

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

$$y = x^2$$

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

→ now very easily we can code this expression

X	$\frac{dY}{dX}$
2	4

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

3	6
4	8
5	10

```
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 we'll do in that case?

i.e. for example:

$$y = x^2$$

$$z = \sin(y)$$

$$\text{we want } \Rightarrow \frac{dz}{dx}$$

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.

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

$$\cos(y) \times (2x)$$

$$2x \cdot \cos(y)$$

↓
x²
↓

$$\frac{dz}{dx} = 2x \cdot \cos(x^2)$$

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' wrt '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 $\frac{dy}{dx}$

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

}

whatever the final expression will get, then we have to code that.

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 fun \rightarrow Complex \rightarrow derivative \rightarrow difficult
 \rightarrow code \rightarrow difficult }

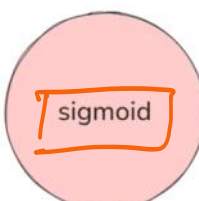
So we are doing this because 'nested function' and finding their derivatives are very much closely related to deep learning.

dataset with 1 feature & label
cgpa placed

cgpa	placed
9.11	1
8.9	1
7	0
6.56	1
4.56	0

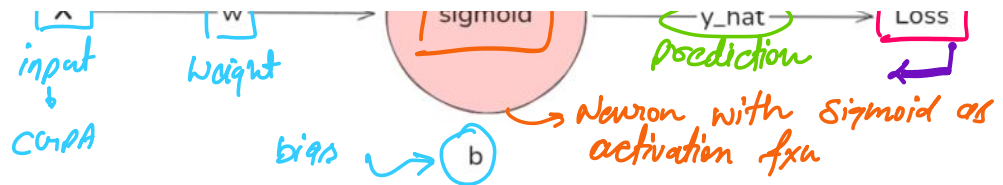
X
input

w
weight



Simplest Possible NN
cuz we used only one neuron

7	0
6.56	1
4.56	0



Training process

1. Forward pass - Compute the output of the network given an input.
2. Calculate loss - Calculate the loss function to quantify the error.
3. Backward pass - Compute gradients of the loss with respect to the parameters.
4. Update gradients - Adjust the parameters using an optimization algorithm (e.g., gradient descent).

Forward Pass Computation

1. Linear Transformation:

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

2. Activation (Sigmoid Function):

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

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

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

After this step we'll come to how much wrong our NN is going to predict the true placed value.

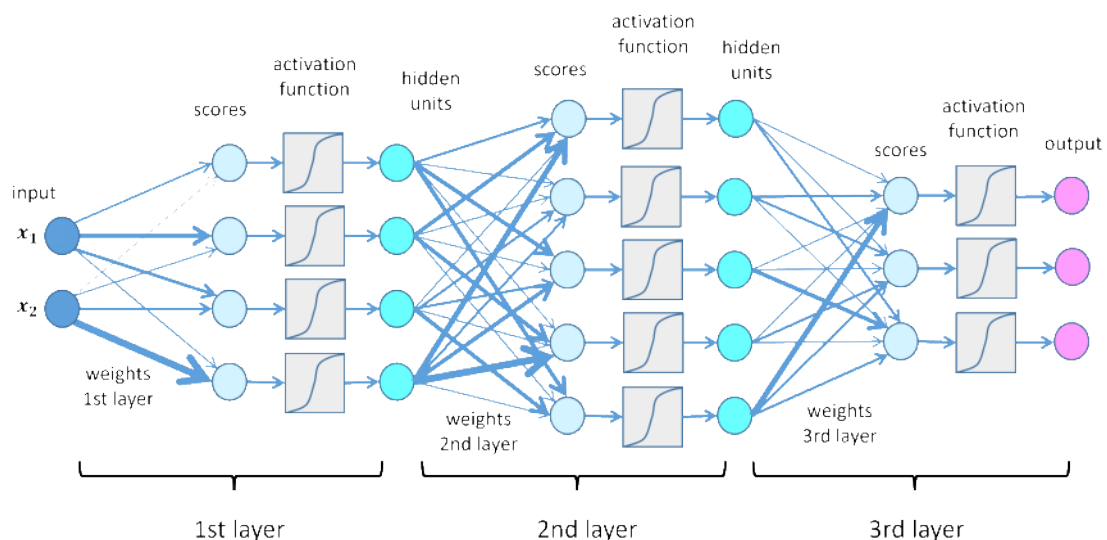
It means now we have to compute $\frac{\partial L}{\partial w}$ & $\frac{\partial L}{\partial b}$

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

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

Again here we have to apply the chain rule to get the derivatives.

Since we have this NN with only 1 neuron.



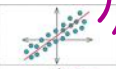


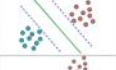

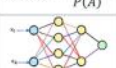



Even if this net is bigger with more

Why if this NN is bigger with many neurons?

→ Imagine the kind of complexity
 this issue is resolved by 'auto-grad'

→ Looking at a top level view of this NN. we see that this is a nested function, where we are first finding 'z' 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			
	Linear Regression	Mean Squared Error	$\mathcal{L}_{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2$
	Logistic Regression	Cross-Entropy Loss	$\mathcal{L}_{MAE} = \frac{1}{N} \sum_{i=1}^N y_i - f(x_i) $
	Decision Tree Classifier	Information Gain or Gini impurity	$\mathcal{L}_{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2$
	Decision Tree Regressor	Mean Squared Error	$\mathcal{L}_{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2}$
	Random Forest Classifier	Information Gain or Gini impurity	
	Random Forest Regressor	Mean Squared Error	
	Support Vector Machines (SVMs)	Hinge Loss	$\mathcal{L}_{Huber} = \begin{cases} \frac{1}{2}(y_i - f(x_i))^2 & : y_i - f(x_i) \leq \delta \\ \delta(y_i - f(x_i) - \frac{1}{2}\delta) & : \text{otherwise} \end{cases}$
	k-Nearest Neighbors	No loss function	$\mathcal{L}_{LogCosh} = \frac{1}{N} \sum_{i=1}^N \log(\cosh(f(x_i) - y_i))$
$P(B A) = \frac{P(B \cap A)}{P(A)}$	Naive Bayes	No loss function	
	Neural Networks	Regression: Mean Squared Error Classification: Cross-Entropy Loss	
	AdaBoost	Exponential loss	
	Gradient Boosting LightGBM CatBoost XGBoost	Regression: Mean Squared Error Classification: Cross-Entropy Loss	
	KMeans Clustering	Within Class Sum of Squares (WCSS)	
Classification Losses (Binary + Multi-class)			
	Binary Cross Entropy (BCE)	Loss function for binary classification tasks.	$\mathcal{L}_{BCE} = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(x_i)) + (1 - y_i) \cdot \log(1 - p(x_i))$
	Hinge Loss	Penalizes wrong and right (but less confident) predictions. Commonly used in SVMs.	$\mathcal{L}_{Hinge} = \max(0, 1 - (f(x) \cdot y))$
	Cross Entropy Loss	Extension of BCE loss to multi-class classification.	$\mathcal{L}_{CE} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(f(x_{ij}))$ N: samples; M: classes
	KL Divergence	Minimizes the divergence between predicted and true probability distribution	$\mathcal{L}_{KL} = \sum_{i=1}^N y_i \cdot \log\left(\frac{y_i}{f(x_i)}\right)$

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:**

- Classifier: Gini impurity or information gain.
- 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 uses the maximum likelihood estimation principle, which, under the assumption of conditional independence, 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.*