

Week 3b - Batch Normalization

Friday, August 14, 2020 11:30 PM

① Normalizing activations in a network

- In earlier weeks, we normalized the inputs ($x = \frac{x - \mu}{\sigma}$)
- In deeper neural networks, normalizing hidden layer outputs also help to train model better
- Implementing Batch Norm

→ Given some intermediate values in NN $\underbrace{z^{(1)} \dots z^{(m)}}_{z^{(l)}(i)}$

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

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

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

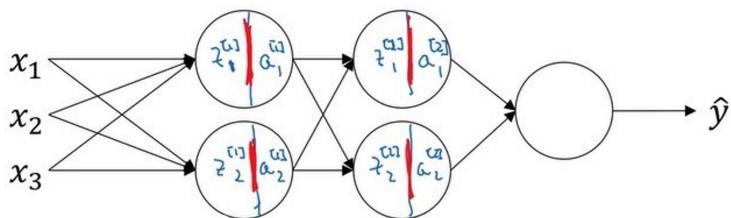
$$\rightarrow z^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta \quad \begin{matrix} \gamma, \beta \text{ are learnable} \\ \text{parameters of model} \\ \text{that can change} \\ \text{mean & variance} \\ \text{to whatever} \\ \text{we want} \end{matrix}$$

$$(z^{(i)} = z^{(i)} \text{ if } \gamma = \sqrt{\sigma^2 + \epsilon})$$

$\beta = \mu$

→ use $\hat{z}^{(l)(i)}$ instead of $z^{(l)(i)}$ in model

② Fitting Batch Norm into a Neural Network



- Batch Norm occurs between calc. of $z^{[l]}$ & $a^{[l]}$
 $\Rightarrow X \xrightarrow{w^{[l]}, b^{[l]}} z^{[l]} \xrightarrow{\text{BN}} \hat{z}^{[l]} \rightarrow a^{[l]} = g(\hat{z}^{[l]}) \xrightarrow{w^{[l]}, b^{[l]}} z^{[l+1]} \xrightarrow{\text{BN}} \hat{z}^{[l+1]} \rightarrow a^{[l+1]} \dots$

- Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}$ {compute
 $\beta^{[l]}, \gamma^{[l]}, \beta^{[2]}, \beta^{[1]}, \dots, \beta^{[L]}, \gamma^{[L]}$ } $\frac{db}{d\beta}, \frac{dw}{d\beta}, \frac{d\gamma}{d\beta}$
 also update $\beta^{[l]} = \beta^{[l]} - \alpha d\beta^{[l]}$
 $\gamma^{[l]} = \gamma^{[l]} - \alpha d\gamma^{[l]}$

- Working with mini-batches

- Working with mini-batches

$$X^{\{1\}} \xrightarrow{w^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \rightarrow g^{[1]}(\tilde{Z}^{[1]}) = \alpha^{[1]} \xrightarrow{w^{[2]}, b^{[2]}} Z^{[2]} \rightarrow \dots$$

$$X^{\{2\}} \xrightarrow{w^{[2]}, b^{[2]}} Z^{[1]} \xrightarrow{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \rightarrow g^{[1]}(\tilde{Z}^{[1]}) = \alpha^{[2]} \xrightarrow{w^{[2]}, b^{[2]}} Z^{[2]} \rightarrow \dots$$

$$\vdots \quad \rightarrow \dots$$

Parameters: $w^{[0]}, \cancel{b^{[0]}}, \beta^{[l]}, \gamma^{[l]}$
 $Z^{[l]}_{(n^{[l]}, 1)}$

$$Z^{[l]} = w^{[l]} \alpha^{[l-1]} + \cancel{b^{[l]}}$$

$$Z^{[l]} = w^{[l]} \alpha^{[l-1]}$$

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

$$\tilde{Z}^{[l]} = \gamma^{[l]} Z^{[l]}_{norm} + \beta^{[l]}$$

$\rightarrow b$ isn't required
since β over-
rides b
Hence β is
enough

- Implementing gradient descent

- for $t=1 \dots \text{num_mini-batches}$
 - Compute forward prop on $X^{\{t\}}$
In each hidden layer, use BN to replace $Z^{[l]}$ with $Z^{[l]}_{norm}$
 - Use backprop to compute $d\alpha^{[l]}, d\cancel{b^{[l]}}, d\beta^{[l]}, d\gamma^{[l]}$
 - Update parameters
 $w^{[l]} = w^{[l]} - \alpha d\alpha^{[l]}$
 $\beta^{[l]} = \dots$
 $\gamma^{[l]} = \dots$

→ works w/ momentum, RMSprop, Adam

③ Why does Batch Norm work

→ serves same purpose as normalizing inputs

→ Batch normalization reduces the problem of input values changing (shifting)

→ Has regularization effect (indirectly)

- each mini-batch is scaled by the mean/variance of that batch
- This adds noise to $Z^{[l]}$ in that mini-batch.
- Effect is similar to dropout (adding noise)
- Effect is reduced by increasing batch size
- Not intended to be used for regularization but only to speed up learning

④ Batch Norm at Test time

④ Batch Norm at Test time

Batch Norm at test time

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

$\underline{\mu}, \underline{\sigma^2}$: estimate very exponentially
 weighted average (across mini-batches).
 $x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots$
 \downarrow
 $\mu_1^{(1)}, \mu_2^{(1)}, \mu_3^{(1)}, \dots$
 $\theta_1, \theta_2, \theta_3, \dots$
 $\epsilon_1^{(1)}, \epsilon_2^{(1)}, \epsilon_3^{(1)}, \dots$
 $\hat{z}^{(1)} = \gamma z_{\text{norm}}^{(1)} + \beta$

μ
 σ^2
 Andrew Ng

- During training, mean & variance is computed for each mini-batch
- In testing, we might have to pro