

ADAM OPTIMIZATION ALGORITHM

* ADAM (Adaptation Moment Estimation) = momentum + RMS prop

* compute dW , db using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dW, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) db \quad \rightarrow \text{MOMENTUM}$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dW^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2 \quad \rightarrow \text{RMS PROP}$$

$$V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2^t), \quad S_{db} = S_{db} / (1 - \beta_2^t)$$

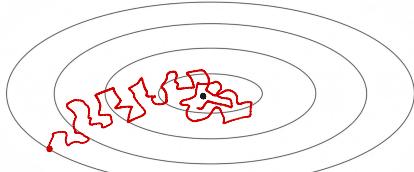
$$W = W - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon}, \quad b = b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon} \quad \left. \right\} \text{UPDATE}$$

* hyperparameters α (learning rate)

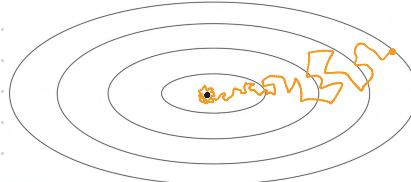


LEARNING RATE DECAY

* reduce learning rate over time



REGULAR α (doesn't converge)



SLOWLY REDUCED α
(you oscillate in a tighter region)

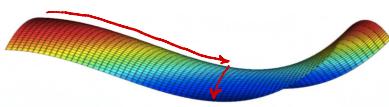
$$\alpha = \frac{1}{1 + \text{decay rate} * \text{epoch}} \alpha_0$$

LEARNING RATE DECAY

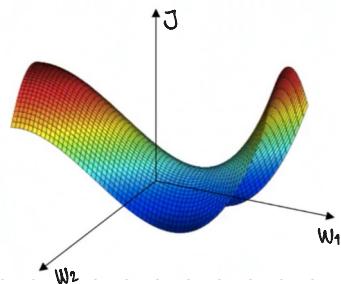
$$\alpha = 0.95^{\text{epoch}} \alpha_0$$

EXPONENTIAL DECAY

THE PROBLEM OF LOCAL OPTIMA



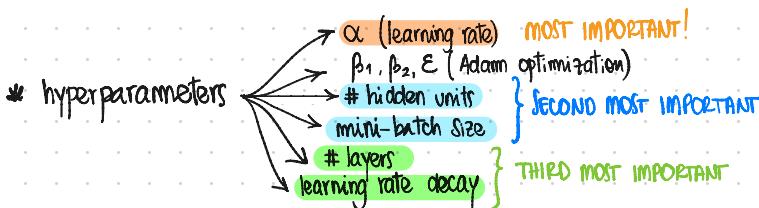
- * the derivative is close to 0 for a long time. They can make learning slow.



SADDLE POINT:

- * minimum along w_1
- * maximum along w_2

TUNING PROCESS

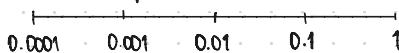


* try random values. Don't use a grid!!

→ coarse-to-fine search process

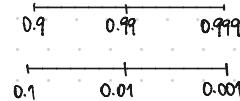
* appropriate scale for hyperparameters:

(#1) learning rate: Search for a value between 0.0001 and 1 within a log scale



(#2) exponentially weighted averages: $\beta = 0.9 \dots 0.999$
(log scale!)

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



(PANDA STRATEGY)

babysitting one model: focus on one model one day at a time

* tuning approach → train many models in parallel

(CAVIAR STRATEGY)

how much computational power you have?

A LOT! → Caviar

↓
NOT much
Panda

BATCH NORMALIZATION

* speeds up learning

* normalize → mean
↓ variance } layer outputs → to make the training of
variance w, b more efficient

before or after
passing through the
activation function?

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

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

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

$$\tilde{z}^{(i)} = \gamma z^{(i)}_{\text{norm}} + \beta \quad \begin{array}{l} \text{LEARNABLE} \\ \text{PARAMETERS} \end{array} \quad \Rightarrow \quad \left. \begin{array}{l} \text{if } \gamma = \sqrt{\sigma^2 + \epsilon} \\ \text{and } \beta = \mu \end{array} \right\} \quad \tilde{z}^{(i)} = z^{(i)}_{\text{norm}}$$

* mostly used with mini-batches

$$x^{[i]} \rightarrow w^{[i]} b^{[i]} \rightarrow z^{[i]} \rightarrow \beta^{[i]} \gamma^{[i]} \rightarrow \tilde{z}^{[i]}$$

* gradient descent + batch normalization:

for $t = 1 \dots$ (# mini-batches):

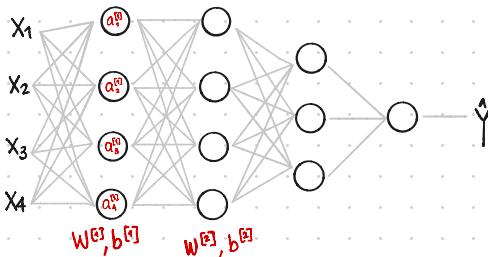
compute forward prop on $X^{[t]}$

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

use back prop to compute $dW^{[l]}$, $db^{[l]}$, $d\beta^{[l]}$, $d\gamma^{[l]}$

update parameters (works with momentum, RMSprop, Adam, etc.).

#1 It makes weights in deeper layers of the NN more robust to changes compared to those in the early layers.



BN LIMITS THE AMOUNT TO WHICH UPDATING THE PARAMETERS IN THE EARLIER LAYERS CAN AFFECT THE DISTRIBUTION OF THE VALUES IN THE UPCOMING LAYERS. IT REDUCES THE PROBLEM OF THE INPUT VALUES CHANGING. THESE VALUES BECOME MORE STABLE.

THE CHANGES IN $W^{[l]}, b^{[l]}$ WILL CHANGE $\alpha_1^{[l]}, \alpha_2^{[l]}, \alpha_3^{[l]},$ AND $\alpha_4^{[l]}$; THAT WILL ALSO CHANGE $W^{[l]}, b^{[l]}$, GENERATING A COVARIATE SHIFT PROBLEM.

→ BUT...

#2 By using a bigger mini-batch size, you reduce the noise together with the regularization effect. But don't use BN as a regularizer.

* at test time, you come up with a separate estimate of μ and σ^2 . You estimate these parameters using exponentially weighted averages across the mini-batches.

$$\begin{array}{c} X^{[1]}, X^{[2]}, X^{[3]} \\ \downarrow \mu^{[1] \text{ EWA}} \quad \downarrow \mu^{[2] \text{ EWA}} \quad \downarrow \mu^{[3] \text{ EWA}} \\ \left. \begin{array}{c} \mu^{[1:3] \text{ EWA}} \\ \sigma^2_{[1:3] \text{ EWA}} \end{array} \right\} \quad \left. \begin{array}{c} \mu, \sigma^2 \Rightarrow z_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}} \end{array} \right\} \end{array}$$

SOFTMAX REGRESSION

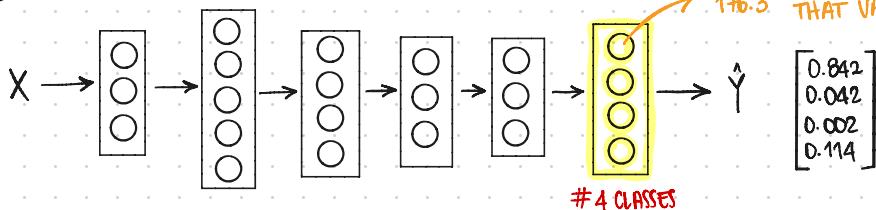
* multiclass classification

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

$$(4.1) \quad a^{[l]} = \frac{e^{z^{[l]}}}{\sum_{i=1}^k t_i} \quad \Rightarrow \quad a_i^{[l]} = \frac{t_i}{\sum_{i=1}^k t_i}$$

ACTIVATION
FUNCTION

example:

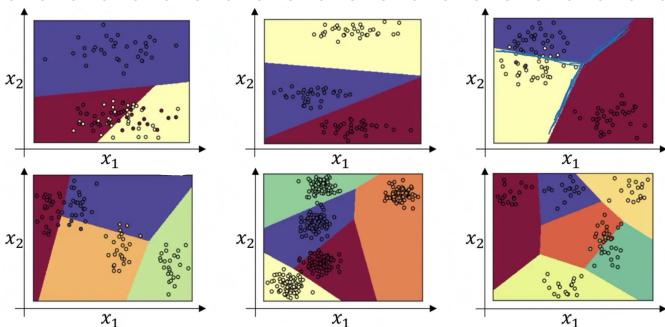


$$\frac{e^5}{176.3} = 0.842 \quad (\text{THE CHANCES OF THAT VALUE BEING } 0 \text{ IS } 84.2\%)$$

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

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

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



* there are examples with no hidden units. As you add layers the NN will be able to classify more complex (non-linear) points.

* Softmax regression generalizes logistic regression to C classes.

* If C=2, softmax reduces to logistic regression.

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

LOSS FUNCTION
(max likelihood estimation)

* $\mathcal{L}(\hat{y}, y)$ small $\rightarrow \hat{y}_j$ big

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

COST FUNCTION

* backpropagation:

$$\frac{\partial J}{\partial z} = dz^{(1)} = \hat{Y} - Y$$