

Improving Deep Neural Networks

WEEK 3

PAGE NO. :
DATE : / /

train / dev / test sets

Data :

training set

holdout test set
cross val/
dev set

for 100 - 1000 - 10000 values

split: 70/30 or 60/20/20

But in big data era with a million values

test data percentage can
be reduced

like 98/1/1

Mismatched train/test dist.

Training set

HD pictures from
Web



dev/test sets

Cat pictures from
your camera / apr

make sure dev and test set come from
the same distribution

$$S:45 \quad \|w^{[e]}\|^2 = \sum_{i=1}^n \sum_{j=1}^n (w_{i,j}^{[e]})^2$$

Bias/Variance

train set error 1% } high variance
 dev set error 11% }

15% } high variance
 30% } high bias

15% } high bias
 16% }

Basic Recipe

High Bias → bigger net neural net,
 (training data more layers
 problems)

High variance → More train data
 (dev set performance) + Regularization
 more appropriate NN arch

Regularization

(for Logistic-

$$\min_{w,b} J(w,b) \quad w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2 + \frac{\gamma}{2m} b^2$$

$$L_2 \text{ regularization} \quad \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

usually omitted

L1 regularization

$$\sum_{2m} \sum_{i=1}^{m_x} |w_i| = \frac{\lambda}{2m} \|w\|_1,$$

w will be sparse

λ is the regularization parameter

* Neural Network

$$J(w^{(l)}, b^{(l)}, w^{(l)}, b^{(l)})$$

$$= \frac{1}{m} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|^2$$

$$\|w^{(l)}\|^2 = \sum_{i=1}^n \sum_{j=1}^{n_l} (w_{ij}^{(l)})^2$$

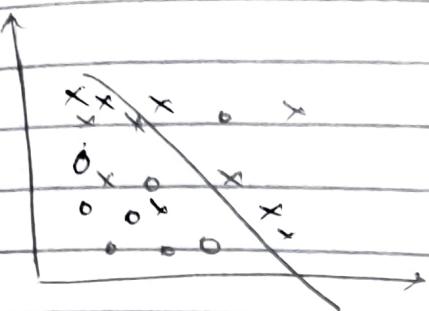
Frobenius norm

Weight updates

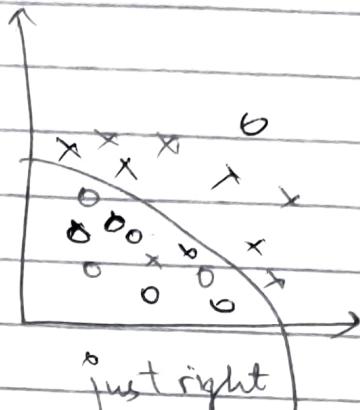
$$w^{(l)} = \partial_w^{(l)} \rightarrow (\text{from backprop} + \frac{\lambda}{m} w^{(l)})$$

also called

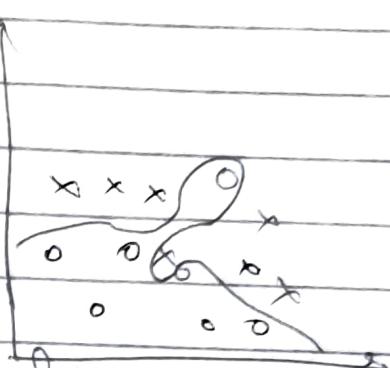
$$-w^{(l)} - \frac{\alpha \lambda}{m} w^{(l)} - \alpha \text{ (from backprop)}$$



high bias



just right



high variance

- build better model

- get more data
- reduce model complexity

regularization also helps getting a more monotonically decreasing cost function

Dropout Regularization

- * Inverted Dropout

$L = 3$ (layers)

$d_3 = \text{np.random.rand}(a_3 \text{ shape}[0], a_3 \text{ .shape}[1]) < \text{keep_prob}$

$a_3 = \text{np.multiply}(a_3, d_3)$

→ creates a boolean matrix with probability of False = 0.2

→ approx 20% of the values will be made 0

$a_3 = \text{Keep prob}$

0.8

dropout is usually used to stop NN from overfitting

$$w^{[1]} = 7 \times 3$$

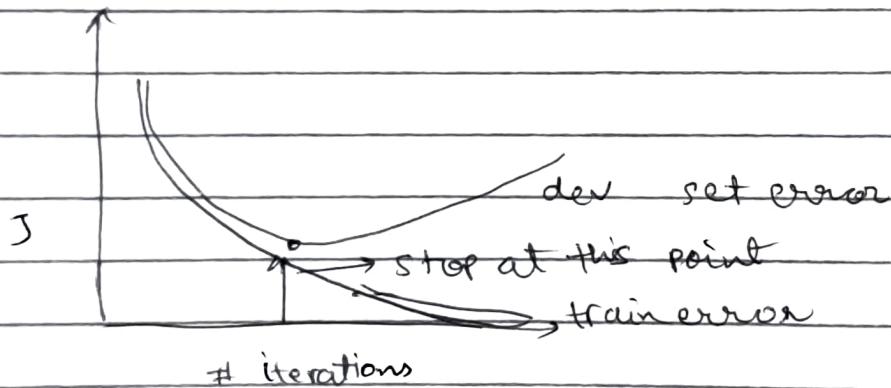
$$w^{[3]} = 3 \times 7$$

Intuition: can't rely on any one feature, so have to spread out weights.

We increase the dropout usually for higher no of input neurons for that layer

More regularization techniques

- Data augmentation: ex for a cv project you can flip images, rotate them, crop them
- Early Stopping:



Orthogonalization: we have two tasks

- optimize w, b to get minimum J
- set up regularized hyperparameters to stop overfitting

#Setting up Normal optimization problem

* Normalizing training sets

mean $x := x - \mu$

var $x := x / \sigma$

* Vanishing / Exploding Gradients

activation values exponentially increase/decrease

When there are many layers. Use proper weight initializations

* Weight Initializations for DNNs

$$z = w_0 + w_1 x_1 + \dots + w_n x_n$$

large $n \rightarrow$ small w_i

$$w^{[l]} = np.random.rand(shape) * np.sqrt(\frac{2}{n^{[l-1]}})$$

ReLU has other variations:

$$\sqrt{\left(\frac{1}{n^{[l-1]}}\right)} \quad \sqrt{\left(\frac{2}{n^{[l-1]} + n^{[l]}}\right)}$$

* Gradient Checking

Optimization Algorithms

* Mini Batch

We have vectorized inputs

$$x = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}_{(n \times m)}$$

$$y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}_{(1, m)}$$

What if $m = 5000000$

divide it into min batches

lets say of 1000 training examples

$$x = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(1000)} & x^{(1001)} & x^{(1002)} & \dots & x^{(2000)} \\ \underbrace{x^{(1)}}_{\sum 1} & \underbrace{x^{(2)}}_{\sum 2} & \dots & \underbrace{x^{(1000)}}_{\sum 1000} & \underbrace{x^{(1001)}}_{\sum 2} & \underbrace{x^{(1002)}}_{\sum 3} & \dots & \underbrace{x^{(2000)}}_{\sum 1000} \end{bmatrix}_{(n \times 1000)}$$

$$y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(1000)} \\ \underbrace{y^{(1)}}_{\sum 1} & \underbrace{y^{(2)}}_{\sum 2} & \dots & \underbrace{y^{(1000)}}_{\sum 1000} \end{bmatrix}$$

We usually shuffle our x & y before partitioning

Alg.

repeat epochs {

for $t=1$ to $t=5000$ {

forward prop on $x^{[+]}$

$$z^{[1]} = w^{[1]} x^{[+]} + b^{[1]}$$

$$A^{[1]} = g^{[1]}(z^{[1]})$$

:

$$A^{[L]} = g^{[L]}(z^{[L]})$$

vectorized
1000 examples

Compute cost $J^{[+]}$ = $\frac{1}{1000} \sum_{i=1}^L f(y^{(i)}, \hat{y}^{(i)})$

$$+ \lambda \sum_{l=1}^L \|w^{[l]}\|^2$$

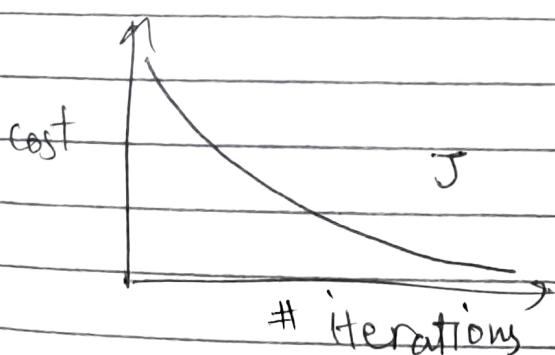
Backprop to compute gradients

$$w^{[l]} := w^{[l]} - \alpha d w^{[l]}, b^{[l]} := b^{[l]} - \alpha d b^{[l]}$$

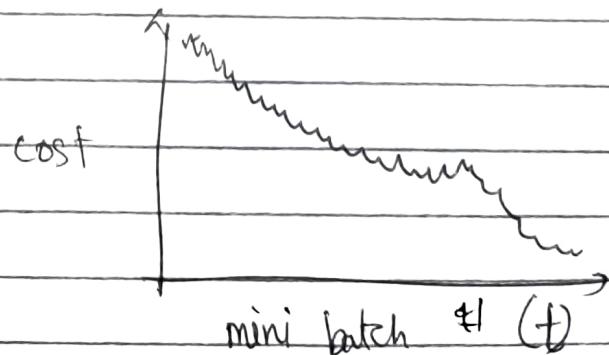
}

}

Batch gradient descent



mini Batch GD



If mini batch size (~~m~~) = m

: we have batch gradient descent
ie $(\bar{x}^{\frac{1}{m}}, \bar{y}^{\frac{1}{m}}) = (\bar{x}, \bar{y})$

If mini batch size = 1

: we have Stochastic GD

every example is its own mini batch
size = 1

Stochastic

mbsize = 1



lose speed from

vectorization

in between

mb size not too

big not too small



fastest learning
~~optimal~~

= m

batch

mbsize = m



too long per

iteration

for large m

If small m : use batch gd

$m \leq 2000$

Typical mini batch sizes: 64, 128, 256, 512

* Exponentially weighted Averages

Moving averages are usually seen in
time series datasets.

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$\beta = 0.9 \Rightarrow$ averaging over last 10 days

$\beta = 0.98 \Rightarrow$ averaging over last 50 days

Curve will be smoother, less sensitive
to immediate changes

$$V_{100} = 0.9V_{99} + 0.1\theta_{100}$$

$$V_{99} = 0.9V_{98} + 0.1\theta_{98}$$

$$V_{98} = 0.9V_{97} + 0.1\theta_{97}$$

$$V_{100} = \underbrace{0.9(0.1\theta_{100} + 0.9(\underbrace{0.1\theta_{99} + 0.9(\dots}_{V_{99}})}_{V_{98}})$$

$$= 0.1\theta_{100} + 0.1 \times (0.9)^1 \theta_{99} + 0.1(0.9)^2 \theta_{98} + 0.1(0.9)^3 \theta_{97} \dots$$

Implementation:

iteration: $V := 0$	$V_0 = 0$ repeat {
$V := \beta V + (1-\beta)\theta_1$	get θ_t
$V := \beta V + (1-\beta)\theta_2$	$V = \beta V + (1-\beta)\theta_t$
...	}

* bias correction

Your predictions lean towards 0 initially
when you start with $v=0$

* GD with momentum - to avoid oscillations in θ
 $V_{dw} = V_{db} = 0$ on iteration t : initial

(compute dW , db on current minibatch)

$$\text{d}V_{dw} = \beta V_{dw} + (1-\beta) dW$$

$$V_{db} = \beta V_{db} + (1-\beta) db$$

$$W := w - \alpha V_{dw}, \quad b := b - \alpha V_{db}$$

$\beta = 0$ implies GD w/o momentum

* RMS prop

On iteration t

compute $d\omega, db$ on current mini batch

$$S_{d\omega} = \beta S_{d\omega} + (1-\beta) d\omega^2 \leftarrow \text{small}$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2 \leftarrow \text{large}$$

$$\omega := \omega - \alpha \frac{d\omega}{\sqrt{S_{d\omega}}} \quad b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$

size

step

mb_size

* Adam optimization Algo

$$V_{d\omega} = 0 \quad S_{d\omega} = 0 \quad V_{db} = 0 \quad S_{db} = 0$$

lose

vector

On iteration t

compute $d\omega, db$ using current mini batch

$$V_{d\omega} = \beta_1 V_{d\omega} + (1-\beta_1) d\omega$$

$$V_{db} = \beta_1 V_{db} + (1-\beta_1) db$$

$$S_{d\omega} = \beta_2 S_{d\omega} + (1-\beta_2) d\omega^2$$

$$S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2$$

$$V_{\text{corrected}} = V_{d\omega} / (1 - \beta_1^t)$$

$d\omega$

$$= V_{d\omega} / (1 - \beta_1^t)$$

$$S_{d\omega}^{\text{corrected}} = S_{d\omega} / (1 - \beta_2^t)$$

$$S_{db} / (1 - \beta_2^t)$$

$$\omega := \omega - \alpha V_{d\omega}^{\text{corrected}}$$

$$\sqrt{S_{d\omega}^{\text{corrected}}} + \epsilon$$

$$b := b - \alpha V_{db}^{\text{corrected}}$$

$$\sqrt{S_{db}^{\text{corrected}}} + \epsilon$$

ϵ = small term usually 10^{-8}

usually,

α needs to be tuned

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

$$\epsilon = 10^{-8}$$

Adam: Adaptive moment estimation

Learning Rate Decay

to tune the α hyperparameter, we use this algorithm

1 epoch = 1 pass through training data

$$\alpha = \frac{1}{1 + \text{decay rate} \times \text{epoch num}} \cdot \alpha_0$$

epoch	α	$\alpha_0 = 0.2$ decay rate = 1
1	0.1	
2	0.067	
3	0.05	
4	0.04	

other methods:

$$\alpha = 0.95^{\text{epoch num}} \cdot \alpha_0$$

$$\alpha = \frac{k}{\sqrt{\text{epoch num}}} \cdot \alpha_0$$

* Problem of Local optima

- typical optima in DL lie at saddle points
- unlikely to stuck in bad local optima
- plateaus can make learning slow
↳ use momentum/rmsprop

Test Practice -

$$\theta_1 = 10$$

$$\theta_2 = 10$$

$$\beta = 0.5$$

$$v_0 = 0$$

$$v_2 = \beta v_1 + (1 - \beta) \theta_2$$

$$\begin{aligned}v_1 &= \beta v_0 + (1 - \beta) \theta_1 \\&= 0.5 \times 10 - 5\end{aligned}$$

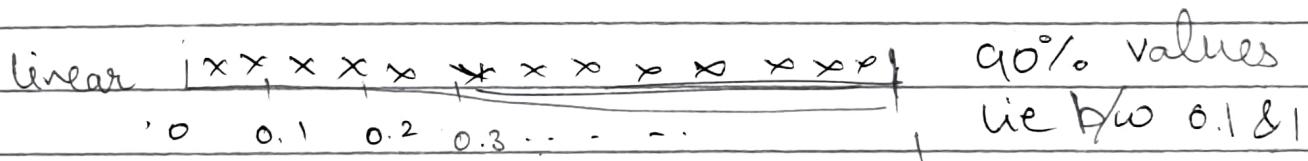
$$\begin{aligned}v_2 &= 0.5 \times 5 + 0.5 \times 10 \\&= 7.5\end{aligned}$$

Hyperparameter Tuning

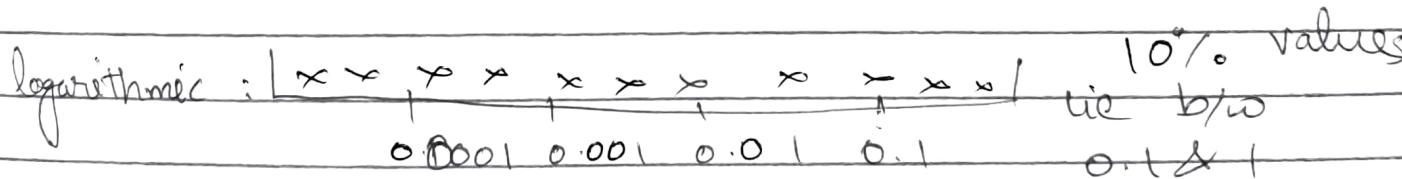
usually in order of importance of tuning:

- α
- β (momentum)
- #hidden units
- mini batch size
- #layers
- learning rate decay
- $\beta_1, \beta_2, \epsilon$ (Adam parameters 0.9, 0.999, 10^{-8}
work for most of the cases)

for use appropriate scale



x: uniform random dist over linear scale



how to apply:

$$r = -4 * np.random.rand()$$

$$\alpha = 10^{**r}$$

↳ usually used for tuning α

Hyperparameters for exponentially weighted averages

$$\beta = 0.9 \dots 0.999 \quad \boxed{1 \times x \times x \times x}$$

last 10 values

last 1000 values

$$1-\beta = 0.1 \dots 0.001$$



this is because

$$0.1 \quad 0.01 \quad 0.001 \\ 10^{-1} \quad 10^{-2} \quad 10^{-3}$$

algorithms are highly sensitive when

$$r \in [-3, -1]$$

$$\beta \approx 1$$

$$1-\beta = 10^r$$

$$\text{ie } \beta = 0.999 \rightarrow \text{last 1000 days} \quad \beta = 1 - 10^r$$

$$0.9995 \rightarrow \text{last 2000 days}$$

$$0.9 \rightarrow \text{last 10}$$

$$0.9005 \rightarrow \text{last } 10^{0.5}$$

notice the difference when
 β is closer to 1

Normalizing Activations Patches

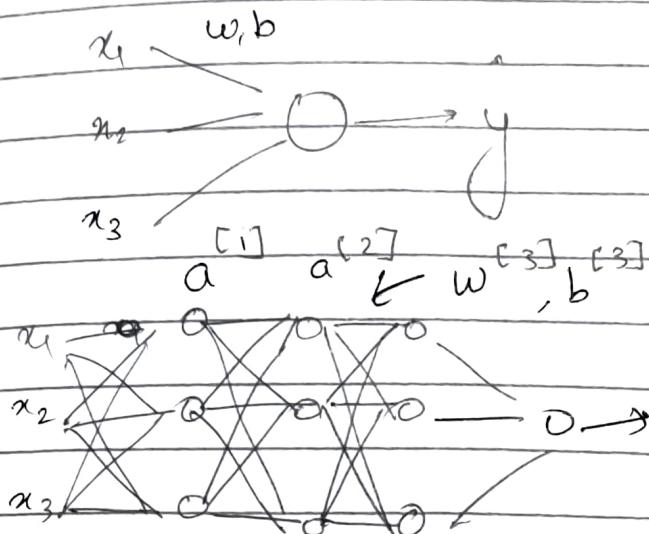
In normalization, we

$$\mu = \frac{1}{m} \sum_i x^{(i)}$$

$$x = x - \mu$$

$$\sigma^2 = \frac{1}{m} \sum_i x^{(i)} - \mu^2$$

$$x = x / \sigma^2$$



we can normalize the inputs to hidden layers too

usually z

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

$$z_{\text{norm}} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

If

$$\gamma = \sqrt{\sigma^2 + \epsilon}$$

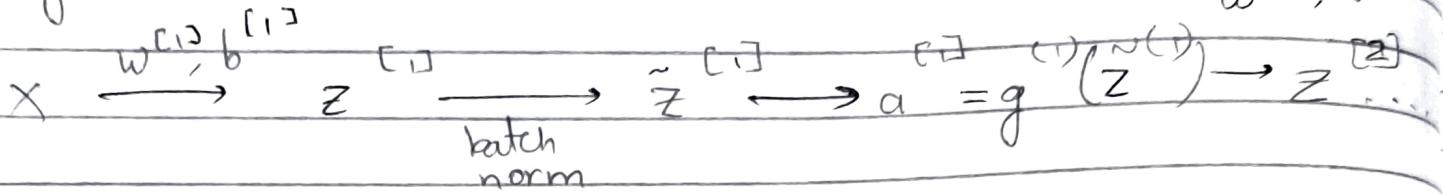
$$\tilde{z}^{(i)} = \gamma z_{\text{norm}} + \beta$$

$$\begin{aligned} \beta &= \mu \\ \tilde{z}^{(i)} &= z^{(i)} \end{aligned}$$

use $\tilde{z}^{(i)}$ instead of $z^{(i)}$

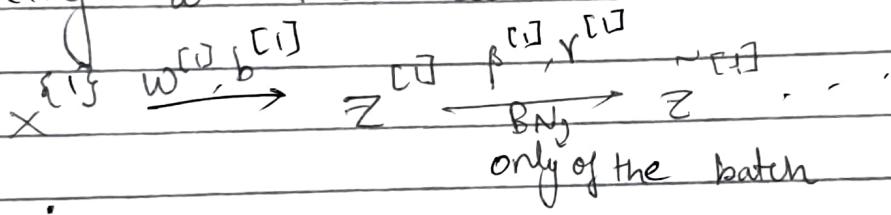
$$U \cdot 3^0 \quad a^{[L]} = \frac{z^L}{\sum_{i=1}^L t_i}$$

Adding Batch normalization to NN



parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}$
 $\beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]}$

Working w mini batches



Multiclass Classification

- Softmax Regression

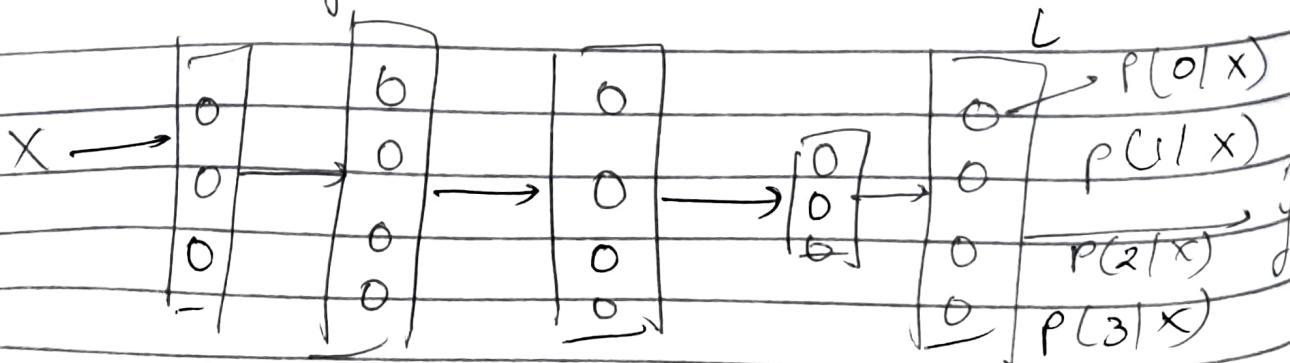
Classes: 1 → cat

2 → dog

3 → chicken

0 → none

C = no of classes = 4 (0, 1, 2, 3)



$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]} \quad (4,1)$$

Activation fn
 $t = e^{\frac{z^{[l]}}{t_i}}$

$$a^{[l]}_{(4,1)} = \frac{e^{z^{[l]}}}{\sum t_i}, \quad a_i^{[l]} = \frac{e^{z_i^{[l]}}}{\sum t_i}$$

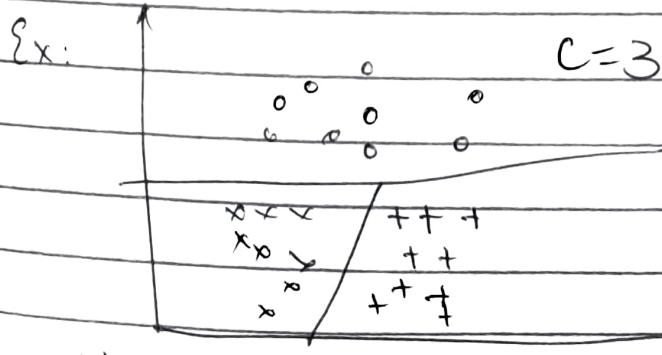
$$z^{[l]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \Rightarrow t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$a^{[l]} = \frac{t}{\sum t_i} = \frac{1}{176.3} \quad \sum t_i = 176.3$$

$$\begin{array}{|c|} \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline \end{array} \rightarrow \frac{e^5}{176.3} = 0.842$$

$$\begin{array}{|c|} \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline 0 \\ \hline \end{array} \rightarrow \frac{e^2}{176.3} = 0.042$$

$$a^{[l]} = g^{[l]}(z^{[l]}) \quad (4,1)$$



Softmax creates linear boundaries (straight lines)

Softmax generalizes logistic to c classes
 if $c=2$: logistic = softmax