

P)

MLE distribution:

$$\text{Suf}(f) = \arg \max_{\theta \in \Theta} P(\theta)$$

If likelihood fn is $P(x|\theta)$, MLE for θ is:

$$\theta_{\text{MLE}} = \arg \max_{\theta} P(x|\theta)$$

$$= \arg \max_{\theta} \prod_i P(x_i|\theta) \rightarrow ①$$

for MAP for θ :

$$\theta_{\text{MAP}} = \arg \max_{\theta} P(x|\theta) \cdot P(\theta)$$

$$= \arg \max_{\theta} \prod_i P(x_i|\theta) \cdot P(\theta) \rightarrow ②$$

if $i = 0, \dots, M-1$.

$$P(\theta) = \frac{1}{M} \text{ in case of equal priors.}$$

from ①,

$$\theta_{\text{MLE}} = \arg \max_{\theta} \left(\sum_i \log P(x_i|\theta) \right)$$

from ②,

$$\theta_{\text{MAP}} = \arg \max_{\theta} \left(\sum_i \log P(x_i|\theta) + \log P(\theta) \right)$$

$$= \arg \max_{\theta} \left(\sum_i \log P(x_i|\theta) + \log \left(\frac{1}{M} \right) \right)$$

$$= \arg \max_{\theta} \left(\sum_i \log P(x_i|\theta) - \log M \right)$$

Now we see that $\log M$ doesn't vary with θ .

$$\text{Hence maximizing } \sum_i \log P(x_i|\theta) - \log M$$

is equivalent to maximizing $\sum_i \log P(x_i|\theta)$ over θ .

Hence MLE and MAP distributions are same in case of uniform prior.

Given N observations x_1, x_2, \dots, x_N , the likelihood of these N observations coming from a Gaussian distribution $N(\mu, \sigma^2)$ is given by:

$$p(x_1, x_2, \dots, x_N | \mu, \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

Using MLE,

$$\begin{aligned} p(x_1, x_2, \dots, x_N | \mu, \sigma^2) &= \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right) \\ &= \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right) \\ &= \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right) \\ &= \sum_{i=1}^N -\frac{1}{2} \log (2\pi\sigma^2) \\ &\quad - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \\ &= -\frac{N}{2} \log (2\pi\sigma^2) - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \end{aligned}$$

To get $\max p(x|\theta)$ we can just take the gradient of $p(x|\theta)$ w.r.t θ

$$\cancel{\nabla p(x|\theta)} \cdot \nabla_\theta p = 0.$$

~~Note~~: $\nabla p(x_1, \dots, x_N | \mu, \sigma^2)$ w.r.t μ : (σ^2 known)

$$\frac{\partial}{\partial \mu} \cdot \left(-\frac{N}{2} \log (2\pi\sigma^2) - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right) = 0.$$

$$\Rightarrow \frac{1}{2\sigma^2} \sum_{i=1}^N 2(x_i - \mu) = 0.$$

$$\sum_{i=1}^N (x_i - \mu) = 0.$$

$$\sum_{i=1}^N x_i - \sum_{i=1}^N \mu = 0.$$

$$\sum_{i=1}^N x_i - \mu N = 0 \quad \therefore \mu = \frac{1}{N} \sum_{i=1}^N x_i$$

Similarly taking gradient w.r.t. σ^2 (μ known)

$$\frac{\partial}{\partial \sigma^2} \left(-\frac{N}{2} \log(2\pi\sigma^2) - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right) = 0.$$

$$-\frac{N}{2} \cdot \frac{2\sigma}{2\pi\sigma^2} - \sum_{i=1}^N \frac{2(x_i - \mu)^2 \cdot 2}{(2\sigma^2)^2} = 0.$$

$$-\frac{N}{2\sigma^2} + \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^4} = 0.$$

$$\sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^4} = \frac{N}{2\sigma^2}$$

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

MLE just gives the mean and variance of the distribution or parameters of the Gaussian distribution and doesn't make any assumptions (unbiased)

i.e) for MAP,

$$\Omega_{MAP} = \arg \max_{\theta} \left(\sum_i \log P(x_i | \theta) + \log P(\theta) \right)$$

$$P(x_1, x_2, \dots, x_N | \mu, \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\frac{\partial}{\partial \mu} \cdot \left[-\frac{N}{2} \log(2\pi\sigma^2) - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} + \log P(\mu) \right] = 0.$$

$$\frac{1}{2\sigma^2} \sum_{i=1}^N 2(x_i - \mu) + \frac{\partial}{\partial \mu} \log P(\mu) = 0.$$

$$\frac{\partial}{\partial \mu} \log P(\mu) = -\frac{1}{\sigma^2} \sum_{i=1}^N x_i + \frac{N\mu}{\sigma^2}$$

$$\mu = \frac{\sigma^2}{N} \sum_{i=1}^N x_i + \frac{\partial}{\partial \mu} \log P(\mu)$$

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i + \frac{\sigma^2}{N} \cdot \frac{\partial}{\partial \mu} \log P(\mu)$$

$$P(\mu) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\mu - \bar{x})^2}{2\sigma^2}}$$

$$\frac{\partial}{\partial \mu} \log P(\mu) = -\frac{1}{2} \log(2\pi\sigma^2) - \frac{\mu^2}{2\sigma^2}$$

$$\frac{\partial}{\partial \mu} \log P(\mu) = -\frac{\mu}{\sigma^2} = -\frac{\mu}{\sigma^2}$$

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i + \frac{\sigma^2}{N} \cdot \frac{\mu}{\sigma^2}$$

$$\mu \left[1 \mapsto \frac{\sum_i x_i}{N y^2} \right] = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\mu = \frac{\left(\frac{1}{N} \sum_{i=1}^N x_i \right) N y^2}{(N y^2 + \sigma^2)} = y^2 \frac{\left(\sum_{i=1}^N x_i \right)}{(N y^2 + \sigma^2)}$$

P-2. From directional derivatives we know

a) $\nabla_{\vec{v}} f = \lim_{h \rightarrow 0} \frac{f(x+h\vec{v}) - f(x)}{h}$

$$\Rightarrow f(x+h\vec{v}) = f(x) + h \nabla f^T \vec{v}$$

$$J(w^k) = \frac{1}{2m} \left[\sum_{i=1}^m \|h_{w^k}(x_i) - y\|_2^2 \right]$$

$$J(w^{k+1}) = J(w^k) + (w^{k+1} - w^k)^T \nabla J$$

$$f(x+h\vec{v}) = f(x) + h \nabla f^T \vec{v} + \frac{h^2}{2} \vec{v}^T H \vec{v} + \dots$$

$$\therefore f(w^{k+1}) = f(w^k) + (w^{k+1} - w^k)^T \nabla f + \frac{(w^{k+1} - w^k)^T H (w^{k+1} - w^k)}{2}$$

$$J(w^k) = \frac{1}{2m} \|h_{w^k}(x) - y\|_2^2$$

$$\therefore J(w^{k+1}) = J(w^k) + \vec{v}^T \nabla f + \frac{1}{2} \vec{v}^T H \vec{v} + \dots$$

where

$$\vec{v} = w^{k+1} - w^k$$

$\nabla f = \text{Jacobian of } v \mapsto w_k$

$H = \text{Hessian of } v \mapsto w_k$

b) Also we know from grad. descent in directional derivatives,

$$\text{if } \vec{v} = -\lambda \nabla f \quad \nabla f = \text{Jacobian}$$

$$\vec{v} = w^{k+1} - w^k$$

$$\therefore w^{k+1} - w^k = -\lambda \nabla f \Rightarrow w^{k+1}$$

$$J(w^k) = \frac{1}{2m} \|h_{w^k}(x) - y\|_2^2 \Rightarrow w^{k+1} = w^k - \lambda \nabla f$$

From eqn ①, $J(w^{k+1}) = J(w^k) - \lambda \nabla f^T \nabla f$
 $= J(w^k) - \lambda \nabla J^T \nabla J$

$\nabla J = \text{Jacobian w.r.t } w_k$

Now we know if ~~and~~ ~~whether~~ $\nabla J(w_k)$ is a
N x 1 vector, $w^T w = \sum_{i=1}^N w_i^2$ which $\|w\|^2$ is a scalar quantity.

$\nabla J(w_k)$ is also a vector of dimension same as w_k , so $\nabla J^T \nabla J$ is a scalar quantity, $\nabla J^T \nabla J = \|J\|_2^2$

$$J(w^{k+1}) = J(w_k) - \lambda \text{ (some scalar quantity)}$$

So J decreases in every iteration of Gradient descent.

c) Optimal learning rate:

$$w^{k+1} = w_k - \lambda \nabla J(w_k)$$

$$J(w^{k+1}) = J(w_k) - \lambda \nabla J(w_k)^T \nabla J(w_k)$$

$$J(w^{k+1}) = J(w_k) - \lambda \nabla J(w_k)^T \nabla J(w_k)$$

$$J(w^{k+1}) = J(w_k) + v^T \nabla J(w_k) + \frac{1}{2} v^T H v \dots \text{from eqn ①.}$$

$$v = w^{k+1} - w_k$$

$$v = -\lambda \nabla J$$

$$J(w^{k+1}) = J(w_k) - \lambda \nabla J(w_k)^T \nabla J(w_k) + \frac{\lambda^2}{2} \nabla J(w_k)^T H \nabla J(w_k)$$

At optimal learning rate we should have reached the global minima ~~in a single step~~, i.e. cost fn won't change any further.

$$J(w^{k+1}) = J(w_k)$$

$$\lambda \nabla J(w_k)^T \nabla J(w_k) = \frac{\lambda^2}{2} \nabla J(w_k)^T H \nabla J(w_k)$$

$$\nabla J^T \nabla J = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$\lambda = \frac{2 \nabla J(w_k)^T \nabla J(w_k)}{\nabla J(w_k)^T H \nabla J(w_k)}$$

λ = optimal learning rate

d) From eqn ①,

$$J(\omega^{k+1}) = J(\omega^k) + v^T \cancel{\frac{\partial f}{\partial \omega}} + \frac{1}{2} v^T H v \quad \cancel{\text{min}}$$
$$v = \omega^{k+1} - \omega^k$$

$$\mathbf{0} = \frac{d}{dv} J(\omega^k) + \frac{1}{2} v^T H$$

$$\cancel{\frac{d}{dv} J(\omega^k)} = -\cancel{\frac{1}{2} v^T H}$$

$$v^T = -\cancel{\frac{d}{dv} J(\omega^k) \cdot H^{-1}}$$

$$\frac{d}{dv} J(\omega^k) = \cancel{\frac{d}{dv} f(\omega^k)} \quad \frac{d}{dv} J(\omega^k) = \nabla J$$

$$\mathbf{0} = \nabla J + v^T \nabla H$$

$$\cancel{v^T \nabla H} \quad v = \cancel{\nabla J} - (\nabla H)^{-1} \nabla J$$

$$\therefore \omega^{k+1} = \omega^k - (\nabla H)^{-1} \nabla J$$

e)

$$\phi_j \leftarrow \mu \phi_j - \frac{\alpha}{N} \sum_{i=1}^N x_j^{(i)} (y^{(i)} - t^{(i)})$$

$$w_j \leftarrow w_j + \phi_j$$

$$x: N \times D \quad y: N \times 1 \quad t: N \times 1$$
$$w: D \times 1 \quad p: D \times 1$$

$$y^{(i)} - t^{(i)} = y - t \quad (N \times 1) -$$

$$x_j: 1 \times D \quad \cancel{\frac{1}{N} \sum_{i=1}^N x_j^{(i)} (y^{(i)} - t^{(i)})} = \cancel{\frac{1}{N} \sum_{i=1}^N x_j^{(i)} (y - t)}$$

for each j

$$y^{(i)} - t^{(i)}: N \times 1$$

$$\phi_j: 1 \times 1$$

$$x_j: N \times 1$$

$$\cancel{\frac{1}{N} \sum_{i=1}^N x_j^{(i)} (y^{(i)} - t^{(i)})} = \cancel{\frac{1}{N} x_j^T (y - t)}$$

$$P = \mu P - \frac{\alpha}{N} x^T (y - t)$$

$$w = w + P$$

P-3

$$a) \quad w^{k+1} = w^k - \alpha \frac{\partial J(w)}{\partial w_k} \quad w^{k+1} = w^k - \alpha \frac{\partial}{\partial w_k} J(w_k)$$

$= w^k$

$$\begin{aligned} \frac{\partial}{\partial w_k} J(w_k) &= \frac{\partial}{\partial w_k} \cdot \frac{1}{2m} \cdot \sum_{i=1}^m (h_\theta(x^i) - y^i)^2 \\ &= \frac{1}{2m} \cdot \frac{\partial}{\partial w_k} \cdot \sum_{i=1}^m \left(\sum_{j=0}^n w_j x_j^i - y^i \right)^2 \\ &= \frac{1}{m} \cdot \sum_{i=1}^m \left(\sum_{j=0}^n w_j x_j^i - y^i \right) \cdot \cancel{\frac{\partial}{\partial w_k}} \cdot \left(\sum_{j=0}^n w_j x_j^i - y^i \right) \\ &= \frac{1}{m} \cdot \sum_{i=1}^m \left(\sum_{j=0}^n w_j x_j^i - y^i \right) x_k^i \end{aligned}$$

$i \rightarrow \text{iterator}$
 for no. of inputs
 $j \rightarrow \text{iterator over features}$

$$w^{k+1} = w^k - \frac{\alpha}{m} \cdot \sum_{i=1}^m \left(\sum_{j=0}^n w_j x_j^i - y^i \right) x_k^i$$

We see that the update of the weights w^k depends on the value of the input x_k . Now suppose that we have

feature 1 ranging from $1 \rightarrow 10$ and another feature 2 ranging from $1 \rightarrow 10^6$. Gradient descent would take larger steps for feature 2 and very small steps for feature 1, which is undesirable, specially if feature 1 and 2 denote similar features with different units. Hence we need to do feature scaling.

$$4.a) . \quad y = x^T \varphi x \quad \text{Consider } N=3.$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \varphi = \begin{bmatrix} \varphi_{11} & \varphi_{12} & \varphi_{13} \\ \varphi_{21} & \varphi_{22} & \varphi_{23} \\ \varphi_{31} & \varphi_{32} & \varphi_{33} \end{bmatrix}$$

$$y = [x_1 \ x_2 \ x_3] \begin{bmatrix} \varphi_{11} & \varphi_{12} & \varphi_{13} \\ \varphi_{21} & \varphi_{22} & \varphi_{23} \\ \varphi_{31} & \varphi_{32} & \varphi_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$= x_1 (\varphi_{11} x_1 + \varphi_{12} x_2 + \varphi_{13} x_3) + x_2 (\varphi_{21} x_1 + \varphi_{22} x_2 + \varphi_{23} x_3) + x_3 (\varphi_{31} x_1 + \varphi_{32} x_2 + \varphi_{33} x_3)$$

$\frac{\partial y}{\partial \varphi}$ should have dimensions of φ .

$$\frac{\partial y}{\partial \varphi} = \begin{bmatrix} \frac{\partial y}{\partial \varphi_{11}} & \frac{\partial y}{\partial \varphi_{12}} & \frac{\partial y}{\partial \varphi_{13}} \\ \frac{\partial y}{\partial \varphi_{21}} & \frac{\partial y}{\partial \varphi_{22}} & \frac{\partial y}{\partial \varphi_{23}} \\ \frac{\partial y}{\partial \varphi_{31}} & \frac{\partial y}{\partial \varphi_{32}} & \frac{\partial y}{\partial \varphi_{33}} \end{bmatrix} = \begin{bmatrix} x_1^2 & x_1 x_2 & x_1 x_3 \\ x_2 x_1 & x_2^2 & x_2 x_3 \\ x_3 x_1 & x_3 x_2 & x_3^2 \end{bmatrix}$$

$$= [x_1 \ x_2 \ x_3] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}^T = x x^T$$

$$= \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix}_{3 \times 1}^T = x x^T$$

$\frac{\partial y}{\partial x}$ has the same dimensions as x .

$$\begin{aligned} \frac{\partial y}{\partial x} &= \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \frac{\partial y}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 2q_1x_1 + q_2x_2 + q_3x_3 + x_2q_{12} + x_3q_{13} \\ q_1x_1 + 2q_2x_2 + q_3x_3 + x_1q_{21} + x_3q_{23} \\ x_1q_{13} + x_2q_{12} + q_{11}x_1 + q_{12}x_2 + q_{13}x_3 \end{bmatrix} \\ &= \begin{bmatrix} 2q_1x_1 + x_2(q_{12} + q_{21}) + x_3(q_{13} + q_{31}) \\ x_1(q_{11} + q_{21}) + x_2(2q_{22} + q_{32} + q_{12}) \\ x_1(q_{13} + q_{31}) + x_2(q_{12} + q_{23}) + 2q_{33}x_3 \end{bmatrix} \\ Q &= \begin{bmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{bmatrix} \quad Q^T = \begin{bmatrix} q_{11} & q_{21} & q_{31} \\ q_{12} & q_{22} & q_{32} \\ q_{13} & q_{23} & q_{33} \end{bmatrix} \end{aligned}$$

$$Q^T Q = \begin{bmatrix} 2q_{11} \\ q_{12} \\ q_{13} \end{bmatrix}$$

$$\begin{aligned} \frac{\partial y}{\partial x} &= \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} 2q_{11} & q_{12} + q_{21} & q_{13} + q_{31} \\ q_{12} + q_{21} & 2q_{22} & q_{23} + q_{32} \\ q_{13} + q_{31} & q_{23} + q_{32} & 2q_{33} \end{bmatrix} \\ &= \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \left(\begin{bmatrix} q_{11} & q_{12} & q_{13} \\ q_{21} & q_{22} & q_{23} \\ q_{31} & q_{32} & q_{33} \end{bmatrix} + \begin{bmatrix} q_{11} & q_{21} & q_{31} \\ q_{12} & q_{22} & q_{32} \\ q_{13} & q_{23} & q_{33} \end{bmatrix} \right) \\ &= x^T (Q + Q^T) \end{aligned}$$

4.b)

$$h_\theta(x) = \theta^T x$$

Using linear reg. with L2 norm cost fn.

$$\begin{aligned} J(\theta) &= \frac{1}{2m} \cdot (x\theta - y)^T (x\theta - y) \\ &= \frac{1}{2m} \cdot (\theta^T x^T - y^T)(x\theta - y) \\ &= \frac{1}{2m} \cdot (\theta^T x^T x\theta - y^T x\theta + y^T y - \theta^T x^T y) \end{aligned}$$

Using commutativity of dot products

$$y^T (x\theta) = (\theta^T x^T) y$$

$$J(\theta) = \frac{1}{2m} (\theta^T X^T X \theta - 2(X\theta)^T y + y^T y)$$

$$\frac{\partial J(\theta)}{\partial \theta}$$

Now ~~$X^T X$ is a scalar~~. $X^T X = \mathbf{I}$. As shown in part 4a)

$$\frac{\partial}{\partial \theta} (\theta^T \mathbf{I} \theta) = \theta^T (\theta + \theta^T)$$

$$= \theta^T (X^T X + X^T X) = 2 \theta^T X^T X$$

$$\text{Let us consider } \frac{\partial}{\partial \theta} (X\theta)^T. \text{ i.e } \frac{\partial}{\partial \theta} \cdot \theta^T X^T = \frac{\partial}{\partial \theta} (X\theta)$$

We know if $\begin{matrix} N \neq W^T X \\ \frac{\partial N}{\partial \theta} = W^T \cancel{\theta} \end{matrix}$

$$X = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \\ \theta_{31} & \theta_{32} & \theta_{33} \end{bmatrix}$$

$$\begin{aligned} \theta^T X^T &= \begin{bmatrix} \theta_{11}/\theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22}/\theta_{23} \\ \theta_{31} & \theta_{32} \end{bmatrix}^T \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \\ &= \begin{bmatrix} \theta_{11} & \theta_{21} & \theta_{31} \\ \theta_{12} & \theta_{22} & \theta_{32} \\ \theta_{13} & \theta_{23} & \theta_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \\ &= \theta_{11}(x_1) + \theta_{21}x_2 + \theta_{31}x_3 + \theta_{12}x_1 + \theta_{22}x_2 + \theta_{32}x_3 \\ &= \begin{bmatrix} \theta_{11}x_1 + \theta_{21}x_2 + \theta_{31}x_3 \\ \theta_{12}x_1 + \theta_{22}x_2 + \theta_{32}x_3 \\ \theta_{13}x_1 + \theta_{23}x_2 + \theta_{33}x_3 \end{bmatrix} \end{aligned}$$

$$\frac{\partial}{\partial \theta} [\theta^T X^T] = \begin{bmatrix} x_1 & x_2 & x_3 \\ x_1 & x_2 & x_3 \\ x_1 & x_2 & x_3 \end{bmatrix}$$

$$\frac{\partial}{\partial \theta} [\theta^T X^T] = \begin{bmatrix} / & / & / \\ / & / & / \\ / & / & / \end{bmatrix}$$

We know

$$\begin{aligned} \frac{\partial}{\partial x} x^T A x &\rightarrow Z \quad Z = x^T (A + A^T) \\ \frac{\partial}{\partial x} x^T Z &= x^T (Z x^{-1} + x^{-T} Z^T) \\ &= x^T Z x^{-1} + Z^T \\ \text{Q.E.D.} &= x^T A + Z^T \end{aligned}$$

$$\frac{\partial}{\partial \theta} (\theta^T X^T) = \frac{\partial}{\partial \theta} (X\theta) \quad (\text{by Commutativity of dot product})$$

$$\frac{\partial}{\partial \theta} (X\theta) = X^T$$

$$\frac{\partial J}{\partial \theta} = \frac{1}{m} \cdot 2 \theta^T X^T X - \frac{1}{m} \cdot 2 \frac{\partial}{\partial \theta} (\theta^T y)$$

$$= \frac{1}{m} \cdot [\theta^T X^T X - X^T y]$$

$$\frac{\partial J}{\partial \theta} = 0$$

$$\theta^T X^T X - X^T y = 0$$

$$X^T y = \theta^T X^T X \Rightarrow X^T y = (X^T X)^T \theta$$

~~you can't do this~~

(using commutativity
of dot prod).

$$\theta^T = (X^T X)^{-1} X^T y$$

$$X^T y = X^T X \cdot \theta$$

$$\therefore \theta = (X^T X)^{-1} X^T y$$

1c). $f_i = f(s_{y_i}) = \frac{e^{s_{y_i}}}{\sum_j e^{s_{y_j}}} \quad (\text{softmax fn})$

$$\frac{\partial f_i}{\partial s_{y_j}} = \frac{\partial}{\partial s_{y_j}} \left(\frac{e^{s_{y_i}}}{\sum_k e^{s_{y_k}}} \right) = \frac{h_i \frac{\partial}{\partial s_{y_j}} e^{s_{y_i}} - e^{s_{y_i}} \frac{\partial}{\partial s_{y_j}} h_i}{h_i^2}$$

$$= \frac{\partial}{\partial s_{y_j}} / \cancel{h_i}$$

Let $g_i = e^{s_{y_i}}$ $h_i = \sum_k e^{s_{y_k}}$

~~$\frac{\partial g_i}{\partial s_{y_j}} = \frac{\partial}{\partial s_{y_j}} e^{s_{y_i}} = e^{s_{y_i}} \frac{\partial}{\partial s_{y_j}} s_{y_i}$~~

$$\frac{\partial f_i}{\partial s_{y_j}} = \frac{h_i e^{s_{y_i}} \frac{\partial}{\partial s_{y_j}} s_{y_i} - e^{s_{y_i}} \frac{\partial}{\partial s_{y_j}} \sum_k e^{s_{y_k}}}{h_i^2}$$

$$\frac{\partial s_{y_i}}{\partial s_{y_j}} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases} \quad \frac{\partial}{\partial s_{y_j}} \sum_k e^{s_{y_k}} = e^{s_{y_j}}$$

~~$-e^{s_{y_i}} e^{s_{y_j}} = -f_i f_j \quad \text{if } i \neq j$~~

$$\frac{h_i e^{s_{y_i}} - e^{s_{y_i}} e^{s_{y_j}}}{h_i^2} = \frac{e^{s_{y_i}}}{h_i} \left[\frac{h_i - e^{s_{y_j}}}{h_i} \right]$$

if $i=j$:

$$f_i \cdot \left[1 - \frac{e^{sy_j}}{h_0} \right] = f_i [1 - f_j]$$

$$\frac{\partial f_i}{\partial s_{y_j}} = \begin{cases} f_i [1 - f_j] & i=j \\ -f_i f_j & i \neq j \end{cases}$$

P-S.

a) $n_H = J$ ~~$f_H = M$~~ stride = s , padding = p . P .
 $n_W = K$ ~~$f_W = N$~~

We know that

~~$n_H' = n_H + 2p - f_H + 1$~~

$$\frac{n_H + 2p - f_H + 1}{s}$$

$$n_W' = \frac{n_W + 2p - f_W + 1}{s}$$

$$\therefore n_H' = \frac{J + 2p - M}{s} + 1$$

$$n_W' = \frac{K + 2p - N}{s} + 1$$

Size of resulting layer:

$$\left(\frac{J + 2p - M}{s} + 1 \right), \left(\frac{K + 2p - N}{s} + 1 \right)$$

b) filter size = $F \times F$
 stride = s

$$n_H = J \quad n_W = K$$

~~$n_H' = \frac{n_H + 2p - F}{s} + 1$~~

$$n_W' = \frac{n_W + 2p - F}{s} + 1$$

If all the dimensions are valid for the convolution, $p=0$ (Assumption)

$$n_H' = \frac{n_H - F}{s} + 1$$

$$n_W' = \frac{n_W - F}{s} + 1$$

c) Yet fully connected layers can be represented as convolutional layers.

Suppose a FC layer takes an input of size

~~$n_H \times n_W \times n_c$~~ and outputs n_f outputs.

We can use a CONV layer of filter size $(n_H \times n_W)$ and n_f filters of such shape with stride = 1, no padding.

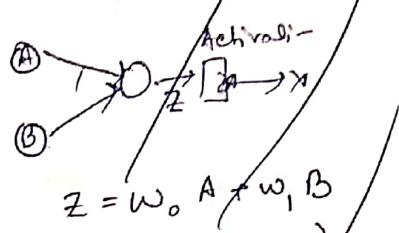
$$n_H' = n_H - n_H + 1 \approx 1$$

$$n_W' = n_W - n_W + 1 \approx 1$$

∴ size of output of conv layer = $1 \times 1 \times n_f$

P-5. Neural Nets & Backpropagation

a)



$$z = w_0 A + w_1 B$$

$$x = \max(0, z)$$

$$\cancel{w_0 \leq 0, w_1 \leq 0}$$

$$x = \max(0, wb) = 0$$

$$\underline{w_0 \leq 0}$$

$$\text{Similarly } \underline{w_1 \leq 0} \dots$$

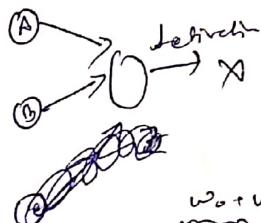
Ans

A	B	x
0	0	0
0	1	0
1	0	0
1	1	1

$$w_0 + w_1 = 1 \rightarrow \underline{w_0 = 1 - w_1}$$

P-5

a)



$$\cancel{w_0 + w_1 + b = z} \rightarrow \underline{z = w_0 A + w_1 B + b}$$

$$x = \text{sigmoid}(z)$$

A	B	x
0	0	0
0	1	0
1	0	0
1	1	1

-20

$$\cancel{w_0 + w_1 + b = z} \rightarrow \underline{z = w_0 A + w_1 B + b}$$

$$z > 4.6$$

$$b = z$$

$$b < -4.6 \rightarrow \underline{2}$$

$$w_0 + b < -4.6 \rightarrow \underline{3}$$

$$w_1 + b < -4.6 \rightarrow \underline{4}$$

if $b = -20$

$$w_0 = 14$$

$$w_1 = 8$$

$$15.4$$

$$-4.6 < w_0$$

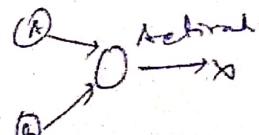
$$\text{If } b = -20, \quad w_0 < \cancel{15.4}$$

$$w_1 < 15.4$$

$$w_0 + w_1 > 24.6$$

$w_0 = 14$
$w_1 = 15$
$b = -20$

b)

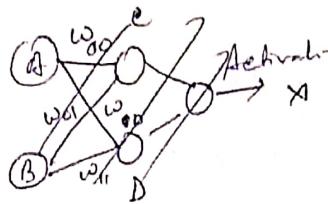
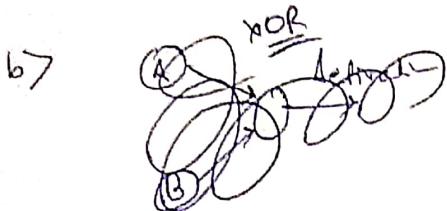


$$z = w_0 A + w_1 B + b$$

A	B	x
0	0	0
0	1	1
1	0	1
1	1	1

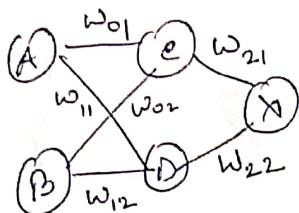
$$\begin{aligned} b &< -4.6 \\ w_1 + b &> 4.6 \\ w_0 + b &> 4.6 \\ w_0 + w_1 + b &> 4.6 \end{aligned}$$

If $b = -10$
 $w_1 = 15$
 $w_2 = 20$
satisfies this condition.



A	B	X
0	0	0
0	1	1
1	0	1
1	1	0

using sigmoid activation fn.



$$X = \bar{A} \cdot B + A \cdot \bar{B}$$

$$C = \bar{A} \cdot B \quad D = A \cdot \bar{B}$$

A	B	C	D
0	0	0	0
0	1	1	0
1	0	0	1
1	1	0	0

$$C = w_{01} A + w_{02} B + b_0$$

$$w_{02} + b_0 > 4.6$$

$$b_0 < -4.6$$

$$w_{01} + b_0 < -4.6$$

$$w_{01} + w_{02} + b_0 < -4.6$$

$$\text{if } b_0 = -10$$

$$w_{02} = 20$$

$$w_{01} = -25$$

$$D = w_{11} A + w_{12} B + b_1$$

$$w_{11} + b_1 > 4.6$$

$$w_{12} + b_1 < -4.6$$

$$b_1 < -4.6$$

$$w_{11} + w_{12} + b_1 < -4.6$$

$$b_1 = -10 - 15$$

$$b_1 = 20$$

$$w_{11} = -5$$

$$w_{12} = -5$$

$$X = w_{21} C + w_{22} D + b_2$$

$$w_{21} = 15 \quad b_2 = -10 \quad (\text{using OR or above})$$

$$w_{22} = 20$$

Similarly we can do so for XOR using a 2-layer network.

c)

$$f(x, w) = \|w \cdot x\|^2 = \sum_{i=1}^n (w \cdot x)_i^2 = \|y\|^2 = \sum_{i=1}^n y_i^2$$

$$x \xrightarrow{\oplus} L_2 \xrightarrow{f(x, w)}$$

f is a scalar fn.

$$y = w \cdot x$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial y} \cdot \frac{\partial y}{\partial x} \frac{\partial f}{\partial y}$$

$$= w^T \frac{\partial f}{\partial y} = 2w^T y = 2w^T w \cdot x \text{ (Ans)}$$

$$\frac{\partial}{\partial y_i} \sum_{i=1}^n y_i^2 = 2y_i$$

$$\therefore \frac{\partial f}{\partial y} = 2y$$

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial y} \cdot \frac{\partial y}{\partial w}$$

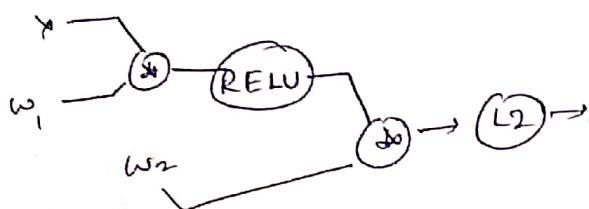
$$y = w^T x$$

$$\begin{aligned}\frac{\partial y}{\partial w} &= \frac{\partial f}{\partial y} \cdot x^T \\ &= 2y \cdot x^T \\ &= 2w^T x \cdot x^T\end{aligned}$$

d) $z_1 = x^T w_1$
 $h_1 = \max(0, z_1)$

$$\hat{y} = h_1 w_2$$

$$L = \|\hat{y}\|_2^2$$



$$\begin{aligned}\frac{\partial L}{\partial w_2} &= \cancel{\frac{\partial L}{\partial \hat{y}}} \cdot \frac{\partial \hat{y}}{\partial w_2} \cdot \frac{\partial L}{\partial y} \\ &= h_1^T \cdot \frac{\partial L}{\partial y} \\ &= h_1^T \cdot 2y \\ &= 2h_1^T y \\ &= 2\max(0, x^T w_1)^T y\end{aligned}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial h_1} \cdot \frac{\partial h_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial L}{\partial \hat{y}} = 2\hat{y}$$

$$\frac{\partial \hat{y}}{\partial h_1} = w_2^T$$

$$\frac{\partial h_1}{\partial z_1} = \begin{cases} 1 & z_1 \geq 0 \\ 0 & z_1 \leq 0 \end{cases}$$

$$\frac{\partial z_1}{\partial w_1} = x^T$$

if $z_1 \geq 0$, $\frac{\partial L}{\partial w_1} = 2\hat{y} \cdot w_2^T \cdot x^T$

if $z_1 \leq 0$, $\frac{\partial L}{\partial w_1} = 0$

where $z_1 = x^T w_1$

where
 $\hat{y} = \max(0, z_1)$

5-e)

$$z = w_0 + w_1 x + w_2 x^2$$

$$y = 1 + e^z$$

$$L = \frac{1}{2} (\log y - \log t)^2$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w_2}$$

$$\frac{\partial L}{\partial y} = \frac{1}{2y} (\log y - \log t)^2$$

$$= \frac{1}{2} \cdot \cancel{2} \cdot (\log y - \log t) \cdot \frac{1}{y}$$

$$= (\log y - \log t) \cdot \frac{1}{y}$$

$$\therefore \frac{\partial L}{\partial y} = (\log y - \log t) \cdot \frac{1}{y}$$

$$\cancel{\frac{\partial y}{\partial z}} \quad \frac{\partial z}{\partial w_2} = 2w_2 x$$

$$\frac{\partial y}{\partial z} = e^z$$

$$\frac{\partial L}{\partial w_2} = (\log y - \log t) \cdot \frac{1}{y} \cdot e^z \cdot 2w_2 x$$

where $z = w_0 + w_1 x + w_2 x^2$

$$y = 1 + e^z$$