

Week 3 :

Hyperparameter Tuning :

1) Tuning process

α

$\rho \sim 0.9$

$\rho_1, \rho_2, \epsilon \sim 0.9, 0.999, 10^{-8}$

layers

hidden units

learning rate decay

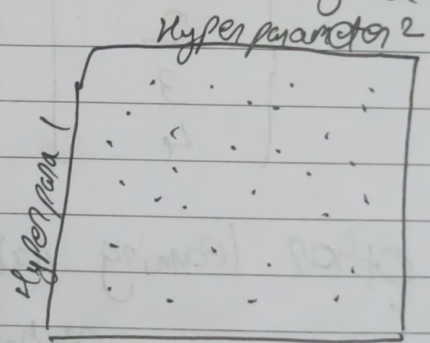
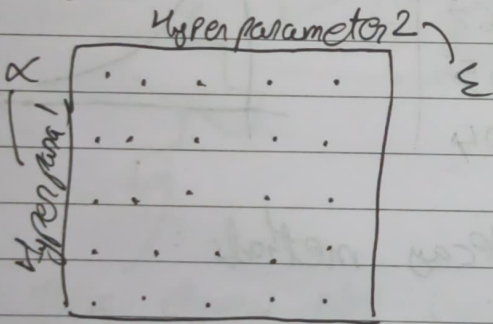
mini-batch size

- most important α

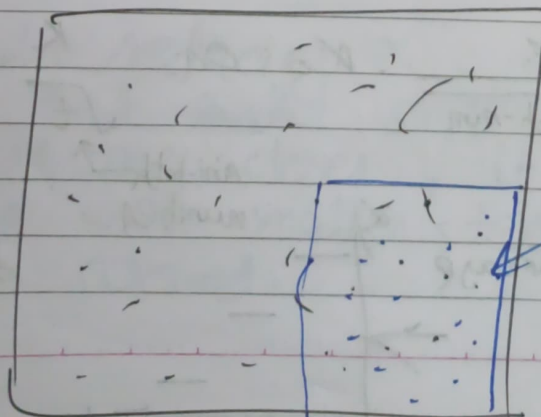
- 2nd

- 3rd

Try random values : Don't use a grid



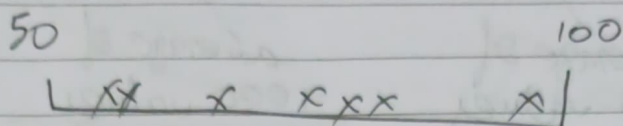
Coarse to fine sampling scheme



Sample more densely in this smaller square

2) Using an appropriate scale to pick hyperparameters

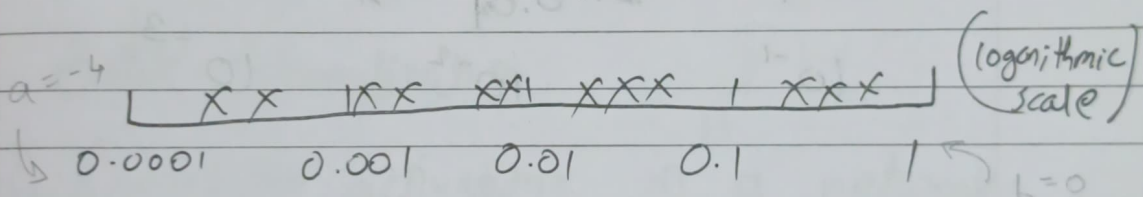
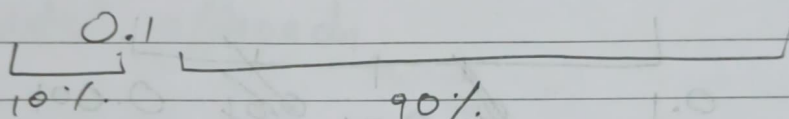
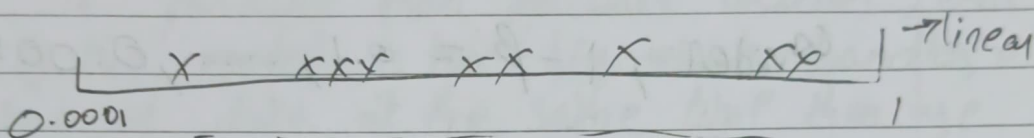
$$n^{[L]} = 50, \dots, 100$$



layers L : 2 - 4

2, 3, 4

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



→ this will be a sampling uniformly at random

In python:

$$r = -4 * \text{random.randn}() \Rightarrow r \in [-4, 0]$$

$$\alpha = 10^r$$

$$\Rightarrow \alpha \in [0.0001, 1]$$

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

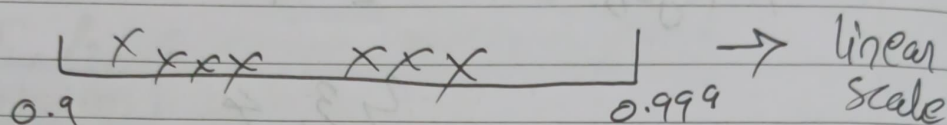
$$r \in [a, b] \text{ uniformly at random}$$

$$\alpha = 10^r$$

Hyparparameters for exponentially weighted averages

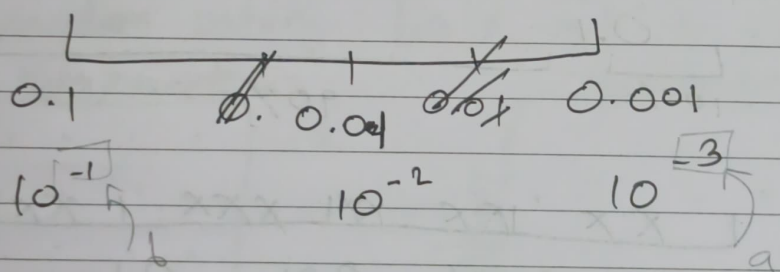
$$\frac{1}{1-\beta} \quad \beta = 0.9, \dots, 0.999$$

\uparrow average of 10 values $\quad \quad \quad \uparrow$ average of 1000 values
 \uparrow $\quad \quad \quad \uparrow$



does not make sense for uniform sampling hence avoid.

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

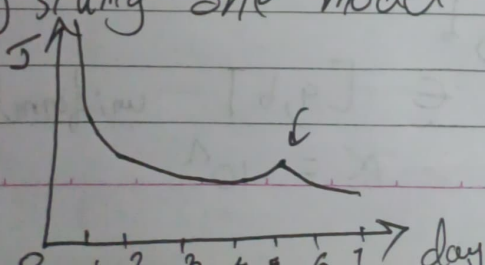


$$\eta \in [-3, -1]$$

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

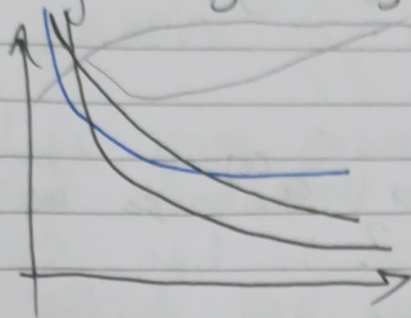
$$\beta = 1 - 10^\eta$$

3) Baby sitting one model:



(panda approach)

Training many models in parallel



→ Pick one that works the best

(Cavier strategy)

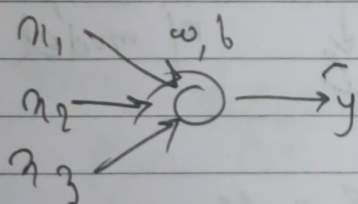
The approach/strategy that we pick depends on the computational resources we have.

If enough computers to train a lot of models in parallel then go with cavier, else if the model is too big which handles a lot of data at the same time then use pandas approach.

Batch Normalization:

1) Normalizing activations in a network:

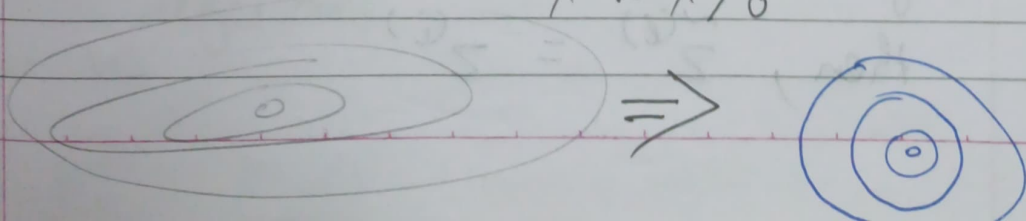
Normalizing inputs to speed up learning

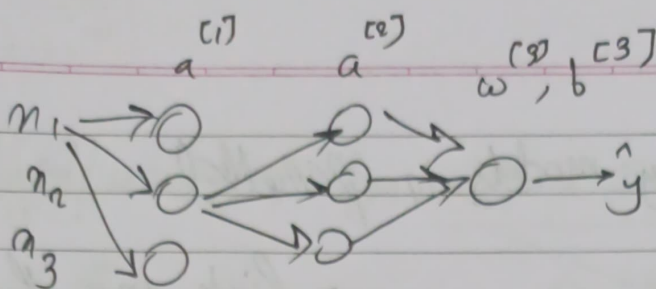


$$\mu = \frac{1}{m} \sum_i x^{(i)}$$

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

$$\tilde{x} = (x - \mu) / \sigma$$





Can we normalize $a^{(2)}$ so as to train faster?
 $w^{(2)}, b^{(2)}$

Normalize $z^{(2)}$

Implementing Batch Norm

Given some intermediate values in NN

$z^{(1)}, \dots, z^{(m)}$
 of some layer l $z^{(l)}(i)$

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

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

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

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

learnable parameters of model

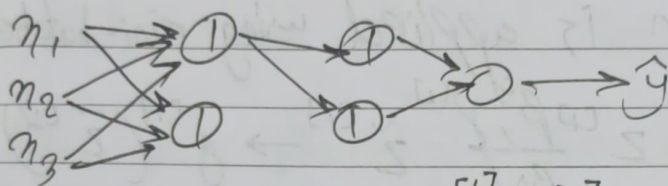
If $\gamma = \sqrt{\sigma^2 + \epsilon}$, $\beta = \mu$

then, $\tilde{z}^{(i)} = z^{(i)}$

Use $\hat{z}^{(i)}[l]$ instead of $z^{(i)}[l]$ for later computations in NN.

Batch norm ensures that the hidden units have standardised mean and variance to 0 and 1 on some other values, which are controlled by β & γ

2) Fitting Batch Norm into a NN



$$\begin{aligned}
 & \mathbf{x} \xrightarrow{w^{[1]}, b^{[1]}} \mathbf{z}^{[1]} \xrightarrow[\text{Batch Norm (BN)}]{\beta^{[1]}, \gamma^{[1]}} \hat{\mathbf{z}}^{[1]} \xrightarrow{w^{[2]}, b^{[2]}} \mathbf{a}^{[1]} = g^{[1]}(\hat{\mathbf{z}}^{[1]}) \\
 & \dots \leftarrow \mathbf{a}^{[2]} = g^{[2]}(\hat{\mathbf{z}}^{[2]}) \leftarrow \hat{\mathbf{z}}^{[2]} \xrightarrow[\text{BN}]{\beta^{[2]}, \gamma^{[2]}} \mathbf{z}^{[2]} \xrightarrow{w^{[3]}, b^{[3]}} \mathbf{a}^{[3]} = g^{[3]}(\hat{\mathbf{z}}^{[3]})
 \end{aligned}$$

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

Note: These values of β are different than the β we had as an hyperparameter in momentum, RMS prop and Adam algorithms.

$d\beta^{[l]}$ can be computed to update β as:

$$\beta^{[l]} = \beta^{[l]} - \alpha d\beta^{[l]} \text{ using gradient descent}$$

We can use other algos like momentum, RMS prop or Adam also to update β .

Using tensorflow we can implement BN as:

→ tf.nn.batch-normalization

Batch-norm is applied using mini-batches also:

$$X^{[l]} \xrightarrow{w^{[l]}, b^{[l]}} z^{[l]} \xrightarrow[\text{BN}]{\beta^{[l]}, \gamma^{[l]}} \tilde{z}^{[l]} \rightarrow g^{[l]}(\tilde{z}^{[l]}) \rightarrow \dots$$

$$X^{[l]} \rightarrow z^{[l]} \xrightarrow[\text{BN}]{} \tilde{z}^{[l]} \rightarrow \dots$$

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

dim = $(n^{[l]}, 1)$ like $z^{[l]}$

since it gets cancelled when we rescale $z^{[l]}$ / on setting it to ≈ 0 permanently.

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

Implementing gradient descent

for $t = 1 \dots \text{num_Mini-batches}$

compute forward prop on $X^{[l]}$

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

Use backprop to compute $dw^{[l]}, db^{[l]}, d\beta^{[l]}, d\gamma^{[l]}$

Update parameters :

$$\left. \begin{aligned} w^{[1]} &:= w^{[1]} - \alpha dw^{[1]} \\ \beta^{[1]} &:= \beta^{[1]} - \alpha d\beta^{[1]} \\ \gamma^{[1]} &:= \gamma^{[1]} - \alpha d\gamma^{[1]} \end{aligned} \right\}$$

Note: Also works with other algos like momentum, RMS prop, adam.

4) Learning on shifting input distribution

5) Batch Norm at test time

Batch norm processed the train data in mini-batches but in ~~test~~ during test phase we will need to process the examples one at a time.

During test time :

Come up with separate μ , σ^2

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

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

$$\theta_1, \theta_2, \theta_3$$

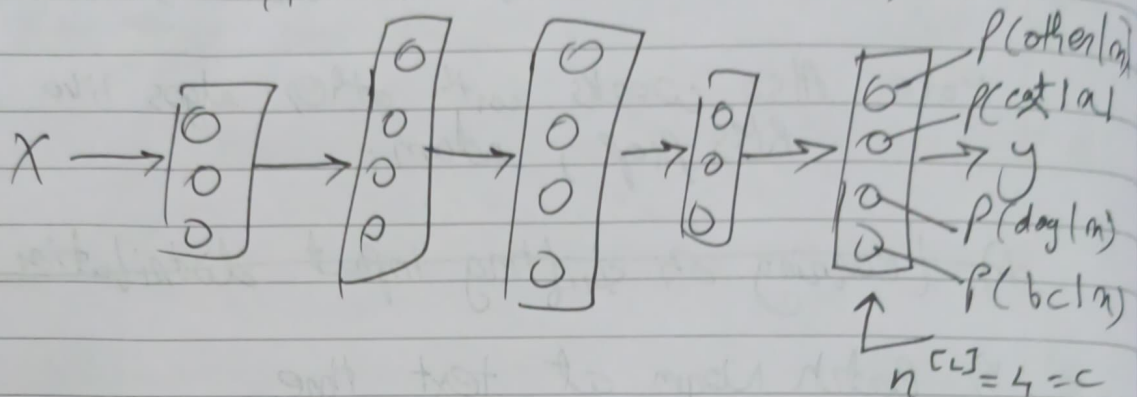
$$z_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}}, \quad \tilde{z} = \gamma z_{\text{norm}} + \beta$$

Multi-class Classification

1) Softmax Regression

eg: Recognizing cats, dogs, and baby chicks

$C = \# \text{ classes } = 4$ (0, 1, 2, 3)



\hat{y} is (4, 1) because it gives us 4 probabilities

Softmax layer:

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

Activation function:

$$t = e^{z^{[L]}} \rightarrow (4, 1)$$

$$\hat{y} = a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^4 t} \rightarrow (4, 1)$$

$$\text{i.e. } a_i^{[L]} = \frac{t_i}{\sum_{i=1}^4 t_i}$$

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

$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}, \quad \sum_{j=1}^4 t_j = 176.3$$

$$a^{[L]} = \frac{t}{176.3}$$

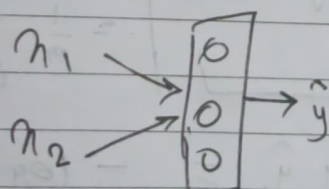
$$\therefore a_1^{[L]} = \frac{e^5}{176.3} = 0.842$$

$$a_2^{[L]} = \frac{e^2}{176.3} = 0.042$$

$$a_3^{[L]} = 0.002, \quad a_4^{[L]} = 0.114$$

softmax activation f^n : $a^{[L]} = g^{[L]}(z^{[L]})$
 $(3, 1)$ $(4, 1)$

softmax examples



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

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

$C=4 \Rightarrow 4$ decision boundaries across the plot when there are no hidden layers.

2) Training a softmax classifier

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \rightarrow t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} \rightarrow a^{[L]} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

Softmax regression generalizes logistic regression to c classes.

If $c = 2$, then softmax reduces to logistic reg.

Note: In contrast to softmax, there is hardmax which gives an output matrix generated from $z^{[L]}$. It outputs 1 for the highest element in $z^{[L]}$ and 0 for the rest of the elements.

Loss function:

eg: $y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $a^{[L]} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$, $c = 4$

$$L(\hat{y}, y) = - \sum_{j=1}^c y_j \log \hat{y}_j \quad \left| \quad J = \frac{1}{n} \sum_{i=1}^n J(w^{[1]}, b^{[1]}, \dots) \right.$$

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

single example:

$$-y_2 \log \hat{y}_2 = -\log \hat{y}_2$$

make \hat{y}_2 as big as possible.

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(n)} \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & \dots & 1 \\ 1 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \end{bmatrix} \leftarrow (1, 1)$$

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

$$= \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \leftarrow (L, m)$$

Gradient descent with softmax:

Back prop: $\delta z_{(L,1)}^{[L]} = \hat{y}^{(1)} - y^{(1)}$

$\nwarrow a^{(L)}$

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