

$y = f(x)$

↳ features

↳ class

$f \rightarrow$ KNN
 \rightarrow Softmax
 \rightarrow Neural Network

LINEAR REGRESSION:

$$y = ax + b$$

$a, b \rightarrow$ parameters

ERROR:

$$\hat{y}^{(i)} = f(x^{(i)})$$

$$\varepsilon_i = (y^{(i)} - \hat{y}^{(i)})^2$$

LOSS / COST FUNCTION:

$$\sum_i (y^{(i)} - f(x^{(i)}))^2$$

LEAST SQUARES SOLUTION:

$$\hat{y} = \theta^T \hat{x} \quad \theta = \begin{bmatrix} a \\ b \end{bmatrix}$$

$$\hat{x} = \begin{bmatrix} x \\ 1 \end{bmatrix}$$

$$L(\theta) = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

$$= \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T \hat{x}^{(i)})^2$$

$$= \frac{1}{2} \sum_{i=1}^n \underbrace{(y^{(i)} - \theta^T \hat{x}^{(i)})^T}_{\text{scalar}} (y^{(i)} - \theta^T \hat{x}^{(i)})$$

$$= \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \hat{x}^{(i)\top} \theta)^T (y^{(i)} - \hat{x}^{(i)\top} \theta)$$

$$= \frac{1}{2} \left(\underbrace{\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}}_{\mathbb{R}^n} - \underbrace{\begin{bmatrix} \hat{x}^{(1)\top} \\ \hat{x}^{(2)\top} \\ \vdots \\ \hat{x}^{(n)\top} \end{bmatrix}}_{\mathbb{R}^{n \times 2}} \theta \right)^T \underbrace{\hat{x}}_{\mathbb{R}^{2 \times 1}}$$

$\underbrace{\qquad\qquad\qquad}_{\mathbb{R}^{1 \times n}}$

$$\left(\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix} - \begin{bmatrix} \hat{x}^{(1)\top} \\ \hat{x}^{(2)\top} \\ \vdots \\ \hat{x}^{(n)\top} \end{bmatrix} \theta \right)$$

$\underbrace{\hspace{10em}}$

$\mathbb{R}^{n \times 1}$

$$= \frac{1}{2} (y - x\theta)^\top (y - x\theta)$$

$$= \frac{1}{2} (y^\top - \theta^\top x^\top) (y - x\theta)$$

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

$$L(\theta) = \frac{1}{2} \left(y^\top y - \underbrace{y^\top x\theta}_{1 \times 1} - \theta^\top x^\top y + \theta^\top x^\top x\theta \right)$$

$$= \frac{1}{2} \left(y^\top y - \underbrace{2y^\top x\theta}_{(1 \times n)(n \times 2)(2 \times 1)} + \underbrace{\theta^\top x^\top x\theta}_{x^\top A x} \right)$$

$$\frac{\delta L}{\delta \theta} = \frac{1}{2} \left(-(2y^\top x)^\top + \left[(x^\top x) + (x^\top x)^\top \right] \theta \right)$$

$$= \frac{1}{2} (-2x^\top y + x^\top x\theta + x^\top x\theta)$$

<p>note: $\hat{y} = \theta^\top x$</p> $\frac{\partial \hat{y}}{\partial x} = \theta$ <p>$\hat{y} = x^\top A x$</p> $\frac{\partial \hat{y}}{\partial x} = (A + A^\top)x$

$$\frac{\partial L}{\partial \theta} = (-x^T y + x^T x \theta)$$

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

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

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

VECTOR AND MATRIX DERIVATIVES:

$$* \quad \frac{\partial \underline{y}}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix} \quad \vec{x} \in \mathbb{R}^n$$

$\frac{\partial \underline{y}}{\partial x} \in \mathbb{R}^n$

$$\Delta y \approx \frac{\partial y}{\partial x_1} \cdot \Delta x_1 + \frac{\partial y}{\partial x_2} \cdot \Delta x_2 + \dots$$

$$+ \frac{\partial y}{\partial x_n} \cdot \Delta x_n$$

$$\Delta y \approx (\nabla_{\vec{x}} \underline{y}^T) \Delta \vec{x}$$

$$* \quad y = \vec{\theta}^T \vec{x} \quad \vec{\theta} \in \mathbb{R}^n \quad \vec{x} \in \mathbb{R}^n$$

$$\frac{\partial \underline{y}}{\partial x} \in \mathbb{R}^n$$

$\hookrightarrow (n \times 1) \times (1 \times n)$

$$y = \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\frac{\partial y}{\partial x_1} = \theta_1, \quad \frac{\partial y}{\partial x_2} = \theta_2$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} = \theta$$

* $y = \vec{x}^T A \vec{x}$

\vec{x}	\downarrow	A	\downarrow	$\frac{\partial y}{\partial x}$
$(1 \times n)$	$(n \times n)$	$(n \times 1)$		$(n \times 1) \times (1 \times n) : (n \times 1)$

$$[x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ A_{21} & \ddots & & \\ \vdots & & \ddots & \\ A_{n1} & & & A_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$= \left[\sum_{i=1}^n A_{i1} x_i \quad \sum_{i=1}^n A_{i2} x_i \quad \dots \quad \sum_{i=1}^n A_{in} x_i \right] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$= \sum_{j=1}^n \sum_{i=1}^n A_{ij} x_i x_j$$

$$\frac{\partial L}{\partial x_i} = 2A_{11}x_1 + \sum_{j=2}^n A_{1j}x_j + \sum_{i=2}^n A_{i1}x_i$$

$(i=j=1)$ $(i=1, j \neq 1)$ $(i \neq 1, j=1)$

$$= \sum_{i=1}^n A_{ii}x_i + \sum_{j=1}^n A_{ij}x_j$$

$$\left[\begin{array}{c} (A_{11}x_1 + A_{12}x_2 + \dots + A_{1n}x_n) \\ (A_{21}x_1 + A_{22}x_2 + \dots + A_{2n}x_n) \\ \vdots \\ (A_{n1}x_1 + A_{n2}x_2 + \dots + A_{nn}x_n) \end{array} \right] \quad \left[\begin{array}{c} (A_{11}x_1 + A_{12}x_2 + \dots + A_{1n}x_n) \\ (A_{21}x_1 + A_{22}x_2 + \dots + A_{2n}x_n) \\ \vdots \\ (A_{n1}x_1 + A_{n2}x_2 + \dots + A_{nn}x_n) \end{array} \right]$$

$$= \begin{bmatrix} A_{11} & A_{21} & \dots & A_{n1} \\ A_{12} & & \ddots & \\ \vdots & \ddots & & \\ A_{1n} & \dots & & A_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ A_{21} & & \ddots & \\ \vdots & & & \\ A_{n1} & \dots & & A_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$= A^T x \quad = \quad A \quad x$$

$$\frac{\partial f}{\partial x} = (A + A^T)x$$

* $f = \vec{y}^T A \vec{x}$ $\frac{\partial f}{\partial A}$

$(1 \times m)$ $(m \times n)$ $(n \times 1)$

$$\begin{bmatrix} y_1 & y_2 & \dots & y_m \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ \vdots & \ddots & & \vdots \\ A_{m1} & \dots & \dots & A_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\sum_{i=1}^m \sum_{j=1}^n A_{ij} y_i x_j$$

$$\frac{\partial f}{\partial A_{11}} = y_1 x_1$$

$$\frac{\partial f}{\partial A} = \begin{bmatrix} y_1 x_1 & y_1 x_2 & \dots & y_1 x_n \\ \vdots & \ddots & & \vdots \\ y_m x_1 & \dots & \dots & y_m x_n \end{bmatrix}$$

$$= \begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_m \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}$$

$m \times 1$

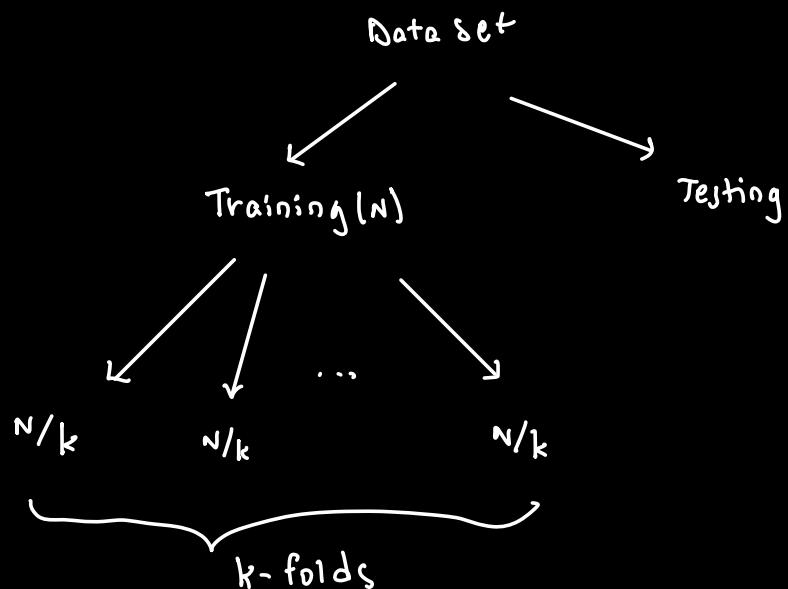
$$= \begin{matrix} \gamma \\ (m \times 1) \end{matrix} \begin{matrix} x^T \\ (1 \times n) \end{matrix}$$

REDUCE OVERFITTING:

→ Regularization

→ More data

K-FOLD CROSS VALIDATION:



→ Train on $(k-1)$ folds

→ Validation on 1 fold (Validation fold)

MAXIMUM LIKELIHOOD ESTIMATION:

Data : H T H H T T H T

Model : $x^{(i)} :$ $\begin{cases} T, & \text{with } p = (1-\theta) \\ H, & \text{with } p = \theta \end{cases}$

Models : Model 1 : $\theta = 1$

Model 2 : $\theta = 0.75$

Model 3 : $\theta = 0.5$

$$MLE = P(\text{Observing the data} | \text{Model}) = \theta^4 (1-\theta)^4$$

Model 1 :

$$MLE = 1 \times 0 \times 1 \times 1 \times 0 \times 0 \times 1 \times 0 = 0$$

Model 2 :

$$MLE = 0.75^4 \times 0.25^4$$

Model 3 :

$$MLE = 0.5^4 \times 0.5^4 \quad \checkmark \quad \text{highest}$$

We can find this also by

$$\frac{\partial (MLE)}{\partial \theta} = 0$$

$$\log \text{MLE} = 4\log \theta + 4\log(1-\theta)$$



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

$$\frac{4}{\theta} + \frac{4}{1-\theta} = 0$$

$$\frac{1}{\theta} = \frac{1}{\theta - 1}$$

$$\theta = 0.5 //$$

MLE - CLASSIFICATION:

We will try to model each class as Gaussian and optimize the parameters for these Gaussian models.

① Assume each class ~ Gaussian

② Find best possible Gaussian model

Class i:

$$x^{(i)} | y^{(i)} = i \sim N(\mu_i, \Sigma_i)$$

N examples
n classes

$\mu_i, \Sigma_i \rightarrow$ PARAMETERS

n classes \Rightarrow 2n parameters.

$$\boldsymbol{\xi}_i \rightarrow \begin{bmatrix} \text{var(dim}_1) & \text{cov} \\ \text{cov(dim}_1, & \text{var(dim}_2) \\ \text{dim}_2) \end{bmatrix}$$

Assume:

$$P(y^{(i)}) = 1/k \quad \} \text{equally likely classes.}$$

$$\text{MLE: } L = \prod_{i=1}^n P(x^{(i)}, y^{(i)} | \theta)$$

$$\theta = (\mu_1, \Sigma_1, \mu_2, \Sigma_2, \dots, \mu_k, \Sigma_k)$$

$$\log L = \sum_{i=1}^n \log P(x^{(i)}, y^{(i)} | \theta)$$

$$= \sum_{i=1}^n \log \left(\underbrace{P(y^{(i)} | \theta)}_{1/k} \cdot P(x^{(i)} | y^{(i)}, \theta) \right)$$

$$= \sum_{i=1}^n \log P(x^{(i)} | y^{(i)}, \theta) + \underbrace{\sum_{i=1}^n \log (P(y^{(i)} | \theta))}_{k_i}$$

$$\log P(x^{(i)} | y^{(i)} = j) = \log \frac{1}{\sqrt{(2\pi)^n |\Sigma_j|}} e^{-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)}$$

$$= \underbrace{\frac{-n}{2} \log 2\pi}_{k_2} - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)$$

$$\log L = \sum_{j=1}^n -\frac{1}{2} \log |\Sigma_{y^{(i)}}| - \frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \Sigma_{y^{(i)}}^{-1} (x^{(i)} - \mu_{y^{(i)}}) + k$$

$$\mu_j = \frac{1}{n} \sum_{i: y^{(i)} = j} x^{(i)}$$

Find μ_i, Σ_i for all $i: 1 \text{ to } n$

TESTING:

Find class of x

Find j such that $\Pr(y=j|x)$ is maximum.

$$\underset{j}{\operatorname{argmax}} \Pr(y=j|x) = \underset{j}{\operatorname{argmax}} \frac{\Pr(y=j, x)}{\Pr(x)}$$

$$= \underset{j}{\operatorname{argmax}} \frac{P(j) \overset{y_j}{\overbrace{P(x|j)}}}{P(x)}$$

↳ independent of j

$$= \underset{j}{\operatorname{argmax}} P(x|j)$$

K- NEAREST NEIGHBORS:

x : Image $32 \times 32 \times 3$ \mathbb{R}^{3072}

y : classes 1 - 10

TRAINING: Store all $(x^{(i)}, y^{(i)})$

PRO:

- * Simple, fast

CONS:

- * Memory-intensive

TESTING: → Find k nearest neighbor of x

→ Choose most common y .

PRO:

- * Simple

CONS:

- * Scale with amount of training data

HYPERPARAMETERS :

→ k

→ distance metric.

OTHER CONS:

→ distance (Say shift image) → not accurate

→ Curse of dimensionality → Distance doesn't make sense.

LINEAR CLASSIFICATION:

x : Image $32 \times 32 \times 3$ \mathbb{R}^{3072} $n = 3072$

y : score to be in each class (1 to 10)
 $c = 10$

$$x \in \mathbb{R}^n$$

$$y \in \mathbb{R}^c$$

$$y = w x + b$$

$$w \in \mathbb{R}^{c \times n}$$

$$b \in \mathbb{R}^c$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_c \end{bmatrix}_{c \times 1} = \begin{bmatrix} w_1^T x + b_1 \\ w_2^T x + b_2 \\ \vdots \\ w_c^T x + b_c \end{bmatrix}_{c \times n} = \begin{bmatrix} -w_1^T - \\ -w_2^T - \\ \vdots \\ -w_c^T - \end{bmatrix}_{c \times n} \begin{bmatrix} x \end{bmatrix}_{n \times 1} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_c \end{bmatrix}_{c \times 1}$$

So each class i has

weight $\rightarrow w_i$

bias $\rightarrow b_i$

Answer: $\arg\max_i (y_i)$

Why not use loss function (least squares) like regression?

$$\text{if } \hat{f} = 7 \quad \text{but } \hat{f} = 6$$

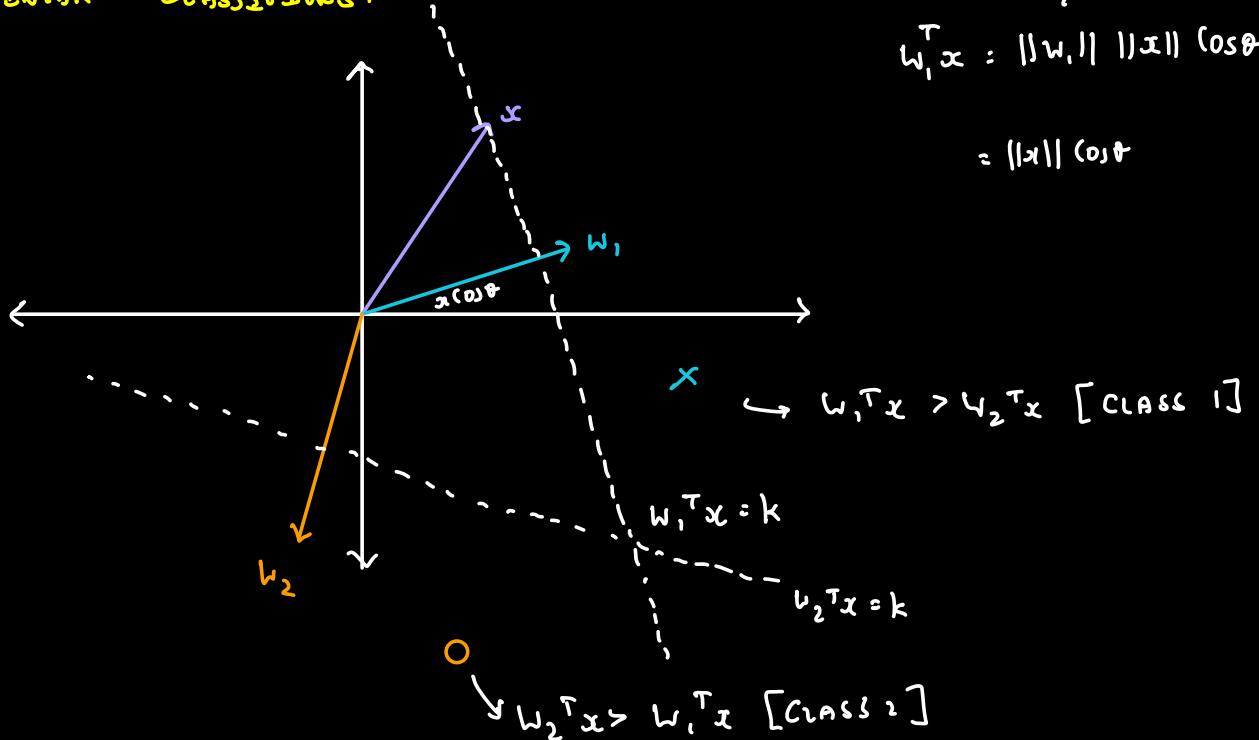
$$\text{error: } (6-7)^2 = 1$$

But if $\hat{f} = 1$

$$\text{error: } (1-7)^2 = 36$$

but both are equally wrong!

LINEAR CLASSIFIERS:



Cons:

→ Fails on
not linearly
separable data.

$\begin{array}{c c} 0 & x \\ \hline x & 0 \end{array}$	or	$\begin{array}{c c c} 0 & 0 & 0 \\ \hline 0 & x & x \\ \hline 0 & 0 & 0 \end{array}$
--	----	--

SOFTMAX CLASSIFIER:

We have $\mathbf{y} = [y_1 \ y_2 \dots \ y_c]^T$

y_i is a score

But to make more intuitive sense, we can convert score to probability. [It is normalized now too].

$$y_i = a_i(x) = w_i^T x + b_i$$

$$\text{softmax}_i(x) = \frac{e^{a_i(x)}}{\sum_{j=1}^c e^{a_j(x)}}$$

$$\text{Prob}(y^{(j)} = i \mid x^{(j)}, \theta) = \text{softmax}_i(x^{(j)})$$

$$P(x, y | \theta) = \prod_{i=1}^m P(x^{(i)}, y^{(i)} | \theta)$$

$$= \prod_{i=1}^m P(x^{(i)} | \theta) P(y^{(i)} | x^{(i)}, \theta)$$

$m \rightarrow$ number of examples
 $c \rightarrow$ classes
 $n \rightarrow$ number of features

$$\theta = \{w, b\}$$

MLE on this:

$$\hat{\theta} = \arg \max_{\theta} \prod_{i=1}^m \underbrace{p(x^{(i)} | \theta) p(y^{(i)} | x^{(i)}, \theta)}_{\text{independent of } \theta}$$

$$= \arg \max_{\theta} \prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}, \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^m \log \frac{\text{softmax}(x^{(i)})}{y^{(i)}}$$

$$= \arg \max_{\theta} \sum_{i=1}^m \log \left[\frac{e^{a_{y^{(i)}}(x^{(i)})}}{\sum_{j=1}^c e^{a_j(x^{(i)})}} \right]$$

$$= \arg \max_{\theta} \frac{1}{m} \sum_{i=1}^m \left[a_{y^{(i)}}(x^{(i)}) - \log \sum_{j=1}^c e^{a_j(x^{(i)})} \right]$$

$$= \arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m \left[\log \sum_j e^{a_j(x^{(i)})} - a_{y^{(i)}}(x^{(i)}) \right]$$

Loss Function:

$$L = \arg \min_{\theta} \frac{1}{m} \sum_{i=1}^m \left[\log \sum_j e^{a_j(x^{(i)})} - a_{y^{(i)}}(x^{(i)}) \right]$$

For one example $x^{(i)} - y^{(i)}$

$$\text{MLE} = \log p_{\pi}(y^{(i)} | x^{(i)}) = a_{y^{(i)}}(x^{(i)}) - \log \sum_j e^{a_j(x^{(i)})}$$

→ If $a_{y^{(i)}}(x^{(i)})$ is the largest score

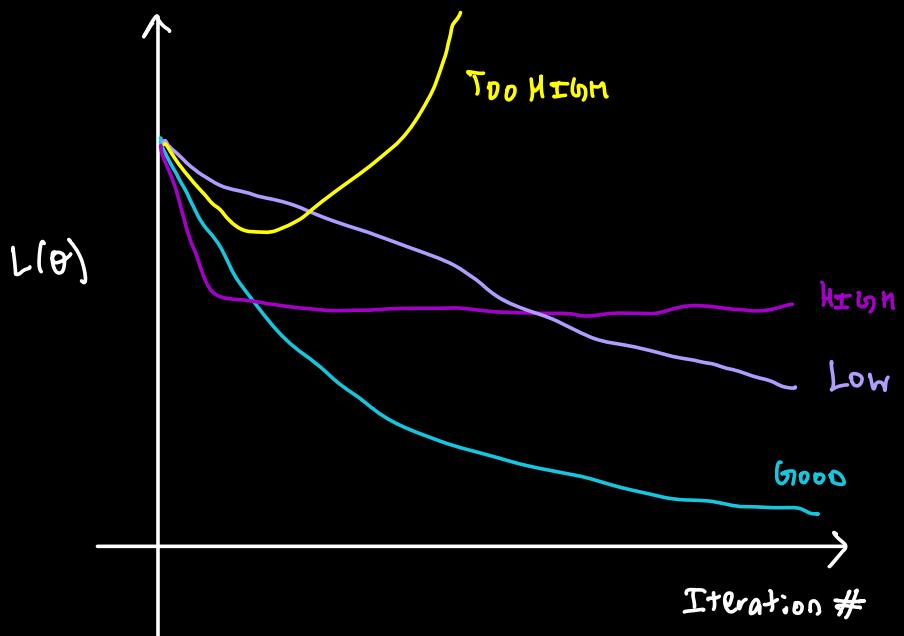
then $\text{MLE} \approx 0$ (optimum)

Loss ≈ 0

→ Else $\text{MLE} < 0$, loss > 0 .

GRADIENT DESCENT:

Interpreting cost function to find learning rate



Why use analytical gradient over numerical?

$$\rightarrow \frac{\partial f}{\partial x}$$

$$\rightarrow \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

- * # of parameters high
- * f may be expensive.

BATCH GRADIENT DESCENT:

Each data provides a noisy estimate
of the gradient at that point.

So, Take all m examples, average the $\frac{\partial f}{\partial w}$ and
apply to w , for each iteration. (Least noisy
estimate).

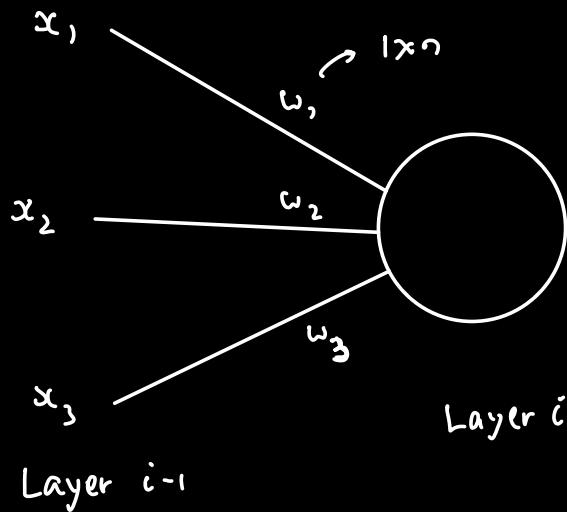
STOCHASTIC / MINI-BATCH GRADIENT DESCENT:

Noisier estimate of gradient with fewer examples.

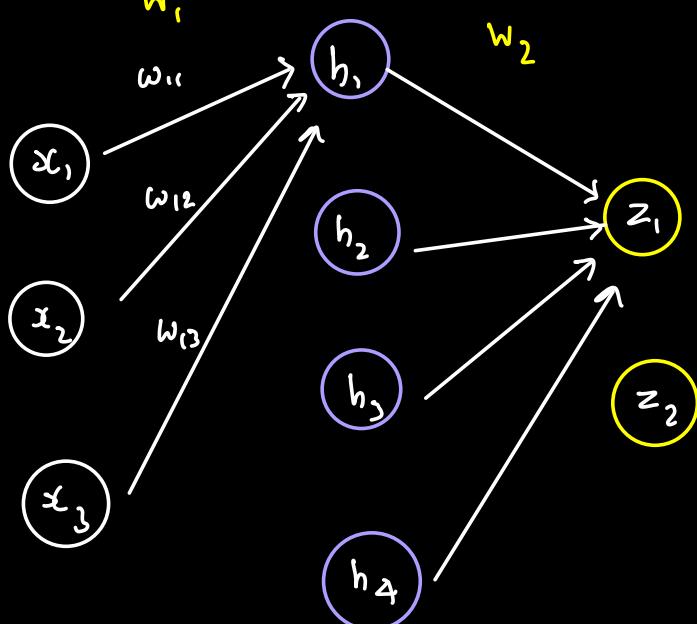
\Rightarrow Regularizes

\Rightarrow Faster

NEURAL NETWORK:



$$f(w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + b)$$



LAYERS = 2

NEURONS = $4 + 2 = 6$

WEIGHTS = $3 \times 4 + 4 \times 2 = 20$

BIASES = $4 + 2 = 6$

PARAMETERS = $20 + 6 = 26$

$w_{i3} \rightarrow$ weight from 3rd neuron

in $i-1^{\text{th}}$ layer to 1st

neuron in i^{th} layer

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \\ h_4 \end{bmatrix} = f \left(\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ \vdots & \ddots & \vdots \\ w_{41} & \dots & w_{43} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_4 \end{bmatrix} \right)$$

$$h = f(w_x + b)$$

$(n_h \times n_x)$
 $(n_h \times 1)$ \downarrow $\downarrow (n_h \times 1)$
 $(n_x \times 1)$

$$z = w_2 h + b_2$$

$(n_b \times 1)$ \downarrow $\downarrow (n_b \times 1)$
 $(n_h \times 1)$

AUTO ENCODER:

$$\begin{array}{ccc} x & & z \\ 0 & h & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ & 0 & 0 \end{array}$$

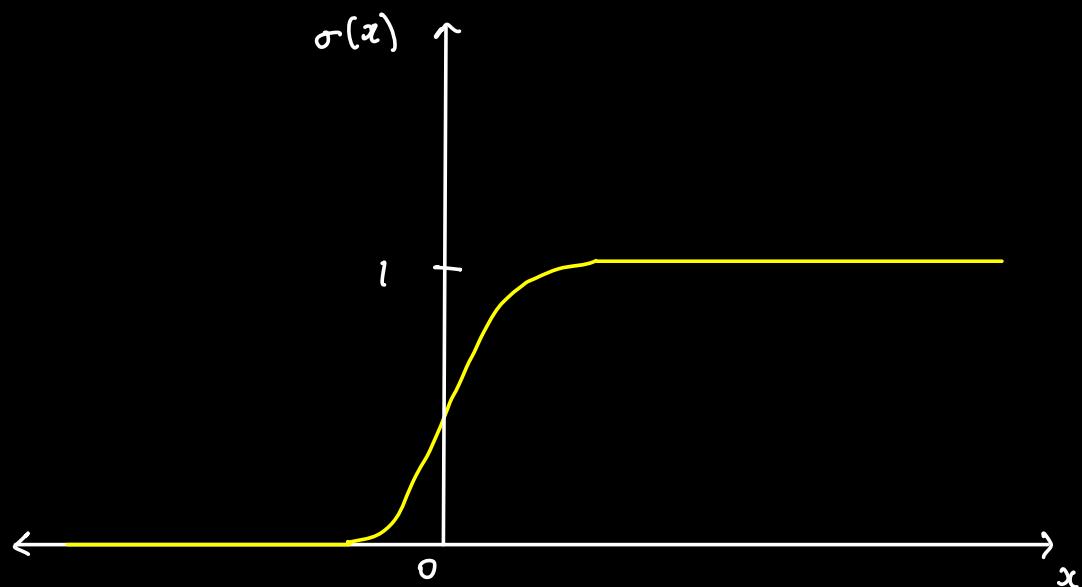
$$L : \|z - x\|^2$$

Dimensionality Reduction.

ACTIVATION FUNCTIONS:

→ SIGMOID ACTIVATION:

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad [\text{Softmax for 2 classes}]$$



$$\begin{aligned}\frac{\partial}{\partial x} \sigma(x) &= \frac{-1}{(1+e^{-x})^2} (-e^{-x}) \\ &= \frac{1}{1+e^{-x}} \times \left(\frac{e^{-x}}{1+e^{-x}} \right)^{-1} \\ &= \sigma(x) (1 - \sigma(x))\end{aligned}$$

DISADVANTAGE:

$$h_{ii} = \sigma(\omega) = \sigma(\omega^T x + b) \quad a = \omega^T x + b$$

$$\frac{\partial L(\omega)}{\partial \omega} = \underbrace{\frac{\partial \sigma(\omega)}{\partial \omega}}_{\text{ }} \cdot \frac{\partial L(\omega)}{\partial \sigma(\omega)}$$

$$\underbrace{\frac{\partial \sigma(\omega)}{\partial a}}_{\text{ }} \cdot \frac{\partial a}{\partial \omega} \xrightarrow{x}$$

$$\sigma(a)(1 - \sigma(a))$$

If $\sigma(\omega^T x + b)$ is close to 0 or 1,

then $\frac{\partial \sigma(\omega)}{\partial a} \approx 0$ and hence

$$\frac{\partial L(\omega)}{\partial \omega} \rightarrow \text{very small number}$$

$$\omega \leftarrow \omega - \varepsilon \underbrace{\frac{\partial L(\omega)}{\partial \omega}}_{\approx 0} \quad \begin{matrix} \omega \text{ won't update} \\ \Rightarrow \text{no learning} \end{matrix}$$

Pros:

→ At $x=0$, linear $\frac{\partial \sigma}{\partial x} \approx 1$ (large).

→ Differentiable everywhere

Cons:

→ At extremes, $\frac{\partial \sigma}{\partial x} \approx 0$ No learning

→ Output not zero-centred (0.5)

↳ activation function ≥ 0

"ZIN ZADING GRADIENT DESCENT"

$x \xrightarrow{w, \sigma} h_1 \xrightarrow{w, \sigma} h_2 \xrightarrow{w} \text{OUTPUT}$

$$h_2 = f(w) = \sigma(w^T h_1 + b) = \sigma(z)$$

where $z = w^T h_1 + b$

$$\frac{\partial h_2}{\partial w} = \frac{\partial z}{\partial w} \cdot \frac{\partial h_2}{\partial z}$$
$$h_1 \quad \downarrow \quad \sigma(z) (1 - \sigma(z))$$
$$\downarrow \quad \geq 0 \quad \geq 0$$

$$\sigma(z) \\ \geq 0$$

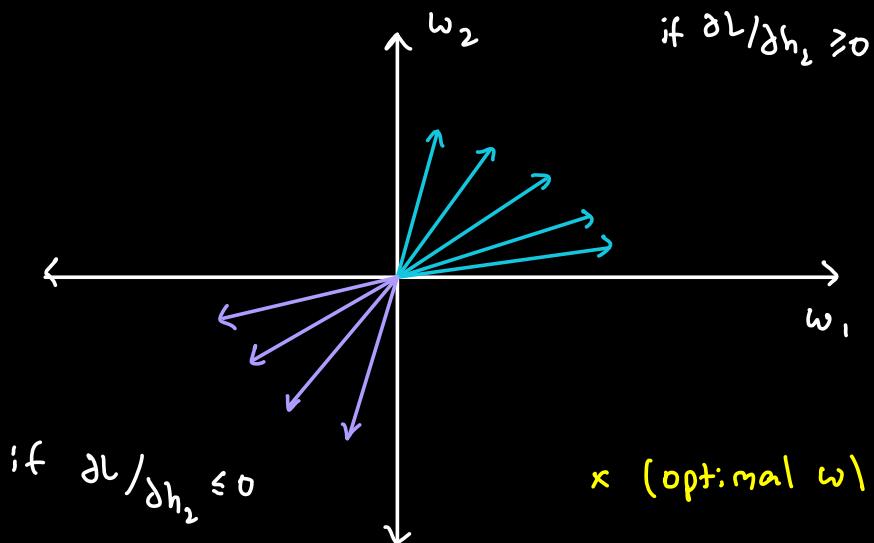
$$\frac{\partial h_2}{\partial w} \geq 0$$

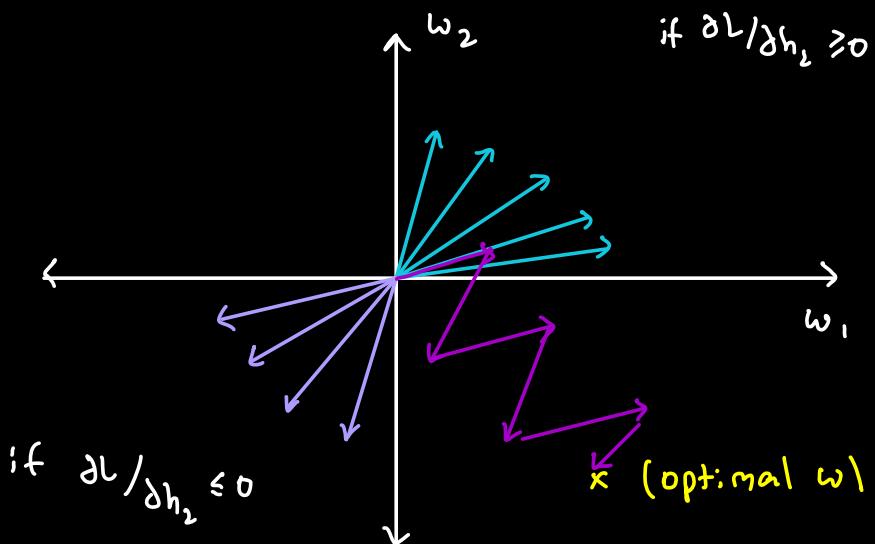
$$\frac{\partial L}{\partial w} = \underbrace{\frac{\partial h_2}{\partial w}}_{\geq 0} \cdot \underbrace{\frac{\partial L}{\partial h_2}}_{\text{Scalar}}$$

so if $\partial L / \partial h_2 \geq 0$ all $\partial L / \partial w \geq 0$

≤ 0 all $\partial L / \partial w \leq 0$

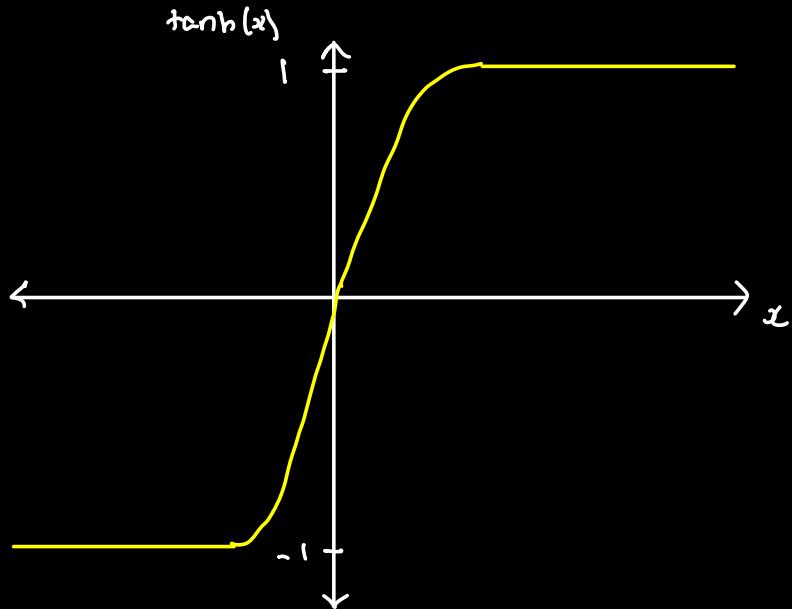
Assume 2 w 's





→ HYPERBOLIC TANGENT

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} - 1$$



$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} - 1$$

$$= \frac{2 - 1 - e^{-x}}{1 + e^{-x}}$$

$$= \frac{1 - e^{-x}}{1 + e^{-x}}$$

$$\frac{\partial \tanh(x)}{\partial x} = \frac{(1 + e^{-x})(e^{-x}) - (1 - e^{-x})(-e^{-x})}{(1 + e^{-x})^2}$$

$$= \frac{e^{-x} + e^{-2x} + e^{-x} - e^{-2x}}{(1+e^{-x})^2}$$

$$= \frac{1+e^{-x} - (1-e^{-x})}{(1+e^{-x})^2}$$

$$= 1 - \tanh^2(x)$$

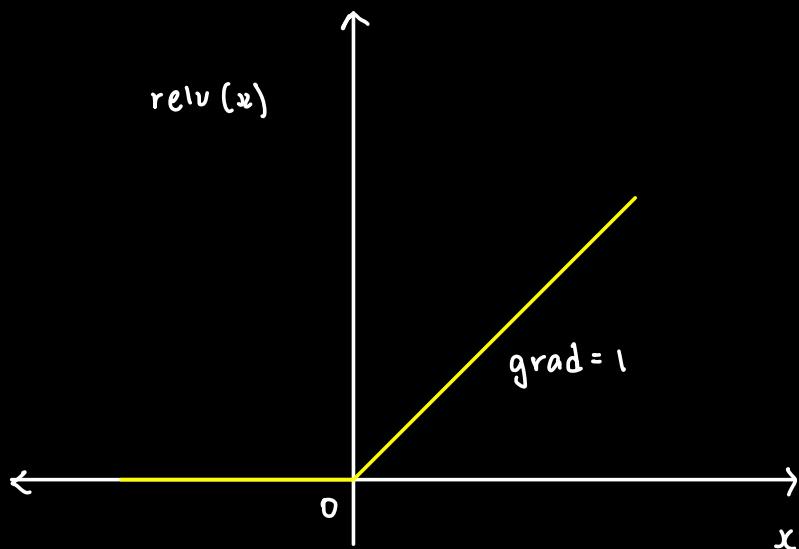
PROS :

- $x=0$ Linear
- Differentiable
- Zero-centered

CONS :

- saturates at large / small values \rightarrow no learning.

→ ReLU UNIT:



$$\text{relu}(x) := \max(0, x)$$

$$\frac{\partial \text{relu}(x)}{\partial x} = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

Pros:

- Converges Faster
- Linear when active
 - ↳ derivative large
- No saturation if $x > 0$

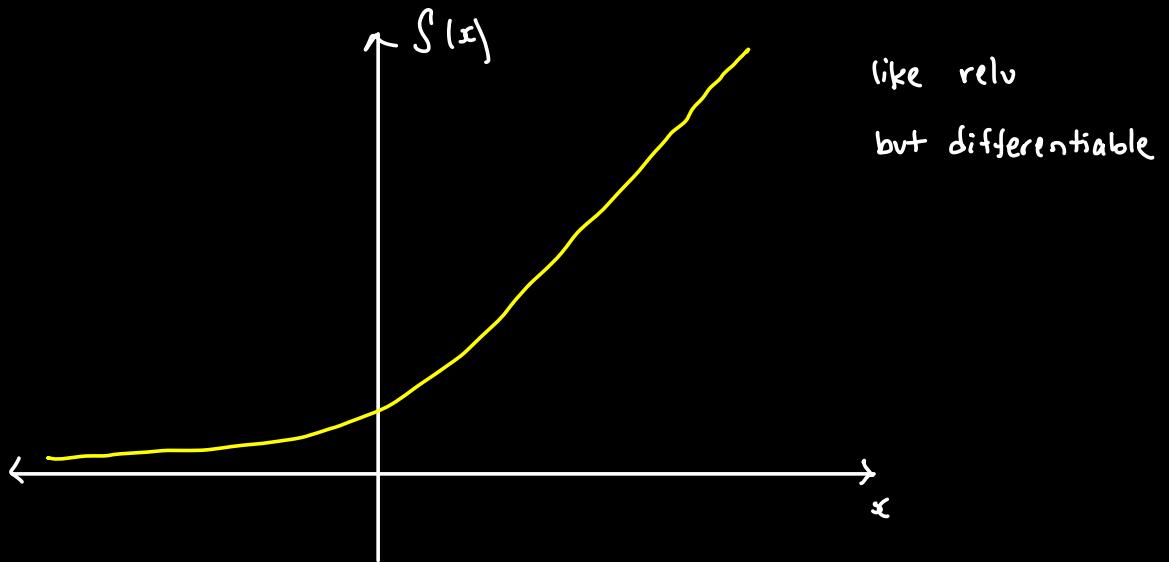
Cons :

- not zero-centric (zig-zag)
- not differentiable at $x=0$
- zero activation \rightarrow no learning.

→ Softplus Unit:

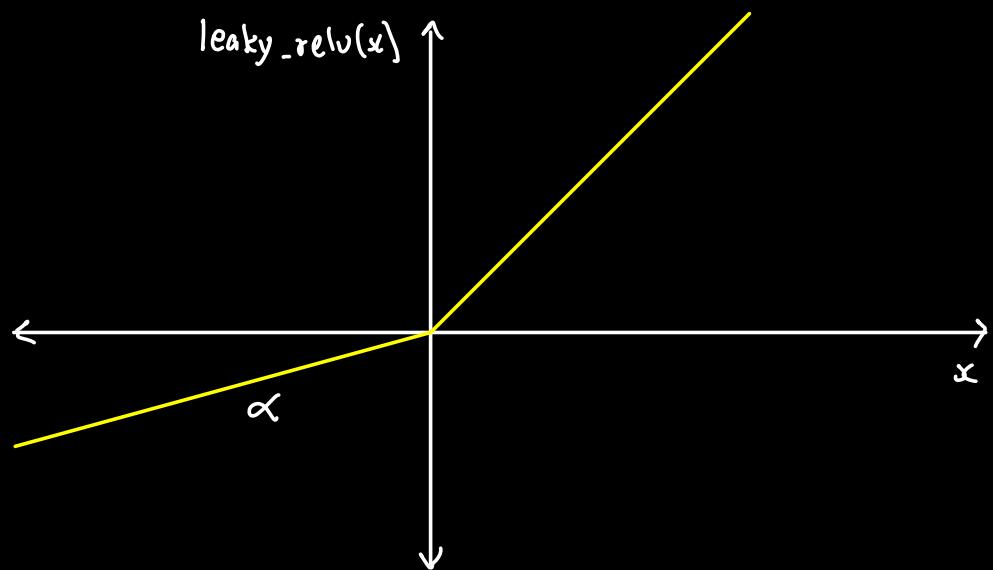
$$S(x) = \log(1 + e^x)$$

$$\frac{\partial}{\partial x} S(x) = \sigma(x)$$



→ Leaky ReLU / PReLU

$$f(x) = \max(\alpha x, x)$$



if $\alpha \rightarrow$ parameter PReLU

\rightarrow hyperparameter leaky ReLU

\rightarrow ELU

$$f(x) = \max(\alpha(e^x - 1), x)$$



Con: compute exponential
(expensive).

Output Activations:

Assume 2 class classification

→ Apply Sigmoid

$$\hat{y} = \sigma(z)$$

which error →

$$\text{MSE} = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \sigma(z^{(i)}))^2$$

$$\text{CE} = - \sum_{i=1}^n \left[y^{(i)} \log \sigma(z^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(z^{(i)})) \right]$$

↓

MLE ~ softmax - 2 classes

↓

$$P(x^{(i)}, y^{(i)} | \theta) = \prod_{i=1}^n P(x^{(i)} | \theta) \cdot P(y^{(i)} | x^{(i)}, \theta)$$

$$P(\text{class 1} | x^{(i)}) = \frac{1}{1 + e^{-x^{(i)}}} = \sigma(x^{(i)})$$

$$P(\text{class 2} | x^{(i)}) = 1 - \sigma(x^{(i)})$$

$$P(y^{(i)} | x^{(i)}, \theta) = \sigma(x^{(i)})^{y^{(i)}} \cdot (1 - \sigma(x^{(i)}))^{(1 - y^{(i)})}$$

MLE

$$\log P = \log \prod_{i=1}^n P(y^{(i)} | x^{(i)}, \theta)$$

$$= \sum_{i=1}^n y^{(i)} \cdot \sigma(x^{(i)}) + (1 - y^{(i)}) (1 - \sigma(x^{(i)}))$$

which error \rightarrow

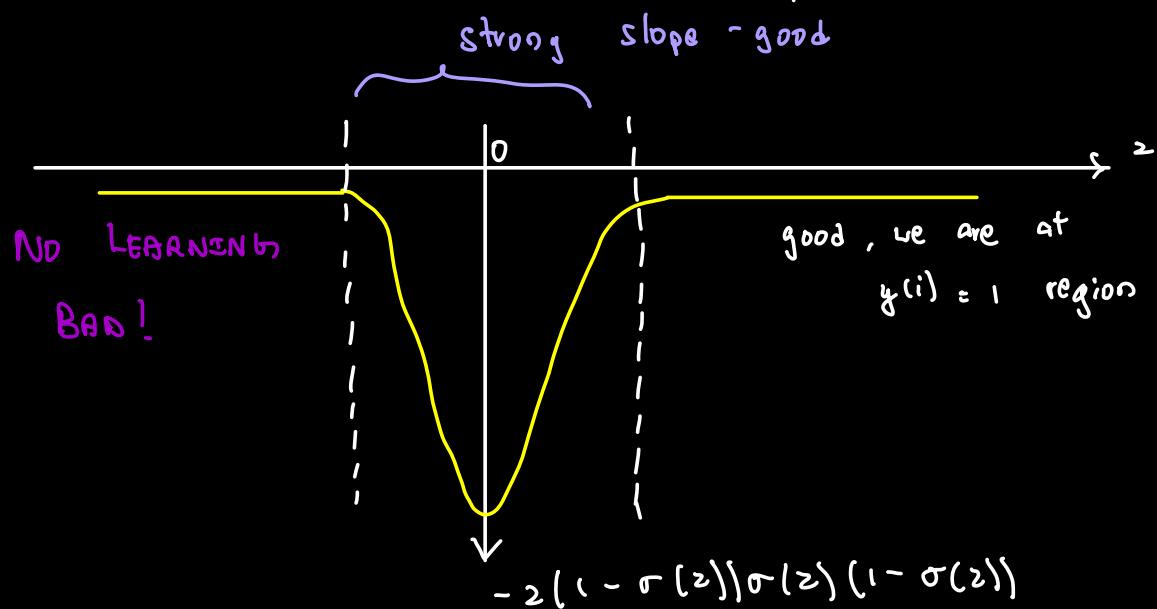
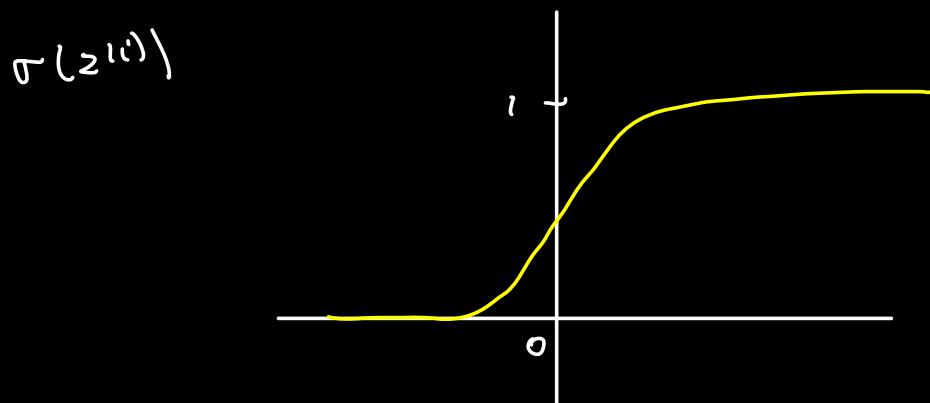
$$\begin{cases} MSE = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \sigma(z^{(i)}))^2 \\ CE = - \sum_{i=1}^n \left[y^{(i)} \log \sigma(z^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(z^{(i)})) \right] \end{cases}$$

(Cross-Entropy)

MSE :

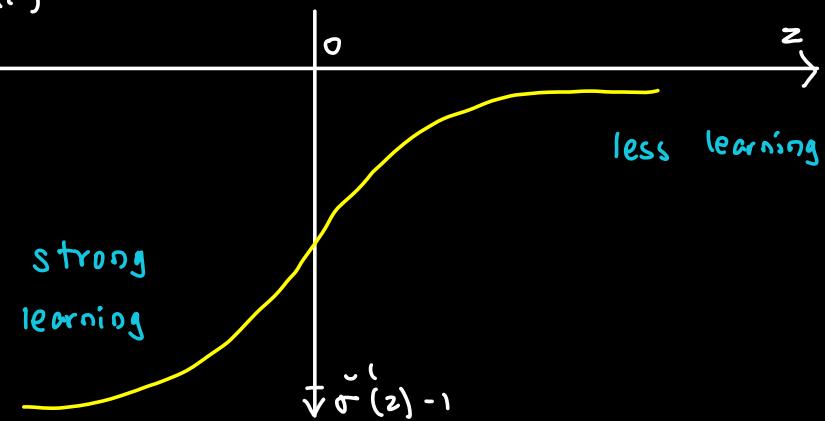
$$\frac{\partial MSE}{\partial z^{(i)}} = -2(y^{(i)} - \underbrace{\sigma(z^{(i)})}_{\partial MSE / \partial \sigma}) \underbrace{(\sigma(z^{(i)}) (1 - \sigma(z^{(i)})))}_{\partial \sigma / \partial z}$$

if $y^{(i)} = 1$



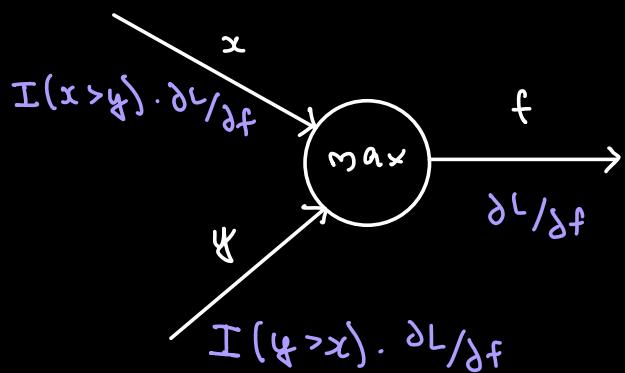
CE:

$$\frac{\delta CE}{\delta z^{(i)}} = \sigma(z^{(i)}) - 1$$



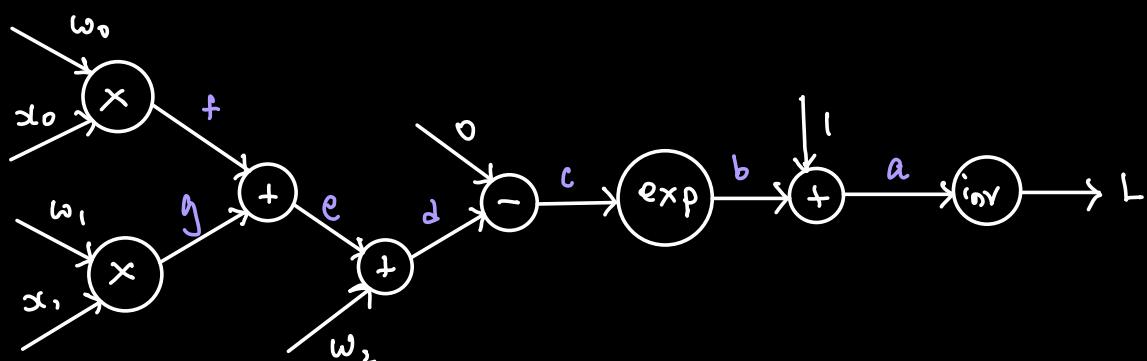
BACK PROPAGATION:

$$\max(x, y)$$



$$I(x > y) = \begin{cases} 1, & x > y \\ 0, & y > x \end{cases}$$

$$L = \frac{1}{1 + e^{-(\omega_0 x_0 + \omega_1 x_1 + \omega_2)}}$$



$$1. \quad L = 1/a$$

$$2. \quad a = 1+b$$

$$3. \quad b = e^c$$

$$\frac{\partial L}{\partial a} = -1/a^2$$

$$\frac{\partial L}{\partial b} = -1/a^2$$

$$\frac{\partial b}{\partial c} = e^c$$

$$\frac{\partial L}{\partial c} = e^c \left(-\frac{1}{a^2} \right)$$

$$4. c = -d$$

$$\frac{\partial L}{\partial d} = e^c/a^2$$

$$5. d = e + \omega_2$$

$$\frac{\partial L}{\partial e} = e^c/a^2$$

$$6. e = f + g$$

$$\frac{\partial L}{\partial f} = \frac{\partial L}{\partial g}$$

$$\frac{\partial L}{\partial \omega_2} = e^c/a^2$$

$$= e^c/a^2$$

$$7. f = \omega_0 x_0$$

$$\frac{\partial f}{\partial \omega_0} = x_0$$

$$8. g = \omega_1 x_1$$

$$\frac{\partial L}{\partial \omega_1} = \frac{e^c}{a^2} \cdot x_1$$

$$\frac{\partial L}{\partial \omega_0} = \frac{\partial f}{\partial \omega_0} \cdot \frac{\partial L}{\partial f} = \frac{e^c x_0}{a^2}$$

$$\frac{\partial L}{\partial x_1} = \frac{e^c \cdot \omega_1}{a^2}$$

$$\frac{\partial L}{\partial x_0} = \frac{e^c \omega_0}{a^2}$$

$$c = -(\omega_0 x_0 + \omega_1 x_1 + \omega_2)$$

$$Q = 1 + e^{-(\omega_0 x_0 + \omega_1 x_1 + \omega_2)}$$

DERIVATIVES OF VECTOR / MATRIX:

$$\vec{y} = \vec{w} \vec{x}$$

↙
↓
→ m x)
 $n \times 1$ $n \times m$

$$\frac{\partial \underline{y}}{\partial x} = \nabla_{\underline{x}} \underline{y} = \begin{bmatrix} | & | & | \\ \partial y_1 / \partial x & \partial y_2 / \partial x & \dots & \partial y_n / \partial x \\ | & | & | \end{bmatrix}$$

$$= \begin{bmatrix} \partial y_1 / \partial x_1 & \partial y_2 / \partial x_1 & \dots & \partial y_n / \partial x_1 \\ \partial y_1 / \partial x_2 & \partial y_2 / \partial x_2 & \dots & \partial y_n / \partial x_2 \\ \vdots \\ \partial y_1 / \partial x_m & \partial y_2 / \partial x_m & \dots & \partial y_n / \partial x_m \end{bmatrix}$$

$$\underline{y} = \underline{w} \underline{x}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1m} \\ \vdots & \ddots & & \vdots \\ w_{n1} & \dots & & w_{nm} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

$$y_1 = w_{11}x_1 + w_{12}x_2 + \dots + w_{1m}x_m = \sum w_{1i}x_i$$

$$y_2 = \sum w_{2i}x_i$$

:

$$y_n = \sum w_{ni}x_i$$

$$\frac{\partial y_1}{\partial x_1} = w_{11}$$

$$\frac{\partial y_1}{\partial x_2} = w_{12}$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_n}{\partial x_2} \\ \vdots & \vdots & & \vdots \\ \frac{\partial y_1}{\partial x_m} & \frac{\partial y_2}{\partial x_m} & \dots & \frac{\partial y_n}{\partial x_m} \end{bmatrix}$$

$$= \begin{bmatrix} w_{11} & w_{21} & \dots & w_{n1} \\ w_{12} & w_{22} & \dots & w_{n2} \\ \vdots & \vdots & & \vdots \\ w_{1m} & w_{2m} & \dots & w_{nm} \end{bmatrix} = w^T \mathbb{Y}$$

TENSOR DERIVATIVES:

$$\mathbf{y} = \mathbf{w}\mathbf{x}$$

\$m \times 1\$
\$n \times 1\$

\$m \times n\$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{w}} \in \mathbb{R}^{m \times n \times m}$$

$$L = \frac{1}{2} \sum_{i=1}^n \| \mathbf{y}^{(i)} - \mathbf{w}\mathbf{x}^{(i)} \|^2$$

$$\text{Find } \frac{\partial L}{\partial w}$$

$$L = \frac{1}{2} \sum_{i=1}^n (\mathbf{y}^{(i)} - \mathbf{w}\mathbf{x}^{(i)})^T (\mathbf{y}^{(i)} - \mathbf{w}\mathbf{x}^{(i)})$$

Let $\mathbf{g}^{(i)} = \mathbf{y}^{(i)} - \mathbf{w}\mathbf{x}^{(i)}$

$$L = \frac{1}{2} \sum_{i=1}^n \mathbf{g}^{(i) T} \mathbf{g}^{(i)}$$

$$\frac{\partial L}{\partial w} = \frac{\partial \mathbf{g}}{\partial w} \cdot \frac{\partial L}{\partial \mathbf{g}}$$

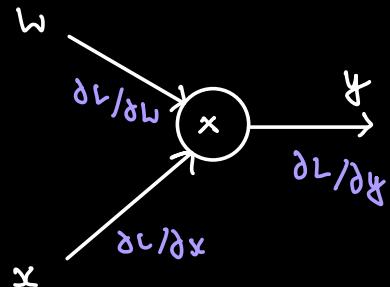
$$m \times n \quad m \times n \times m \quad m \times 1$$

$$\hookrightarrow \mathbf{g}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \mathbf{g}} \mathbf{g}^T$$

$$\frac{\partial}{\partial x} (x^T y) = y$$

$$\frac{\partial}{\partial x} (w x) = w^T$$



$$y = wx$$

↗ ↘ m x 1
n x 1 ↗ n x m

$$\frac{\partial L}{\partial x} = w^T \frac{\partial L}{\partial y}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \cdot x^T$$

$$y = wx$$

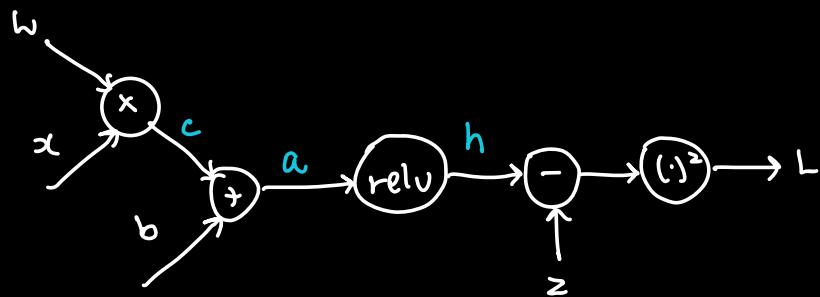
✓ ↘ n x m ↗ m x p
n x p

$$\frac{\partial L}{\partial x} = w^T \frac{\partial L}{\partial y}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \cdot x^T$$

$$L = (h - z)^2$$

$$h = \text{ReLU}(w_x + b)$$



$$\frac{\partial L}{\partial h} = 2(h - z)$$

$$h = \text{ReLU}(a)$$

$$\frac{\partial h}{\partial a} = \mathbb{I}(a > 0)$$

$$\frac{\partial L}{\partial a} = \frac{\partial h}{\partial a} \frac{\partial L}{\partial h} = \mathbb{I}(a > 0) \cdot 2(h - z)$$

$$a = b + c$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial c} = \mathbb{I}(a > 0) \cdot 2(h - z)$$

$$c = w_x$$

$$\frac{\partial L}{\partial x} = w^\top \frac{\partial L}{\partial c} = w^\top \mathbb{I}(a > 0) \cdot 2(h - z)$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial c} \cdot x^\top = \mathbb{I}(a > 0) \cdot 2(h - z) x^\top$$

How TO TRAIN NEURAL NETWORKS BETTER?

- Weight Initialization
- Batch Normalization
- Regularizations
 - L1 + L2 Norm
- Dataset Augmentation
- Model Ensembles
- Dropout.

WEIGHT INITIALIZATION:

- All zeroes
 - All the neurons in hidden unit will be Symmetric
- Random small values

$$h_1 = \text{relu}(w_1, x)$$

$h_2 = \text{relu}(w_2 h_1)$ As w_i are small, activations decay to zero.

Assume $z = w_{10} h_0$

$$\frac{\partial L}{\partial w_{10}} = \frac{\partial L}{\partial z} \cdot h_0^\top \xrightarrow{o's} \frac{\partial L}{\partial w_{10}} = 0$$

No learning.

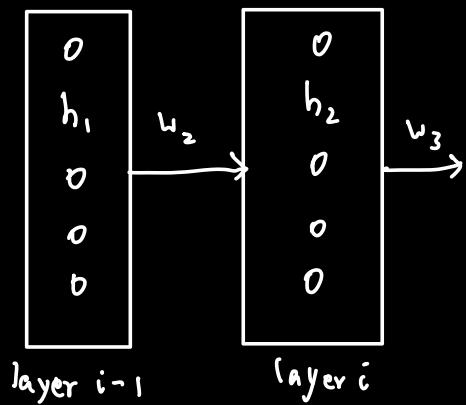
→ Random large values

Activations explode.

$$\frac{\delta L}{\delta w_{10}} = \frac{\delta L}{\delta z} \cdot h_0^T$$

↑ huge , outside the Space-

XAVIER INITIALIZATION:



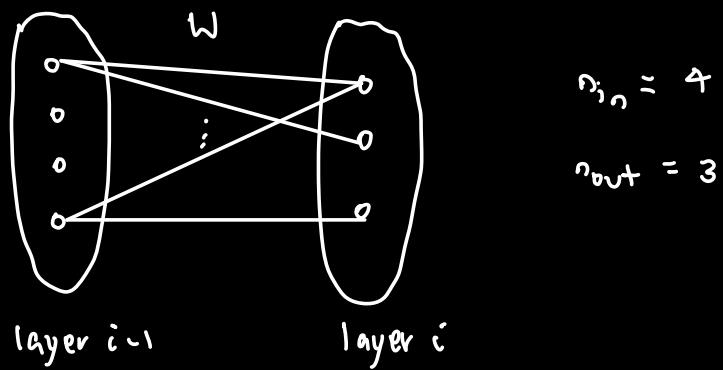
CONDITIONS:

$$\Rightarrow \text{var}(h_1) \approx \text{var}(h_2) \approx \dots \approx \text{var}(h_i)$$

$\Rightarrow \text{var}(\nabla_{w_i} L) \approx \text{var}(\nabla_{w_j} L)$

only for
the first
step.

$$\text{var}(w_{ij}) = \frac{1}{n_{in}}$$
 or $\frac{1}{n_{out}}$ or $\frac{2}{n_{in} + n_{out}}$



$$\text{var}(w_{ij}) = 1/4 \quad \text{or} \quad 1/3 \quad \text{or} \quad 2/7$$

Weights in layer i :

$$\mathcal{N}\left(0, \frac{2}{n_{in} + n_{out}}\right)$$

Pros:

- * Weights not 0 for few layer networks

Cons:

- * Large number of layers \rightarrow weights die down
- * Doesn't work for ReLU
 ↳ ReLU sets all negative to 0, so it kills half the units. So "He" initialization says

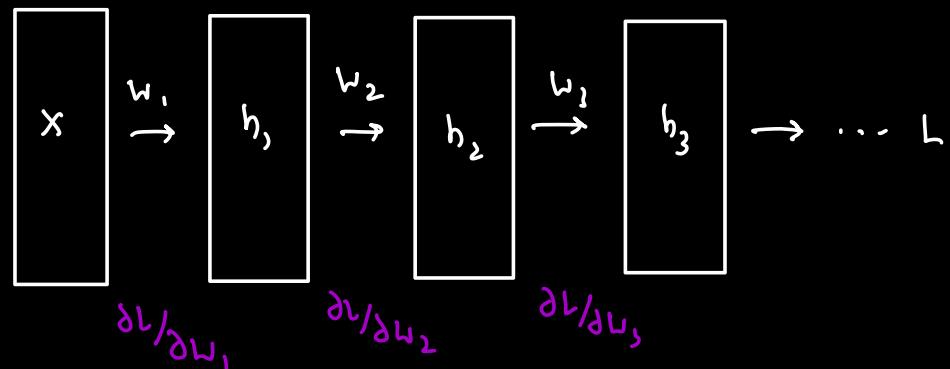
use $\text{var}(w_{ij}) = 2/n_{in} //$

[MSRA]

BATCH NORMALIZATION:

What are we solving?

[Internal Covariate Shift,
Loss Surface Smoother]



We find $\delta L / \delta w_2$ that best decreases final loss,

given the current h_1 . But in gradient descent,

we update w_1 and w_2 simultaneously.

So the new value of w_2 which was best

for the old h_1 may not be best for
the new h_1 that will be calculated

with the updated w_1 .

e.g. h_1 had mean 100. w_2 felt that $w_2 h_1$
was low and to decrease loss, it learns

to increase w_2 . But may be at the same

time w_1 learns and increases too. Now

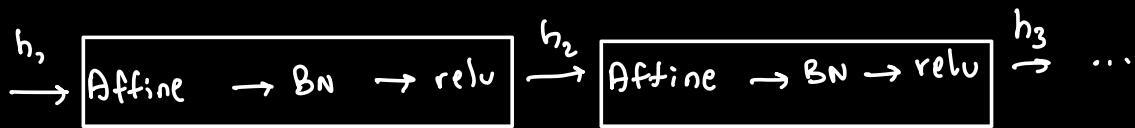
h_1 is large and so is w_2 . This might

have increased the overall loss!

So, we can ensure the statistics of hidden layers are same by normalizing them!

$$E(h_i) = 0$$

$$\text{var}(h_i) = 1$$



PROBLEM:

→ All units have same statistics. But some features might need a different mean and variance.

⇒ Scale and Shift

METHOD:

$$1. \hat{x}_i := \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

$$\mu_i := \frac{1}{m} \sum_{j=1}^m x_i^{(j)}$$

$$\sigma_i^2 := \frac{1}{m} \sum_{j=1}^m (x_i^{(j)} - \mu_i)^2$$

$$2. y_i = \gamma_i \hat{x}_i + \beta_i$$

m - training examples

$$E(y_i) = \beta_i$$

in mini-batch.

$$\text{var}(y_i) = \gamma_i^2$$

So, \rightarrow Bias the solution to make features
with mean=0, variance = 1

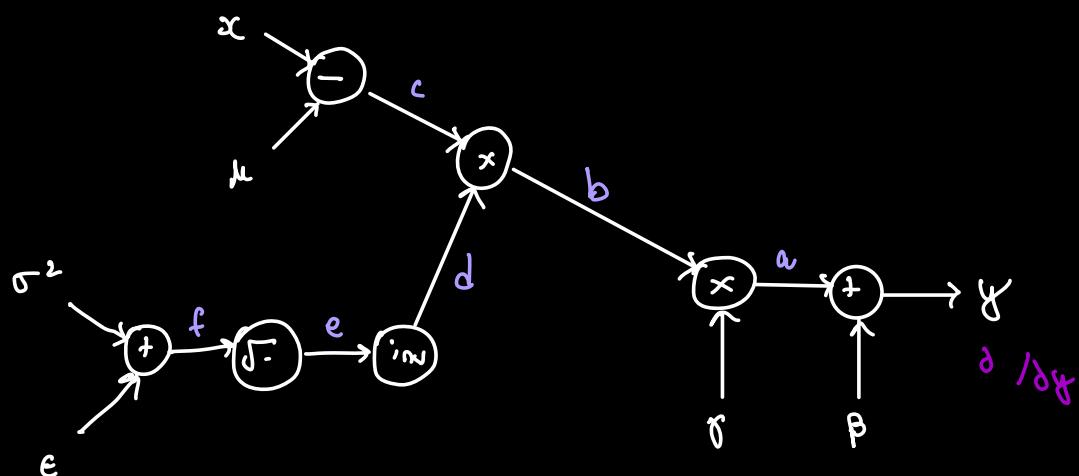
\rightarrow But also allow it to overcome this if
the constraint is too much.

$\beta_i, \delta_i \rightarrow$ learnable parameters.

\rightarrow less sensitive to weight initialization.

\rightarrow Applied separately to EACH NEURON.

BACKPROPAGATION:



$$y = a + \beta$$

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial \beta} = \frac{\partial L}{\partial y}$$

$$a = \delta b$$

$$\frac{\partial a}{\partial \gamma} = b \quad \frac{\partial a}{\partial b} = \gamma$$

$$\frac{\partial L}{\partial b} = \frac{\partial a}{\partial b} \frac{\partial L}{\partial a} = \gamma \frac{\partial L}{\partial y}$$

$$\frac{\partial L}{\partial \gamma} = b \frac{\partial L}{\partial y}$$

$$d = \frac{1}{e}$$

$$\frac{\partial d}{\partial e} = \frac{-1}{e^2}$$

$$\frac{\partial L}{\partial c} = \frac{-1}{e^2} \cdot c \delta \frac{\partial L}{\partial y}$$

$$e = \sqrt{f}$$

$$\frac{\partial e}{\partial f} = \frac{1}{2\sqrt{f}}$$

$$\frac{\partial L}{\partial f} = \frac{1}{2\sqrt{f}} \left(\frac{-1}{e^2} \right) c \delta \frac{\partial L}{\partial y}$$

Each training data returns this gradient. Total is sum across them (Law of total gradient).

$$b = cd$$

$$\frac{\partial L}{\partial c} = d \delta \frac{\partial L}{\partial y}$$

$$\frac{\partial L}{\partial d} = c \delta \frac{\partial L}{\partial y}$$

$$c = x - \mu$$

$$\frac{\partial L}{\partial x} = d \delta \frac{\partial L}{\partial y}$$

$$\frac{\partial L}{\partial \mu} = -d \delta \frac{\partial L}{\partial y}$$

$$f = \sigma^2 + \epsilon$$

$$\begin{aligned} \frac{\partial L}{\partial \sigma^2} &= \frac{1}{2\sqrt{f}} \left(\frac{-1}{e^2} \right) c \delta \frac{\partial L}{\partial y} \\ &= \frac{\partial L}{\partial \epsilon} \end{aligned}$$

$$\frac{\partial L}{\partial \beta} = \sum_{i=1}^m \frac{\partial L}{\partial y^{(i)}}$$

μ contributes to σ^2

$$\frac{\partial L}{\partial \mu} = -\delta \sum_{j=1}^m \frac{\partial L}{\partial y^{(j)}} + \frac{\partial L}{\partial \sigma^2} \cdot \frac{\partial \sigma^2}{\partial \mu}$$

$$\sigma^2 := \frac{1}{m} \sum_{j=1}^m (x^{(j)} - \mu)^2$$

$$\frac{\partial \sigma^2}{\partial \mu} = \frac{-2}{m} \sum_{j=1}^m (x^{(j)} - \mu)$$

$$\frac{\partial L}{\partial \mu} = \frac{-1}{\sqrt{\sigma^2 + \epsilon}} \delta \sum_{j=1}^m \frac{\partial L}{\partial y^{(j)}} + \frac{1}{m(\sigma^2 + \epsilon)^{3/2}} \sum_{j=1}^m (x - \mu) \frac{\partial \delta}{\partial y^{(j)}}$$

$$\frac{\partial L}{\partial x^{(j)}} = \delta \sum_{j=1}^m \frac{\partial L}{\partial y^{(j)}} + \frac{\partial \sigma^2}{\partial x^{(j)}} \cdot \frac{\partial L}{\partial \sigma^2} \rightarrow \frac{\partial \mu}{\partial x^{(j)}} \cdot \frac{\partial L}{\partial \mu}$$

↓

$$2(x^{(j)} - \mu)/m$$

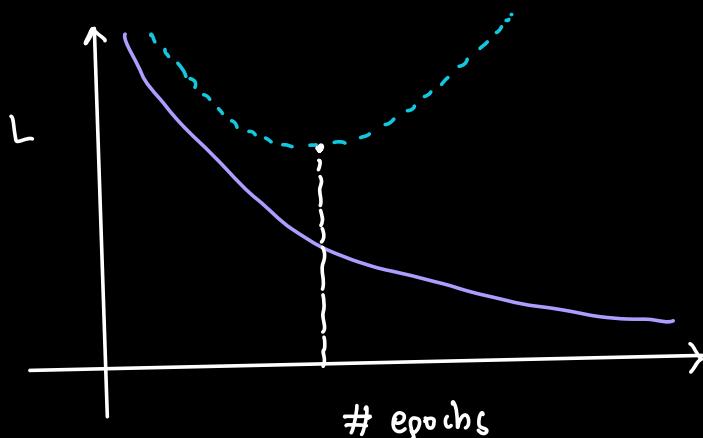
PROS OF BATCH NORMALIZATION:

- Allows higher learning rates (not sure update of other weights would impact, so smaller steps were taken)
- Reduce strong dependencies on initialization.

REGULARIZATION:

Modification to reduce generalization error
but not training error.

EARLY STOPPING:



PARAMETER NORM PENALTIES:

$$J(\theta; x, y) + \alpha \Omega(\theta)$$

→ L₂ REGULARIZATION:

$$\Omega(\theta) = \frac{1}{2} w^T w \quad \frac{1}{2} \|w\|_F^2$$

Ridge Regression / Tikhonov Regularization

$$\frac{\partial \Omega(\theta)}{\partial w} = w$$

$$\text{Cost} : \tilde{J}(w; x, y) = J(w; x, y) + \frac{\alpha}{2} w^T w$$

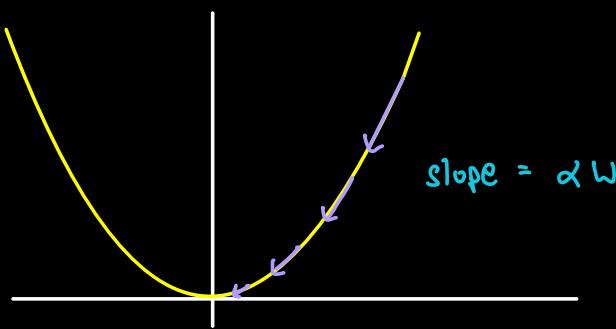
$$\nabla_w \tilde{J}(w; x, y) = \alpha w + \nabla_w J(w; x, y)$$

$$w \leftarrow \underbrace{(1 - \epsilon \alpha)}_{\text{Decreases weights}} w - \epsilon \nabla_w J(w; x, y)$$

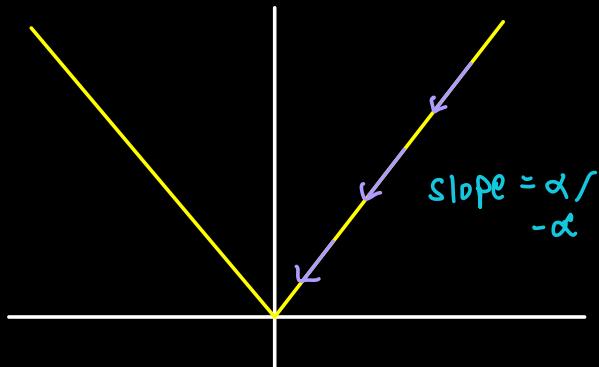
→ L1 REGULARIZATION:

$$\Omega(w) = \|w\|_1 = \sum_i |w_i|$$

→ Tends to push more weights to 0.



2-NORM



1-NORM

Slope is constant, doesn't decrease with weight.

So,

$$\downarrow \\ w \leftarrow w - \epsilon \alpha w - \epsilon \nabla_w J$$

$$w \leftarrow w - \underbrace{\epsilon \alpha}_{\text{constant}} - \epsilon \nabla_w J$$

→ Used for feature selection.

DATA AUGMENTATION:

→ Add cropped /flipped ... images

→ Inject noise

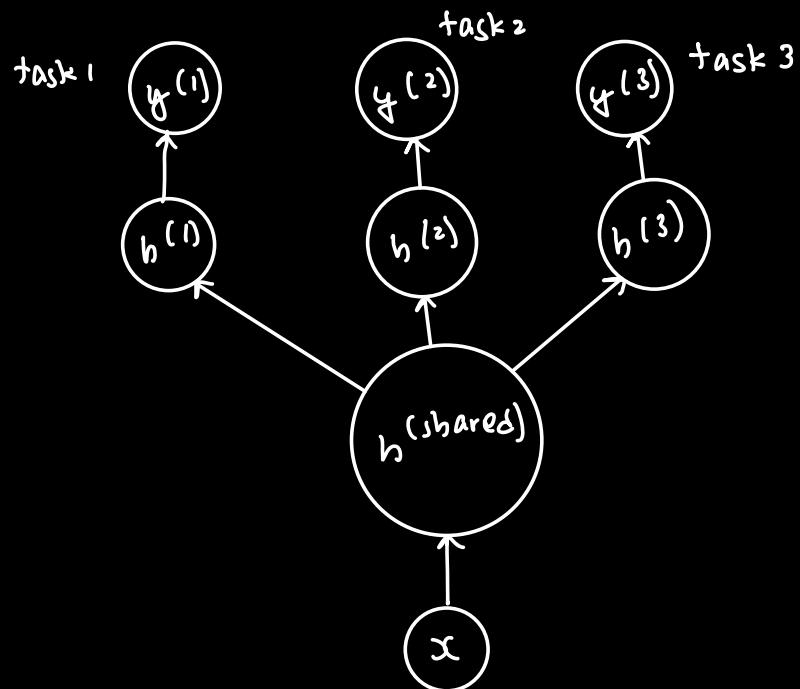
→ Label Smoothing

Instead of $[0.0 \dots 1.0 \dots]$

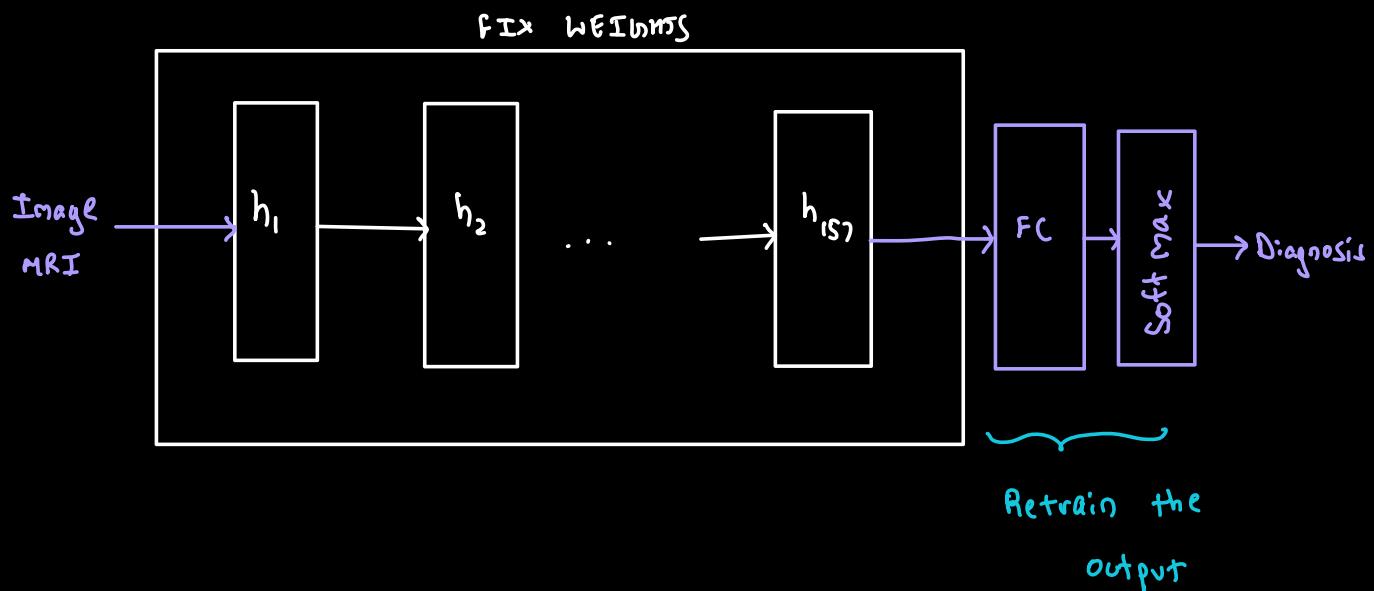
dog

$[0.011 \quad 0.011 \dots \quad 0.9 \quad 0.011 \quad 0.011 \dots 0.011]$

MULTITASK LEARNING:



TRANSFER LEARNING:



ENSEMBLE METHODS:

→ Train multiple different models

→ Average their results together.

k independent models

↳ Best case: $1/k$ decrease in error

↳ Worst case: Same error.

PRO:

→ Neural Networks when repeatedly trained → independent errors

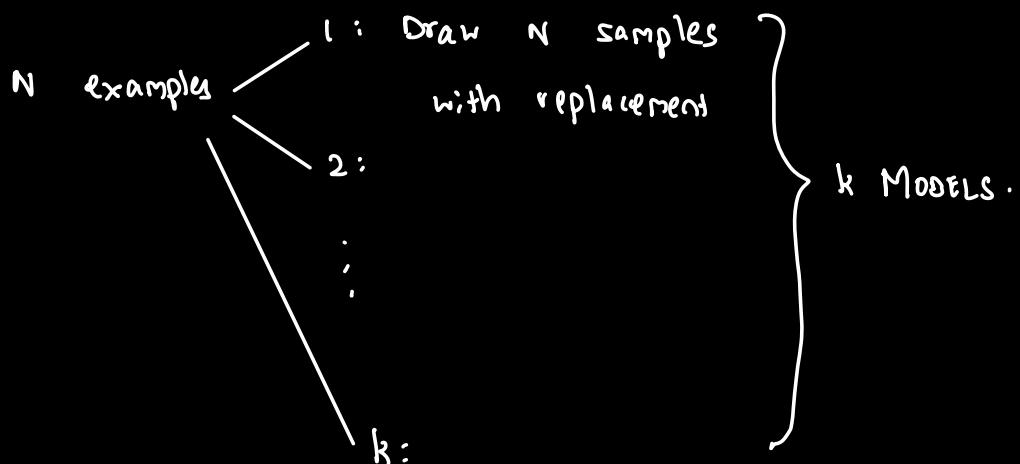
→ Almost always better accuracy.

Con:

→ Too much time to train.

Solution: Take suboptimal models. Set learning rate as cosine function and take each of the dip as a model.

→ Bagging (Bootstrap AGGREGATING):



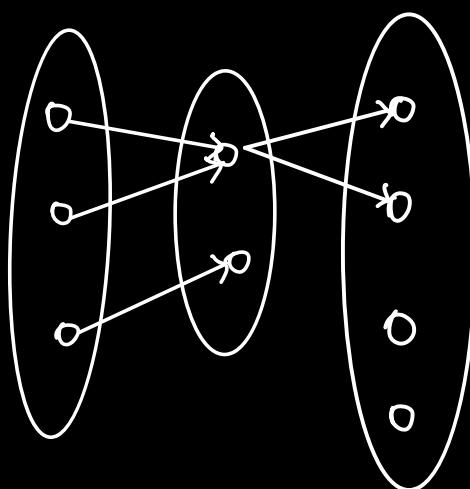
DROPOUT:

→ Approximation of Bagging

Binary Mask:

$p \rightarrow$ probability of seeing a 0.

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \quad \# \text{ of units/neurons.}$$



mask $\rightarrow (q \times 1)$

2^q possible configurations.

Affine \rightarrow relu \rightarrow batchnorm \rightarrow dropout



this count affect batchnorm statistics.

\rightarrow what about test time?

Assume $p=0.5$

What if we use all units?

\rightarrow over activate

Over many iterations, contributions of $w_i h_i$ to h_{out} is actually $(1-p) w_i h_i$.

$$h_{out} = \text{relu}(w_i h_i) * (1-p)$$

Pass p to test!

BETTER: INVERTED DROPOUT:

Scale at training:

$$h_i = (h_i < (1-p)) / (1-p)$$

PROS:

\rightarrow Estimates bagging

\rightarrow Regularizes each hidden unit to work well in different contexts.

\rightarrow Encodes redundant features.

STOCHASTIC GRADIENT DESCENT:

Loop

→ Sample m examples

$$\rightarrow g = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} J(\theta)$$

$$\rightarrow \theta \leftarrow \theta - \varepsilon g$$

Until stopping criterion is met.

MOMENTUM:

→ Maintain running mean of the gradients.

$$v = 0$$

$$\alpha \in [0, 1]$$

Loop

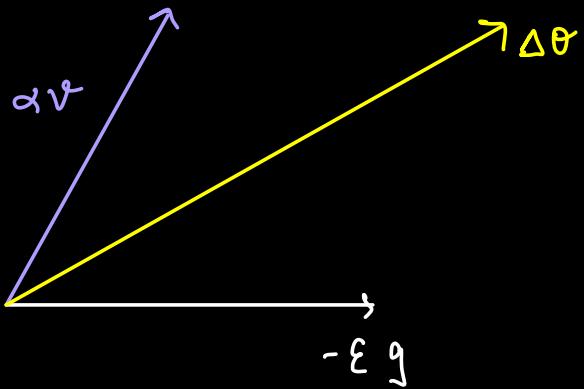
$$\rightarrow \text{Compute } g = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} J(\theta)$$

$$\rightarrow v \leftarrow \alpha v - \varepsilon g$$

$$\rightarrow \theta \leftarrow \theta + v$$

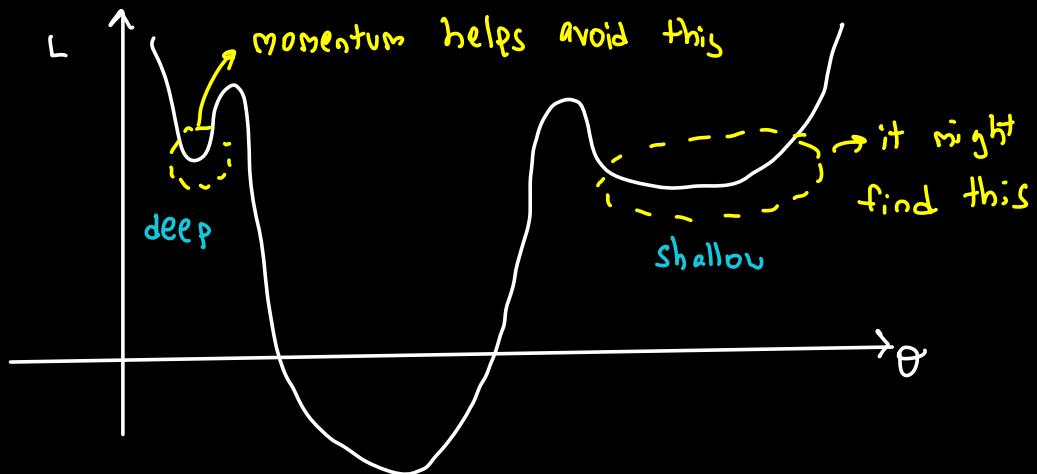
Until stopping criterion is met.

$$\text{So, } v_4 = -(\alpha^3 g_1 + \alpha^2 g_2 + \alpha g_3 + g_4)$$



PRO:

- Solves zig-zagging
- Helps to avoid deep local minima. Momentum pushes it.



Shallow local minima: Small change in θ , doesn't change L much. Test data might be slightly different. So small changes in θ shouldn't affect much. This is Robust!

NESTEROV MOMENTUM:

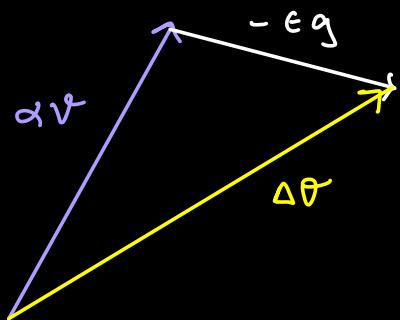
$$v^t = 0$$

$$\alpha \in [0, 1]$$

Loop

- Compute $g = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} J(\theta + \alpha v^t)$
- $v^t \leftarrow \alpha v^t - \epsilon g$
- $\theta \leftarrow \theta + v^t$

until stopping criterion is met.



ADAPTIVE GRADIENT (ADAGRAD):

If we took larger steps in a particular direction,

start taking smaller steps along that direction now.

$$\alpha = 0$$

$$\gamma = 10^{-7}$$

Loop

$$\rightarrow \text{Compute } g = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta_i} J(\theta)$$

$$\rightarrow \alpha \leftarrow \alpha + g \odot g \quad \odot \rightarrow \text{element wise product}$$

$$\rightarrow \theta \leftarrow \theta - \frac{\epsilon}{\sqrt{\alpha} + \gamma} \odot g$$

until stopping criterion is met.

n neurons / unit:

$$g \in \mathbb{R}^m$$

$$g \odot g : \begin{bmatrix} g_1^2 \\ g_2^2 \\ \vdots \\ g_n^2 \end{bmatrix}$$

$$\frac{\varepsilon}{\sqrt{a_i + \gamma}} \odot g = \begin{bmatrix} \frac{\varepsilon}{\sqrt{a_1 + \gamma}} \cdot g_1 \\ \frac{\varepsilon}{\sqrt{a_2 + \gamma}} \cdot g_2 \\ \vdots \\ \frac{\varepsilon}{\sqrt{a_n + \gamma}} \cdot g_n \end{bmatrix}$$

Con:

→ α always increases → step size always decreases.

No way for it to increase. Slower.

minima
in this case, we want it to move
in one direction. But as we have
moved along γ , we will slow down
by the decreasing step size.

RMSProp:

Allows step size to increase.

$$\alpha = 0$$

$$\gamma = 10^{-7}$$

$$\beta = [0, 1]$$

Loop

$$\rightarrow \text{compute } g = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} J(\theta)$$

$$\rightarrow \alpha \leftarrow [\beta \alpha + (1-\beta) g \odot g] \longrightarrow \alpha \text{ can}$$

$$\rightarrow \theta \leftarrow \theta - \frac{\epsilon}{\sqrt{\alpha} + \gamma} \odot g \quad \text{decrease.}$$

Until stopping criterion is met.

Pro:

Faster than Adagrad.

RMSProp + Momentum:

$$v = 0$$

$$\alpha = 0$$

$$\gamma = 10^{-7}$$

$$\beta = [0, 1]$$

$$\alpha = [0, 1]$$

Loop

$$\rightarrow \text{Compute } g = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} J(\theta)$$

$$\rightarrow \alpha \leftarrow \beta \alpha + (1-\beta) g \odot g$$

$$\rightarrow v \leftarrow \alpha v - \frac{\epsilon}{\sqrt{\alpha} + v} \odot g$$

$$\rightarrow \theta \leftarrow \theta + v$$

until stopping criterion is met.

ADAM:

Momentum \rightarrow shows the right direction

Adagrad \rightarrow take more judicious steps in that direction.

Without bias correction:

$v = 0$ "First moment"

$a = 0$ "Second moment"

$\gamma = 10^{-7}$

$\beta_1 = [0, 1]$

$\beta_2 = [0, 1]$

Loop

\rightarrow Compute $g = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} J(\theta)$

$\rightarrow v \leftarrow \beta_1 v + (1 - \beta_1) g$

$\rightarrow a \leftarrow \beta_2 a + (1 - \beta_2) g \odot g$

$\rightarrow \theta \leftarrow \theta - \frac{\epsilon}{\sqrt{a} + \gamma} \odot v$

Until stopping criterion is met.

With bias correction:

To avoid taking larger steps at the starting.

$$v = 0 \quad \text{"First moment"}$$

$$\alpha = 0 \quad \text{"Second moment"}$$

$$\gamma = 10^{-7}$$

$$\beta_1 = [0, 1]$$

$$\beta_2 = [0, 1]$$

$$t = 0$$

Loop

$$\rightarrow \text{Compute } g = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} J(\theta)$$

$$\rightarrow t \leftarrow t + 1$$

$$\rightarrow v \leftarrow \beta_1 v + (1 - \beta_1) g$$

$$\rightarrow a \leftarrow \beta_2 a + (1 - \beta_2) g \odot g$$

$$\rightarrow \tilde{v} = \frac{v}{(1 - \beta_1^t)}$$

$$\rightarrow \tilde{a} = \frac{1}{(1 - \beta_2^t)} a$$

$$\rightarrow \theta \leftarrow \theta - \frac{\epsilon}{\sqrt{\tilde{a}} + \gamma} \odot \tilde{v}$$

until stopping criterion is met.

As $t \rightarrow \infty$, $\tilde{v} \rightarrow v$

PRD:

→ most used, best.

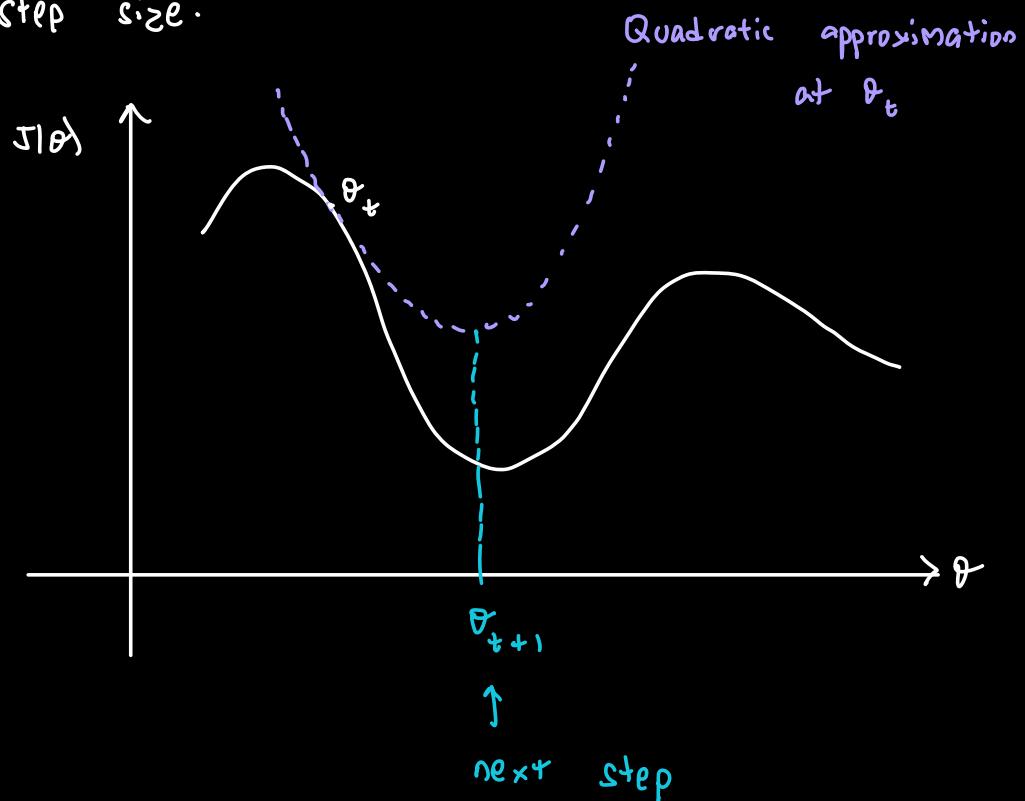
CDN:

→ Might find steep local minima.

SECOND ORDER METHODS:

Use curvature of cost function to decide

step size.



So deep $J(\theta) \rightarrow$ small steps

shallow $J(\theta) \rightarrow$ large steps.

NEWTON'S METHOD:

Loop

$$\rightarrow \text{Compute } g = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} J(\theta)$$

\rightarrow Compute Hessian, H

$$\rightarrow \theta \leftarrow \theta - H^{-1}g$$

until stopping criterion is met.

Conj:

\rightarrow Hard to calculate H^{-1} (Computationally intensive)

\rightarrow Memory to store Hessian.

$$n \approx 1 \times 10^6$$

4 bytes

totally 3.8 TB

\rightarrow H requires large batch size.

CHALLENGES:

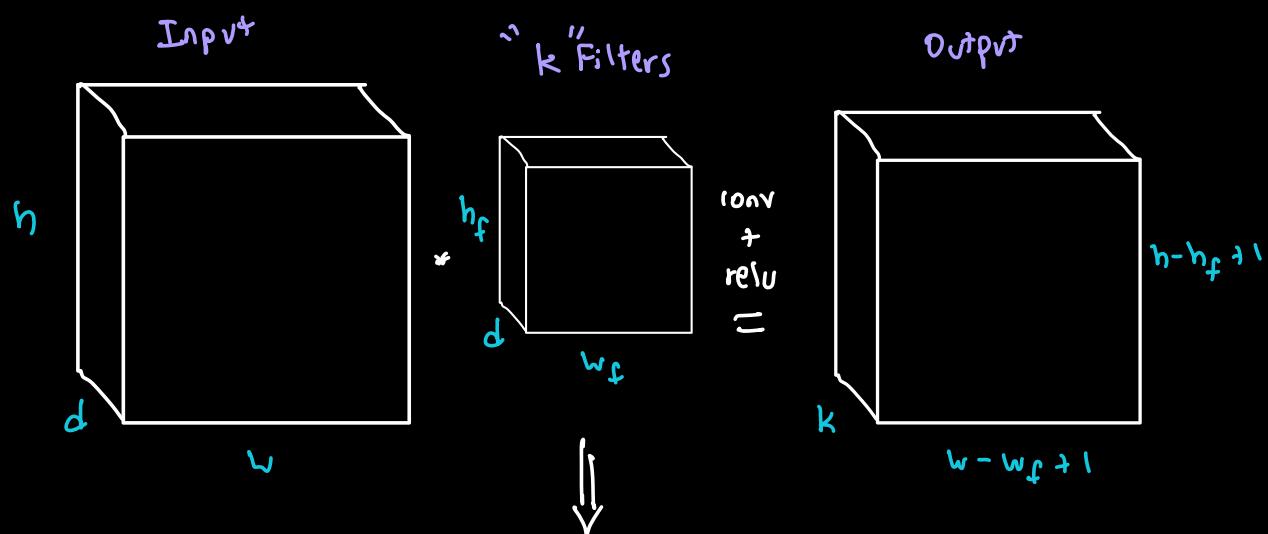
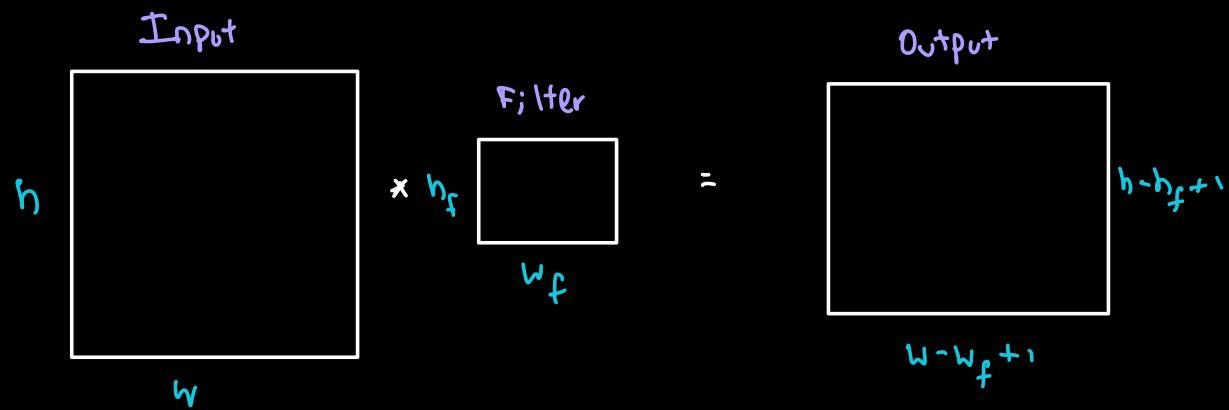
→ EXPLODING GRADIENTS:

$$\|g\| \geq \text{clip}$$

$$g \leftarrow \frac{g}{\|g\|} \cdot \text{clip}.$$

CONVOLUTIONAL NEURAL NETWORK:

→ Fully Connected Network has a lot of parameters.



common for all the units of input

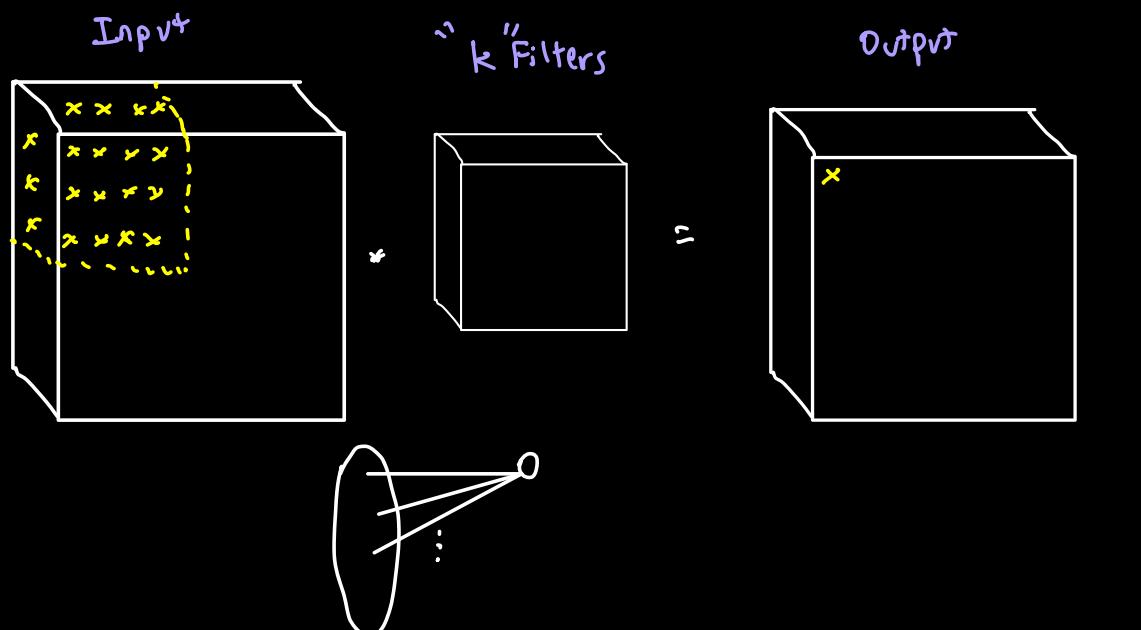
"Shared parameters"

Sparsely Connected Layers:

Unlike Fully Connected Network, the top left of

output depends only on top left of input.

Number of neurons it depends on is decided by the filter.



CNN: $w_f \times h_f \times d$ weights per neuron.

PRO:

→ Less computational memory

→ Less computational time.

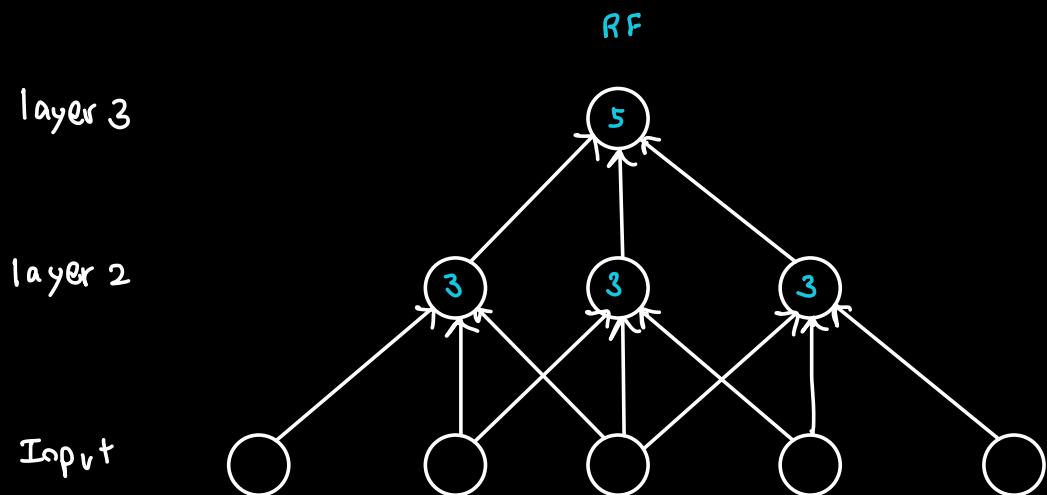
CON:

→ No neuron may have seen all of the input.

Solution: More layers.

RECEPTIVE FIELD:

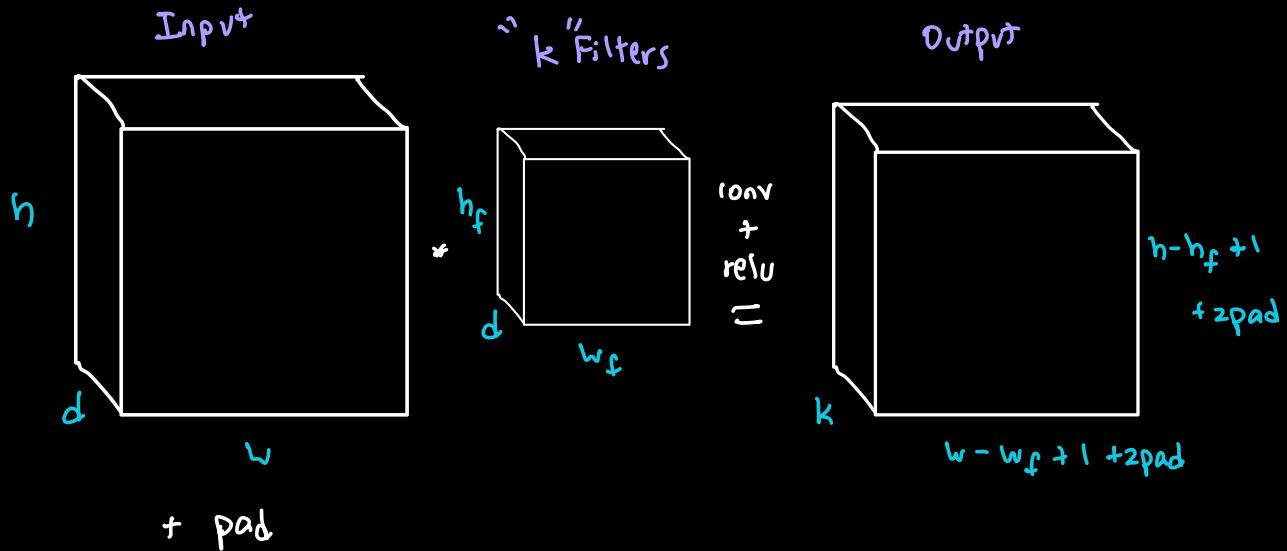
Amount of input the neuron sees.



2 ARCHITECTURES

	FN	CNN
Input:	$32 \times 32 \times 3 = 3072$	$32 \times 32 \times 3$
	\downarrow	\downarrow 100 filters $4 \times 4 \times 3$
Output / feature maps	500 output	$29 \times 29 \times 100 = 84,100$
Parameters :	$3072 \times 500 \approx 1.5 \text{ million}$ + 1 bias	$100 \times (4 \times 4 \times 3 + 1)$ bias $= 4900$

padding:



if $w_f = h_f = 3$, set $\text{pad} = 1$

output and input \rightarrow same h, w .

Stride:

In $7 \times 7 \rightarrow$ we can do stride = 1, 2 or 4.
not 3.

$$\left(\frac{w - w_f + 2\text{pad}}{\text{stride}} + 1, \frac{h - h_f + 2\text{pad}}{\text{stride}} + 1 \right)$$

POOLING LAYER:

(w_p, h_p) filter with a stride.

Operation: $\max()$ / $\text{avg}()$.

$$\left(\frac{w - w_p}{\text{stride}} + 1, \frac{h - h_p}{\text{stride}} + 1, d \right)$$

PRO:

→ Introduce spatial invariances.

If data is say shifted one pixel,

pooling layer's output for most cases

will remain the same.

LENET-5:

Handwriting Recognition

$$22 \times 32 \xrightarrow{6 \text{ } (5 \times 5) \text{ filters}} 28 \times 28 \times 6$$

Number of connections = each output has
 $(5 \times 5 + 1)$ connections
to input

$$= 26 \times 28 \times 28 \times 6$$

$$= 122,304 //$$

ALEXNET:

$$227 \times 227 \times 3$$

ZF NET:

- smaller filters, smaller strides - earlier layers
- more filters in deeper layers.

VGGNET:

Use small filters.

→ Lesser weights

→ More can be stacked to get same receptive field → more activations.

$$1. \frac{\partial}{\partial w} \text{tr}(wA) = A^T$$

$$\frac{\partial}{\partial w} \text{tr}(wAw^T) = wA^T + wA$$

$$\frac{\partial}{\partial w} \text{tr}(ABC) = \frac{\partial}{\partial w} \text{tr}(CBA) = \frac{\partial}{\partial w} \text{tr}(BAC) = \frac{\partial}{\partial w} \text{tr}(ABC)^T$$

$$2. xy = z \quad x = zy^{-1} \quad \text{not} \quad y^{-1}z$$

↑

$$\hookrightarrow xy y^{-1} = zy^{-1}$$

$$3. \frac{\partial}{\partial x} (x^T y) = y$$

$$\frac{\partial}{\partial x} (wx) = w^T$$

$$4. y = wx$$

$$\frac{\partial L}{\partial x} = w^T \frac{\partial L}{\partial y}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \cdot x^T$$

$$5. \quad y = w x$$

$$\begin{matrix} \checkmark & \downarrow & \swarrow \\ n \times p & n \times m & m \times p \end{matrix}$$

$$\frac{\partial L}{\partial x} = w^T \frac{\partial L}{\partial y}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} x^T$$

$$6. \quad y = x^T \theta x$$

$$\frac{\partial y}{\partial x} = (A + A^T)x$$

$$y = \theta^T x$$

$$\frac{\partial y}{\partial x} = \theta$$

$$7. \quad y = x^T A y$$

$$\frac{\partial y}{\partial x} = Ay$$

$$\frac{\partial y}{\partial y} = A^T x$$

$$\frac{\partial y}{\partial A} = x y^T$$

$$8. \quad \frac{\partial \mathbf{g}}{\partial \mathbf{x}^T} = \left(\frac{\partial \mathbf{g}}{\partial \mathbf{x}} \right)^T$$

$$9. \quad \mathbf{g} = |\mathbf{k}|$$

$$\frac{\partial \mathbf{g}}{\partial \mathbf{k}} = |\mathbf{k}| (\mathbf{k}^{-1})^T$$

$$10. \quad \overrightarrow{\frac{\partial (\mathbf{x} \mathbf{x}^T)}{\partial \mathbf{x}}} = \overrightarrow{\frac{\partial (\mathbf{x}^T \mathbf{x})}{\partial \mathbf{x}}} = 2\mathbf{x}$$

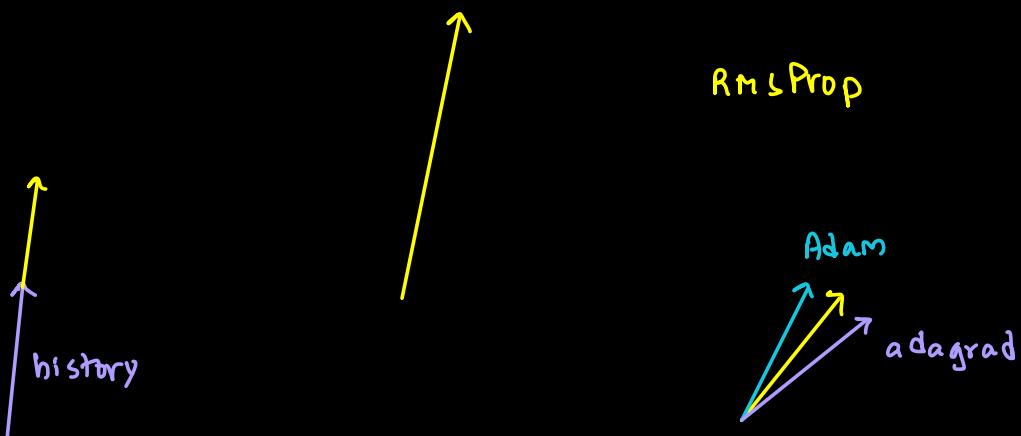
$$\frac{\partial L}{\partial \mathbf{x}} = 2 \left(\frac{\partial L}{\partial \mathbf{x} \mathbf{x}^T} \right) \mathbf{x}$$

$$11. \quad \overrightarrow{\frac{\partial}{\partial \mathbf{x}}} \text{tr}(\mathbf{A}) = \text{tr} \left(\overrightarrow{\frac{\partial \mathbf{A}}{\partial \mathbf{x}}} \right)$$

$$12. \quad \frac{\partial}{\partial \mathbf{P}} \|\mathbf{P}\|_2 = \frac{\mathbf{P}}{\|\mathbf{P}\|_2}$$

SGD

Momentum



RMSProp

Adam

adagrad

RMSProp

↳ bias towards 0

in the starting-

Hw1:

$$\hat{y} = w \cdot x$$

$$L = \frac{1}{2} \sum_{i=1}^n (\hat{y}^{(i)} - w \cdot x^{(i)})^T (\hat{y}^{(i)} - w \cdot x^{(i)})$$

$$y - w \cdot x = \begin{bmatrix} | & | & | \\ \vdots & \vdots & \vdots \\ \hat{y}^{(1)} & \hat{y}^{(2)} & \dots & \hat{y}^{(n)} \\ | & | & | & | \end{bmatrix} - w \begin{bmatrix} | & | & | \\ \vdots & \vdots & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(n)} \\ | & | & | & | \end{bmatrix}$$

$$= \begin{bmatrix} | & | & | \\ \vdots & \vdots & \vdots \\ (\hat{y}^{(1)} - w \cdot x^{(1)}) & (\hat{y}^{(2)} - w \cdot x^{(2)}) & \dots & (\hat{y}^{(n)} - w \cdot x^{(n)}) \\ | & | & | & | \end{bmatrix}$$

$$(y - w \cdot x)^T (y - w \cdot x) = \begin{bmatrix} (\hat{y}_1 - w \cdot x_1)^T (\hat{y}_1 - w \cdot x_1) & (\hat{y}_1 - w \cdot x_1)^T (\hat{y}_2 - w \cdot x_2) & \dots \\ (\hat{y}_2 - w \cdot x_2)^T (\hat{y}_1 - w \cdot x_1) & (\hat{y}_2 - w \cdot x_2)^T (\hat{y}_2 - w \cdot x_2) & \dots \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$L = \frac{1}{2} \text{tr} (y - w \cdot x)^T (y - w \cdot x) \approx \frac{1}{2} \text{tr} (y^T - x^T w^T)(y - w \cdot x)$$

$$= \frac{1}{2} \text{tr} \left(y^T y - \underbrace{y^T w x}_{w x y^T} - \underbrace{x^T w^T y}_{x^T w^T w x} + \underbrace{x^T w^T w x}_{w x x^T} \right)$$

$$\frac{\partial L}{\partial w} = 0$$

$$\frac{\partial}{\partial w} \text{tr}(y^T w x) = \frac{\partial}{\partial w} \text{tr}(\underbrace{w x y^T}_{w x y^T})$$

$$- y x^T$$

$$\frac{\partial}{\partial w} \text{tr}(x^T w^T y) = \frac{\partial}{\partial w} \text{tr}(\underbrace{y^T w x}_{w x y^T})$$

$$- y x^T$$

$$\frac{\partial}{\partial w} \text{tr}(x^T w^T w x)$$

$$x x^T w^T w$$

$$\underbrace{w x x^T}_{w x x^T} w^T$$

$$w x x^T + w x x^T$$

$$\Rightarrow w x x^T - y x^T = 0$$

$$y x^\top = W x x^\top$$

$$W = (y x^\top) (x x^\top)^{-1}$$

Now:

$$\Pr(y^{(j)} : i | x^{(j)}, \theta) \approx \text{softmax}_i(x^{(j)})$$

$$\begin{aligned} \text{MLE} &= P(x^{(j)}, y^{(j)} | \theta) \\ &= P(x^{(j)} | \theta) \cdot P(y^{(j)} | x^{(j)} | \theta) \end{aligned}$$

$$\begin{aligned} \log L &= \log P(y^{(j)} | x^{(j)} | \theta) \\ &\approx (\mathbf{w}_{y^{(j)}}^\top \mathbf{x} + b_{y^{(j)}}) - \log \sum_{k=1}^C (e^{\mathbf{w}_k^\top \mathbf{x} + b_k}) \end{aligned}$$

$$\text{if } i = y^{(j)}$$

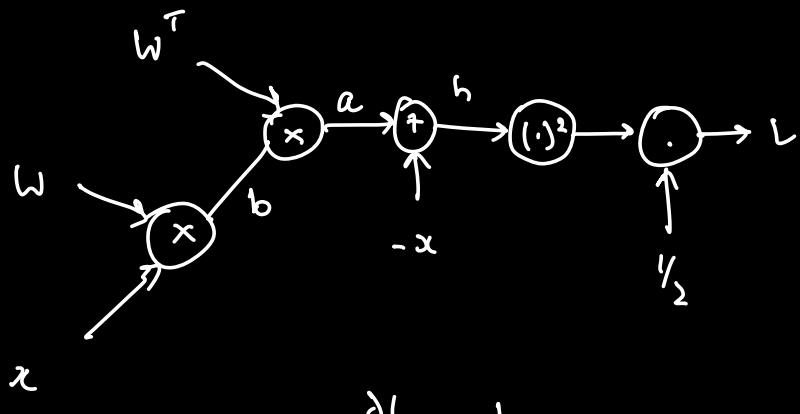
$$\frac{\partial L}{\partial w_i} = \mathbf{x} - \frac{1}{\sum_{k=1}^C e^{\mathbf{w}_k^\top \mathbf{x}}} e^{\mathbf{w}_i^\top \mathbf{x}} \cdot \mathbf{x}$$

else

$$\frac{\partial L}{\partial w_i} = \frac{-1}{\sum_{k=1}^C e^{\mathbf{w}_k^\top \mathbf{x}}} e^{\mathbf{w}_i^\top \mathbf{x}} \cdot \mathbf{x}$$

HW3:

$$1. \quad L = \frac{1}{2} \| w^T w x - x \|^2$$



$$\frac{\partial L}{\partial h} = b$$

$$\frac{\partial L}{\partial x} = - \left(\frac{\partial L}{\partial h} \right) = -b$$

$$\frac{\partial L}{\partial a} = b \quad a = w^T b$$

$$\frac{\partial L}{\partial b} = w \frac{\partial L}{\partial a} \quad \frac{\partial L}{\partial w^T} = \frac{\partial L}{\partial a} \cdot b^T$$

$$= w b \quad \vdots b b^T$$

$$\frac{\partial L}{\partial w} = b b^T$$

$$\hat{y} = w \cdot x$$

$$\frac{\partial L}{\partial z} = w^T \frac{\partial L}{\partial \hat{y}}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}} \cdot x^T$$

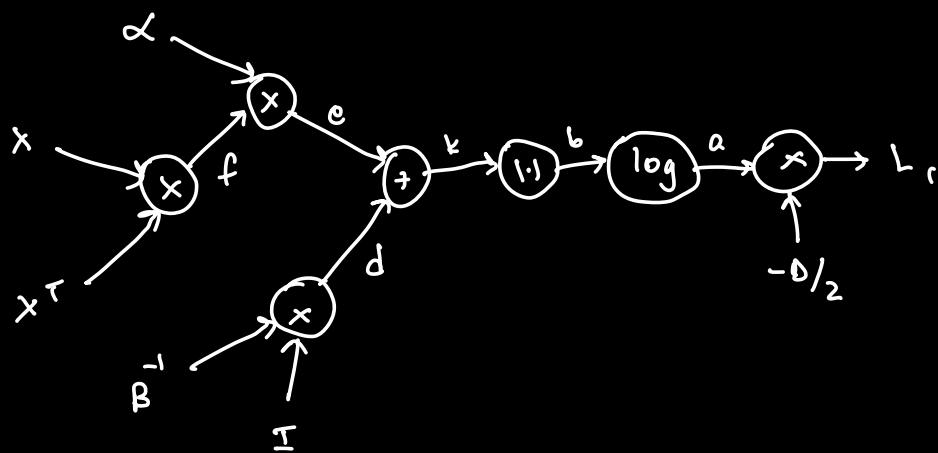
$$b = w \cdot x$$

$$\frac{\partial L}{\partial x} = w^T w h$$

$$\frac{\partial L}{\partial w} = w h x^T$$

$$\delta L / \delta w = w (w^T w x \cdot x) x^T + w x (w^T w x - x)^T$$

2)



$$\frac{\partial L_1}{\partial a} = -\frac{D}{2}$$

$$a = \log b$$

$$\frac{\partial a}{\partial b} \approx \frac{1}{b}$$

$$\frac{\partial L_1}{\partial b} = \frac{\partial a}{\partial b} \quad \frac{\partial L_1}{\partial a} = \frac{1}{b} \cdot \left(-\frac{D}{2} \right)$$

$$b = |k|$$

$$\frac{\partial b}{\partial k} = |k| (k^{-1})^T$$

$$\frac{\partial L}{\partial k} = \frac{\partial b}{\partial k} \frac{\partial L}{\partial b}$$

$$= |k| (k^{-1})^T \left(-\frac{D}{2b} \right)$$

$$\frac{\partial L}{\partial e} = \frac{\partial L}{\partial d} = \frac{\partial L}{\partial k}$$

$$\frac{\partial L}{\partial f} = \alpha \frac{\partial L}{\partial k}$$

$$f = xx^T$$

$y = w \alpha$
$\frac{\partial L}{\partial \alpha} = w^T \frac{\partial L}{\partial y}$
$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \cdot \alpha^T$

$$\frac{\partial L}{\partial x^T} = x^T \frac{\partial L}{\partial f}$$

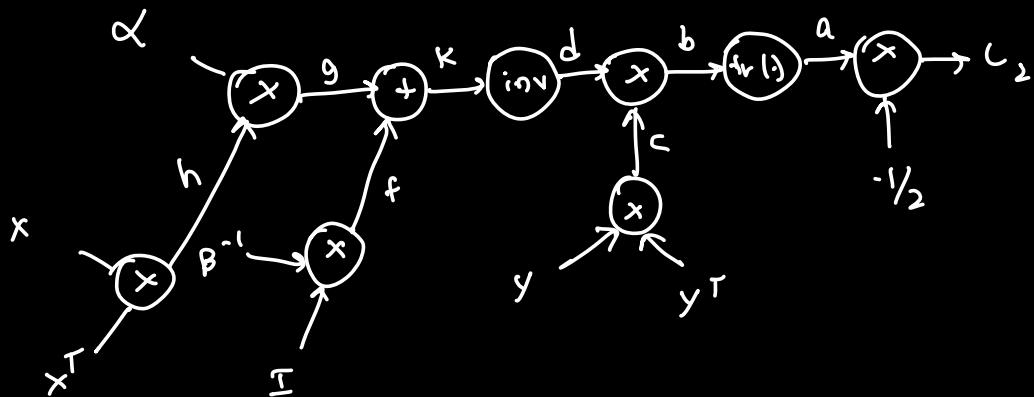
$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial f} \cdot x$$

$$\frac{\partial L}{\partial x^T} = x^T \alpha |k| (k^{-1})^T \left(-\frac{D}{2b} \right)$$

$$= \left(-\frac{\alpha D}{2b} \right) |k| x^T (k^{-1})^T$$

$$\frac{\partial L}{\partial x} = \left(-\frac{\alpha \Delta}{2b} \right) |k| (k^{-c})^r x$$

$$+ \left(-\frac{\alpha \Delta}{2b} \right) |k| (k^{-c})^r x$$



$$\alpha = \text{tr}(b)$$

$$\frac{\partial \alpha}{\partial b} = \text{tr} \left(\frac{\partial b}{\partial b} \right) = I$$

$$\frac{\partial L_2}{\partial b} = -I/2$$

$$\alpha b$$

$$\frac{\partial L_2}{\partial d} = \gamma y^T \left(-\frac{I}{2} \right)$$

$$\frac{\partial L}{\partial h} = \alpha \frac{\partial L}{\partial k}$$

$$\frac{\partial h}{\partial x} = 2x$$

$$\frac{\partial L_2}{\partial a} = -1/2 \quad d = e^{-1}$$

$$\frac{\partial L}{\partial k} = -k^{-r} \frac{\partial L}{\partial d} k^{-r}$$

$$= -k^{-r} y y^T \left(-\frac{I}{2} \right) k^{-r}$$

$$\frac{\partial L}{\partial g} = \frac{\partial L}{\partial k}$$

$$\frac{\partial L}{\partial h} = \alpha \frac{\partial L}{\partial k}$$

$$\frac{\partial L}{\partial x} = \frac{\partial h}{\partial x} \frac{\partial L}{\partial h} = 2\alpha \left(+ K^{-T} x y^T \frac{I}{2} K^{-T} \right)$$

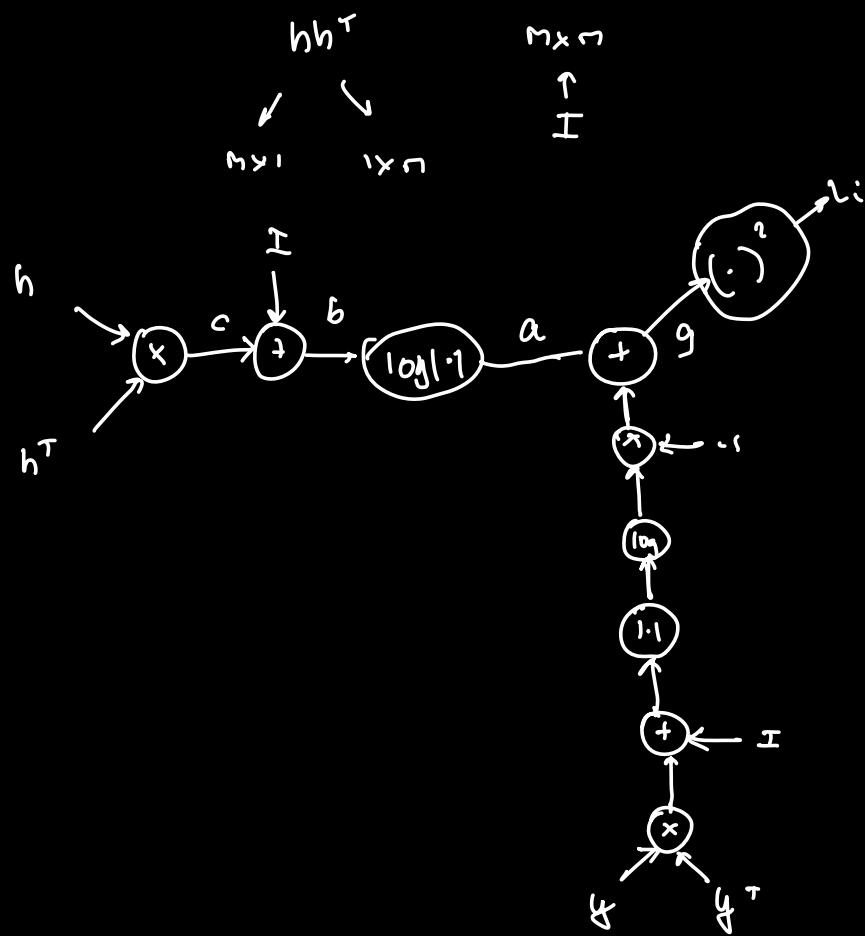
$$= \alpha K^{-T} x y^T K^{-T} x$$

$$L_1 \rightarrow 1 \times 1 \quad \quad \quad \overset{m \times n}{\curvearrowright}$$

$$h = b x + b$$

$$\downarrow \quad \quad \quad \downarrow \quad \quad \quad \downarrow$$

$$m \times 1 \quad n \times 1 \quad m \times 1$$



$$\frac{\partial L_i}{\partial a} = 1$$

$$a = \log b$$

$$\frac{\partial L_i}{\partial b} = \frac{\partial a}{\partial b} \cdot \frac{\partial L_i}{\partial a} = (b^T)^{-1}$$

$$\frac{\partial L_i}{\partial c} = (b^T)^{-1}$$

$$c = hh^T$$

$$\frac{\partial L_i}{\partial h} = 2 \frac{\partial L_i}{\partial c} h = 2(b^T)^{-1} h$$

$$h = \text{ReLU}\left(\frac{a}{w^T x + b}\right)$$

$$h = \text{ReLU}(a)$$

$$\frac{\partial h}{\partial a} = 1 \quad (a > 0)$$

$$I(a > 0) \quad \frac{\partial L_i}{\partial a} = I(a > 0) \cdot 2(b^T)^{-1} h \times 2g$$

$$a = w^T x + b$$

$$\frac{\partial L_i}{\partial x} = w^T I(a > 0) \cdot 2(b^T)^{-1} h \times 2g$$

$$\frac{\partial \ell_i}{\partial w} = I(a > b) \cdot 2(b^T)^{-1} h(x^T) \times 2g$$

$$b) \quad a) \quad 27 \times 27 \times 9b \\ 3 \times 3 \times 9b \times 384$$

$$3 \times 3 \times 384 \times 384$$

$$3 \times 3 \times 384 \times 256$$

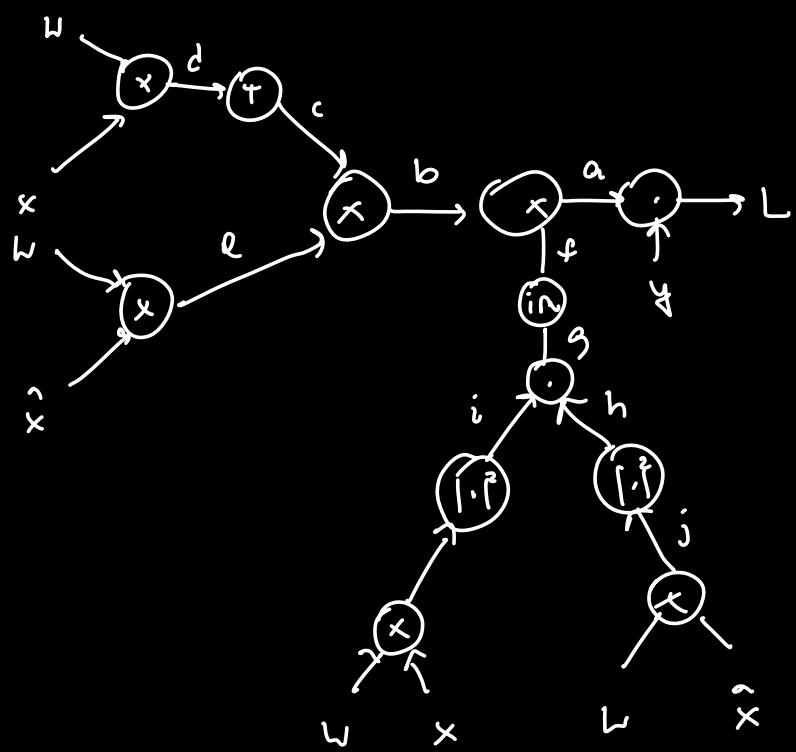
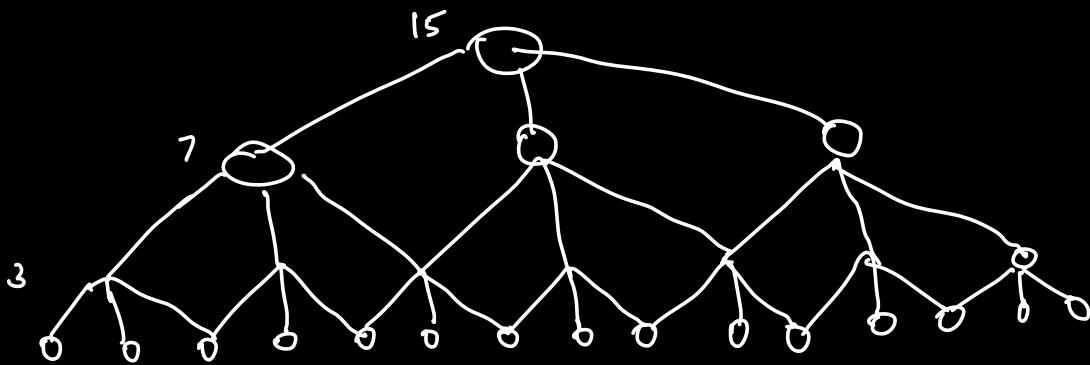
$$27 \times 27 \times 256$$

$$6 \times 6 \times 256 \times 4096$$

$$27 \times 27 \times 256 \times 4$$

$$6 \times 6 \times 256 \times 4$$

$$4096 \times 4$$



$$g = \frac{1}{\sqrt{1}}$$

$$\frac{\partial L}{\partial f} = by$$

$$\frac{\partial L}{\partial g} = \left(\frac{-1}{f^2} \right) by$$

$$\frac{\partial L}{\partial h} = \left(\frac{i}{f^2} \right) by$$

$$h = |ij|^2$$

$$\frac{\partial L}{\partial a} = g$$

$$g = bf$$

$$\frac{\partial a}{\partial b} = f$$

$$\frac{\partial L}{\partial b} = \frac{\partial g}{\partial b} \frac{\partial L}{\partial a} = fy$$

$$\frac{\partial h}{\partial j} = 2j$$

$$b \in \mathbb{C}$$

$$\frac{\partial L}{\partial c} = \frac{\partial L}{\partial b} \cdot e^T$$

$$= f_y e^T$$

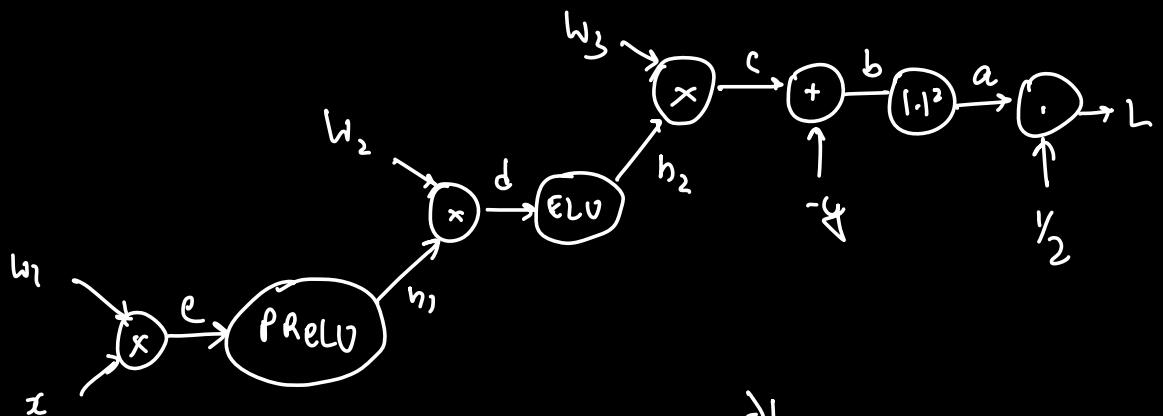
$$\frac{\partial L}{\partial e} = c^T \frac{\partial L}{\partial b}$$

$$= c^T f_y$$

$$\Phi = Wx$$

$$\frac{\partial L}{\partial x} = w^T \frac{\partial L}{\partial \Phi}$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \Phi} \cdot x^T$$



$$\frac{\partial L}{\partial a} = \frac{1}{2}$$

$$\frac{\partial L}{\partial b} = b$$

$$\frac{\partial L}{\partial c} = b$$

$$C = w_3 h_2$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial c} \cdot h_2^\top$$

$$= b h_2^\top$$

$$\underbrace{\frac{\partial L}{\partial h_2}}_{\frac{\partial L}{\partial b}} = w_3^\top \frac{\partial L}{\partial c}$$

$$= w_3^\top b$$

$$\frac{\partial L}{\partial d} = \frac{\partial h_2}{\partial d} \frac{\partial L}{\partial h_2}$$

$$h_2 = \text{ELU}(d)$$

$$\frac{\partial h_2}{\partial d} = \begin{cases} \alpha e^d & , d < 0 \\ 1 & , d \geq 0 \end{cases}$$

$$\frac{\partial h_2}{\partial d} = I(d \geq 0) + I(d < 0) \alpha e^d$$

$$\frac{\partial L}{\partial d} = \frac{\partial h_2}{\partial d} \cdot w_3^\top b$$

$$d = w_2 h_1$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial d} \cdot h_1^\top \quad \frac{\partial L}{\partial h_1} = w_2^\top \frac{\partial L}{\partial d} = w_2^\top \frac{\partial h_2}{\partial d} w_3^\top b$$

$$h_1 = \text{ReLU}(e)$$

$$\frac{\partial h_1}{\partial e} = \begin{cases} 0, & e < 0 \\ 1, & e \geq 0 \end{cases}$$

$$\frac{\partial L}{\partial e} = (I(e < 0) \alpha + I(e \geq 0)) w_2^\top (I(d \geq 0) + I(d < 0) \alpha e^d) w_3^\top b$$

→ (c-4)

$$e = \omega, x$$

$$\frac{\delta L}{\delta \omega_i} = \frac{\delta L}{\delta e} \cdot x^T \quad \frac{\delta L}{\delta x} = \omega_i^T \frac{\delta L}{\delta e}$$