (\phi(z) + x) \cdot \frac{d}{dw} (wx + b)$$

$$= \frac{dl}{dy} \cdot \phi'(z) \cdot x$$

$$\text{① } \frac{dl}{dw} = \frac{dl}{dy} \cdot \frac{dy}{dv} \cdot \frac{dv}{dh} \cdot \frac{dh}{dz} \cdot \frac{dz}{dw}$$

$$= \frac{dl}{dy} \cdot \frac{d}{dv} (\phi(v) + h) \cdot \frac{d}{dh} (vh + c) \cdot \frac{d}{dz} (\phi(z) + x) \cdot \frac{d}{dw} (wx + b)$$

$$= \frac{dl}{dy} \cdot \phi'(v) \cdot v \cdot \phi'(z) \cdot x$$

$$\text{So, } \frac{dl}{dw} = \text{①} + \text{②} = \boxed{\frac{dl}{dy} \cdot \phi'(z) \cdot x [\phi'(v) \cdot v + \mathbf{I}] = \frac{dl}{dw}}$$

\mathbf{I} is identity matrix.

Now, $\frac{dl}{db} =$ gradient through ① wt. layer + gradient through skip connection. ②

$$\text{① } \frac{dl}{db} = \frac{dl}{dy} \cdot \frac{dy}{dv} \cdot \frac{dv}{dh} \cdot \frac{dh}{dz} \cdot \frac{dz}{db}$$

$$= \frac{dl}{dy} \cdot \frac{d}{dv} (\phi(v) + h) \cdot \frac{d}{dh} (vh + c) \cdot \frac{d}{dz} (\phi(z) + x) \cdot \frac{d}{db} (wx + b)$$

$$= \frac{dl}{dy} \cdot \phi'(v) \cdot v \cdot \phi'(z) \cdot 1$$

$$\textcircled{2} \quad \frac{dl}{db} = \frac{dl}{dy} \cdot \frac{dy}{dz} \cdot \frac{dz}{dh} \cdot \frac{dh}{db}$$

$$= \frac{dl}{dy} \cdot \phi'(z)$$

So,

$$\frac{dl}{db} = \textcircled{1} + \textcircled{2} = \left[\frac{dl}{dy} \cdot \phi'(z) \cdot [\phi'(x) \cdot V + I] \right] = \frac{dl}{db}$$

3. EM for Probabilistic PCA.

(a) z (latent code vector) drawn from std. Gaussian $N(0, I)$
 we will consider z to be the 1-D projection.

$$z \sim N(0, 1)$$

and,

$$x|z \sim N(u \cdot z, \sigma^2 I)$$

I changed the mean from $z \cdot u$ to $u \cdot z$ for the following reason.

Here is my approach.

$x \in \mathbb{R}^{J \times N}$ (there are 'N' samples of x , each of which has 'J' dimensions).

each sample $x_i \in \mathbb{R}^J$ where $i = 1$ to N .

u is the principal component / parameter vector. i.e., when $x^{(i)}$ is projected on u we get $z^{(i)}$ which is the latent representation of $x^{(i)}$ in 1-D space.

So, $u \in \mathbb{R}^{J \times 1}$ (transform J dimension to 1)

and, $z = u^T x$
 so, $z \in \mathbb{R}^{1 \times N}$ ('N' samples each of which is 1-D).

so, $z^{(i)} \in \mathbb{R}$ (scalar).

So, when projecting back z to a higher dimension to get approximation of x we get $x = u \cdot z$.

So, $X = U \cdot Z$.

$(J \times 1) \cdot (1 \times N) = (J \times N)$ which is the dimension for X .

Hence, my reason to change $Z \cdot U$ to $U \cdot Z$.

So, we have

$$Z^{(i)} \sim N(0, 1)$$

$$X^{(i)} | Z^{(i)} \sim N(U \cdot Z^{(i)}, \sigma^2)$$

Now, from the relations given in the Appendix, we can infer the values of $p(X)$ and $p(Z|X)$

$$p(X) \sim N(A\mu + b, A\Sigma A^T + S)$$

where, $A = U$ $\mu = 0$ $b = 0$ $\Sigma = I$ $S = \sigma^2 I$.

So, we get

$$p(X^{(i)}) \sim N(0, U \cdot U^T + \sigma^2 I)$$

and finally, $p(Z|X) = N(C(A^T S^{-1}(X - b) + \Sigma^{-1}\mu), C)$

where $C = (\Sigma^{-1} + A^T S^{-1} A)^{-1}$

Replacing the values we get - $p(Z^{(i)} | X^{(i)}) =$

$$C(A^T S^{-1}(X - b) + \Sigma^{-1}\mu)$$

$$= C(U^T (\sigma^2 I)^{-1} X)$$

$$= C \cdot \left(\frac{U^T X}{\sigma^2} \right)$$

$$C = (1 + U^T (\sigma^2)^{-1} U)^{-1}$$

$$C = \left(1 + \frac{U^T U}{\sigma^2} \right)^{-1}$$

$$C = \left(\frac{\sigma^2 + U^T U}{\sigma^2} \right)^{-1}$$

$$C = \frac{\sigma^2}{U^T U + \sigma^2}$$

The reason we can divide this is

$U^T U$ and σ^2 are both scalars.

$$\text{So, } \mu = \frac{\sigma^2}{\sigma^2 + u^T u} \cdot \frac{u^T x}{\sigma^2} = \frac{u^T x}{\sigma^2 + u^T u}.$$

$$c = \frac{\sigma^2}{\sigma^2 + u^T u}.$$

$$\text{So, } p(z^{(i)} | x^{(i)}) \sim N \left(\frac{u^T x}{\sigma^2 + u^T u}, \frac{\sigma^2}{\sigma^2 + u^T u} \right).$$

* $\sigma^2 + u^T u$ is a scalar, so, we can divide the term by them. if, they were vectors we would multiply the numerators by $(\sigma^2 + u^T u)^{-1}$.

$$\text{Now, } m = E[z | x] = \mu = (\sigma^2 I + u^T u)^{-1} \cdot u^T x$$

$$\text{and, } s = E[z^2 | x]$$

$$\text{we know, } \text{var}(z | x) = E[z^2 | x] - E[z | x]^2.$$

Substituting the values we get.

$$(\sigma^2 + u^T u)^{-1} \cdot \sigma^2 = E[z^2 | x] - ((\sigma^2 + u^T u)^{-1} \cdot u^T x)^2$$

$$s = (\sigma^2 + u^T u)^{-1} \cdot \sigma^2 I + ((\sigma^2 + u^T u)^{-1} \cdot u^T x)^2$$

$$m^{(i)} = \frac{u^T x^{(i)}}{\sigma^2 + u^T u}$$

$$s^{(i)} = \frac{\sigma^2}{\sigma^2 + u^T u} + \frac{(u^T x^{(i)})^2}{(\sigma^2 + u^T u)^2}$$

Since, $u^T x^{(i)} = z^{(i)}$ (scalar) & $\sigma^2 + u^T u$ is scalar, $m^{(i)}$ & $s^{(i)}$ are scalars.

(b) M-Step

$$u_{\text{new}} \leftarrow \underset{u}{\text{argmax}} \sum_{i=1}^N E_{q(z^{(i)})} [\log p(z^{(i)}, x^{(i)})].$$

$E_{q(z^{(i)})}$ is the expected value of $q(z)$ for all values $z^{(i)}$.

Lets start with the $\log p(z^{(i)}, x^{(i)})$ term.

$$\log [p(x^{(i)}, z^{(i)})] = \log [p(x^{(i)} | z^{(i)}) \cdot p(z^{(i)})]$$

$$= \log p(x^{(i)} | z^{(i)}) + \log p(z^{(i)})$$

Now, we know, $x^{(i)} | z^{(i)}$ is multivariate Gaussian as x has J dims.

$$x^{(i)} | z^{(i)} \sim N(u \cdot z^{(i)}, \sigma^2)$$

we know, for multivariate Gaussian $p(x; \mu, \Sigma) =$

$$\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \cdot e^{-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)}$$

putting the values we get.

$$p(x^{(i)} | z^{(i)}) = \frac{1}{(2\pi)^{J/2} \sigma} \cdot e^{-\frac{1}{2\sigma^2} (x^{(i)} - u \cdot z^{(i)})^T \cdot (x^{(i)} - u \cdot z^{(i)})}$$

σ^2 is scalar.

$$\text{So, } \log p(x^{(i)} | z^{(i)}) = \log \frac{1}{(2\pi)^{J/2} \sigma} - \frac{1}{2\sigma^2} (x^{(i)} - u \cdot z^{(i)})^T \cdot (x^{(i)} - u \cdot z^{(i)})$$

Now, $p(z^{(i)})$ is a univariate dist. as $z^{(i)}$ is a scalar.

$$p(z^{(i)}) = \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{1}{2} (z^{(i)})^2}$$

$$\text{So, } \log p(z^{(i)}) = \log \frac{1}{\sqrt{2\pi}} - \frac{1}{2} (z^{(i)})^2$$

$$\text{So, } \log [p(x^{(i)}, z^{(i)})] = \log \frac{1}{(2\pi)^{J/2} \sigma} - \frac{1}{2\sigma^2} (x^{(i)} - u \cdot z^{(i)})^T \cdot (x^{(i)} - u \cdot z^{(i)})$$

$$+ \log \frac{1}{\sqrt{2\pi}} - \frac{1}{2} (z^{(i)})^2$$

Now, since, we will optimize w.r.t. u , we can remove the terms that are not dependent on u .

So, we have

$$-\frac{1}{2\sigma^2} (x^{(i)} - u z^{(i)})^T (x^{(i)} - u z^{(i)})$$

All other terms would yield '0' if differentiated w.r.t. u , so, we drop them.

So, we have after expanding

$$-\frac{1}{2\sigma^2} \left\{ x^{(i)T} x^{(i)} - 2x^{(i)T} u z^{(i)} + u^T u (z^{(i)})^2 \right\}$$

Now, putting back in our original equation, and equating to 0

$$\frac{\partial}{\partial u} \sum_{i=1}^N E_{q(z^{(i)})} \left[-\frac{1}{2\sigma^2} (x^{(i)T} x^{(i)} - 2x^{(i)T} u z^{(i)} + u^T u (z^{(i)})^2) \right] = 0$$

$$-\frac{1}{2\sigma^2} \frac{\partial}{\partial u} \left[\sum_{i=1}^N E_{q(z^{(i)})} (x^{(i)T} x^{(i)} - 2x^{(i)T} u z^{(i)} + u^T u (z^{(i)})^2) \right] = 0$$

Now, using property of expectation.

$$\frac{\partial}{\partial u} \sum_{i=1}^N x^{(i)T} x^{(i)} - 2x^{(i)T} u E[z^{(i)}] + u^T u E[(z^{(i)})^2] = 0$$

Now, $E[z^{(i)}] = m^{(i)}$ & $E[(z^{(i)})^2] = s^{(i)}$

$$\sum_{i=1}^N \frac{\partial}{\partial u} (x^{(i)T} x^{(i)} - 2x^{(i)T} u m^{(i)} + u^T u s^{(i)}) = 0$$

$$\Rightarrow \sum_{i=1}^N -2x^{(i)T} m^{(i)} + 2u s^{(i)} = 0$$

$$\Rightarrow \sum_{i=1}^N x^{(i)T} m^{(i)} = \sum_{i=1}^N u s^{(i)}$$

$$\Rightarrow u = \frac{\sum_{i=1}^N x^{(i)T} m^{(i)}}{\sum_{i=1}^N s^{(i)}} \quad \leftarrow \text{Ans}$$