

## What is Internal Covariate Shift?

**Internal covariate shift** refers to the change in the distribution of a neural network's layer inputs during training, as the parameters of the preceding layers are updated. This phenomenon occurs because the input to each layer depends on the outputs of all previous layers, and as these layers' parameters (weights and biases) are updated via backpropagation, the distribution of inputs to subsequent layers shifts. This can make training deep neural networks slower and more unstable, as each layer must continuously adapt to these changing input distributions.

The term was introduced in the context of **Batch Normalization** (Ioffe & Szegedy, 2015), which aims to mitigate this issue by normalizing layer inputs to maintain a consistent distribution during training.

## Why is Internal Covariate Shift a Problem?

1. **Training Instability:** As the input distribution to a layer changes, the layer's parameters (e.g., weights) may no longer be optimal, requiring constant readjustment. This can lead to slower convergence or even training divergence.
2. **Vanishing/Exploding Gradients:** Shifting distributions can exacerbate gradient issues, especially in deep networks, making it harder for gradients to propagate effectively.
3. **Hyperparameter Sensitivity:** Models become more sensitive to learning rates and other hyperparameters, as they must compensate for the shifting distributions.

## Example to Illustrate Internal Covariate Shift

Let's consider a simple example in the context of training a deep neural network, such as a multi-layer perceptron (MLP) or a Transformer-based large language model (LLM).

### Scenario: Training a Deep Neural Network

- **Network Setup:** Suppose you have a 10-layer MLP trained to classify text

sentiment (positive/negative) using a dataset of movie reviews. Each layer applies a linear transformation followed by a ReLU activation:  $y = \text{ReLU}(Wx + b)$ .

- **Training Process:** During training, the network uses stochastic gradient descent (SGD) to update weights and biases based on mini-batches of data.

## Internal Covariate Shift in Action

- **Initial State:** At the start of training, the weights in the first layer ( $W_1$ ) produce outputs that feed into the second layer. Let's say the outputs of the first layer (inputs to the second layer) have a mean of 0 and a standard deviation of 1 after the ReLU activation.
- **After Several Updates:** As training progresses, the weights ( $W_1$ ) are updated, causing the outputs of the first layer to shift. For example, the mean of the outputs might increase to 2, and the standard deviation might grow to 3 due to changes in the weight distribution or input data variance in different mini-batches.
- **Impact on Later Layers:** The second layer, which was trained assuming inputs with mean 0 and standard deviation 1, now receives inputs with a different distribution (mean 2, standard deviation 3). This shift forces the second layer's weights to adapt to this new distribution, even if the actual task (sentiment classification) hasn't changed. This process repeats across all layers, causing a cascading effect.

## Consequences

- **Slower Convergence:** The second layer (and subsequent layers) must continuously adjust to these shifting inputs, slowing down the optimization process.
- **Loss Spikes:** If the distribution shift is significant (e.g., due to a large learning rate or a noisy mini-batch), the loss function may spike, as seen in training LLMs when the loss curve exhibits sharp increases.
- **Example in LLMs:** In a Transformer model trained for language modeling, the output of one attention layer feeds into the next. If the attention weights change significantly, the input distribution to the next layer (e.g., token embeddings or attention scores) may shift, causing the subsequent layer to receive inputs with a different mean or variance. This can lead to unstable

training, especially in deep models with many layers.

## How Batch Normalization Addresses Internal Covariate Shift

Batch Normalization (BatchNorm) mitigates internal covariate shift by normalizing the inputs to each layer to have a fixed mean (typically 0) and variance (typically 1) within each mini-batch. The process is:

1. Compute the mini-batch mean ( $\mu_B$ ) and variance ( $\sigma_B^2$ ).
2. Normalize the layer inputs:  $\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$ .
3. Scale and shift the normalized values using learnable parameters ( $\gamma$ ) and ( $\beta$ ):  $y_i = \gamma \hat{x}_i + \beta$ .

This ensures that each layer sees inputs with a consistent distribution, reducing the impact of internal covariate shift.

### Example with BatchNorm

- **Without BatchNorm:** In the MLP example, the second layer's inputs might shift from a mean of 0 to 2 across epochs, requiring the layer to relearn its weights.
- **With BatchNorm:** The second layer's inputs are normalized to have mean 0 and variance 1 for each mini-batch, stabilizing the training process. The learnable parameters ( $\gamma$ ) and ( $\beta$ ) allow the layer to adapt the normalized inputs to the task without suffering from distribution shifts.

### Real-World Analogy

Imagine a chef (the neural network layer) cooking a dish (processing inputs) using ingredients (layer inputs) that change unpredictably each day (e.g., sometimes sweet, sometimes salty). The chef must constantly adjust the recipe, slowing down the cooking process. BatchNorm is like standardizing the ingredients (e.g., always using 1 tsp of salt), so the chef can focus on perfecting the recipe without worrying about ingredient variability.

### Connection to Loss Spikes in LLM Training

Internal covariate shift can contribute to loss spikes in LLM training (as

discussed in your previous question). For example:

- A sudden change in the distribution of token embeddings or attention outputs in a Transformer layer can cause a downstream layer to produce erroneous predictions, leading to a spike in the loss curve.
- Techniques like Layer Normalization (common in LLMs) or adaptive optimizers like AdamW (with bias correction, as discussed earlier) help mitigate these shifts, stabilizing training.

## **Summary**

Internal covariate shift is the change in the distribution of layer inputs during training due to updates in preceding layers' parameters. It slows convergence and can cause instability, such as loss spikes in LLMs. For example, in a deep MLP or Transformer, shifting input distributions (e.g., mean shifting from 0 to 2) force layers to adapt continuously. Batch Normalization or Layer Normalization reduces this issue by standardizing inputs, improving training stability.

If you'd like a visualization (e.g., a chart showing how layer input distributions shift over time or the effect of BatchNorm), or a more specific example (e.g., in a Transformer model), let me know!