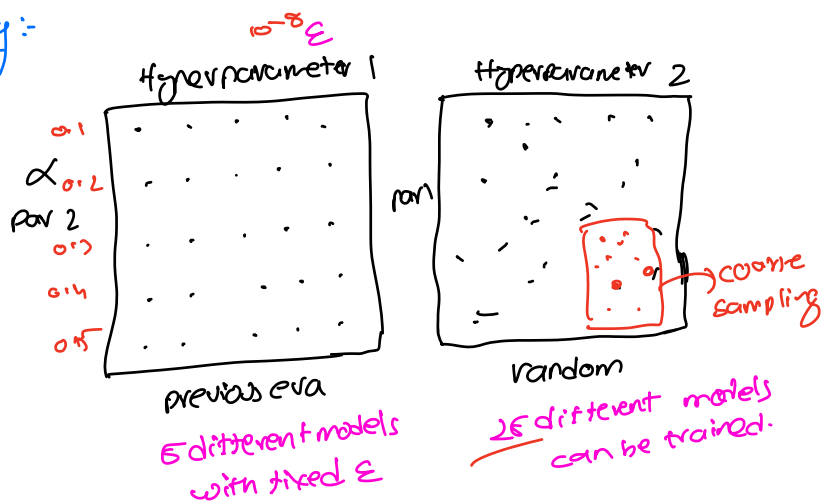


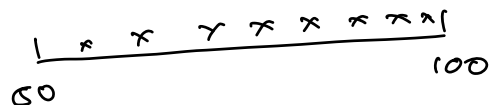
lec1: Hyperparameter Tuning:-

- α
- β 0.9
- $\beta_1, \beta_2, \epsilon$
- layers
- hidden units
- learning rate decay
- mini batch size



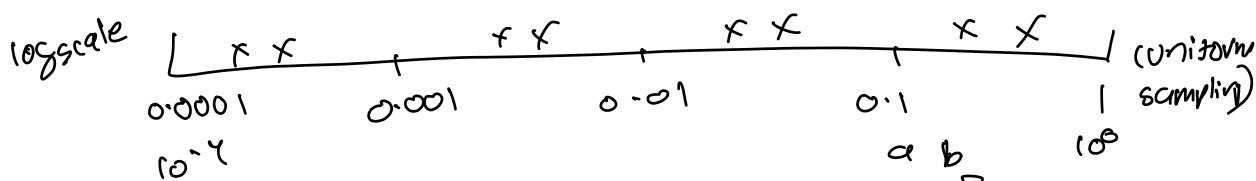
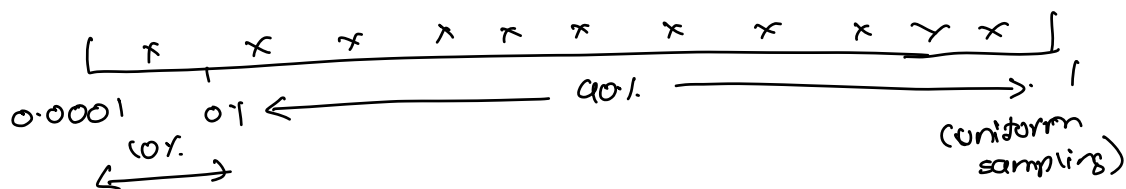
lec2: scale to pick hyperparameters:-

$$n^{[2]} = 50, \dots, 100 \text{ (uniformly random)}$$



$$\# \text{layers: } L: 2-4 \text{ (2, 3, 4) (uniformly random sampling)}$$

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



$$v = -4 * \text{np.random.rand}() \quad v \in [-4, 0]$$

$$\alpha = 10^v \quad v \in [-4, 0]$$

$$10^a \dots 10^b$$

$\beta = 0.9, \dots, 0.999$
 \downarrow
 avg 10 values

\downarrow
 avg 1000 values

$| \times \times \times \times \times \times \times \times |$
 $0.9 \quad \text{uniform sampling} \quad 0.999$
 not good.

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

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

$\gamma \in [-3, -1] \rightarrow \text{uniformly sampling str.}$

$1-\beta = 10^r$

$\beta: 0.9000 \rightarrow 0.4005 \rightarrow 10$

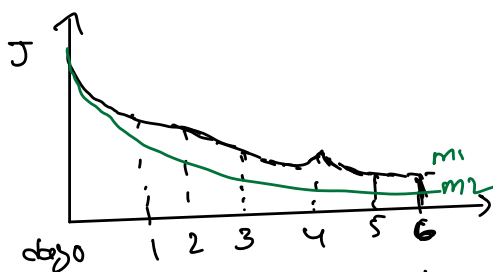
$\beta = 1-10^r$

$\beta: 0.999 \rightarrow 0.9995$
 $1000 \quad \quad \quad -1000$

dec3: Hyperparameter search process:

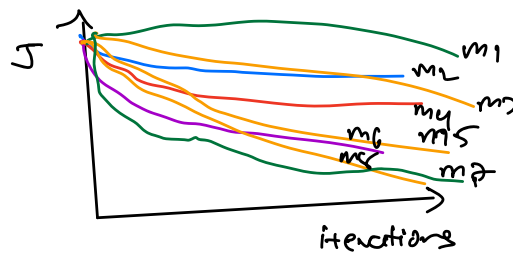
\rightarrow reevaluate hyperparameters occasionally
 train many models in parallel

Babysitting one model



change the parameters
 as the model train.

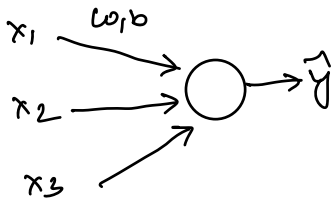
Panda approach



cava approach. (mt)

lec4: Batch Normalization

Normalizing inputs to speed up learning



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

$$x = x - \mu$$

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

$$x = x / \sigma$$

Given some intermediate values in NN $z^{(1)}, \dots, z^{(m)}$ for $\tilde{z}^{(i)}$

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

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

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

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

learnable parameters of model.

if $\gamma = \sqrt{\sigma^2 + \epsilon}$ then $\hat{\tilde{z}}^{(i)} = \tilde{z}^{(i)}$

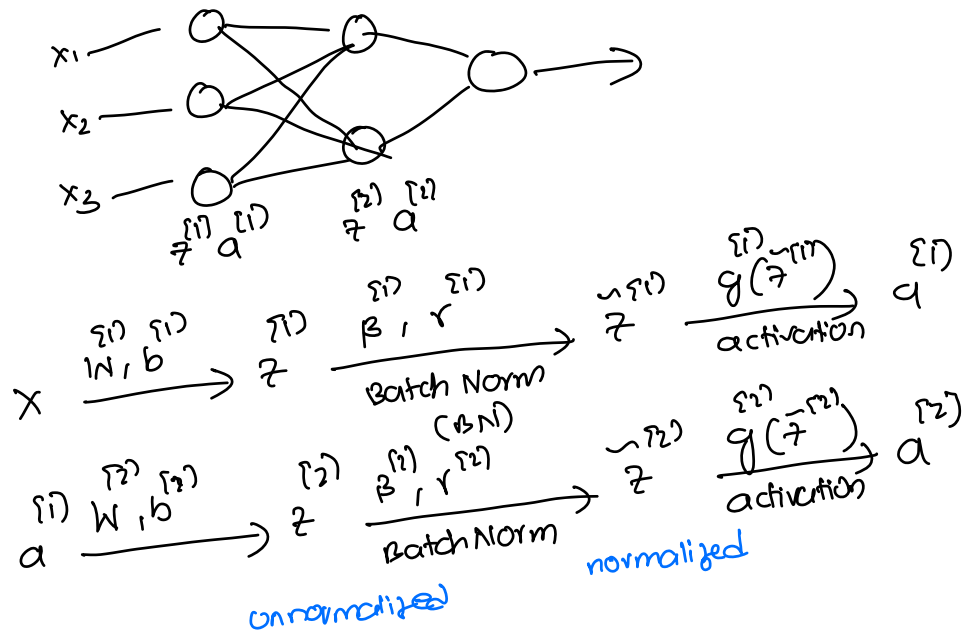
$$\beta = \mu$$

use $\hat{\tilde{z}}^{(i)}$ instead of $\tilde{z}^{(i)}$ for next process.



(γ, β)
we can
vary

lec 6: Batch Norm to NN



parameters:-

$$\begin{aligned} & W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)} \dots W^{(L)}, b^{(L)} \\ & \beta^{(1)}, \gamma^{(1)}, \beta^{(2)}, \gamma^{(2)} \dots \beta^{(L)}, \gamma^{(L)} \end{aligned}$$

$$\begin{aligned} \text{compute } d\beta^{(L)} & \quad \beta^{(L)} = \beta^{(L)} - \alpha d\beta^{(L)} \\ d\gamma^{(L)} & \quad \gamma^{(L)} = \gamma^{(L)} - \alpha d\gamma^{(L)} \end{aligned} \Bigg/ \text{update.}$$

→ nn. batch-normalization.

BN with minibatch:-

for $t=1$, to minibatches.

$$\begin{array}{l}
 x^{(1)} \xrightarrow{w^{(1)}, b^{(1)}} z^{(1)} \xrightarrow[\text{BN}]{\beta^{(1)}, \gamma^{(1)}} \tilde{z}^{(1)} \xrightarrow{g^{(1)}} a^{(1)} \\
 a^{(1)} \xrightarrow{w^{(2)}, b^{(2)}} z^{(2)} \xrightarrow[\text{BN}]{\beta^{(2)}, \gamma^{(2)}} \tilde{z}^{(2)} \xrightarrow{g^{(2)}} a^{(2)}
 \end{array}$$

parameters:- $w^{(1)}, \cancel{b^{(1)}}, \beta^{(1)}, \gamma^{(1)}$

$$z^{(1)} = w^{(1)} a^{(1)} + \cancel{b^{(1)}}$$

BN makes $z^{(1)}$ - mean 0, $\sigma^2 = 1$, then rescale to $\beta^{(1)}, \gamma^{(1)}$

Hence we can omit $b^{(1)}$

$$\tilde{z}_{\text{norm}}^{(1)} = \gamma^{(1)} z_{\text{norm}}^{(1)} + \boxed{\beta^{(1)}} \rightarrow \text{mean.}$$

$$\begin{aligned}
 z^{(1)} &= (n^{(1)}, 1) & b^{(1)} &= (n^{(1)}, 1) \\
 \therefore \beta^{(1)} &= (n^{(1)}, 1) & \gamma^{(1)} &= (n^{(1)}, 1)
 \end{aligned}$$

Implementing GD:

for $t = 1, \dots, \text{num minibatches}$

1. compute forward prop on $x^{(t)}$
 - in each hidden layer, use BN to replace $z^{(l)}$ with $\tilde{z}^{(l)}$
2. use backprop to compute $dw^{(l)}$, ~~$db^{(l)}$~~ , $d\beta^{(l)}$, $d\gamma^{(l)}$

3. update parameters

$$w^{(l)} = w^{(l)} - \alpha dw^{(l)}$$

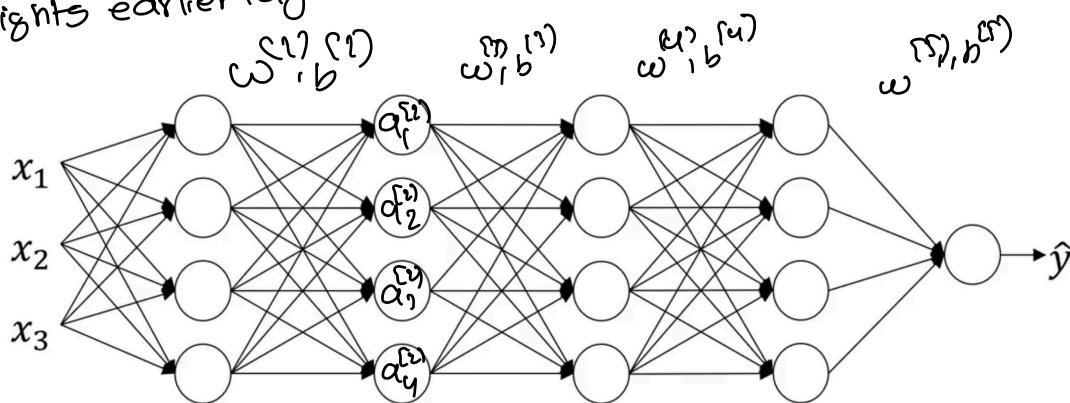
$$\beta^{(l)} = \beta^{(l)} - \alpha d\beta^{(l)}$$

$$\gamma^{(l)} = \gamma^{(l)} - \alpha d\gamma^{(l)}$$

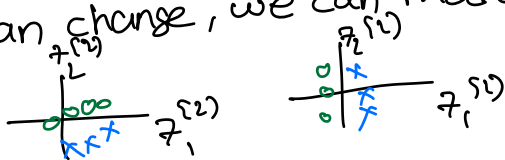
deco: Why BN really works:-

reason:

It makes weights, later or deeper more robust to changes to weights earlier layers of neural network, say in layer one.



BN helps model to learn distribution of hidden layer, so even though $a_i^{(l)}$ can change, we can model that.



BN reduces the problem of input value changing.

BN as regularization:

→ each mini batch $x^{(t)}$ is scaled by mean/variance computed on just that mini batch.

→ This adds some noise to value $z^{(l)}$ within that mini batch. so, similar to dropout, it adds some noise to each hidden layers activation.

→ This has slight regularization effect.

mini batch: 64 512
 more noise less noise
 more regularization less regularization

dec 7: BN at test time:

→ BN process one mini batch at a time during training

→ But during test time, we may need to process one example at a time.

training:

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

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

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

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

Test:

μ, σ^2 : estimate using exponentially weighted average across mini batches.

in layer l :

$$\begin{array}{ccccccc} x^{(1)}, & x^{(2)}, & x^{(3)} & \dots & \dots & \dots & \dots \\ \downarrow & \downarrow & \downarrow & & & & \\ \mu & \mu & \mu & & & & \\ \sigma^2 & \sigma^2 & \sigma^2 & & & & \end{array}$$

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

↑ weighted avg. ↓ weighted avg.

$$\tilde{z} = \sqrt{z_{\text{norm}}} + \beta$$

decg: SoftMax Regression:-

- cats, dogs, baby chick, other
1, 2, 3, 0

$C = \text{number of classes} = 4 \in (0, 1, 2 \dots C-1)$

$$n^{[L]} = C = 4$$

$$n^{[L]} = 4$$

$$\vec{y} = (n^{[L]}, 1) = (4, 1)$$

$$O \quad P(\text{other} | x)$$

$$O \quad P(\text{cat} | x)$$

$$O \quad P(\text{dog} | x)$$

$$O \quad P(\text{bc} | x)$$

softmax layers:-

$$z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]}$$

Activation function:

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

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^4 t_j}$$

$$a_i = \frac{t_i}{\sum_{i=1}^4 t_i}$$

example:-

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

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

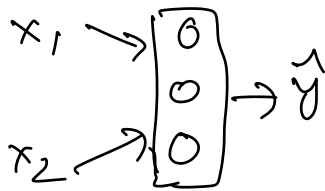
$$\sum_{i=1}^4 t_i = 176.3$$

$$z^{[L]} \rightarrow a^{[L]}$$

$$\vec{y} = a^{[L]} = \frac{t}{176.3} = \begin{bmatrix} 148.4/176.3 \\ 7.4/176.3 \\ 0.4/176.3 \\ 20.1/176.3 \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

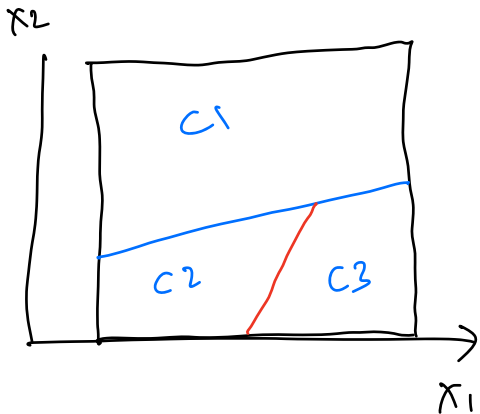
$$a^{[L]} = g \left(z^{[L]} \right)$$

softmax examples



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

$$a^{[1]} = \hat{y} = g(z^{[1]}) \quad (\text{softmax layer})$$



lec 9: Training a softmax classifier:

$$z^{[1]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \rightarrow a^{[1]} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

softmax hardmax

if $C=2$, softmax \rightarrow logistic regression

$$a^{[1]} = \begin{bmatrix} 0.842 \\ -0.158 \end{bmatrix}$$

Loss function:

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$\downarrow \alpha(\hat{y}, y) = - \sum_{j=1}^4 y_j \log \hat{y}_j = - y_2 \log \hat{y}_2 = - \log \hat{y}_2 \downarrow \text{small}$$

$\Rightarrow \hat{y}_2 \uparrow$ big as possible.

$$J(w^{(0)}, b^{(0)}, \dots) = \frac{1}{m} \sum_{i=1}^m \alpha(\hat{y}^{(i)}, y^{(i)})$$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]_{(4, m)}$$

$$\hat{Y} = [\hat{y}^{(1)} \ \hat{y}^{(2)} \ \dots \ \hat{y}^{(m)}]_{(4, m)}$$

GD:-

$$a^{[L]} = \sigma^{[L]}(z^{[L]}) = \hat{y} \rightarrow \alpha(\hat{y}, y)$$

backprop:-

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

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

dec 10:- Deep learning frameworks:

Tensorflow:

$$J(\omega) = \omega^2 - 10\omega + 25$$

$$J(\omega) = (\omega - 5)^2$$

$\omega = 5$ minimizes $J(\omega) = 0$

$$x(0)\omega^2 + x(1)\omega + x(2)$$

