



Batch Normalization

Deep Neural Network

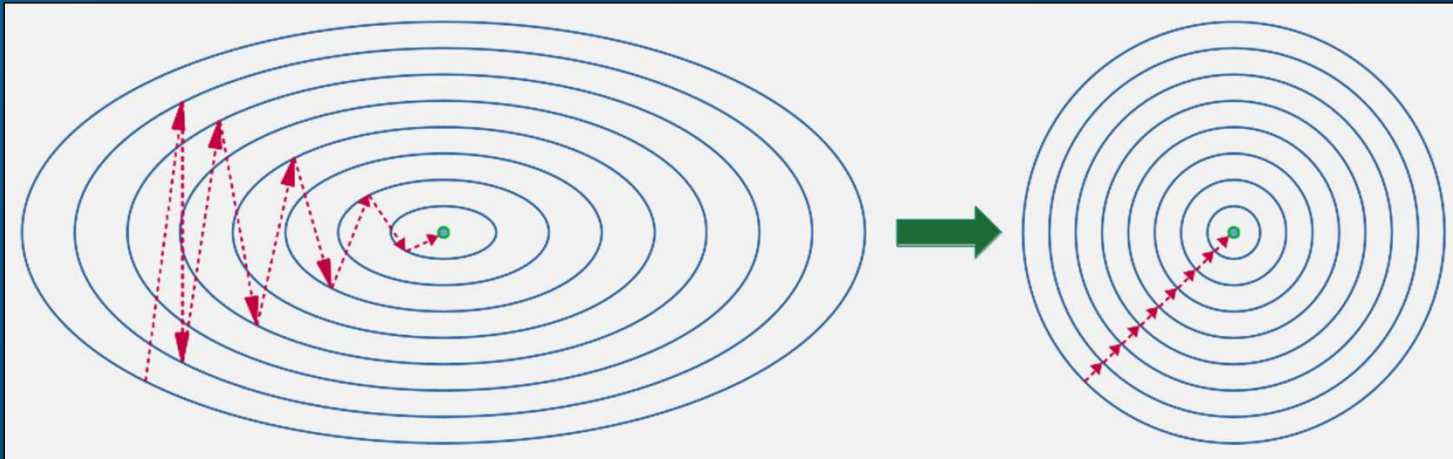
Session 13

Pramod Sharma

pramod.sharma@prasami.com

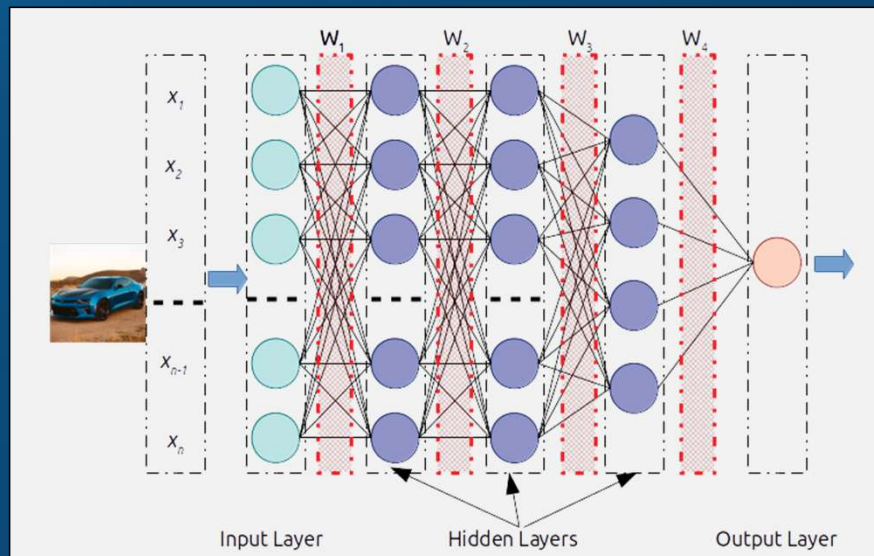
Batch Normalization

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



Batch Normalization

- ❑ 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?

Batch Normalization

- ❑ Batch normalization (also known as batch norm) [by Sergey Ioffe 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...
 - ❖ Higher Learning Rates: Can use much larger learning rates without training diverging
 - ❖ Less Sensitive to Weight Initialization: Networks are more robust to poor initializations
 - ❖ Reduced Need for Regularization: The inherent noise in batch statistics acts as a regularizer
 - ❖ Faster Convergence: All of the above leads to needing fewer epochs to get good results
- ❑ 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

Batch Normalization

- In General, any Z^i can be normalized

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

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

- $Z^i \text{ Norm} = \frac{Z^i - \mu}{\sqrt{\sigma^2}}$

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

- where γ and β are *parameters we can train*

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

Batch Normalization

- In General, any Z^i can be normalized

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

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

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

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

- where γ and β are parameters, we can **Train**

- if $\gamma = \frac{1}{\sqrt{\sigma^2}}$ and $\beta = \frac{\mu}{\sqrt{\sigma^2}}$; $z^i_{\text{Norm}} = \hat{z}$

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

Batch Normalization

□ In General, any Z^i can be normalized

□ mean $\mu = \frac{\sum Z^i}{m}$

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

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

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

□ where γ and β are parameters, we can train

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

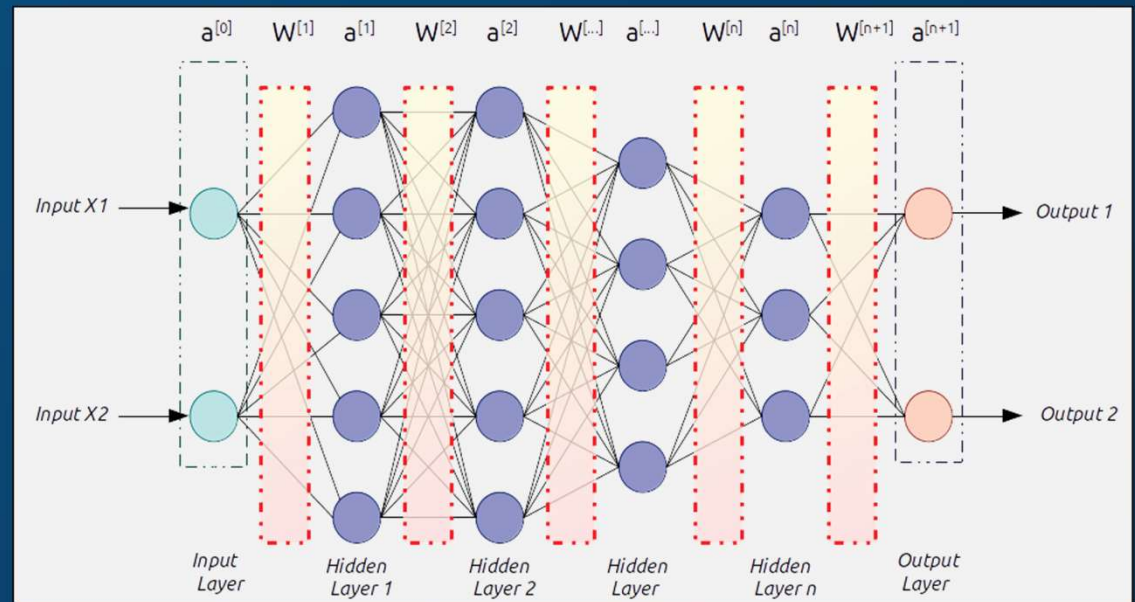
Lets add a small ε to prevent zero divide error...

Batch Normalization

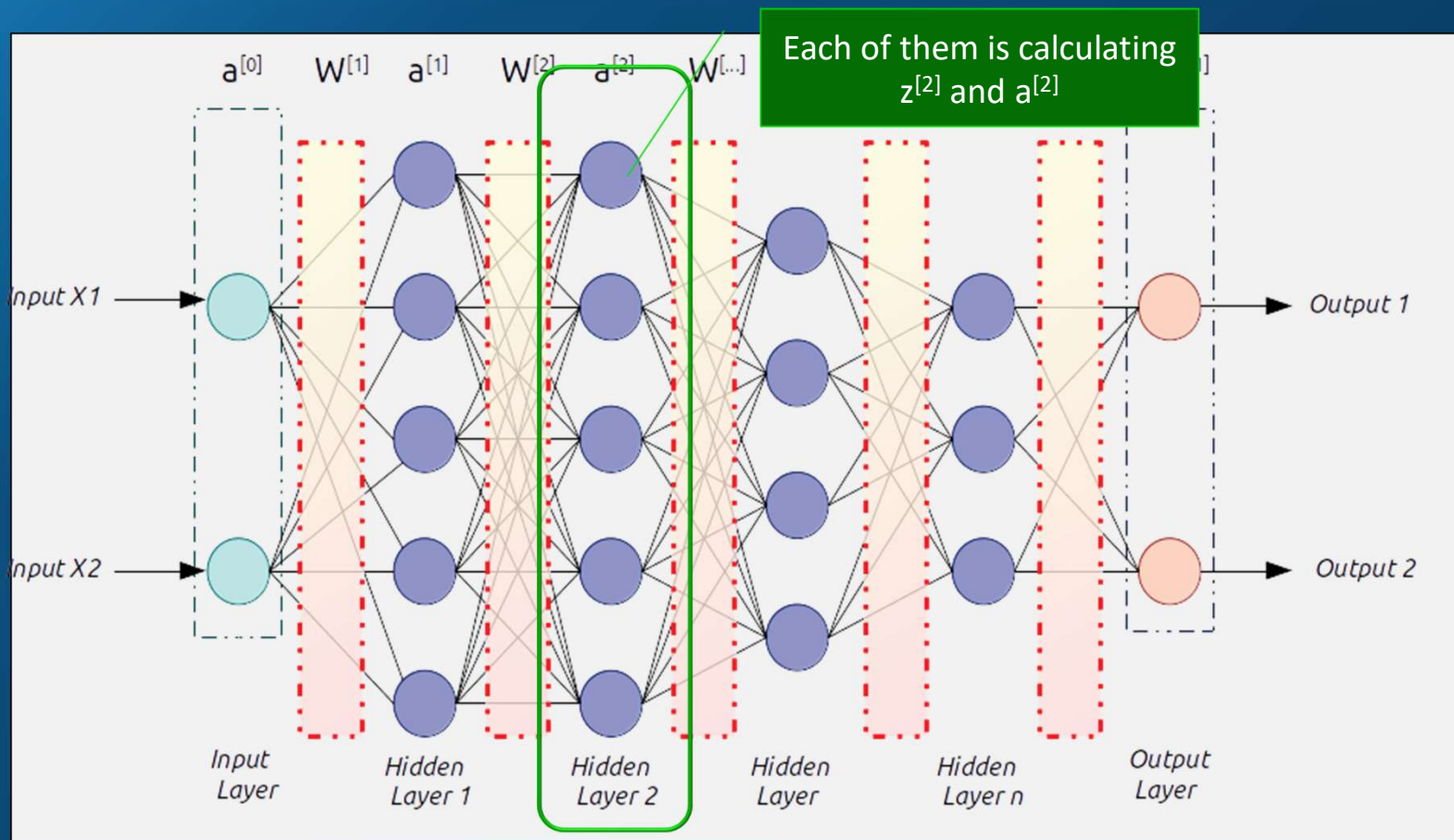
□ Notes:

- ❖ 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?

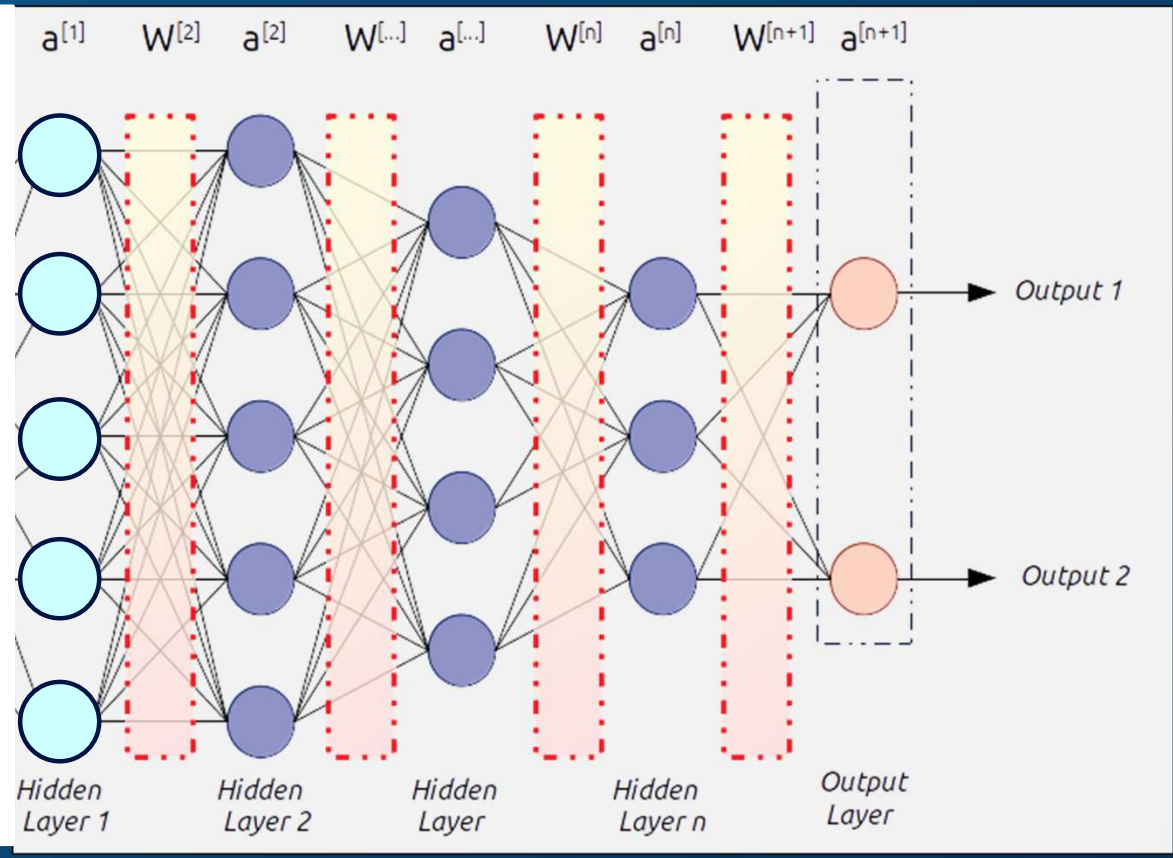


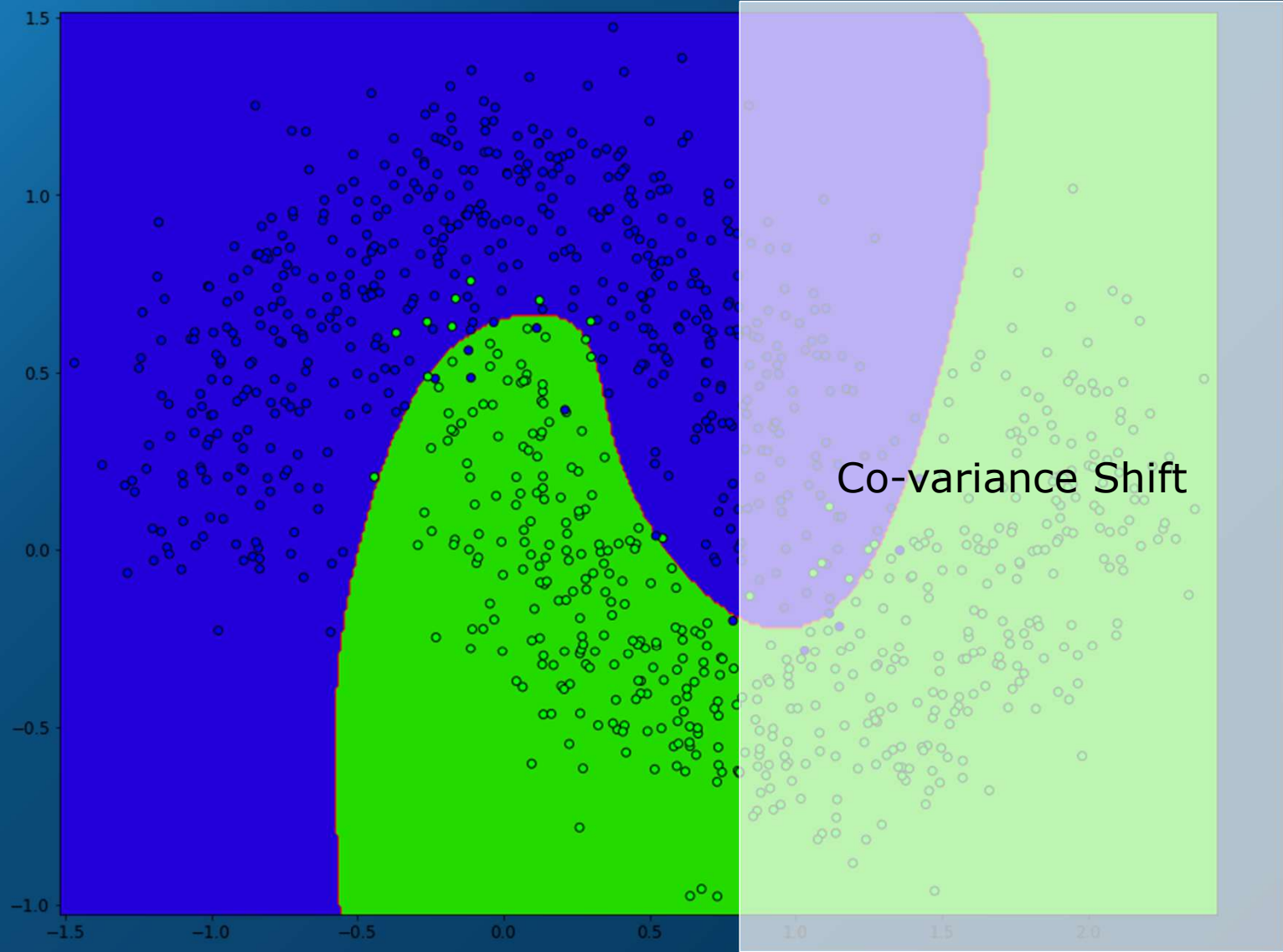
Batch Normalization



Batch Normalization

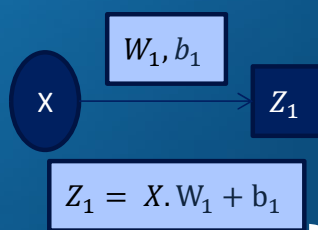
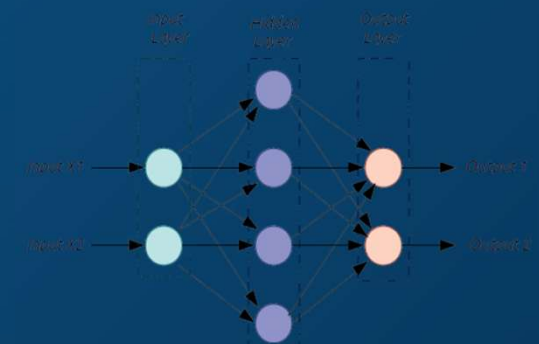
For $a^{[2]}$ all $a^{[1]}$ are acting as input features





Batch Normalization

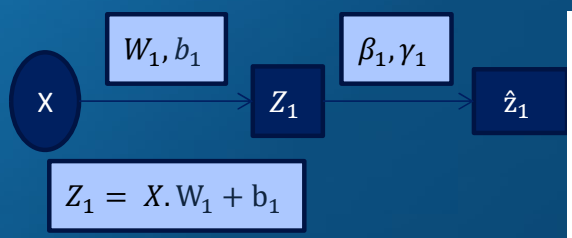
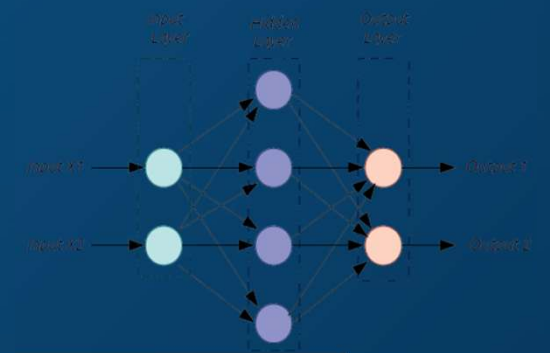
Forward and back propagation with batch norm:



Our standard equation to calculate z_1 .

Batch Normalization

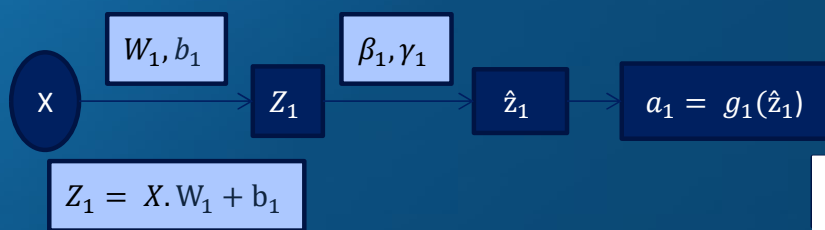
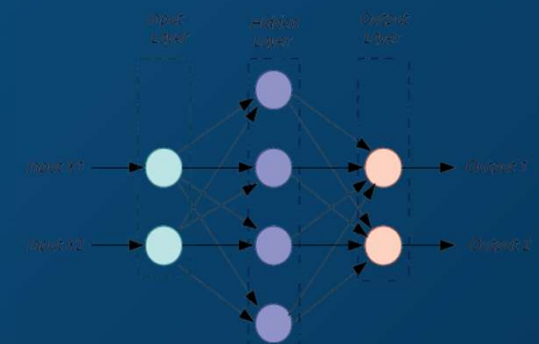
Forward and back propagation with batch norm:



Calculate \hat{z}_1 , based on β_1, γ_1

Batch Normalization

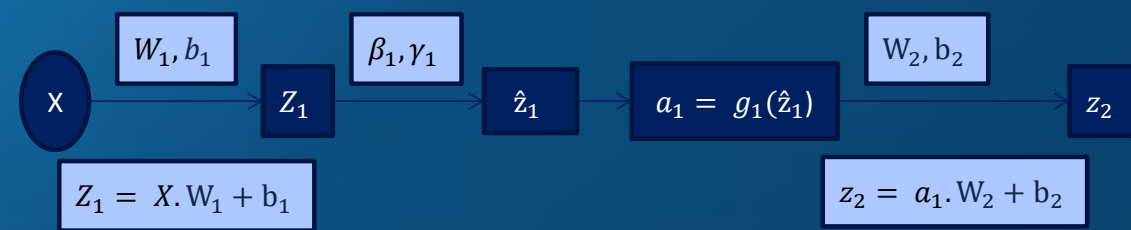
Forward and back propagation with batch norm:



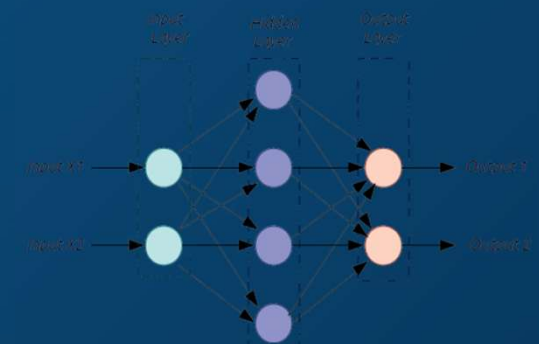
Apply activation function $g_1(\hat{z}_1)$

Batch Normalization

Forward and back propagation with batch norm:

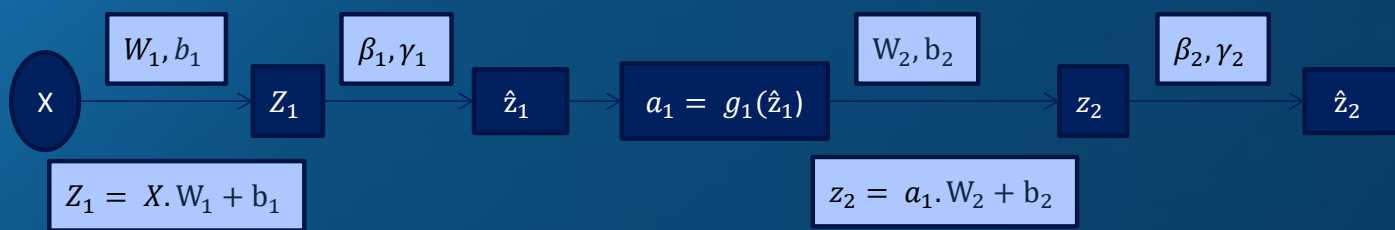
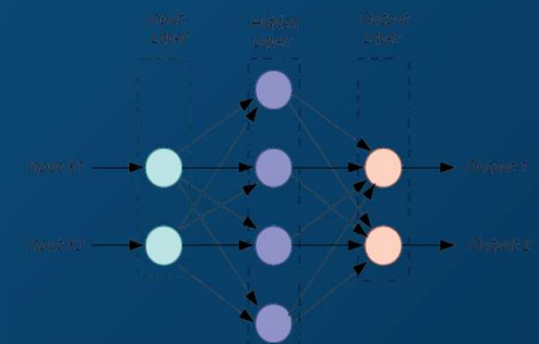


Calculate z_2 as usual



Batch Normalization

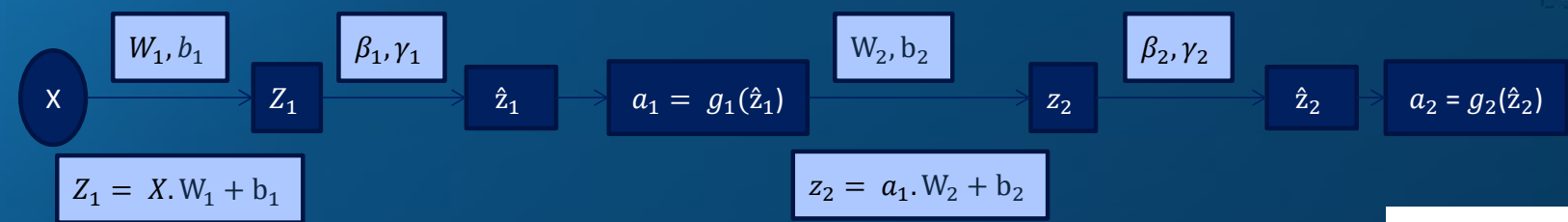
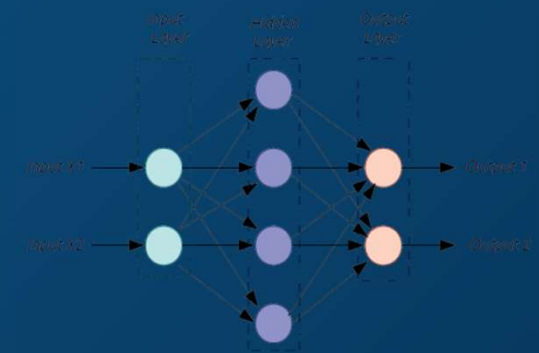
Forward and back propagation with batch norm:



We know how to calculate \hat{z}_2

Batch Normalization

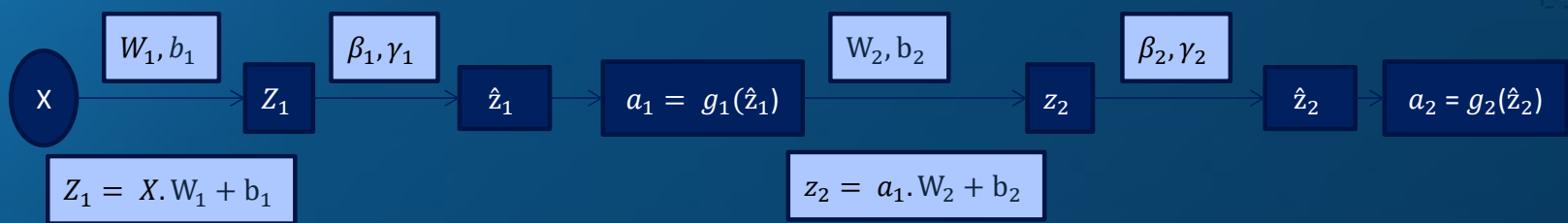
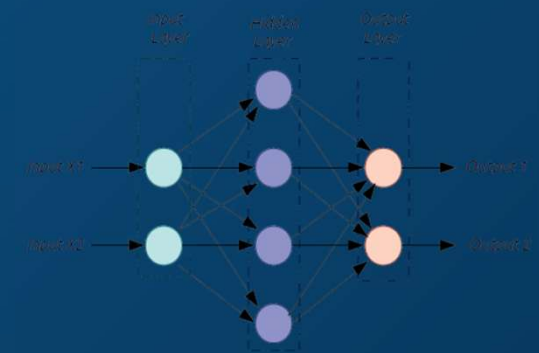
Forward and back propagation with batch norm:



We also know how to calculate a_2

Batch Normalization

Forward and back propagation with batch norm:



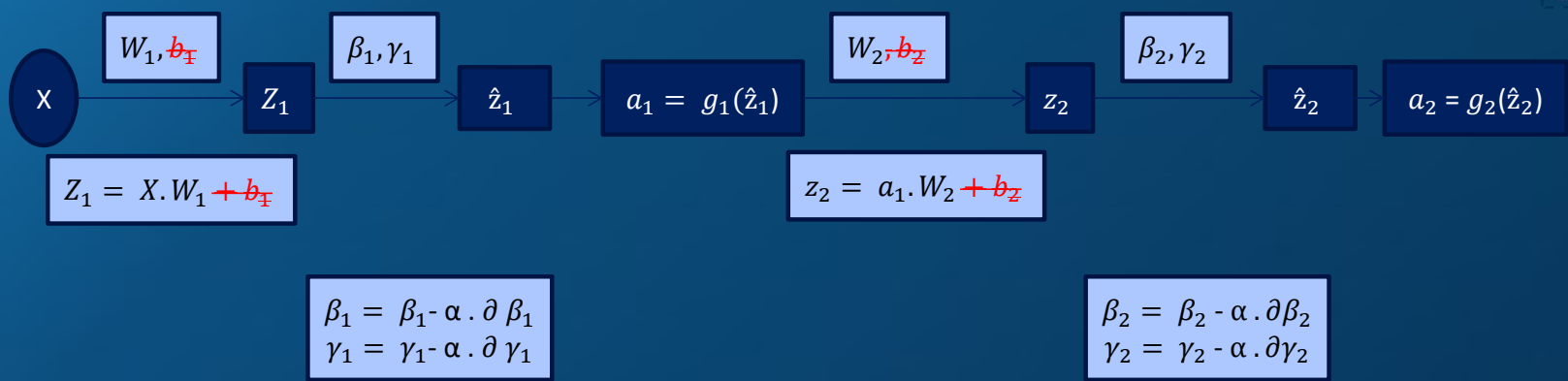
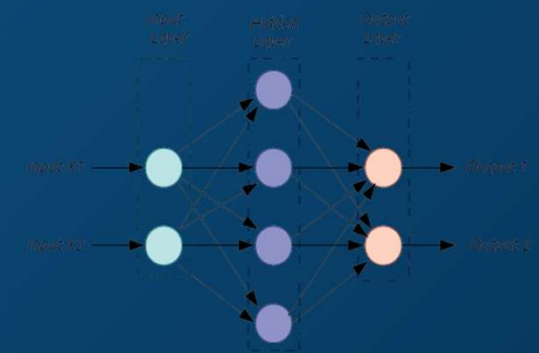
$$\begin{aligned} \beta_1 &= \beta_1 - \alpha \cdot \partial \beta_1 \\ \gamma_1 &= \gamma_1 - \alpha \cdot \partial \gamma_1 \end{aligned}$$

Using the gradient descent,
update β 's, γ 's along with
 W 's and b 's

$$\begin{aligned} \beta_2 &= \beta_2 - \alpha \cdot \partial \beta_2 \\ \gamma_2 &= \gamma_2 - \alpha \cdot \partial \gamma_2 \end{aligned}$$

Batch Normalization

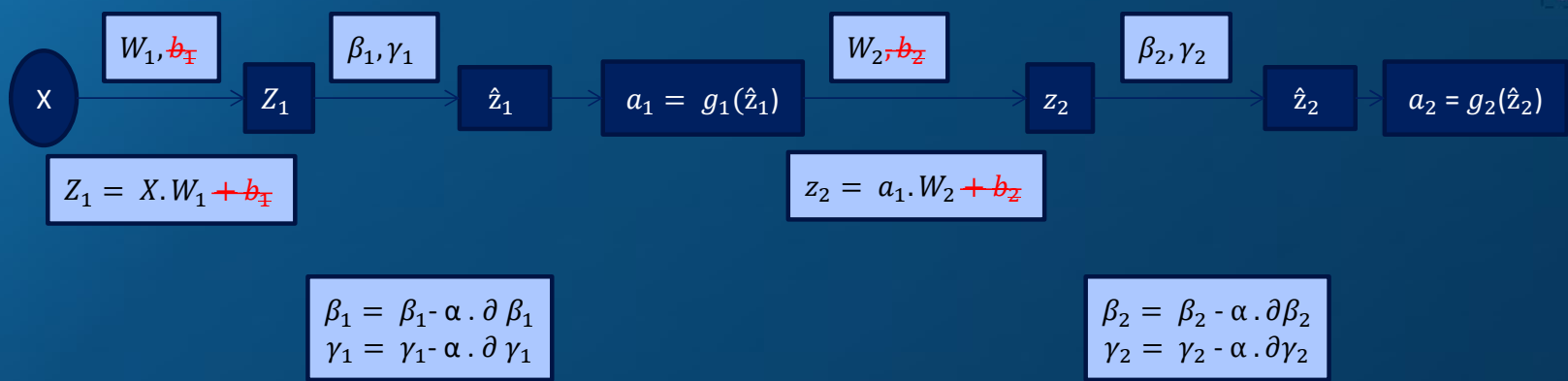
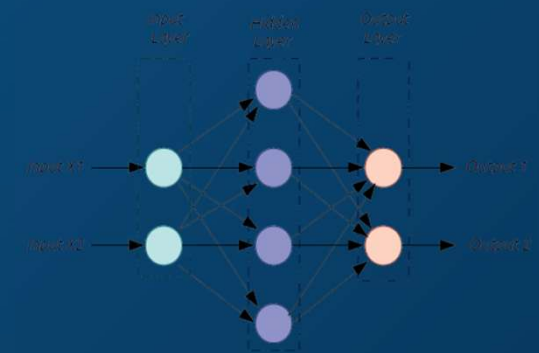
Forward and back propagation with batch norm:



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!!

Batch Normalization

Forward and back propagation with batch norm:

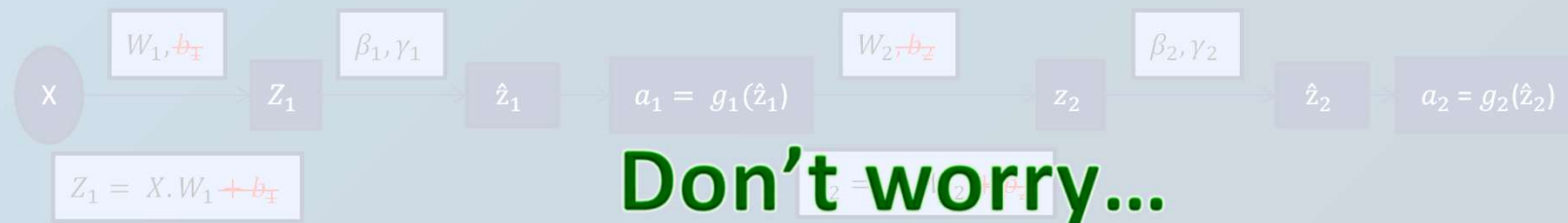


And at test/validation time using a exponentially weighted average!
 So while training do not forget to save exponentially weighted values or simply running average!!

Batch Normalization

- Forward and back propagation with batch norm:

Too many calculation steps...



Don't worry...

Most frameworks have one line code to do it.

Batch Normalization – Code Sample

```
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...

Batch Normalization



Remember β, γ are parameters you train!

At test time, we use these fixed running_mean and running_var values for normalization.

Reflect...

- ❑ 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

Next Session... Recurrent Neural Networks