batchnorm1

# Batch Normalization

## Normalizing activations in a network

deeplearning.ai

---

## Normalizing inputs to speed up learning

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

$X = X - \mu$ ← element-wise

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

$X = X / \sigma^2$

won't it be nice to normalize $a^{[2]}$

Can we normalize $a^{[2]}$ so as to train $w^{[3]}, b^{[3]}$ faster

Normalize $\boxed{\frac{z^{[2]}}{\uparrow}}$

from elongated to ··

flow about a deeper model

This is what batch norm (like

→ There's debate about whether you should normalize $a^{[2]}$ or $z^{[2]}$.  Andrew Ng prefers $\boxed{z^{[2]}}$

Andrew Ng

---

## Implementing Batch Norm

Given some intermediate values in NN   $z^{(i)}, \ldots, z^{(m)}$   $z^{[l](i)}$

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

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

$z^{(i)}_{norm} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$   ← in case $\sigma^2 = 0$

Gamma, beta

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

If $\gamma = \sqrt{\sigma^2 + \epsilon}$ ←

$\beta = \mu$ ←

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

learnable parameters of model.

$X \leftarrow$

$z^{(i)} \leftarrow$

for some hidden layer
[l] omitted in this slide.

You don't always want to want your hidden layer value forced to mean 0.

how to rescale?

e.g. if you have ... activation

$$\Rightarrow \tilde{z}^{(i)} = \gamma \, z_{norm}^{(i)} + \beta \qquad$$ learnable parameters of model.

Use $\tilde{z}^{[l](i)}$ instead of $z^{[l](i)}$

Andrew Ng

how to. rescale?

forced ...
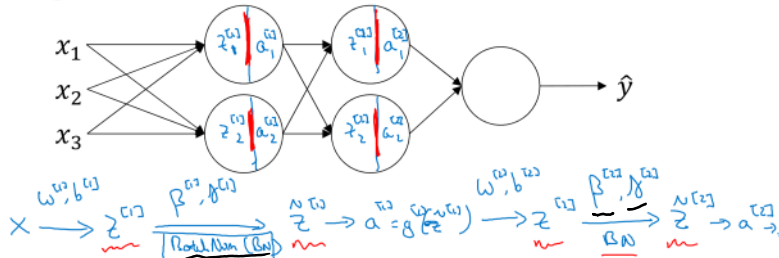
e.g. if you have sigmoid activation

adjust $\gamma \cdot \beta$

---

batchnorm2

# Batch Normalization

## Fitting Batch Norm into a neural network

deeplearning.ai

---

# Adding Batch Norm to a network

$$X \longrightarrow z^{[1]} \xrightarrow[\text{Batch Norm (BN)}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \to a^{[1]} = g^{[1]}(\tilde{z}^{[1]}) \longrightarrow z^{[2]} \xrightarrow[\text{BN}]{\beta^{[2]}, \gamma^{[2]}} \tilde{z}^{[2]} \to a^{[2]} \to \dots$$

$W^{[1]}, b^{[1]}$     $\beta^{[1]}, \gamma^{[1]}$     $W^{[2]}, b^{[2]}$     $\beta^{[2]}, \gamma^{[2]}$

(Parameters): $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots, W^{[l]}, b^{[l]}$

$\to \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[l]}, \gamma^{[l]}$

$\to \beta$
$\to \sigma \to$ distinguish from the hyperparam of momentum etc.

$d\beta^{[2]}$     $\beta^{[2]} = \beta^{[2]} - \alpha \, d\beta^{[2]} \longleftarrow$ these can be updated ..

tf.nn.batch-normalization $\Longleftarrow$     tensorflow

Andrew Ng

Andrew Ng Deep Learning. Page 2

# Working with mini-batches

$X^{\{1\}} \xrightarrow{W^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow[\beta N]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \to g^{[1]}(\tilde{Z}^{[1]}) = a^{[1]} \xrightarrow{W^{[2]}, b^{[2]}} Z^{[2]} \to \dots$

$X^{\{2\}} \to Z^{[1]} \xrightarrow[\boxed{BN}]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \to \dots$

Just
this
mini-batch

$X^{\{3\}} \to \dots$

Parmters: $W^{[l]}, \not{b^{[l]}}, \beta^{[l]}, \gamma^{[l]}.$

$\quad\quad\quad\quad (n^{[l]},1) \quad (n^{[l]},1) \quad (n^{[l]},1)$

$Z^{[l]}$
$(n^{[l]},1)$

# of layers
$\beta^{[3]}\gamma^{[2]}$ normalization
$a^{[3](1)} \; a^{[3](1)} \; a^{[3](3)}$

$\to z^{[l]} = W^{[l]} a^{[l-1]} + \not{b^{[l]}}$

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

$z^{[l]}_{norm}$

$\to \tilde{z}^{[l]} = \gamma^{[l]} z^{[l]}_{norm} + \boxed{\beta^{[l]}}$

important →

win batch normalization
$b^{[l]}$ not necessary
⇒ will be canceled out by normalization
or set $b^{[l]} = 0$
batch norm't zero out the means of $z^{[l]}$
⇒ So set $b^{[l]} = 0.$

Andrew Ng

# Implementing gradient descent

for $t = 1 \dots$ num MiniBatches
$\quad$ Compute forward prop on $X^{\{t\}}.$
$\quad\quad$ In each hidden layer, use BN to repa $z^{[l]}$ with $\tilde{z}^{[l]}.$
$\quad$ Use backprop to compute $dW^{[l]}, \not{db^{[l]}}, d\beta^{[l]}, d\gamma^{[l]}$
$\quad$ Update parms $W^{[l]} := W^{[l]} - \alpha \, dW^{[l]}$
$\quad\quad\quad\quad\quad\quad \beta^{[l]} := \beta^{[l]} - \alpha \, d\beta^{[l]}$
$\quad\quad\quad\quad\quad\quad \gamma^{[l]} := \dots$

Works w/ momentm, RMSprop, Adam.

BH = batch normalization

Andrew Ng

batchnorm3

# Batch Normalization

## Why does Batch Norm work?

deeplearning.ai

## Learning on shifting input distribution

$x_1$
$x_2$  → $\hat{y}$
$x_3$

Cat $y = 1$    Non-Cat $y = 0$

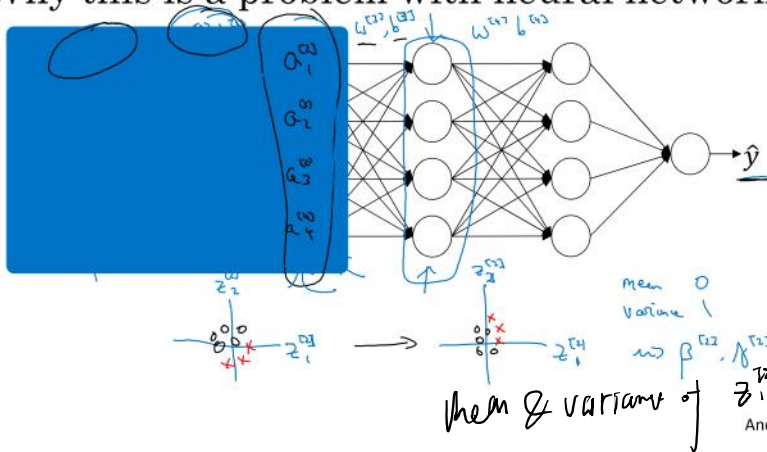$y = 1$    $y = 0$

"Covariate shift"

$x \longrightarrow y$

make weights robust to changes of
change of weights on previous layers

model of cat detection
trained on black cat picture

✳⟶ test on colored cats.
this is an example of a changed
data distribution

## Why this is a problem with neural networks?



$a_1^{[2]}$
$a_2^{[2]}$
$a_3^{[2]}$
$a_4^{[2]}$

$z_2^{[2]}$
$z_1^{[2]}$

Mean 0
Variance 1

$\beta^{[2]}, \gamma^{[2]}$

mean & variance of $z_1^{[2]}, z_2^{[2]}$

Andrew Ng

Batch norm reduces the amount that the hidden-layer values change

Batch Norm 的意义

叫以慢一点儿/�可以

won't change.

Stablise reduce the effect of input value on later layers

$\Rightarrow$ speed up learning

$\Delta$ So does it help reusing pre-trained model

## 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 $z^{[l]}$ within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.
- This has a slight regularization effect.

$\tilde{z}^{[l]}$   64, 128   $z^{[2]}$

$\mu, \sigma^2$

数据增强
data augmentation
了?

Mini-batch: 64 $\longrightarrow$ 512

Andrew Ng

② regularization.

Add noise

Next: what to do at test time?

batchnorm4

**8.** Which of the following statements about $\gamma$ and $\beta$ in Batch Norm are true?

- ☑ They can be learned using Adam, Gradient descent with momentum, or RMSprop, not just with gradient descent.

  Correct

- ☐ $\beta$ and $\gamma$ are hyperparameters of the algorithm, which we tune via random sampling.

  Un-selected is correct

- ☑ The optimal values are $\gamma = \sqrt{\sigma^2 + \varepsilon}$, and $\beta = \mu$.

  This should not be selected

- ☑ They set the mean and variance of the linear variable $z^{[l]}_i$ of a given layer.

  Correct

- ☐ There is one global value of $\gamma \in \mathbb{R}$ and one global value of $\beta \in \mathbb{R}$ for each layer, and applies to all the hidden units in that layer.

  Un-selected is correct

---

) computed at mini batches

# Batch Norm at test time

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

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

$$z^{(i)}_{norm} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \varepsilon}}$$ *estimated*

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

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

$X^{\{1\}}, X^{\{2\}}, X^{\{3\}}, \ldots$

$\mu^{\{1\}[l]}$  $\mu^{\{2\}[l]}$  $\mu^{\{3\}[l]}$ → $\mu$

$\theta_1$  $\theta_2$  $\theta_3$   $\sigma^2$

$\sigma^{2\{1\}[l]}$  $\sigma^{2\{2\}[l]}$

$z_{norm} = \frac{z - \mu}{\sqrt{\sigma^2 + \varepsilon}}$   $\tilde{z} = \gamma z_{norm} + \beta$

$\beta - \gamma, \mu, \sigma^2$

estimate $\mu, \sigma^2$

Use running average to estimate $\mu, \sigma^2$

keep a running average

use the exponentiated average

Andrew Ng