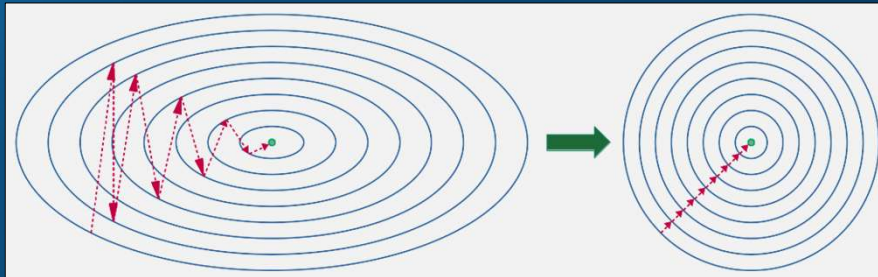# Batch Normalization

Deep Neural Network
Session 13
Pramod Sharma
pramod.sharma@prasami.com

---

# Batch Normalization

2

❑ It definitely helps to normalize input data
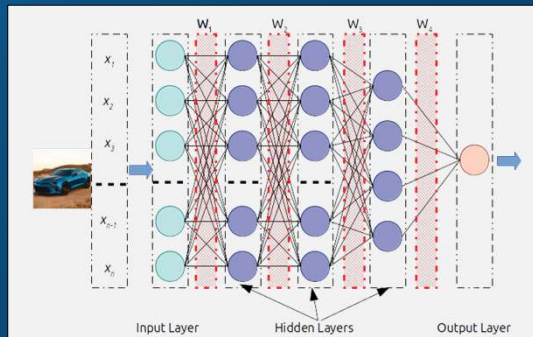
❑ Gradient converges faster

pra-sami

# Batch Normalization

3

□ What about hidden layer?

□ After all activations from previous layer are inputs for current layer…



□ Will it help if we normalize the hidden layers too?

pra-sami

---

# Batch Normalization

4

□ Batch normalization (also known as batch norm) [by Sergey Loffe and Christian Szegedy in 2015]
  ❖ Make artificial neural networks faster
  ❖ More stable through normalization of the input layer by re-centering and re-scaling
  ❖ Wider choices of hyper- parameter…

□ In theory, its normalizing activation values of the respective layers

□ In practice, it works better if we normalize 'z'
  ❖ Look at the documentation for details

pra-sami

## Batch Normalization

5

- In General, any $Z^i$ can be normalized

$$\text{mean } \mu = \frac{\sum Z^i}{m}$$

$$\text{std } \sigma^2 = \frac{1}{m} \Sigma (Z^i - \mu)^2$$

- $Z^i Norm = \frac{Z^i - \mu}{\sqrt{\sigma^2}}$
- $\hat{z} = \gamma . Z^i \text{ Norm} + \beta$
- where $\gamma$ and $\beta$ are paramters we can tune
- if $\gamma = \frac{1}{\sqrt{(\sigma^2)}}$ and $\beta = \frac{\mu}{\sqrt{\sigma^2}}$ ; $Z^i$ Norm $= \hat{z}$

pra-sami

---

## Batch Normalization

6

- In General, any $Z^i$ can be normalized

$$\text{mean } \mu = \frac{\sum Z^i}{m}$$

$$\text{std } \sigma^2 = \frac{1}{m} \Sigma (Z^i - \mu)^2$$

$$z^i_{Norm} = \frac{z^i - \mu}{\sqrt{\sigma^2}}$$

$$\hat{z} = \gamma . z^i_{Norm} + \beta$$

- where $\gamma$ and $\beta$ are parameters, we can **Train**
- if $\gamma = \frac{1}{\sqrt{\sigma^2}}$ and $\beta = \frac{\mu}{\sqrt{\sigma^2}}$ ; $Z^i_{Norm} = \hat{z}$

Instead of using $z^i_{Norm,}$ researchers realized that its better to derive $\hat{z}$ with two trainable parameters.

Intuition is that by normalizing z, we are introducing bias in the system. Hence it makes sense to train these parameters

pra-sami

## Batch Normalization

7

□ In General, any $Z^i$ can be normalized

□       mean μ = $\frac{\sum Z^i}{m}$

□       std $\sigma^2 = \frac{1}{m}\Sigma(Z^i - \mu)^2$

     ❖     $Z^i$ Norm $= \frac{Z^i - \mu}{\sqrt{\sigma^2 + \varepsilon}}$

     ❖     $\hat{z} = \gamma \cdot Z^i$ Norm $+ \beta$

□ where $\gamma$ and $\beta$ are parameters, we can train

□ if $\gamma = \frac{1}{\sqrt{\sigma^2 + \varepsilon}}$ and $\beta = \frac{\mu}{\sqrt{\sigma^2 + \varepsilon}}$ ; $Z^i$ Norm $= \hat{z}$

> Lets add a small $\varepsilon$ to prevent zero divide error…

pra-sami

---

## Batch Normalization

8

□ In General, any $Z^i$ can be normalized

□       mean μ = $\frac{\sum Z^i}{m}$

□       std $\sigma^2 = \frac{1}{m}\Sigma(Z^i - \mu)^2$

     ❖   Z^i Norm $= \frac{z^i - \mu}{\sqrt{\sigma2 + \varepsilon}}$

     ❖   $\hat{z} = \gamma \cdot zi$ Norm $+ \beta$

□ where $\gamma$ and $\beta$ are parameters, we can train

□ if $\gamma = \frac{1}{\sqrt{\sigma2 + \varepsilon}}$ and $\beta = \frac{\mu}{\sqrt{\sigma2 + \varepsilon}}$ ; zi Norm $= \hat{z}i$
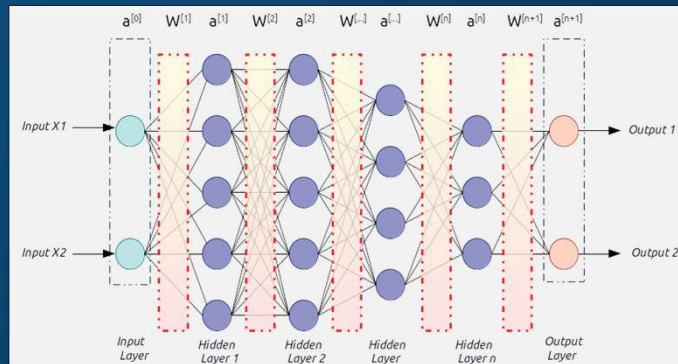
pra-sami

# Batch Normalization

9

- ❑ Notes:
  - ❖ Batch norm is used along with mini batches
  - ❖ Batch norm is applied to the batch under consideration only irrespective of other mini batches
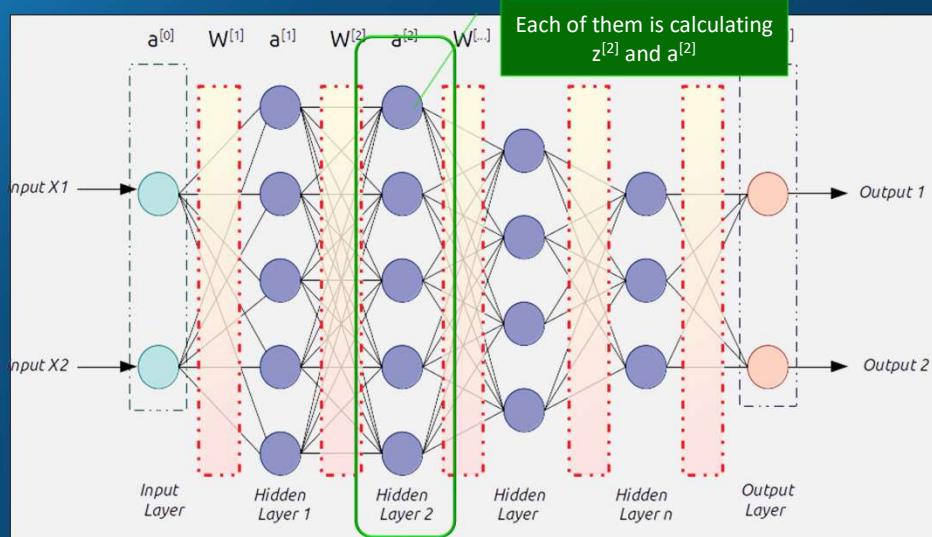
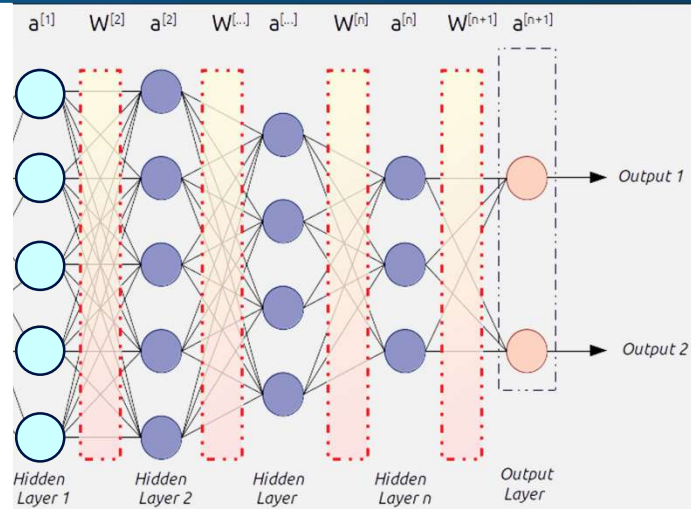- ❑ Where does it fit in overall scheme?



# Batch Normalization

10

Each of them is calculating $z^{[2]}$ and $a^{[2]}$

## Batch Normalization



For $a^{[2]}$ all $a^{[1]}$ are acting as input features



Co-variance Shift

# Batch Normalization

13

❑ Forward and back propagation with batch norm:

$W_1, b_1$

X

$Z_1$

$Z_1 = X.W_1 + b_1$

Our standard equation to calculate $z_1$.

pra-sami

---

# Batch Normalization

14

❑ Forward and back propagation with batch norm:

$W_1, b_1$

X

$Z_1$

$\beta_1, \gamma_1$

$\hat{z}_1$

$Z_1 = X.W_1 + b_1$

Calculate $\hat{z}_1$, based on $\beta_1, \gamma_1$

pra-sami

## Batch Normalization

15

❑ Forward and back propagation with batch norm:

$$X \xrightarrow{W_1, b_1} Z_1 \xrightarrow{\beta_1, \gamma_1} \hat{z}_1 \rightarrow a_1 = g_1(\hat{z}_1)$$

$$Z_1 = X . W_1 + b_1$$

Apply activation function $g_1(\hat{z}_1)$

pra-sami

## Batch Normalization

16

❑ Forward and back propagation with batch norm:

$$X \xrightarrow{W_1, b_1} Z_1 \xrightarrow{\beta_1, \gamma_1} \hat{z}_1 \rightarrow a_1 = g_1(\hat{z}_1) \xrightarrow{W_2, b_2} z_2$$

$$Z_1 = X . W_1 + b_1 \qquad z_2 = a_1 . W_2 + b_2$$

Calculate $z_2$ as usual

pra-sami

## Batch Normalization

17

❑ Forward and back propagation with batch norm:

$$X \xrightarrow{W_1, b_1} Z_1 \xrightarrow{\beta_1, \gamma_1} \hat{z}_1 \rightarrow a_1 = g_1(\hat{z}_1) \xrightarrow{W_2, b_2} z_2 \xrightarrow{\beta_2, \gamma_2} \hat{z}_2$$

$$Z_1 = X.W_1 + b_1$$

$$z_2 = a_1.W_2 + b_2$$

We know how to calculate $\hat{z}_2$

11/25/2024

pra-sami

---

## Batch Normalization

18

❑ Forward and back propagation with batch norm:

$$X \xrightarrow{W_1, b_1} Z_1 \xrightarrow{\beta_1, \gamma_1} \hat{z}_1 \rightarrow a_1 = g_1(\hat{z}_1) \xrightarrow{W_2, b_2} z_2 \xrightarrow{\beta_2, \gamma_2} \hat{z}_2 \rightarrow a_2 = g_2(\hat{z}_2)$$

$$Z_1 = X.W_1 + b_1$$

$$z_2 = a_1.W_2 + b_2$$

We also know how to calculate $a_2$

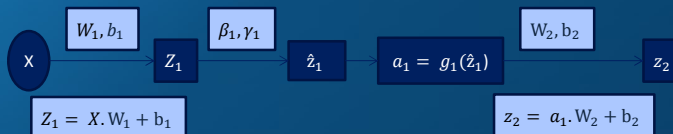11/25/2024

pra-sami

## Batch Normalization

19

❑ Forward and back propagation with batch norm:

$W_1, b_1$ → $Z_1$ → $\beta_1, \gamma_1$ → $\hat{z}_1$ → $a_1 = g_1(\hat{z}_1)$ → $W_2, b_2$ → $z_2$ → $\beta_2, \gamma_2$ → $\hat{z}_2$ → $a_2 = g_2(\hat{z}_2)$

X

$Z_1 = X.W_1 + b_1$

$z_2 = a_1.W_2 + b_2$

$\beta_1 = \beta_1 - \alpha . \partial \beta_1$
$\gamma_1 = \gamma_1 - \alpha . \partial \gamma_1$

Using the gradient descent, update $\beta$'s, $\gamma$'s along with W's and b's

$\beta_2 = \beta_2 - \alpha . \partial \beta_2$
$\gamma_2 = \gamma_2 - \alpha . \partial \gamma_2$

pra-sami

---

## Batch Normalization

20

❑ Forward and back propagation with batch norm:

$W_1, \cancel{b_1}$ → $Z_1$ → $\beta_1, \gamma_1$ → $\hat{z}_1$ → $a_1 = g_1(\hat{z}_1)$ → $W_2, \cancel{b_2}$ → $z_2$ → $\beta_2, \gamma_2$ → $\hat{z}_2$ → $a_2 = g_2(\hat{z}_2)$

X

$Z_1 = X.W_1 \cancel{+ b_1}$

$z_2 = a_1.W_2 \cancel{+ b_2}$

$\beta_1 = \beta_1 - \alpha . \partial \beta_1$
$\gamma_1 = \gamma_1 - \alpha . \partial \gamma_1$

$\beta_2 = \beta_2 - \alpha . \partial \beta_2$
$\gamma_2 = \gamma_2 - \alpha . \partial \gamma_2$

One more thing, since we are normalizing our Z's, keeping b's in the equation does not make any sense now.
Being the constant it will get eliminated!!

pra-sami

## Batch Normalization

21

❑ Forward and back propagation with batch norm:

$W_1, b_1$ | $\beta_1, \gamma_1$ | $W_2, b_2$ | $\beta_2, \gamma_2$

X → $Z_1$ → $\hat{z}_1$ → $a_1 = g_1(\hat{z}_1)$ → $z_2$ → $\hat{z}_2$ → $a_2 = g_2(\hat{z}_2)$

$Z_1 = X.W_1 + b_1$

$z_2 = a_1.W_2 + b_2$

$$\beta_1 = \beta_1 - \alpha . \partial \beta_1$$
$$\gamma_1 = \gamma_1 - \alpha . \partial \gamma_1$$

$$\beta_2 = \beta_2 - \alpha . \partial \beta_2$$
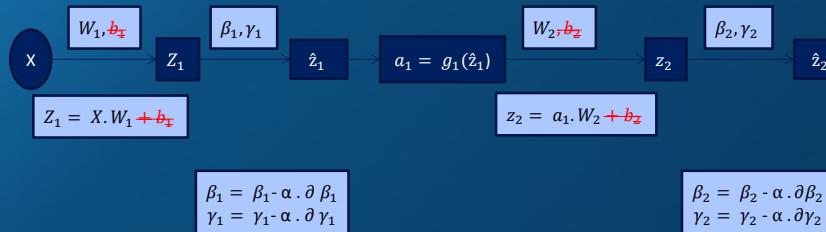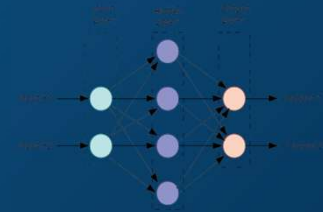$$\gamma_2 = \gamma_2 - \alpha . \partial \gamma_2$$

And at test/validation time using a exponentially weighted average!
So while training do not forget to save exponentially weighted values or simply running average!!

11/25/2024

pra-sami

---

## Batch Normalization

22

❑ Forward and back propagation with batch norm:

# Too many calculation steps...

$W_1, b_1$ | $\beta_1, \gamma_1$ | $W_2, b_2$ | $\beta_2, \gamma_2$

X → $Z_1$ → $\hat{z}_1$ → $a_1 = g_1(\hat{z}_1)$ → $z_2$ → $\hat{z}_2$ → $a_2 = g_2(\hat{z}_2)$

$Z_1 = X.W_1 + b_1$

# Don't worry...
# Most frameworks have one line code to do it.

$\gamma_1 = \gamma_1 - \alpha . \partial \gamma_1$

$\gamma_2 = \gamma_2 - \alpha . \partial \gamma_2$

11/25/2024

pra-sami

## Batch Normalization – Code Sample

23

```
model = tf.keras.models.Sequential(
    [
        tf.keras.layers.RNN( keras.layers.LSTMCell(units), input_shape=(None, input_dim) ),
        tf.keras.layers.BatchNormalization(),
        tf.keras.layers.Dense(output_size),
    ]
)


class Net(nn.Module):
        def __init__(self):
        super(Net, self).__init__()
        self.dense1 = nn.Linear(in_features=320, out_features=50)
        self.dense1_bn = nn.BatchNorm1d(50)
        self.dense2 = nn.Linear(50, 10)
```

❑ And it is applied to mini batches only….

❑ Batch Norm can be updated using any of the optimization functions…

11/25/2024

pra-sami

---

## Batch Normalization

24

Remember $\beta, \gamma$ are parameters you train!

11/25/2024

pra-sami

## Reflect…

25

- What is the primary purpose of Batch Normalization in deep learning?
  - A) To prevent overfitting
  - B) To reduce the number of parameters in the model
  - C) To accelerate training and reduce internal covariate shift
  - D) To increase the depth of the neural network

- Answer: C) To accelerate training and reduce internal covariate shift

- At which stage is Batch Normalization applied in a neural network?
  - A) After the input layer
  - B) After the activation function
  - C) Before the loss calculation
  - D) Before or after the activation function, depending on the implementation

- Answer: D) Before or after the activation function, depending on the implementation

- Which of the following is a key step in Batch Normalization?
  - A) Normalizing the gradient updates
  - B) Normalizing the activations by subtracting the batch mean and dividing by the batch standard deviation
  - C) Initializing weights to zero
  - D) Adding noise to the input data

- Answer: B) Normalizing the activations by subtracting the batch mean and dividing by the batch standard deviation

- What are the two learnable parameters introduced in Batch Normalization?
  - A) Gamma and Beta
  - B) Alpha and Beta
  - C) Theta and Gamma
  - D) Sigma and Mu

- Answer: A) Gamma and Beta

11/25/2024

pra-sami

---

26

Next Session… Recurrent Neural Networks

11/25/2024

pra-sami