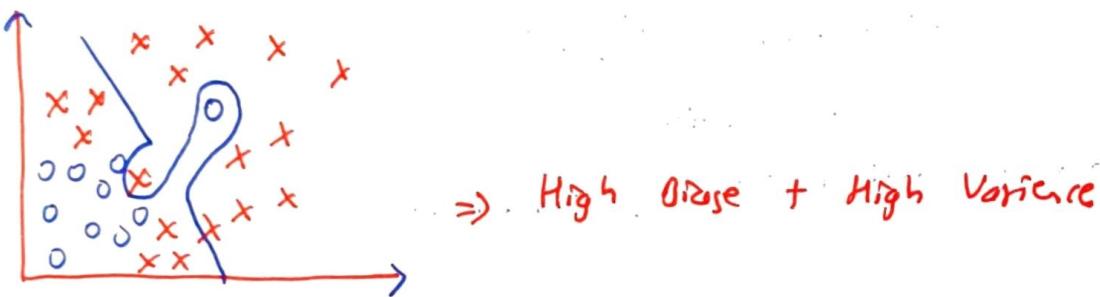


IMPROVING NEURAL NETWORKS: HYPERPARAMETER TUNNING

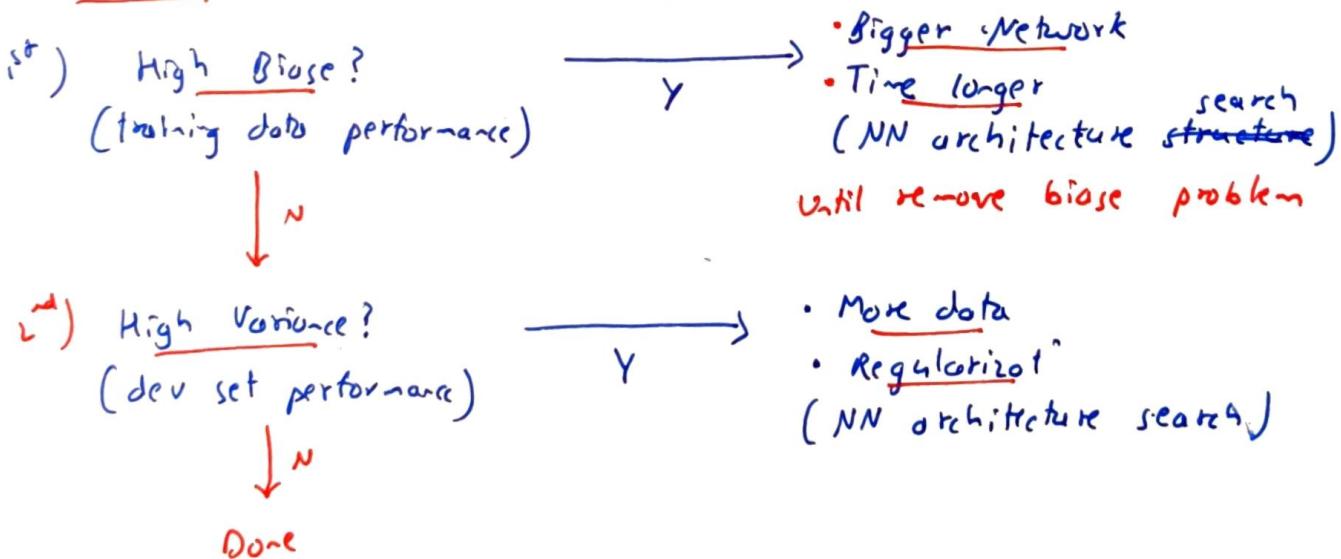
BIASE & VARIANCE

Train set error:	10%	15%	15%	0.5%
Dev set error:	11%	16%	30%	1%
	\Rightarrow High Variance	\Rightarrow High bias	\Rightarrow high variance \Rightarrow high bias	\Rightarrow low bias \Rightarrow low variance

Optimal error (Bayes) error: $\approx 0\%$



Basic Recipe for Machine Learning:



"Bias Variance trade off" \rightarrow earlier Deep Learning era
 { "not much methods to correct one without hurting other"
 - But now it can be done in Deep Learning

REGULARIZATION

- to prevent overfitting or reduce variance

logistic regression: $w \in \mathbb{R}^n$, $b \in \mathbb{R}$, $\lambda = \text{regularization parameter}$

$$\min J(w, b)$$

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

$$\underline{L_2 \text{ regularization}}: \|w\|_2^2 = \sum_{j=1}^n w_j^T w$$

$$\underline{L_1 \text{ Regularization}}: \frac{\lambda}{m} \sum_{j=1}^n |w_j| = \frac{\lambda}{2m} \|w\|_1$$

w will be sparse \because a lot of zeros

Omit
 $\because b$ is only one parameter, whereas w is high-dimensional parameter vector, which is more beneficial to be regularized

Neural Network:

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w_l^{[l]}\|_F^2$$

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{n^{[l]}} (w_{ij}^{[l]})^2 : \text{"Frobenius norm"}$$

$$w: (n^{[L-1]}, n^{[L]})$$

$$\boxed{dw^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]}} \quad \text{"Weight decay"}$$

$$w^{[l]} = w^{[l]} - \alpha dw^{[l]}$$

$$\Rightarrow w^{[l]} = w^{[l]} - \alpha \left((\text{from backprop}) + \frac{\lambda}{m} w^{[l]} \right)$$

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

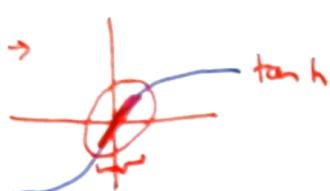
$$\text{or } \boxed{w^{[l]} = w^{[l]} \left(1 - \frac{\alpha \lambda}{m} \right) - \alpha (\text{from backprop})}$$

slightly less than 1

How does regularization prevent overfitting?

→ By reducing weights by smaller & smaller unit, $\omega^{[1]} \approx 0$, which may turn the problem of high variance or overfitting to a problem of high bias or underfitting.

But in the process, there must be an intermediary point where the fitting is "just right".



$$\text{If } \lambda \uparrow \Rightarrow \omega^{[1]} \downarrow \therefore z^{[1]} = \omega^{[1]} a^{[1-1]} + b^{[1]}$$

$$\therefore \text{if } \omega^{[1]} \downarrow \Rightarrow z^{[1]} \downarrow \downarrow$$

it will be in small range (shown in graph)

∴ $g(z)$ will be linear

Every layer ≈ linear

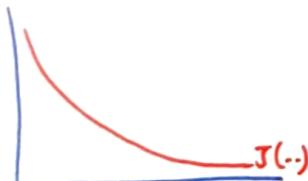
Preventing it from taking complex form like:



⇒ Reducing Overfitting

Implementation tip:

$$J(\cdot) = \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_j \| \omega^{(j)} \|_F^2$$



⇒ Now new $J(\cdot)$ has this additional term

so while doing gradient descent, it must be included in the function

Dropout Regularization:

- Go through each layer of NN and set some probability of dropping some neuron from the Network.
- ⇒ Removing ingoing and outgoing nodes from it
- ⇒ Reduced Network ⇒ smaller examples to be trained ⇒ Reducing Overfitting

1) Inverted Dropout:

Illustrate with layer $l=3$, keep-prob = 0.8

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

$$a_3' = \text{np.multiply}(a_3, d_3) \quad \left\{ \begin{array}{l} \text{chances of elements of } d_3 \text{ to be zero, this makes} \\ \text{activations for that elements to be zero} \end{array} \right.$$

$a_3' / \text{keep_prob} \rightarrow \text{Inverted Dropout}$

Eg.: 50 units in $l=3$, if 0.8 probability of keeping them ⇒ 10 units shut off

$$z^{[4]} = w^{[4]} \cdot o^{[3]} + b^{[4]}$$

↳ reduced by 20%

$$l=0.8$$

⇒ Making test time easier

Making predict @ test time:

$$a^{[0]} = x$$

No drop out

$$z^{[1]} = w^{[1]} a^{[0]} + b^{[1]}$$

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

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \dots$$

$$\downarrow \nearrow y$$

Why does drop out work?

One Neuron

~ can't rely on any one feature, so have to spread out weights on each of previous units \Rightarrow shrinking weights (squared norm of weights)

$\Rightarrow \sim L_2$ Regularizatⁿ

\rightarrow can have keep-prob diff. for different layer

layers with
like n more parameters can have
higher or lower keep-prob
layer with less parameter can
have keep-prob close to 1.0

Downside:

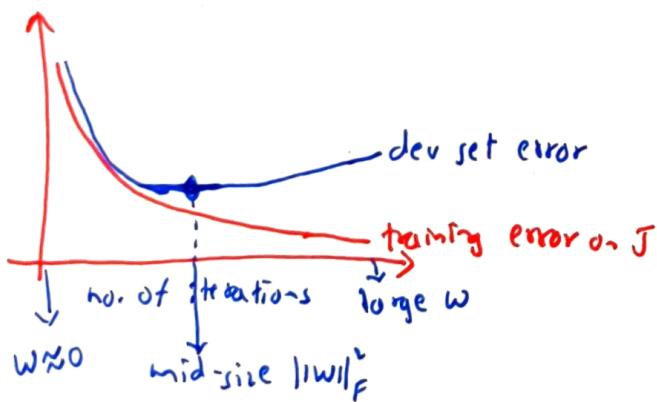
- cost f^* $J(\cdot)$ is not well defined

Other Regularization Techniques:

i) Data augmentation:

Like flipping image; Random crops of image; random distortions & transforms

ii) Early stopping:



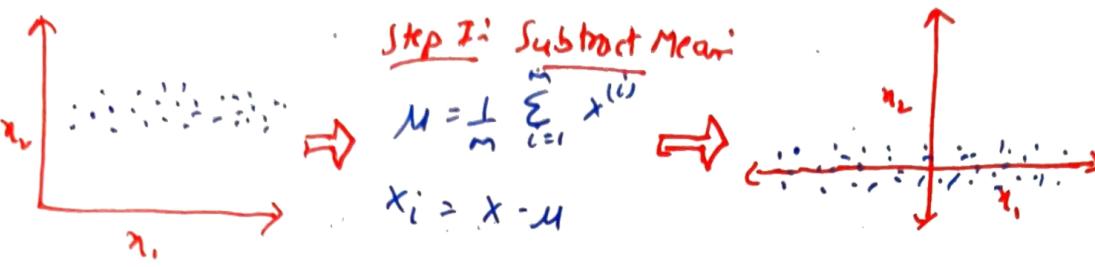
In ML: Downside!

- optimize cost f^* J
- Not overfit

- Early stopping both these tasks and can't do one at a time, so can't try much complicated tasks

Orthogonalizatⁿ

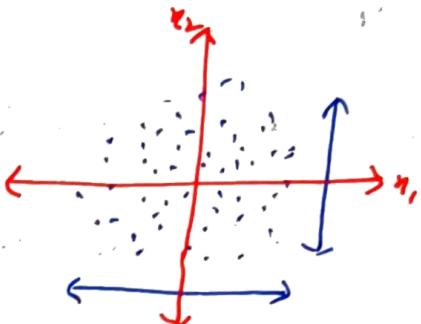
Normalizing Inputs



Step II: Normalize Variance:

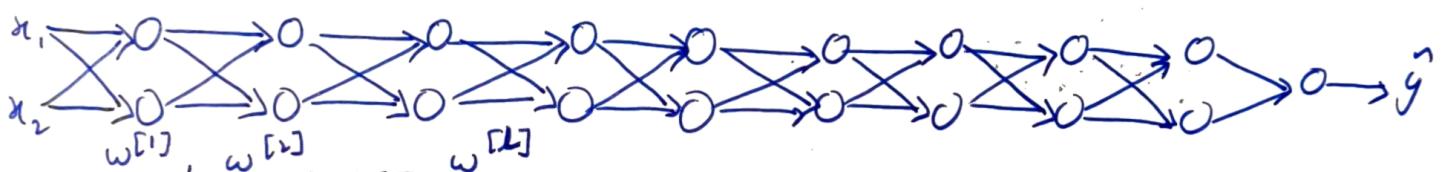
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \bar{x})^2 \quad \text{element-wise square}$$

$$x_i = \frac{x^{(i)} - \bar{x}}{\sigma}$$



Note: Use same \bar{x} & σ^2 to normalize train & test set

VANISHING / EXPLODING GRADIENTS



$$g(z) = z, \quad b^{[L]} = 0$$

$$y^* = w^{[L]} w^{[L-1]} w^{[L-2]} \dots w^{[2]} w^{[1]} \underbrace{w^{[1]} x}_z$$

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

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

$$a^{[2]} = g(z^{[2]}) = g(w^{[2]} a^{[1]})$$

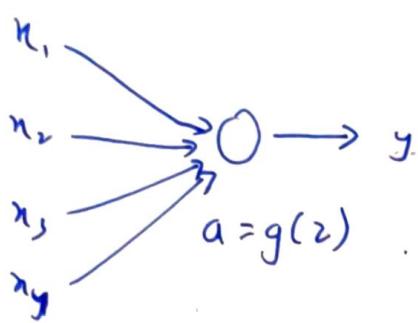
$$w^{[1]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \quad \hat{y} = w^{[L]} \underbrace{\begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x}_{\rightarrow 1.5^{L-1} x}$$

$$\text{If } w^{[1]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \quad \dots \quad \underbrace{\begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^{L-1} x}_{\rightarrow 0.5^{L-1} x}$$

$\Rightarrow w^{[1]} > 1 \Rightarrow$ Activations can explode

$\text{If } w^{[1]} < 1 \Rightarrow$ Activations can decrease exponentially

Weight Initialization for Deep Networks



$$z = \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_n x_n + b$$

let, $b > 0$

larger the n is smaller ω_i

one way: $\text{Var}(\omega_i) = \frac{1}{n}$

$$\omega^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{1}{n^{(l-1)}}\right)$$

If ReLU: $\text{Var}(\omega_i) = \frac{2}{n}$, $g^{[l]}(z) = \text{ReLU}(z)$

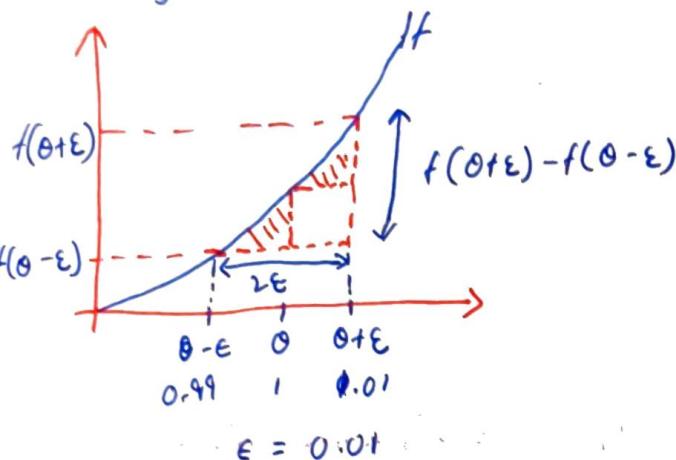
$$\boxed{\omega^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{2}{n^{(l-1)}}\right)}$$

If tanh: $\text{var}(\omega_i) = \sqrt{\frac{1}{n^{(l-1)}}}$ } - Xavier Initialization

$$\sqrt{\frac{2}{n^{(l-1)} * n^{[l]}}}$$

Numerical approximation of gradients

Checking your derivative computation:



$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx g(\theta)$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001$$

$$g(\theta) = 3\theta^2 = 3$$

approx error: 0.0001

(pre slide: 3.0301 error: 0.03)

$$\because f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \xrightarrow{\text{approx}} O(\epsilon^0) = (0.01) = 0.0001$$

GRADIENT CHECKING

Take $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}$ & reshape into a big vector θ

$\underbrace{w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}}_{\text{(Concatenate)}}$

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$$

only do with $d\theta^{[1]}, d\theta^{[2]}, \dots, d\theta^{[L]}$ & reshape into big vector $d\theta$

$\underbrace{d\theta^{[1]}, d\theta^{[2]}, \dots, d\theta^{[L]}}_{\text{Concatenate}}$

Is do the gradient of $J(\theta)$?

\Rightarrow Gradient checking (Grad Check)

for eg

Grad Check:

for each i :

$$d\theta_{approx}^{[i]} = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\approx d\theta^{[i]} = \frac{\partial J}{\partial \theta_i}$$

check it

$$d\theta_{approx} \approx d\theta$$

To do so:

check

$$\boxed{\frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2}} \rightarrow \text{Euclidean distance}$$

In practice: $\epsilon = 10^{-7}$; then $\approx 10^{-7}$ - great
 $\approx 10^{-5}$ - need to be checked
 $\approx 10^{-3}$ - wrong

Implementation of Grad Check:

- Don't use in training - only to debug
- If algo fails grad check, look @ θ -points to try to identify bug.
try to find from where they arise $\begin{smallmatrix} dw \\ db \end{smallmatrix}$
- Remember Regularization

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum \|\theta^{[j]}\|_F^2$$

$d\theta$ = gradient of J wrt. θ [includes regularization term]
- Doesn't work with dropout
 - to work with check - turn off dropout - $\text{keep_prob} = 1.0$
- Run @ random initialization; perhaps again after some training

OPTIMIZATION ALGORITHMS

Mini-Batch Gradient Descent:

Vectorization allows you to efficiently compute on m examples.

$$X = \begin{bmatrix} x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)} \end{bmatrix}_{(n_x, m)}$$

$$Y = \begin{bmatrix} y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(m)} \end{bmatrix}_{(1, m)}$$

What if $m = 5,000,000$? - training still be very slow

∴ Mini-batches of 1,000 each

$$X = \begin{bmatrix} x^{(1)}, x^{(2)}, \dots, x^{(1,000)}, x^{(1,001)}, \dots, x^{(2,000)}, \dots, x^{(m)} \end{bmatrix}_{(n_x, m)}$$

$x^{\{1\}}_{(n_x, 1000)}$ $x^{\{2\}}_{(n_x, 1000)}$ \vdots $x^{\{5,000\}}_{(n_x, 1000)}$

$$Y = \begin{bmatrix} y^{(1)}, y^{(2)}, \dots, y^{(1,000)}, y^{(1,001)}, \dots, y^{(2,000)}, \dots, y^{(m)} \end{bmatrix}_{(1, m)}$$

$y^{\{1\}}_{(1, 1000)}$ $y^{\{2\}}_{(1, 1000)}$ \vdots $y^{\{5,000\}}_{(1, 1000)}$

For $t = 1, \dots, 5,000$ ↗ 1 step of gradient descent using $x^{\{t\}}, y^{\{t\}}$

Forward propagation $X^{\{t\}}$:

$$\left. \begin{aligned} Z^{[1]} &= w^{[1]} X^{\{t\}} + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) \end{aligned} \right\} \text{Vectorized implementation (1,000 example)}$$

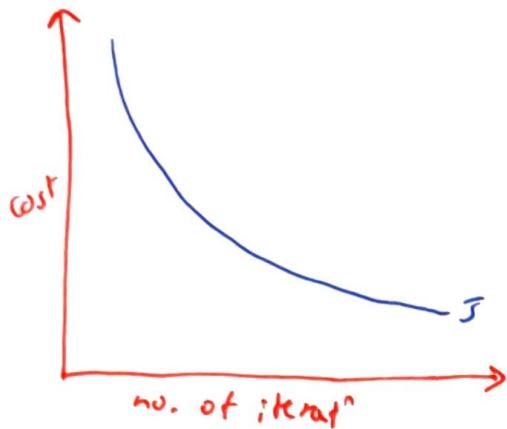
Compute cost $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^L \mathcal{L}(g^{(i)}, y^{(i)}) + \frac{\lambda}{2 \times 1000} \sum_k \|w^{[k]}\|_F^2$

Backprop to compute gradients using $J^{\{t\}} \frac{\partial}{\partial w^{[k]}, b^{[k]}} x^{\{t\}}, y^{\{t\}}$

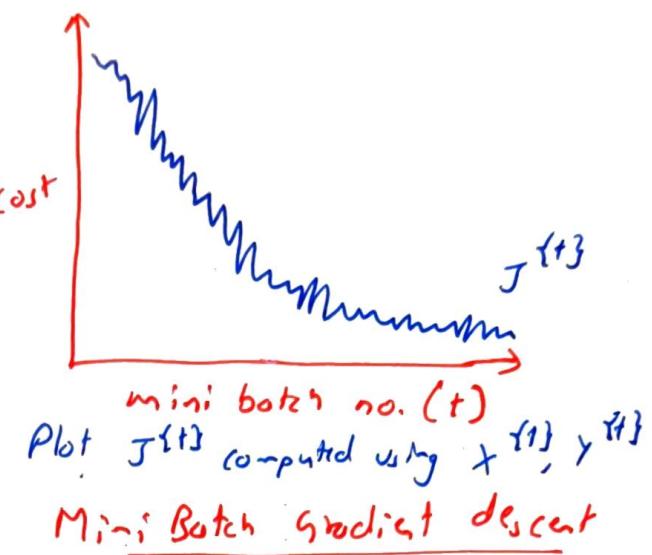
$$w^{[k]} = w^{[k]} - \alpha d w^{[k]} ; b^{[k]} = b^{[k]} - \alpha d b^{[k]}$$

"epoch" → pass through training sets

Training with mini batch gradient descent:



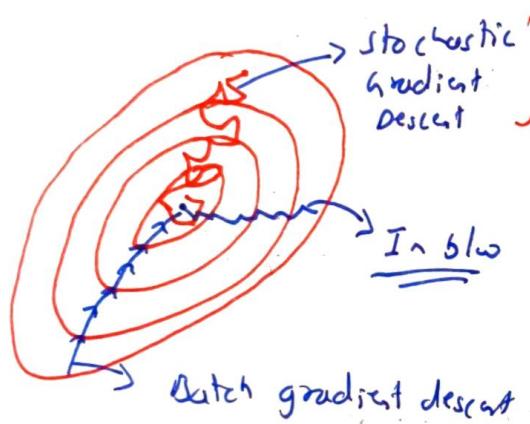
Batch Gradient Descent



Size of mini-batch:

- If mini batch size = m : Batch gradient descent $(x^{(t)}, y^{(t)}) = (x, y)$

If mini batch size = 1 : Stochastic gradient descent - every eg. is own mini batch



In practice: somewhere in b/w

$\frac{1}{d} m$

Faster Learning
• Vectorizatⁿ

• Make progress without
needed to wait
for processing entire
training set

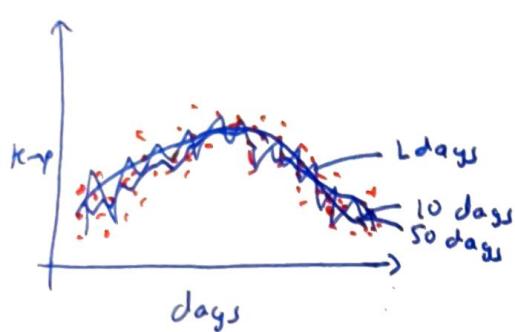
Note: - If small training set: - Use batch gradient descent.
 $(n \leq 2000)$

- Typical mini-batch sizes: $\frac{64}{2^6}, \frac{128}{2^7}, \frac{256}{2^8}, \frac{512}{2^9}$

- Mini-batch size: mini-batch fits in CPU/GPU memory
 $x^{(t)}, y^{(t)}$

Exponentially Weighted Averages:

Temperature in London:



$$\begin{aligned}
 V_0 &= 0 \\
 V_1 &= 0.9V_0 + 0.1\theta_1 \\
 V_2 &= 0.9V_1 + 0.1\theta_2 \\
 V_3 &= 0.9V_2 + 0.1\theta_3 \\
 &\vdots \\
 V_t &= 0.9V_{t-1} + 0.1\theta_t
 \end{aligned}$$

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

$$\beta = 0.9 \approx 10 \text{ days' temp}$$

$$\beta = 0.98 \approx 50 \text{ days}$$

$$\beta = 0.5 \approx 2 \text{ days}$$

V_t : approximately average over $\approx \frac{1}{1-\beta}$ days' temp.

$$V_{100} = 0.99V_{99} + 0.1\theta_{100} \rightarrow 0.1\theta_{98} + 0.9V_{97}$$

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

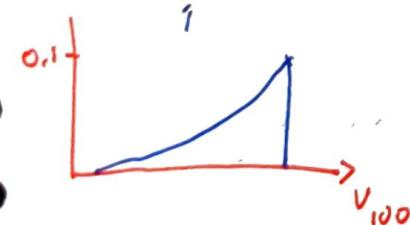
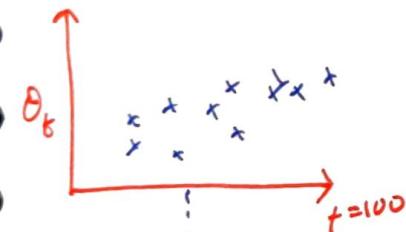
$$= 0.1\theta_{100} + 0.1 \times 0.9\theta_{99} + 0.1(0.9)^2\theta_{98} + 0.1(0.9)^3\theta_{97} + \dots$$

$$0.9^{\frac{10}{10}} \approx 0.35 \approx \frac{1}{e}$$

$$(1-\varepsilon)^{\frac{1}{\varepsilon}} = \frac{1}{e}$$

$\left\{ \begin{array}{l} \varepsilon = 0.1 \text{ in this case} \end{array} \right.$

$$\varepsilon = 1 - \beta$$



Implementing exponentially weighted avg.:

$$V_0 = 0$$

$$V_1 = \beta V_0 + (1-\beta) \theta_1$$

$$V_2 = \beta V_1 + (1-\beta) \theta_2$$

:

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

$$V_\theta = 0$$

$$V_\theta := \beta v + (1-\beta) \theta,$$

$$V_\theta := \beta v + (1-\beta) \theta_2$$

$$\vdots$$

$$\overline{V_\theta = 0}$$

Repeat {

Get Next θ_t

$$V_\theta := \beta V_\theta + (1-\beta) \theta_t$$

}

Bias Correction:

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t, \text{ let } \beta = 0.98$$

$$V_0 = 0$$

$$V_1 = 0.98 V_0 + 0.02 \theta_1 = 0.02 \theta_1$$

$$V_2 = 0.98 V_1 + 0.02 \theta_2 = 0.98(0.02) \theta_1 + 0.02 \theta_2 \\ = 0.0196 \theta_1 + 0.02 \theta_2$$

$$\frac{V_t}{1-\beta^t}$$

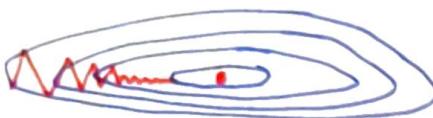
$$\text{When } t=2: 1-\beta^t = 1-(0.98)^2 = 0.0396$$

$$\frac{V_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396} \} \text{ weighted average}$$

⇒ removes bias

⇒ as t increases the bias also approaches 0 ($\because \theta$ approaches 0)

Gradient Descent with Momentum



- i) It prevents from using higher learning rate, as when used, gradient descent will overshoot
- ii) want less slower learning on vertical axis but want faster learning on horizontal axis

\Rightarrow Gradient descent with momentum:

Momentum:

On iterat' t:

Compute dW, db on current mini-batch

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

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

Update weights:

$$w = w - \alpha V_{dw} ; b = b - \alpha V_{db}$$

What it does: is Average out the oscillation both axis \Rightarrow approx. 0 on vertical axis

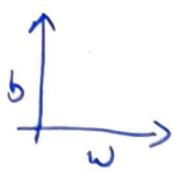
bigger on horizontal axis

\Rightarrow 2 Hyperparameter α , β
 ↓
 learning rate

Momentum
 standard $\beta = 0.9$

\approx avg. over last ≈ 10 gradients

RMSprop (Root Mean Square Prop):



- want to slow down learning on b axis
- " " fast the " " w axis

On iterth t:

Compute $d\omega, db$ on current mini-batch

$$s_{d\omega} = \beta s_{d\omega} + (1-\beta) \frac{dw^2}{\sqrt{s_{d\omega} + \epsilon}} \xrightarrow{\text{element wise square}}$$

$$s_{db} = \beta s_{db} + (1-\beta) \frac{db^2}{\sqrt{s_{db} + \epsilon}} \xrightarrow{\text{large}}$$

$$\omega := \omega - \alpha \frac{dw}{\sqrt{s_{d\omega} + \epsilon}} \xrightarrow{\text{small}}$$

$$b := b - \alpha \frac{db}{\sqrt{s_{db} + \epsilon}} \xrightarrow{\text{large}}$$

$\Rightarrow w$ updated fast

$\Rightarrow b$ update will be slow

Adam Optimization Algorithm:

(Momentum + RMSprop)

$$v_{d\omega} = 0, s_{d\omega} = 0, v_{db} = 0, s_{db} = 0$$

On iterth t:

Compute $d\omega, db$ using current mini-batch

$$v_{d\omega} = \beta_1 v_{d\omega} + (1-\beta_1) dw, v_{db} = \beta_1 v_{db} + (1-\beta_1) db \leftarrow \text{momentum}$$

$$s_{d\omega} = \beta_2 s_{d\omega} + (1-\beta_2) dw^2, s_{db} = \beta_2 s_{db} + (1-\beta_2) db^2 \leftarrow \text{RMSprop}$$

$$v_{d\omega}^{\text{corrected}} = \frac{v_{d\omega}}{(1-\beta_1^t)}, v_{db}^{\text{corrected}} = \frac{v_{db}}{(1-\beta_1^t)}$$

$$s_{d\omega}^{\text{corrected}} = \frac{s_{d\omega}}{(1-\beta_2^t)}, s_{db}^{\text{corrected}} = \frac{s_{db}}{(1-\beta_2^t)}$$

Updating:

$$\omega := \omega - \alpha \frac{v_{d\omega}^{\text{corrected}}}{\sqrt{s_{d\omega}^{\text{corrected}}} + \epsilon}$$

$$b := b - \alpha \frac{v_{db}^{\text{corrected}}}{\sqrt{s_{db}^{\text{corrected}}} + \epsilon}$$

Hyperparameters choices:

- α : needs to be tuned
- β_1 : 0.9 (d w)
- β_2 : 0.999 (d w^2)
- ϵ : 10^{-8}

Usually - fix use std. values of β_1, β_2
 ϵ & try tuning α

Learning Rate Decay:

1 epoch = 1 pass through the data

$$\alpha = \frac{1}{1 + \text{decay rate} * \text{epoch}} \alpha_0$$

Epoch	α
1	0.1
2	0.67
3	0.5
4	0.4
:	:

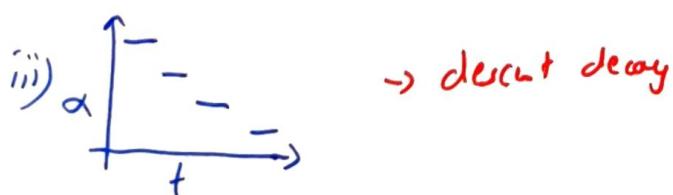
$$\alpha_0 = 0.2$$

decay rate = 1

other learning rate decay methods:

i) $\alpha = 0.95^{\text{epoch-no.}} \cdot \alpha_0$ - exponential decay

ii) $\alpha = \frac{k}{\sqrt{\text{epoch-no.}}} \cdot \alpha_0$



iv) Manual decay

HYPERPARAMETERS TUNING

$\alpha \rightarrow$ most imp.

$\beta \approx 0.9$

$\beta_1, \beta_2, \epsilon$

no. of layers

no. of hidden units

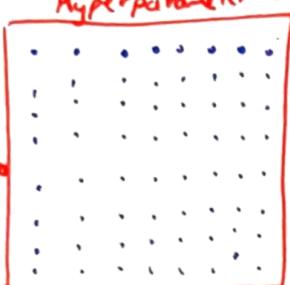
learning rate decay

mini-batch size

2nd most imp.

Hyperparameter 2 $\rightarrow \epsilon$

α
Hyperparameters 1



Grid Method

[useful only when very few hyperparameters]

Hyperparameter 2 $\rightarrow \epsilon$

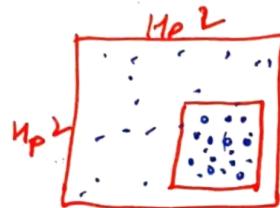
α
Hyperparam-
eter 1



Random

[Better for more no. hyperparameters]

Coarse to fine scheme \rightarrow



- start by exploring bigger region
- then focusing on smaller region where you are getting better results
- then do more sampling in that region

Using an appropriate scale to pick hyperparameters:

i) Picking hyperparameters @ random:

$$n^{[1]} = 50, \dots, 100$$

$\xrightarrow{50} \xrightarrow{\dots} \xrightarrow{100}$ } picking w. @ random

no. of layers : 2-4 } Sampling layer @ random

ii) Appropriate scale for hyperparameters:

$$\alpha = 0.0001, \dots, 1$$

$$0.0001 \xrightarrow{\dots} \xrightarrow{\dots} \xrightarrow{\dots} \xrightarrow{\dots} \xrightarrow{\dots}, \quad \left\{ \text{Linear scale not to be used} \right\}$$

$$0.0001 \xrightarrow{\dots} 0.001 \xrightarrow{\dots} 0.01 \xrightarrow{\dots} 0.1 \xrightarrow{\dots} 1, \quad \left\{ \text{Logarithmic scale} \right\}$$

$$r = -4 * np.random.rand() \leftarrow r \in [-4, 0]$$

$$\alpha = 10^r \leftarrow 10^{-4} \dots 10^0$$

$$10^a \approx 10^b \quad \& \quad a \in [a, b] \quad \alpha = 10^r$$

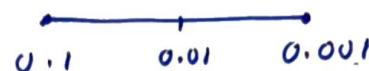
In this eg. $a = \log_{10}(0.0001) = -4$, $b = \log_{10}(1) = 0$

iii) Hyperparameters for exponentially weighted averages:

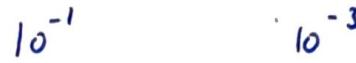
$$\beta = 0.9 \dots 0.999$$



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



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



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

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

NOTE:

If β is close to 1, the sensitivity of results gets changes even with very small changes in β

If $\beta: 0.9000 \rightarrow 0.9005$ } \rightarrow might not have much effect on the changes

$\beta: 0.999 \rightarrow 0.9995$ } ~ 10 values

\rightarrow will have a greater effect

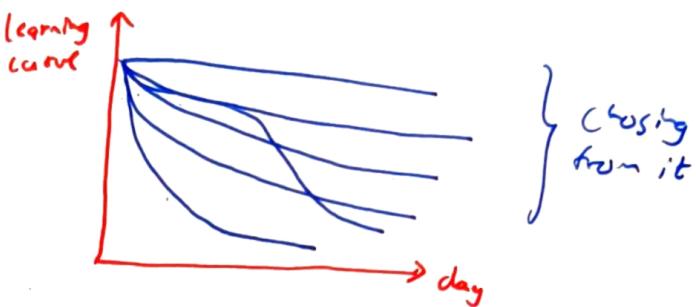
i. Logarithmic scale should be used ~~instead~~ as it give proper weightage to the scale, rather than linear scale.

Hyperparameters tuning in practice: Pandas vs. Caviar

1) Babysitting One Model:



2) Training Many models in parallel:



Panda Approach

{ "the were very few children,
but when have, take care of them only" }

Caviar Approach

{ "when fish reproduce, the reproduced
many eggs and see if any of
them survive" }

BATCH NORMALIZATION

Normalizing inputs to speed up learning

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

$$x_{\text{norm}} = x - \mu$$

$$\sigma^2 = \frac{1}{n} \sum_i x^{(i)2}$$

$$x = \frac{x - \mu}{\sigma^2}$$

In deeper Nets: Can we normalise the values of $a^{[L]}$ so as to train $w^{[L+1]}, b^{[L+1]}$ faster?

Normalize $z^{[L]}$ or $a^{[L]}$
↓
more used in practice

Given some intermediate values in NN

$$\underbrace{z^{(1)}, \dots, z^{(n)}}_{z^{(L)}[i]}$$

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

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

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

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

↳ Learnable parameters of model

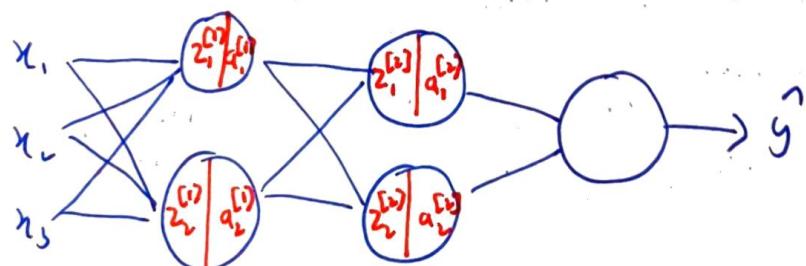
} allows to set mean of $\tilde{z}^{(i)}$ whatever you want it to be

$$\text{If: } \gamma = \sqrt{\sigma^2 + \epsilon} \quad \text{Then: } \tilde{z}^{(i)} = z^{(i)}$$

$$\beta = \mu$$

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

Adding Batch Norm to a Network:



$$X \xrightarrow{\omega^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow[\text{Batch Norm (BN)}]{\gamma^{[1]}, \beta^{[1]}} \tilde{Z}^{[1]} \rightarrow a^{[1]} = g(\tilde{Z}^{[1]}) \xrightarrow{\omega^{[2]}, b^{[2]}} Z^{[2]}$$

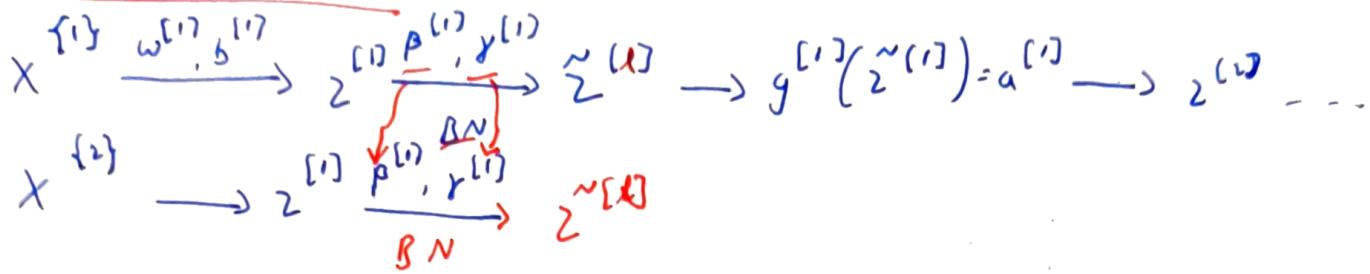
Parameters: $\omega^{[1]}, b^{[1]}, \omega^{[2]}, b^{[2]}, \dots, \omega^{[L]}, b^{[L]}$

$\rightarrow \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]}$

$\beta^{[L]} = \beta^{[L]} - d\beta^{[L]}$

$$a^{[L]} \leftarrow \tilde{Z}^{[L]}$$

Working with mini-batch:



Params: $w^{[l]}, b^{[l]}, \beta^{[l]}, \gamma^{[l]}$

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

$\Rightarrow z^{[l]} = w^{[l]} a^{[l-1]}$

$z_{\text{norm}}^{[l]} = \gamma^{[l]} z_{\text{norm}}^{[l]} + \beta^{[l]}$

Batch Norm 1st scale $z^{[l]}$ to mean 0 and then subtract & rescale it
with β & γ

Whatever is the value of $b^{[l]}$, it will get subtracted out $\therefore b^{[l]}$ is ignored

Mean Subtract step

Implementing gradient descent:

for $t = 1 \dots \text{no. of mini batches}$

Compute forward prop on mini batch $X^{[t]}$

In each hidden layer, use BN to replace $\underline{z}^{[l]}$ with $\tilde{z}^{[l]}$

Use back prop to compute $d\omega^{[l]}, d\beta^{[l]}, d\gamma^{[l]}$

Update parameters $\omega^{[l]} := \omega^{[l]} - \alpha d\omega^{[l]}$

$\beta^{[l]} := \beta^{[l]} - \alpha d\beta^{[l]}$

$\gamma^{[l]} = \dots$

{ no. of b 's removed } \therefore we removed b

Why Batch Norm Works?

- It makes weights deeper in deeper NN more robust to changes as compared to earlier weights in NN

"Covariate Shift" - If we maped $X \rightarrow y$ & trained the model, if the distribution of X changes we have to retrain the model

- It is true even if the σ function does not change
- It become worse if function also changes

Happens in DNN when $a^{[l]}$ or $z^{[l]}$ keeps on changing because of $w^{[l-1]}, b^{[l-1]}$ and this in turn affect the reltn with $w^{[l+1]}, b^{[l+1]}$

- But even if the values of $z^{[l]}$ changes, but the value of mean & variance remain same, thus $\beta^{[l]}, \gamma^{[l]}$
 \Rightarrow allowing one layer to learn somewhat independently from previous layers

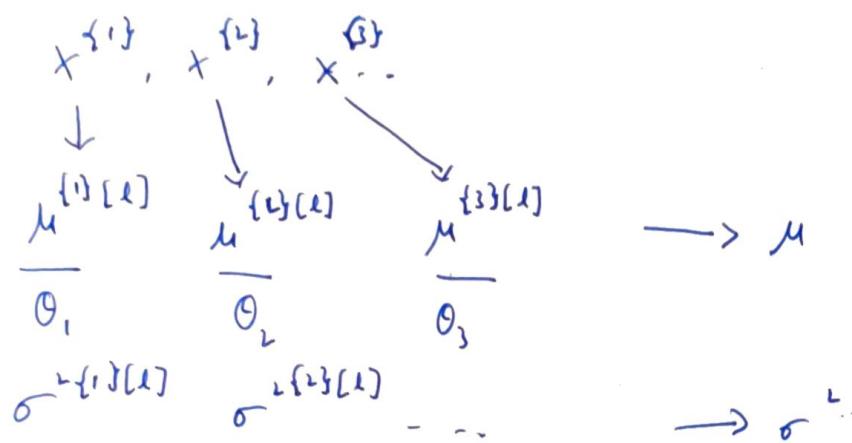
2) Batch Norm as regularization:

- Each mini-batch is scaled by the mean & variance computed on just that mini-batch \Rightarrow little bit noise to $z^{[l]}$
 \Rightarrow making scaling process of $z^{[l]}$ noisier
- So N to dropout, it adds some noise to each hidden layer's activation.
 \Rightarrow slight regularization effect \because noise added is small

If use a bigger mini-batch size \rightarrow reduces this regularization effect

Batch Norm at test time:

μ, σ^2 : estimate using exponentially weighted average (across mini-batch)



At test time:

$$z_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \tilde{z} = \gamma z_{\text{norm}} + \beta$$

MULTIPLE - CLASS CLASSIFICATION

SOFTMAX: on layer L ; C classes = $0, 1, 2, 3$

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

Activation: t^i :

$$t = e^{(z^{[L]})}$$

$$(4,1) \quad a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^4 t_j}; \quad a_i^{[L]} = \frac{t_i}{\sum_{j=1}^4 t_j}$$

$$\text{Let, } z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}; \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}. \quad \sum_{j=1}^4 t_j = 176.3$$

$$a^{[L]} = \frac{t}{176.3} = \frac{1}{176.3} \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

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

Training a Softmax Classifier:

$$(4,1) \quad z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} \quad C=4, \quad g^{[L]} \quad a^{[L]} = g^{[L]} (z^{[L]}) = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

contrast to "hard max"

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Softmax regression generalizes logistic regression to C classes.

Loss function:

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} - \text{cat} \quad a^{[L]} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \quad C = 4$$

$$\boxed{\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^C y_j \log \hat{y}_j}$$

In given e.g. $y_1 = y_3 = y_4 = 0$

$\Rightarrow \mathcal{L} = -y_2 \log \hat{y}_2 = -\log \hat{y}_2 \quad \} \Rightarrow \text{need to make } \log \hat{y}_2 \text{ big}$

Cost:

$$J(\omega^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \cdots & y^{(n)} \end{bmatrix} \\ = \begin{bmatrix} 0 & 0 & 0 & \cdots \\ 1 & 0 & 0 & \cdots \\ 0 & 0 & 1 & \cdots \\ 0 & 1 & 0 & \cdots \end{bmatrix}_{(4, n)}$$

$$\hat{Y} = \begin{bmatrix} \hat{y}^{(1)} \\ \vdots \\ \hat{y}^{(n)} \end{bmatrix} \\ = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}_{(4, n)}$$

Gradient Descent with Softmax:

$$z^{[L]}_{(4,1)} \rightarrow a^{[L]} = \hat{y} \rightarrow \mathcal{L}(\hat{y}, y)$$

Backprop:

$$\frac{dz^{[L]}_{(4,1)}}{dz^{[L]}} = \hat{y} - y$$

$$\frac{\partial J}{\partial z^{[L]}}$$