

## Hyperparameter Tuning

### Tuning Process

Learning Rate  $\alpha$

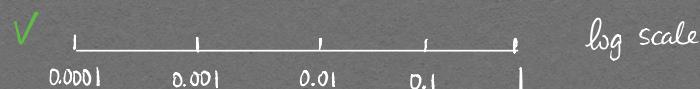
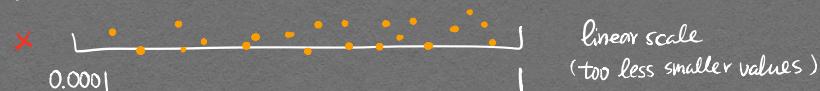
$\beta$ , # hidden units, mini-batch size

# layers, learning rate decay

$\beta_1, \beta_2, \epsilon$  usually just default value

Using an appropriate scale to pick hyperparameters

learning rate



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

$$\alpha = 10^r$$

Similarly for  $\beta = 0.9, \dots, 0.999$

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

$$10^{-1} \quad 10^{-3}$$

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

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

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

$$\frac{10^{-1}}{0.9} \quad \frac{10^{-2}}{0.99}$$

$$\frac{0.1}{0.99} \quad \frac{0.01}{0.999}$$

$$r \in [0, 1)$$

$$10^{-1}$$

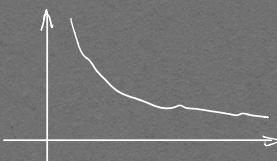
$$10^{-2} \dots$$

$$10^{-1.2}$$

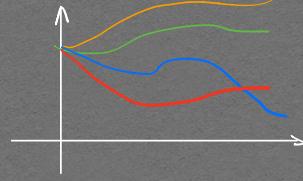
$$10^{-1.2}$$

### In Practice . Pandas vs Caviar

Babysitting one model  
(not enough computational resource)



Train Many Models in Parallel

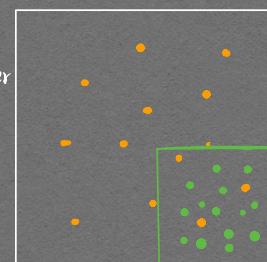


### Batch Normalization

Normalize Inputs to speed up learning

Can we normalize  $a^{[l-1]} / Z^{[l-1]}$  to train  $W^{[l]}, b^{[l]}$  faster?

### Hyperparameter 1



Coarse to Fine  
Densely Search  
in small region

Try randomly sampling  
instead of grid searching

Given some intermediate values in NN.  $\bar{z}^{(1)}, \bar{z}^{(2)}, \dots, \bar{z}^{(m)}$

$$\mu = \frac{1}{m} \sum_{i=1}^m \bar{z}^{(i)} \quad \text{If } \gamma = \sqrt{\sigma^2 + \epsilon}$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (\bar{z}^{(i)} - \mu)^2 \quad \beta = \mu$$

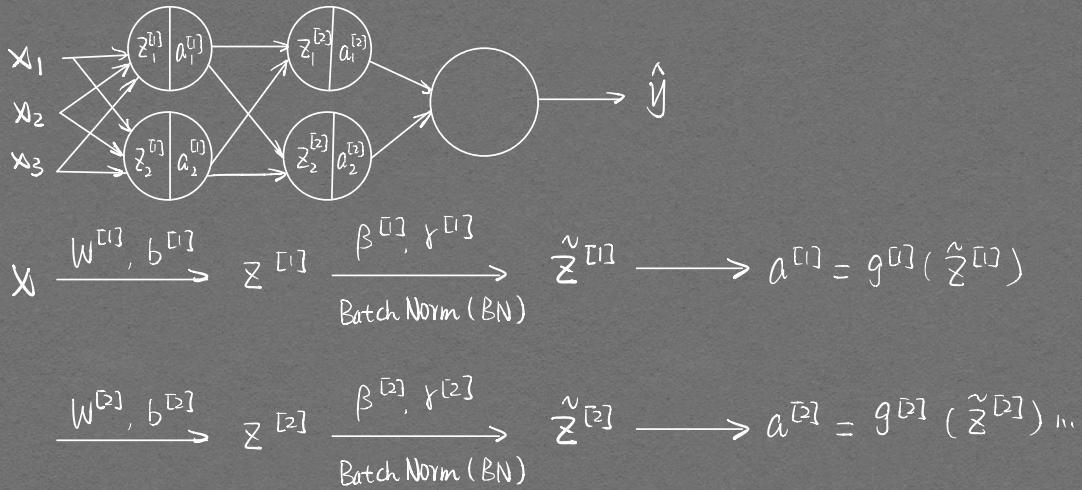
then

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

$$\tilde{z}^{(i)} = \gamma \bar{z}_{\text{norm}}^{(i)} + \beta \quad \text{learnable parameters}$$

Use  $\tilde{z}^{(i)}$  instead of  $\bar{z}_{\text{norm}}^{(i)}$  (mean 0, variance 1 is not always desired)

### Fitting Batch Norm to a Neural Network



Parameters:  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots, W^{[L]}, b^{[L]}$     $b^{[L]}$  is actually eliminated  
 $\beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots$  due to batch normalization

### Implementing Gradient descent

for  $t = 1, 2, \dots, \# \text{mini batches}$  {

Compute Forward Prop on  $x^{[t]}$

In each hidden layer,  $z^{[l]} \leftarrow \tilde{z}^{[l]}$

Use Back Prop to compute  $dW^{[l]}, db^{[l]}, d\beta^{[l]}, d\gamma^{[l]}$

Update Parameters

}

## Why does Batch Norm work?

② Weights in deeper layers are more robust to changes in earlier layers.

"Covariance Shift" If input distribution changes,  
even if the mapping does not change,  
we might still need to retrain our parameters.

Batch Norm keeps hidden units' mean / variance stable.

Earlier hidden units are limited to shift less.

## ③ Batch Norm as Regularization

Each mini-batch is scaled by the mean / variance computed on just that mini-batch.

This adds some noise to the values of  $Z^{[l]}$

Batch Norm could go with Dropout.

Larger minibatch size mitigates the regularization effect.

## Batch Norm At Test Time

$\mu, \sigma^2$ : estimate using exponentially weighted average across mini-batches

$$\begin{aligned} X^{[1]} &\rightarrow \boxed{\mu^{[1][l]}, \sigma^{[1][l]}} \\ X^{[2]} &\rightarrow \boxed{\mu^{[2][l]}, \sigma^{[2][l]}} \\ &\vdots \\ &\boxed{\mu_{\text{AVG}}, \sigma_{\text{AVG}}} \end{aligned}$$

## Multi-class Classification - Softmax Classification

$P_1$
$P_2$
$P_3$
$P_4$

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

Activation Function

$$t = e^{Z^{[l]}} \quad a^{[l]} = \frac{e^{Z^{[l]}}}{\sum_{i=1}^4 t_i} = \frac{t}{\sum_{i=1}^4 t_i}$$

## Train A Softmax Layer

$$y = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \hat{y} = \begin{bmatrix} .3 \\ .2 \\ .1 \\ .4 \end{bmatrix} \quad L(y, \hat{y}) = -\sum_{j=1}^4 y_j \log \hat{y}_j \\ = -\log \hat{y}_2$$

$$J(W, b) = \frac{1}{m} \sum_{i=1}^m L(y_i, \hat{y}_i)$$

make  $\hat{y}_2$  big

$$Y = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}] \\ = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & \dots \\ 0 & 0 & 0 \end{bmatrix} \quad \hat{Y} = \begin{bmatrix} .2 & .3 & .9 \\ .4 & .3 & .1 \\ .2 & .4 & \dots \\ .2 & 0 & 0 \end{bmatrix}$$

$$dZ^{[l]} = \hat{y} - y \text{ for Backprop}$$

## Programming Frameworks

### Deep Learning Frameworks

Pytorch

Tensorflow  
Keras

Ease of Programming

Truly Open Source with Good Governance

Motivating Problem

$$J(w) = w^2 - 10w + 25$$

w = tf.variable for parameters

$$\text{coefficients} = np.array()$$

x = tf.placeholder()

cost = ...

train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)

init = tf.global\_variables\_initializer()

sess = tf.Session

sess.run(init)

for i in range(1000):

sess.run(train, feed\_dict = {x: coefficients})

print(sess.run(w))