# Gradient Clipping Explained

<div align="right">what purpose does gradient clipping serve ?</div>

**Gradient clipping** is a technique used in training deep neural networks to prevent the **exploding gradient problem** — a scenario where gradients during backpropagation become too large, leading to unstable weight updates and failure to converge.

## 🔧 Purpose of Gradient Clipping:

- **Stabilizes Training**: Keeps gradients within a reasonable range, especially in deep networks or RNNs.

- **Prevents NaNs**: Large gradients can cause weights to overflow, resulting in `NaN` values.

- **Improves Convergence**: Avoids erratic updates that derail the learning process.

## 💡 How It Works:

The idea is to **scale down** gradients when their norm exceeds a predefined threshold.

**Example: Clipping by Norm**

```python
torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
```

If the total norm of gradients > 1.0, they are rescaled proportionally so the norm equals 1.0.

## 📍 Common Use Cases:

- Recurrent Neural Networks (RNNs), LSTMs, GRUs

- Deep reinforcement learning

- Very deep feedforward networks

Would you like a visual illustration of how exploding gradients affect weight updates?

<div align="right">explain norm in gradient clipping</div>

In **gradient clipping**, the **norm** refers to the **magnitude of the entire gradient vector** across all parameters. It is used to determine whether the gradients are "too large" and need to be scaled down.

# 🔢 What is the Norm?

The most common norm used is the **L2 norm** (Euclidean norm), which measures the overall size of the gradient vector.

**For a gradient vector g:**

$$\|g\|_2 = \sqrt{g_1^2 + g_2^2 + \cdots + g_n^2}$$

This is simply the square root of the sum of squares of all individual gradients from all parameters.

---

# 🧠 Why Norm Matters in Clipping

When training a model, if:

$$\|g\|_2 > \text{max\_norm}$$

then the gradients are too large, and we clip them like this:

$$g_{\text{clipped}} = g \cdot \frac{\text{max\_norm}}{\|g\|_2}$$

This **scales all gradients proportionally** so their norm becomes exactly `max_norm`.

---

# 📌 Example:

Suppose your model has three parameters with gradients:

- $g_1 = 3$
- $g_2 = 4$
- $g_3 = 12$

L2 norm:

$$\|g\|_2 = \sqrt{3^2 + 4^2 + 12^2} = \sqrt{9 + 16 + 144} = \sqrt{169} = 13$$

If `max_norm = 5`, then clip like this:

$$\text{scale} = \frac{5}{13}$$

$$g_{\text{clipped}} = g \cdot \frac{5}{13} \Rightarrow [3, 4, 12] \cdot \frac{5}{13} = [1.15, 1.54, 4.62]$$

---

Would you like an explanation of other norms like L1 or max norm in this context too?

how is max norm choosen ?

Choosing the **max norm** (`max_norm`) for gradient clipping is a **hyperparameter tuning decision** — there's no universal value, but here's how you can choose it effectively:

---

## ✅ General Guidelines for Choosing `max_norm`

1. **Start Small**:

   - Common starting values are **1.0**, **5.0**, or **10.0**.

   - Example (PyTorch):

     ```python
     torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
     ```

2. **Use Model Feedback**:

   - If training **explodes** or gives `NaN` values: try lowering the norm.

   - If training **slows down or stalls**: try increasing it slightly.

3. **Smaller for RNNs/Transformers**:

   - RNNs, LSTMs, and deep Transformers are more prone to exploding gradients, so `max_norm = 1.0` is common.

4. **Check Gradient Norms During Training**:

   - Log the gradient norm per batch before clipping:

     ```python
     total_norm = torch.norm(torch.stack([torch.norm(p.grad.detach(), 2) for p in
     ```

```
    model.parameters() if p.grad is not None]))
```

- Then decide a reasonable upper bound based on the typical range.

---

## 📌 Rule of Thumb

| Model Type | Recommended `max_norm` |
|---|---|
| Small feedforward | 5–10 |
| Deep CNNs | 1–5 |
| RNNs/LSTMs | 0.25–1 |
| Transformers (NLP) | 0.5–2 |

---

## 🔁 Tune Like Any Other Hyperparameter

Just like learning rate or batch size, try a few values and monitor:

- Loss curves
- Gradient norms
- Stability of training

---

Would you like a code snippet to auto-log gradient norms and visualize them while training?

> does tensorflow or pytorch gradient clipping work exactly as you wrote mathematically ?

Great question — and yes, **both TensorFlow and PyTorch implement gradient clipping in ways that closely match the math**, but there are **minor practical differences** in behavior, especially in **how the norm is computed** and **when clipping is applied**.

---

## ✅ PyTorch: `clip_grad_norm_`

```python
torch.nn.utils.clip_grad_norm_(parameters, max_norm, norm_type=2.0)
```

### Behavior:

- Computes the total norm:

$$\text{total\_norm} = \left( \sum_i \|g_i\|^p \right)^{1/p}$$

  where `g_i` is the gradient of parameter $i$, and `p` is the `norm_type` (default = 2).

- If `total_norm > max_norm`, rescales **all gradients**:

$$g_i \leftarrow g_i \cdot \frac{\text{max\_norm}}{\text{total\_norm} + \epsilon}$$

- Done **in-place** (the `_` suffix).

📌 **Important**: PyTorch adds a small `epsilon` to avoid divide-by-zero.

---

## ✅ TensorFlow: `clip_by_global_norm`

```python
clipped_grads, _ = tf.clip_by_global_norm(gradients, clip_norm)
```

### Behavior:

- Computes:

$$\text{global\_norm} = \sqrt{\sum_i \|g_i\|_2^2}$$

- If `global_norm > clip_norm`, rescales:

$$g_i \leftarrow g_i \cdot \frac{\text{clip\_norm}}{\text{global\_norm}}$$

📌 Similar logic as PyTorch, but the API is more functional (returns new clipped tensors instead of modifying in-place).

## 🔍 Key Differences Summary:

| Aspect | PyTorch | TensorFlow |
|---|---|---|
| Default norm type | L2 (can set L1, ∞, etc.) | Always L2 |
| In-place vs return | Modifies grads in-place (`clip_grad_norm_`) | Returns new clipped grads |
| Logging norm | Returns total norm | Returns both clipped grads and norm |
| Use in optim step | After `.backward()` before `.step()` | After `tape.gradient()` before `apply_gradients()` |

Would you like a side-by-side code example in PyTorch vs TensorFlow?