

Tuning Process.

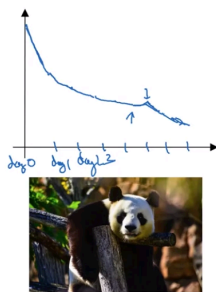
- Try Random values, don't use grid
- Coarse to fine

Pick Scale for hyper parameters

- Log scale
- Uniform scale

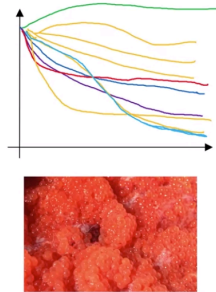
2 Approach

Babysitting one model



Panda

Training many models in parallel



Caviar

Andrew N

Batch Normalization

- More robust
- Easier for hyperparam search
- Much bigger range for hyperparam
- Train deep network

Basically normalize each layer for a multilayer neural network so that each layer train faster. Apply normalization to hidden layers

Implementing Batch Norm

Given some intermediate values in NN

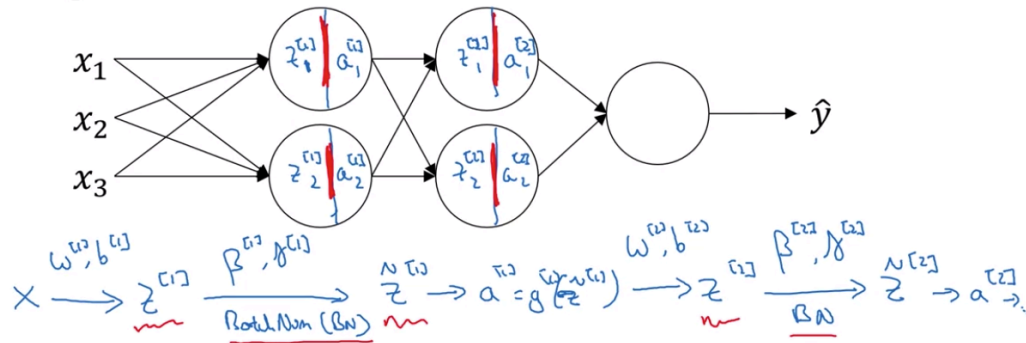
$$\mu = \frac{1}{m} \sum_i z^{(i)}$$
$$\sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$
$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$
$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

If $\gamma = \sqrt{\sigma^2 + \epsilon}$ and $\beta = \mu$ then $\tilde{z}^{(i)} = z^{(i)}$

learnable parameters of model.

Essentially computing the identity function. Use \tilde{z}_i instead of z_i

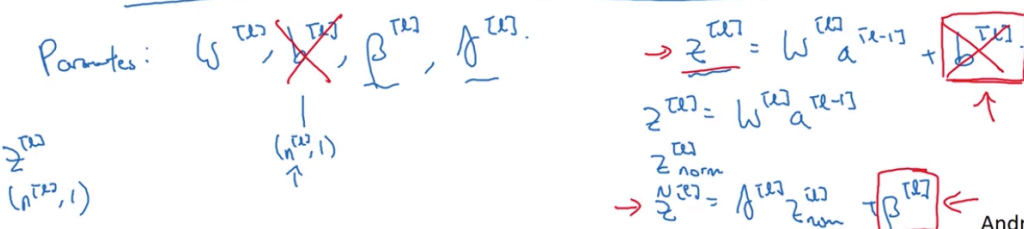
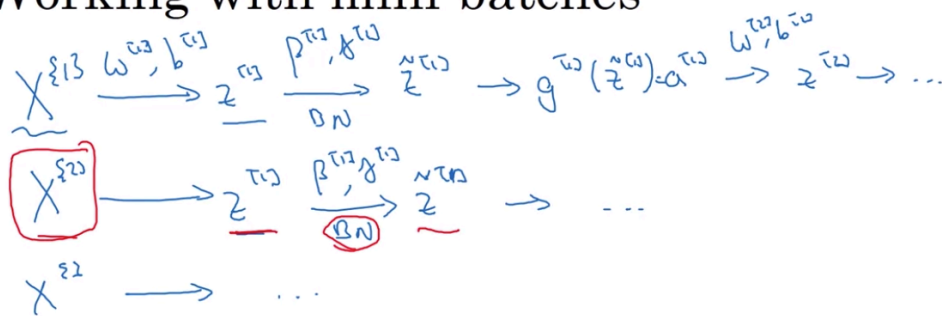
Adding Batch Norm to a network



Adding batch norm before activation.

`tf.nn.batch_normalization`

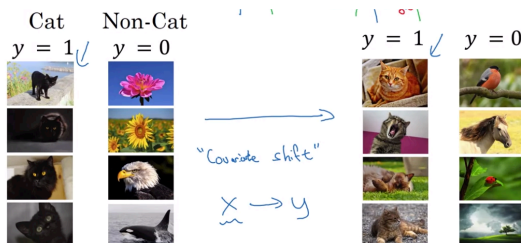
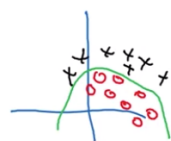
Working with mini-batches



Andrew Ng

This is no point keeping bias term, beta for BatchNorm controls the shift.

Covariance shift -> if you learn some mapping (like with the data on the left, it might not do well on the data on the right even though the ground truth remain the same)



Batch norm reduce the intermediate hidden value changes, reduce covariance shift's effect.
 Limit the amount the input value changes (earlier layer changes) affect the later layer changes. (More stable)
 Make the later layer learn easier.

Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch. X
- This adds some noise to the values $z^{[l]}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations. $\tilde{z}^{[l]}$, μ , σ^2
- This has a slight regularization effect.

Using a bigger mini batch size, reduce regularization effect. Don't rely on batch norm as a regularization as the effect is not as good as dropout.

Batch Norm at test time

$$\begin{aligned} \rightarrow \mu &= \frac{1}{m} \sum_i z^{(i)} \\ \rightarrow \sigma^2 &= \frac{1}{m} \sum_i (z^{(i)} - \mu)^2 \\ \rightarrow z_{\text{norm}}^{(i)} &= \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \\ \rightarrow \tilde{z}^{(i)} &= \gamma z_{\text{norm}}^{(i)} + \beta \end{aligned}$$

Epsilon added to variance for numerical stability

Mean and variance computed on the entire mini batch. But sometimes at test time, maybe only one instance is available.

Solution. Use Exponentially weighted moving average.

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

$x^{[1]}, x^{[2]}, x^{[3]}, \dots$

$\mu^{[1]}, \mu^{[2]}, \mu^{[3]}, \dots \rightarrow \mu$

$\sigma^{[1]}, \sigma^{[2]}, \sigma^{[3]}, \dots \rightarrow \sigma^2$

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

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

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

Andrew Ng

Softmax vs Hardmax

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

"hard max"

- If $C = 2$, softmax reduces to logistic regression
- Softmax is a generalization of logistic regression to more classes

Loss function

Handwritten notes for the loss function:

- Diagram showing a vector $y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ with a circled '1' and a note $y_1 = 1$. A note $y_1 = y_2 = y_3 = 0$ is also present.
- Diagram showing a vector $\hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$ with a circled '0.2' and a note $\hat{y}_2 = 0.2$.
- Equation for the loss function: $\mathcal{L}(\hat{y}, y) = -\sum_{j=1}^n y_j \log \hat{y}_j$. A note $-y_2 \log \hat{y}_2 = -\log \hat{y}_2$ is shown with an arrow pointing to the term.
- Equation for the total loss: $J(w^{(1)}, b^{(1)}, \dots) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$.
- Text: "Make \hat{y}_2 big."

Tensorflow

remember to initialize your variables, create a session and run the operations inside the session.

```
a = tf.constant(2)
b = tf.constant(10)
c = tf.multiply(a,b)
print(c)
```

```
> Tensor("Mul:0", shape=(), dtype=int32)
```

```
sess = tf.Session()
print(sess.run(c))
```

```
> 20
```

Feed data like this, used for feeding in training data

```
x = tf.placeholder(tf.int64, name = 'x')
print(sess.run(2 * x, feed_dict = {x: 3}))
sess.close()
```

```
> 6
```