

3.0 OVERVIEW PyTorch `autograd`

24 August 2025 02:47 AM Avinash Yadav

① **autograd** → PyTorch's automatic differentiation tool.

The Why?

Let's take a mathematical relation as:

$$y = x^2$$

Now let's say we have to write a Python script for a given x we need to find the derivative of y with respect to x :

i.e. for given x we want $\frac{dy}{dx}$

it is simple i.e. we have to find

$$\frac{dy}{dx} = 2x$$

now very easily we can code this expression

X	$\frac{dY}{dX}$
2	4
3	6
4	8
5	10

and whenever we keep any new value of x we can easily get the the value of $\frac{dy}{dx}$

```
def dy_dx(x):  
    return 2*x  
  
dy_dx(3)
```

since this was easier to do, we solve it manually and later coded it as well.

But what if, we have a complex relation/equation which we have to differentiate, then what will do in that case?

i.e. for example:

$$Y = x^2$$

$$Z = \sin(Y)$$

We want $\frac{dz}{dx}$

Now again we have to do the same task, but this time it is little difficult, now we have to chain rule of differentiation in order to solve this.

$$\begin{aligned} \text{i.e. for } \frac{dz}{dx} &= \frac{dz}{dy} \times \frac{dy}{dx} \\ &\downarrow \qquad \downarrow \\ \frac{d(\sin y)}{dy} \times \frac{d(x^2)}{dx} & \\ &\downarrow \qquad \downarrow \\ \cos(y) \times 2x & \\ &\downarrow \\ \frac{dz}{dx} = \cos(y) \times 2x & \end{aligned}$$

Now again we got the formula and for this we can again write the code such that for a given x we will get the derivative of z w.r.t x : i.e. $\frac{dz}{dx}$

```
import math  
  
def dz_dx(x):  
    return 2 * x * math.cos(x**2)  
  
dz_dx(4)
```

Now what if we got one more level of difficulty. let's say we have:

$$Y = x^2$$

$$Z = \sin(Y)$$

$$U = e^Z$$

& we want $\frac{du}{dx}$

Therefore now for $\frac{du}{dx}$ we need: $\frac{du}{dx} = \frac{du}{dz} \times \frac{dz}{dy} \times \frac{dy}{dx}$

So what we saw above that: As the complexity of nested function increases, the difficulty of finding their derivatives and coding them also increases.

{ Nested func → Complex → derivative → difficult } L → code → difficult }

So we are doing this because nested function and finding their derivatives are very much closely related to deep learning.

Why is this NN bigger with many neurons?

→ Imagine the kind of complexity this issue is resolved by autograd.

→ Looking at a top level view of this NN, we see that this is a nested function, where we are first finding in the feeding it to activation function to get out y_{pred} and then computing the loss based on it.

Different Kinds of Loss Functions

Loss Function Name	Description	Function
Regression Losses		
Mean Squared Error	Measures average squared error in prediction. But is only used for training.	$L_{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2$
Root Mean Squared Error	Measures absolute average bias in prediction. Also called L1 Loss.	$L_{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2}$
Huber Loss	A combination of MSE and MAE. It is parameter-free function.	$L_{Huber} = \begin{cases} \frac{1}{2} (y_i - f(x_i))^2 & : y_i - f(x_i) \leq \delta \\ y_i - f(x_i) - \frac{1}{2}\delta & : \text{otherwise} \end{cases}$
Log-Cosh Loss	Similar to Huber Loss + non-parametric, differentiable, expensive.	$L_{LogCosh} = \frac{1}{N} \sum_{i=1}^N \log[\cosh(f(x_i) - y_i)]$
Classification Losses (Binary + Multi-class)		
Binary Cross Entropy (BCE)	Loss function for binary classification tasks	$L_{BCE} = -\sum_{i=1}^N y_i \cdot \log(p(x_i)) + (1 - y_i) \cdot \log(1 - p(x_i))$
Hinge Loss	Penalizes wrong and right (but not intermediate) predictions. Commonly used in SVMs.	$L_{Hinge} = \max(0, 1 - f(x) \cdot y)$
Cross-Entropy Loss	Extension of BCE loss to multi-class classification tasks	$L_{CE} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{C-1} y_{ij} \log(p_{ij})$
KL Divergence	Minimizes the divergence between predicted and true probability distribution	$L_{KL} = \sum_{i=1}^N y_i \cdot \log(\frac{y_i}{f(x_i)})$

1. Linear Regression: Mean Squared Error (MSE). This can be used with and without regularization, depending on the situation.

2. Logistic regression: Cross-entropy loss or Log Loss, with and without regularization.

3. Decision Tree and Random Forest:

1. Classifier: Gini impurity or information gain.

2. Regressor: Mean Squared Error (MSE)

4. Support Vector Machines (SVMs): Hinge loss. It penalizes both wrong and right (but less confident) predictions. Best suited for creating max-margin classifiers, like in SVMs.

5. k-Nearest Neighbors (kNN): No loss function. kNN is a non-parametric lazy learning algorithm. It works by retrieving instances from the training data, and making predictions based on the k nearest neighbors to the test data instance.

6. Naive Bayes: No loss function. Naive Bayes doesn't have an explicit "loss function" in the same way iterative algorithms do because its parameters are determined directly through frequency counts and conditional probabilities, not through iterative optimization. Instead of a traditional training loop, it yields a simple, direct, closed-form solution for its parameters, effectively skipping the need for an explicit loss function during training.

7. Neural Networks: They can use a variety of loss functions depending on the type of problem. The most common ones are:

1. Regression: Mean Squared Error (MSE).

2. Classification: Cross-Entropy Loss.

8. AdaBoost: Exponential loss function. AdaBoost is an ensemble learning algorithm. It combines multiple weak classifiers to form a strong classifier. In each iteration of the algorithm, AdaBoost assigns weights to the misclassified instances from the previous iteration. Next, it trains a new weak classifier and minimizes the weighted exponential loss.

9. Other Boosting Algorithms:

1. Regression: Mean Squared Error (MSE).

2. Classification: Cross-Entropy Loss.

10. Kmeans: The K-Means loss function, known as Inertia or within-cluster sum of squares (WCSS), is the sum of the squared Euclidean distances between each data point and its assigned cluster's centroid.