# What is Batch Normalization?

Batch Normalization (BN) is a technique in deep learning used to improve the training of deep neural networks by normalizing the inputs to each layer. It stabilizes the learning process, reduces the sensitivity to the choice of hyperparameters, and often speeds up convergence.

---

## Why Use Batch Normalization?

1. **Internal Covariate Shift**: During training, the distribution of activations in a network changes as the parameters of the previous layers update. This is called "internal covariate shift," and it slows down training. BN addresses this issue by normalizing the inputs to each layer, keeping their distribution stable.

2. **Improved Training Speed**: By normalizing activations, the network can use higher learning rates without diverging, speeding up convergence.

3. **Regularization**: BN introduces a slight regularization effect, reducing the need for other regularization techniques like dropout.

4. **Reduced Dependence on Initialization**: Makes the training process less sensitive to the initial weights.

---

## How Batch Normalization Works

1. **Input Normalization**: For each mini-batch of data during training, BN normalizes the inputs to a layer to have:

   - **Mean = 0**

   - **Standard deviation = 1**

   This is done using the following formulas:

$$\mu_B = \frac{1}{m} \sum_{i=1}^{m} x_i$$

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_B)^2$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

Here:

- $\mu_B$: Mean of the batch.

- $\sigma_B^2$: Variance of the batch.

- $\epsilon$: Small constant added for numerical stability.

- $\hat{x}_i$: Normalized input.

2. **Learnable Parameters**: Instead of using the normalized values directly, BN introduces two learnable parameters:

   - **Scale ($\gamma$)**: Controls the spread of the normalized values.

   - **Shift ($\beta$)**: Adjusts the center of the normalized values.

$$y_i = \gamma \hat{x}_i + \beta$$

3. **Backpropagation**: The parameters $\gamma$ and $\beta$ are updated during training, allowing the model to learn the optimal scaling and shifting for each layer.

---

## Batch Normalization in Practice

1. **Where to Apply BN?**

   - BN is typically applied after the linear transformation (e.g., dense or convolutional layer) but **before** the activation function.

2. **Usage in Convolutional Layers**:

   - For convolutional layers, BN normalizes over the batch and spatial dimensions (height and width).

---

## Advantages of Batch Normalization

1. **Stabilizes Training**: Normalized activations reduce issues caused by high learning rates or poor initialization.

2. **Speeds Up Convergence**: Allows the use of higher learning rates, leading to faster training.

3. **Improves Generalization**: Adds a slight regularization effect, helping reduce overfitting.

4. **Compatibility**: Works well with various architectures (e.g., CNNs, RNNs) and optimizers.

---

# Implementation in Python (Using Keras)

Batch normalization is easy to apply in deep learning frameworks like TensorFlow/Keras:

### Example in a Fully Connected Network

```python
from tensorflow.keras.models import Sequential from tensorflow.keras.layers import Dense,
BatchNormalization, Activation model = Sequential([ Dense(128, input_shape=(784,)),
BatchNormalization(), Activation('relu'), Dense(64), BatchNormalization(),
Activation('relu'), Dense(10, activation='softmax') ])
```

### Example in a Convolutional Neural Network

```python
from tensorflow.keras.models import Sequential from tensorflow.keras.layers import
Conv2D, BatchNormalization, Activation, Flatten, Dense model = Sequential([ Conv2D(32,
(3, 3), input_shape=(64, 64, 3)), BatchNormalization(), Activation('relu'), Flatten(),
Dense(10, activation='softmax') ])
```

---

# Disadvantages

1. **Extra Computation**: BN introduces additional computations during training.

2. **Dependence on Batch Size**: Small batch sizes can lead to unstable estimates of mean and variance.

3. **Potential Incompatibility**: In some cases, especially with very small batches or certain architectures (like RNNs), alternative normalization techniques (e.g., Layer Normalization) may work better.

---

# Key Points to Remember

- BN reduces sensitivity to hyperparameters like learning rate.

- It may not work well with small batch sizes (use Group Normalization or Layer Normalization instead).

- It is applied differently in training and inference:

- **Training**: Uses batch mean and variance.
- **Inference**: Uses running averages of mean and variance computed during training.