# 3.1 WHAT IS `autograd`

01 September 2025     12:37 AM     ☒ Avinash Yadav

`**autograd**` is a core component of PyTorch that *provides automatic differentiation for tensor operations.* It *enables gradient computation*, which is essential for training machine learning models using optimization algorithms like gradient descent.

↳ with the help of it, we can perform automatic differentiation on any kind of tensoris operation

When training neural networks, the most frequently used algorithm is **back propagation**. In this algorithm, parameters (model weights) are adjusted according to the **gradient** of the loss function with respect to the given parameter. To compute those gradients, PyTorch has a built-in differentiation engine called `torch.autograd`. It supports automatic computation of gradient for any computational graph.

## EXAMPLE 1.)

1.)   $Y = x^2$ $\xrightarrow{\text{To FIND}}$ $\dfrac{dY}{dx}$   for given 'x'

```python
1  import torch
```
✓ 0.0s

*Scalar Tensor*

```python
1  x = torch.tensor(3.0, requires_grad=True)
```
✓ 0.0s

we are setting require gradient attribute as True. by default it is False. By this we are telling PyTorch that we want to calculate derivative of tensor 'x'. so it starts tracking it.

So whenever we want to calculate derivate of a tensor; then while creating, we have to set 'requires_grad' attribute as True.

```python
1  y = x**2
```
✓ 0.0s

'y' calculated in forward dxn

```python
1  print(x)
2
3  print(y)
```
✓ 0.0s

```
tensor(3., requires_grad=True)
tensor(9., grad_fn=<PowBackward0>)
```

So when we set it True, PyTorch make computation graph internally

$Forward \rightarrow Y = x^2$

$X \rightarrow 9x \rightarrow y$

backward

$\dfrac{dy}{dx}$

computed value

```python
1  y.backward()
```
✓ 0.0s

Now going backward, it will compute $\dfrac{dy}{dx}$ i.e. $2x$
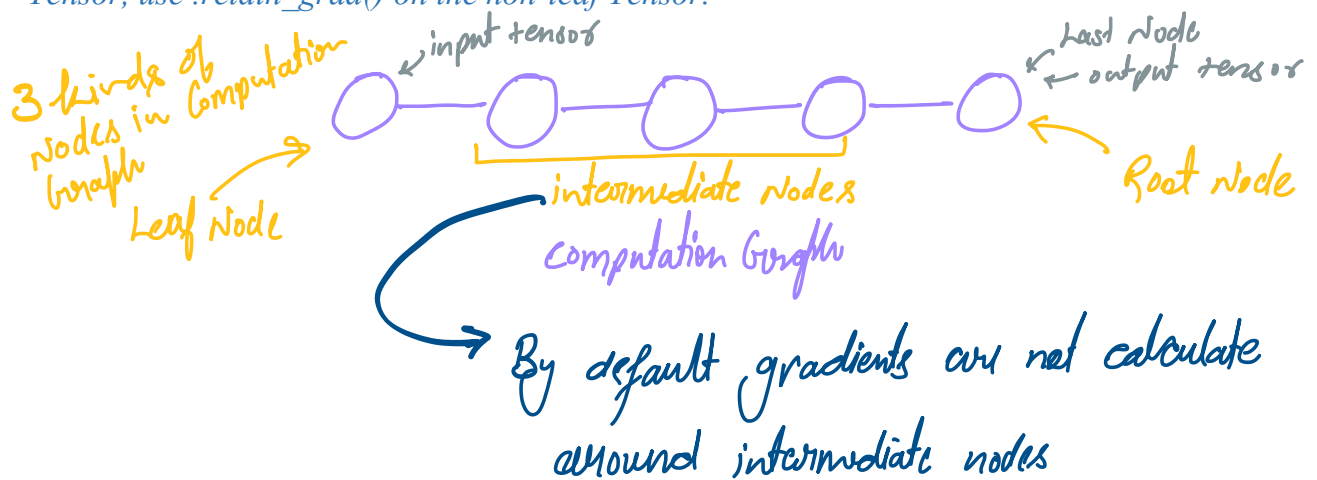
So while going backward PyTorch needs to know what operation was performed here while it went for forward pass. that why it remembers it as 'PowBackward0' for backward pass

```python
1  x.grad
```
✓ 0.0s

```
tensor(6.)
```

this is $\dfrac{dy}{dx}$ = $f'(x) = f'(3)$
= 6

*The .grad attribute of a Tensor that is not a leaf Tensor won't be populated during autograd.backward(). If you indeed want the .grad field to be populated for a non-leaf Tensor, use .retain_grad() on the non-leaf Tensor.*

3 kinds of nodes in Computation graph

input tensor

Last Node ← output tensor

Leaf Node

intermediate nodes

Computation Graph

Root Node

By default gradients are not calculate around intermediate nodes

Conceptually, autograd keeps a record of data (tensors) and all executed operations (along

with the resulting new tensors) in a directed acyclic graph (DAG) consisting of Function objects. In this DAG, ***leaves are the input tensors***, ***roots are the output tensors***. By tracing this graph from roots to leaves, you can automatically compute the gradients using the chain rule.

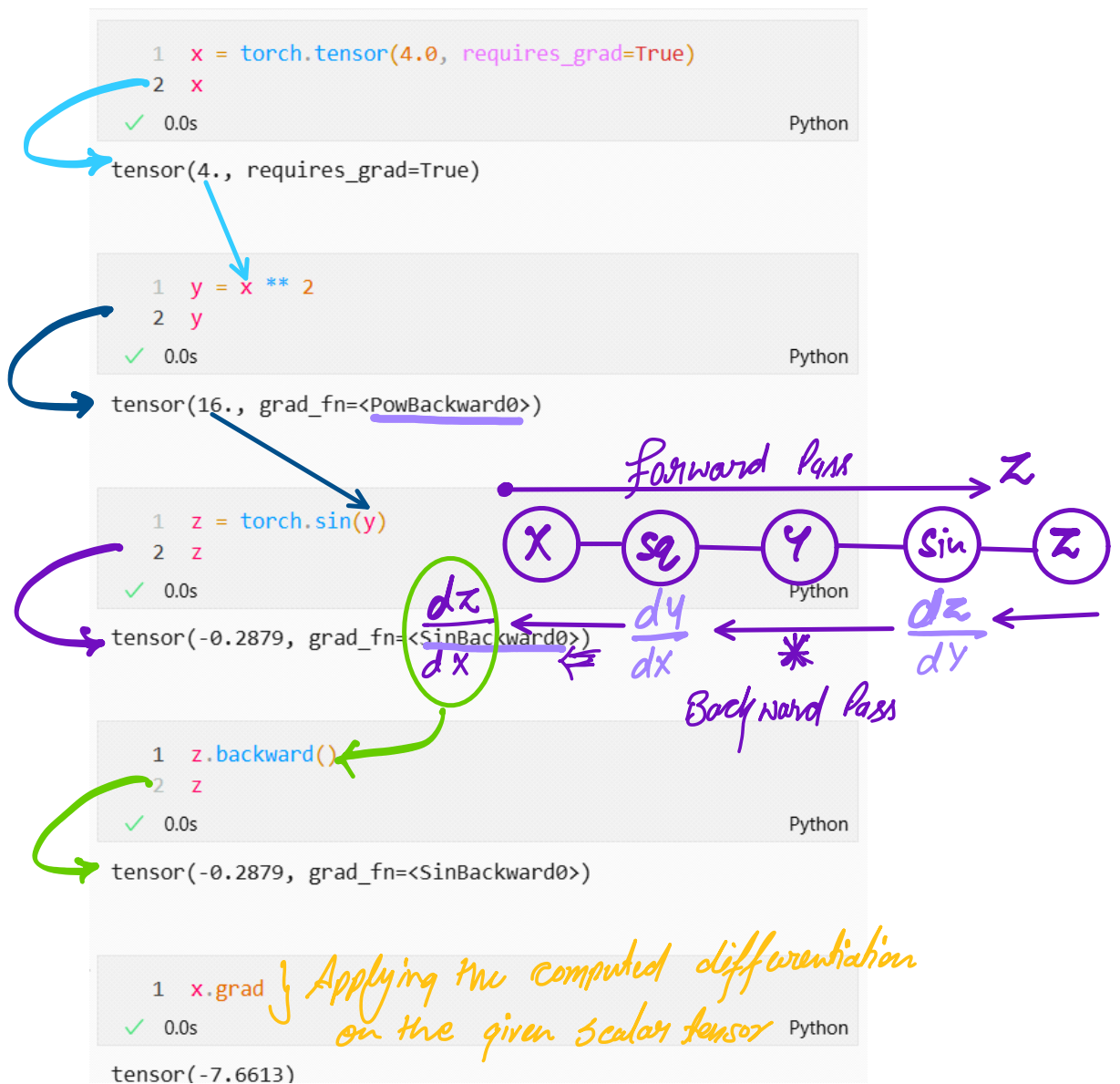In a **forward pass**, autograd does two things simultaneously:
- run the requested operation to compute a resulting tensor
- maintain the operation's *gradient function* in the DAG.

The **backward pass** kicks off when *.**backward()*** is called on the DAG ***root.autograd*** then:
- computes the gradients from each *.**grad_fn***,
- accumulates them in the respective tensor's *.**grad attribute***
- using the chain rule, propagates all the way to the leaf tensors.

## $\underline{\text{Example } \alpha.)}$

$$Y = x^2$$

$$Z = \sin(Y)$$

$\left.\vphantom{\begin{matrix}a\\b\end{matrix}}\right\}$ we need $\dfrac{dZ}{dx}$

```python
1  x = torch.tensor(4.0, requires_grad=True)
2  x
```
✓ 0.0s        Python

tensor(4., requires_grad=True)

```python
1  y = x ** 2
2  y
```
✓ 0.0s        Python

tensor(16., grad_fn=<PowBackward0>)

```python
1  z = torch.sin(y)
2  z
```
✓ 0.0s        Python

tensor(-0.2879, grad_fn=<SinBackward0>)

```python
1  z.backward()
2  z
```
✓ 0.0s        Python

tensor(-0.2879, grad_fn=<SinBackward0>)

*Forward Pass* → $z$

$X$ — $sq$ — $Y$ — $Sin$ — $Z$

$\dfrac{dz}{dx}$ ← $\dfrac{dy}{dx}$ ← $*$ ← $\dfrac{dz}{dy}$ ←

*Backward Pass*

```python
1  x.grad
```
✓ 0.0s        Python

} *Applying the computed differentiation on the given scalar tensor*

tensor(-7.6613)

# How Does autograd works?

INPUT →

| CGPA | PLACED |
|------|--------|
| 6.7  | 0 → NO |

1. **Linear Transformation:**

$$z = w \cdot x + b$$

2. **Activation (Sigmoid Function):**

$$y_{\text{pred}} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3. **Loss Function (Binary Cross-Entropy Loss):**

$$L = -[y_{\text{target}} \cdot \ln(y_{\text{pred}}) + (1 - y_{\text{target}}) \cdot \ln(1 - y_{\text{pred}})]$$

since we have 2 parameters 'w' & 'b'. Therefore we'll need to find 2 derivates.

1·Y derivative of loss wrt 'w' $\Rightarrow \partial L/\partial w$

2·Y derivative of loss wrt 'b' $\Rightarrow \partial L/\partial b$

$$\therefore \quad \frac{\partial L}{\partial w} = \frac{\partial L}{\partial y_{pred}} \times \frac{\partial y_{pred}}{\partial z} \times \frac{\partial z}{\partial w}$$

$$\downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow$$

from eq.3    from eq.2    from eq.1

Similarly;

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial y_{pred}} \times \frac{\partial y_{pred}}{\partial z} \times \frac{\partial z}{\partial b}$$

$$\downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow$$

from eq.3    from eq.2    from eq.1

So now computing derivates:

$$\frac{\partial h}{\partial y_{pred}} = \frac{(y_{pred} - y)}{y_{pred}(1 - y_{pred})} \qquad \frac{\partial z}{\partial w} = x$$

$$\frac{\partial y_{pred}}{\partial z} = y_{pred}(1 - y_{pred}) \qquad \frac{\partial z}{\partial b} = 1$$

$$\therefore \frac{\partial h}{\partial w} = \frac{\partial h}{\partial y_{pred}} \times \frac{\partial y_{pred}}{\partial z} \times \frac{\partial z}{\partial w}$$

$$= \frac{(y_{pred} - y)}{y_{pred}(1 - y_{pred})} * y_{pred}(1 - y_{pred}) * x$$

$$= (y_{pred} - y) * x$$

$$\therefore \frac{\partial h}{\partial w} = x(y_{pred} - y)$$

Similarly

$$\frac{\partial h}{\partial b} = (y_{pred} - y)$$

```python
1  import torch
2
3  # Inputs
4  x = torch.tensor(6.7)   # Input feature          → x
5  y = torch.tensor(0.0)   # True label (binary)    → y
6
7  w = torch.tensor(1.0)   # Weight                 y initial setup with Random value
8  b = torch.tensor(0.0)   # Bias
```
✓  0.0s                                                                Python

*Calculating loss*

$$L = - [y_{\text{target}} \cdot \ln(y_{\text{pred}}) + (1 - y_{\text{target}}) \cdot \ln(1 - y_{\text{pred}})]$$

```python
1  # Binary Cross-Entropy Loss for scalar
2  def binary_cross_entropy_loss(prediction, target):
3      epsilon = 1e-8  # To prevent log(0)
4      prediction = torch.clamp(prediction, epsilon, 1 - epsilon)
5      return -(target * torch.log(prediction) + (1 - target) * torch.log(1 - prediction))
```
✓  0.0s                                                                Python

```python
1  # Forward pass
2  z = w * x + b   # Weighted sum (linear part)      → Eq.1
3
4  y_pred = torch.sigmoid(z)  # Predicted probability → Eq.2
5
6  # Compute binary cross-entropy loss
7  loss = binary_cross_entropy_loss(y_pred, y)        → Eq.3
8  loss
```
✓  0.0s                                                                Python

tensor(6.7012)

**BACK-PROPAGATION**

```python
1   # Derivatives:
2   # 1. dL/d(y_pred): Loss with respect to the prediction (y_pred)
3   dloss_dy_pred = (y_pred - y)/(y_pred*(1-y_pred))
4
5   # 2. dy_pred/dz: Prediction (y_pred) with respect to z (sigmoid derivative)
6   dy_pred_dz = y_pred * (1 - y_pred)
7
8   # 3. dz/dw and dz/db: z with respect to w and b
9   dz_dw = x   # dz/dw = x
10  dz_db = 1   # dz/db = 1 (bias contributes directly to z)
11
12  dL_dw = dloss_dy_pred * dy_pred_dz * dz_dw
13  dL_db = dloss_dy_pred * dy_pred_dz * dz_db
```
✓  0.0s                                                                Python

```python
1  print(f"Manual Gradient of loss w.r.t weight (dw): {dL_dw}")
2  print(f"Manual Gradient of loss w.r.t bias (db): {dL_db}")
```
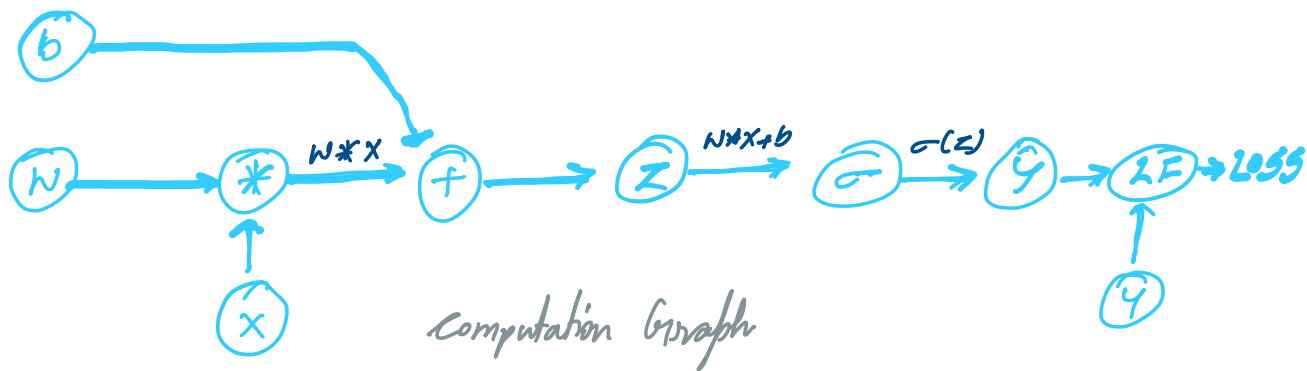✓  0.0s                                                                Python

Manual Gradient of loss w.r.t weight (dw): 6.691762447357178
Manual Gradient of loss w.r.t bias (db): 0.998770534992218

Computation Graph

# Using `autograd`

```python
1  x = torch.tensor(6.7)
2  y = torch.tensor(0.0)
3  print(x)
4  print(y)
```
We are not using `requires_grad = True` because we do not have to calculate the derivatives Wrt to X & Y

✓ 0.0s                                    Python

```
tensor(6.7000)
tensor(0.)
```

```python
1  w = torch.tensor(1.0, requires_grad=True)
2  b = torch.tensor(0.0, requires_grad=True)
3  print(w)
4  print(b)
```

✓ 0.0s                                    Python

```
tensor(1., requires_grad=True)
tensor(0., requires_grad=True)
```

# STARTING FORWARD PASS :

```python
1  z = w*x + b
2  z
```
Eq. 1.

✓ 0.0s                                    Python

```
tensor(6.7000, grad_fn=<AddBackward0>)
```

```python
1  y_pred = torch.sigmoid(z)
2  y_pred
```
Eq 2.

✓ 0.0s                                    Python

```
tensor(0.9988, grad_fn=<SigmoidBackward0>)
```

```python
1  loss = binary_cross_entropy_loss(y_pred, y)
2  loss
```
Eq.3.

✓ 0.0s                                    Python

```
tensor(6.7012, grad_fn=<NegBackward0>)
```

```python
1  loss.backward()
```

✓ 0.0s                                    Python

```python
1  loss.backward()
```
✓ 0.0s                                                                Python

```python
1  print(f"By using autograd function (dw): {w.grad}")
2  print(f"By using autograd function (db): {b.grad}")
```
✓ 0.0s                                                                Python

```
By using autograd function (dw): 6.6917619705200195
By using autograd function (db): 0.9987704753875732
```

— ✗ — ✗ — ✗ — ✗ — ✗ — ✗ —

# Clearing Gradients

The key concept behind `autograd` in PyTorch is **gradient accumulation**. Each time you call `.backward()` on a tensor, the gradients for all tensors with `requires_grad=True` are added (accumulated) to their `.grad` attributes. This means if you run the backward pass multiple times without clearing the gradients, the values in `.grad` will keep increasing, reflecting the sum of all computed gradients.

**Why does this happen?**
This behavior is useful when training neural networks using mini-batches. You might want to accumulate gradients over several batches before updating the model parameters.

**How to manage gradient accumulation:**

- To avoid unwanted accumulation, always clear gradients before a new backward pass using `.zero_()` on the `.grad` attribute:

  ```python
  x.grad.zero_()
  ```

- Alternatively, use `optimizer.zero_grad()` when working with optimizers.

**Summary:**

- `.backward()` accumulates gradients in `.grad`.
- Always clear gradients before a new backward pass unless you intentionally want to accumulate them.

**FIRST RUN**

```python
1  x = torch.tensor(2.0, requires_grad=True)
2  x
```
✓ 0.0s                                           Python

```
tensor(2., requires_grad=True)
```

**SECOND RUN**          **THIRD RUN**

FP
```python
1  y = x ** 2
2  y
```
✓ 0.0s                          Python

```
tensor(4., grad_fn=<PowBackward0>)
```

FP
```python
1  y = x ** 2
2  y
```
✓ 0.0s                          Python

```
tensor(4., grad_fn=<PowBackward0>)
```

FP
```python
1  y = x ** 2
2  y
```
✓ 0.0s                          Python

```
tensor(4., grad_fn=<PowBackward0>)
```

BP
```python
1  y.backward()
```
✓ 0.0s                          Python

BP
```python
1  y.backward()
```
✓ 0.0s                          Python

BP
```python
1  y.backward()
```
✓ 0.0s                          Python

```
1  x = torch.tensor(2.0, requires_grad=True)
2  x
✓ 0.0s                                              Python
```
tensor(2., requires_grad=True)

```
1  y = x ** 2
2  y
✓ 0.0s                                              Python
```
tensor(4., grad_fn=<PowBackward0>)

```
1  y.backward()
✓ 0.0s                                              Python
```

```
1  x.grad
✓ 0.0s                                              Python
```
tensor(4.)

```
1  x.grad.zero_()
✓ 0.0s                                              Python
```
tensor(0.)

# Disabling Gradient Tracking

Sometimes, we may need to perform computations without tracking gradients or calculating derivatives.

**Scenarios Where Disabling Gradient Tracking is Useful**

- **Model Inference:**
  When making predictions with a trained model, gradients are not needed.

- **Model Evaluation:**
  During validation or testing phases, to save memory and computation.

- **Feature Extraction:**
  When using a model to extract features from data without updating weights.

- **Saving/Loading Model Outputs:**
  When storing intermediate results for later use.

- **Visualizations:**
  When plotting or analyzing outputs that do not require gradients.

- **Deployment:**
  In production environments where only forward passes are performed.

In such cases, PyTorch provides several ways to disable gradient tracking:

- **Set `requires_grad` to `False`:**
  You can turn off gradient tracking for a tensor by setting its `requires_grad` attribute to `False` using `requires_grad_(False)`.

- **Detach a tensor:**
  Use `.detach()` to create a new tensor that does not require gradients and is disconnected from the computation graph.

- **Use `torch.no_grad()` context:**
  Wrap your code inside a `with torch.no_grad():` block to temporarily disable gradient tracking for all operations within the block.

Disabling gradient tracking is useful for inference, evaluation, or any scenario where derivatives are not needed, as it reduces memory usage and speeds up computations.