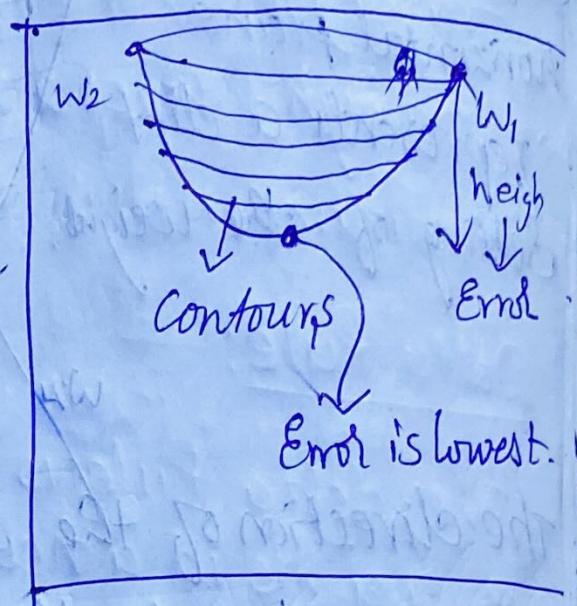


UNIT 1 Chapter 2

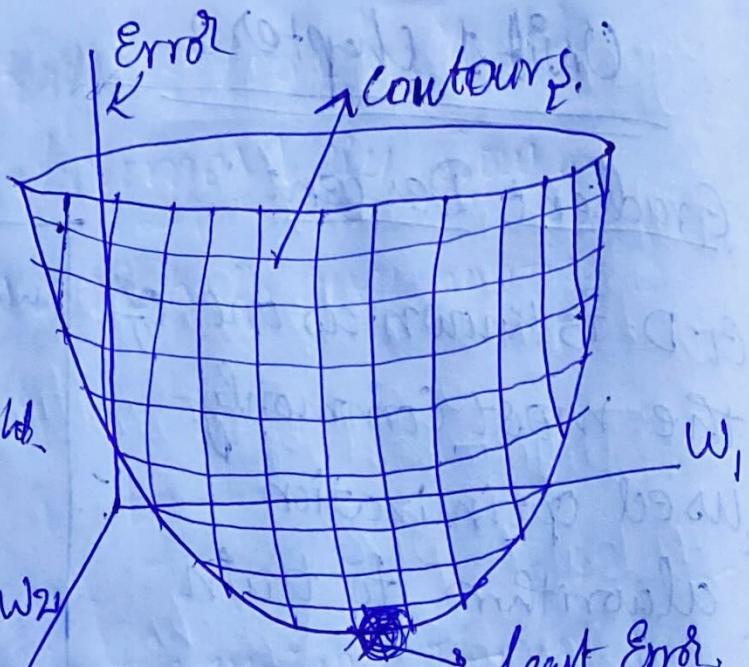
Gradient Descent:

G.D. is known as one of the most commonly used optimization algorithms to train ML models by means of minimizing errors between actual & expected results.



- It helps to find the local minimum of a function.
- The main objective of using a gradient descent algorithm is to minimize the cost function using iteration.
- The 3D space where the horizontal dimensions corresponds to the weights w_1 and w_2 ; and the vertical dimension corresponds to the value of the error function E .
- As shown in figure the circles are elliptical "contours". Where the minimum error is at the center of the ellipse.

* Each point on the horizontal plane represents a cliff setting of the weight.



The direction of the steepest descent is always \perp to these contours.

- We start by initializing the weights (w_1 , w_2) randomly.
- Evaluate the gradient at this starting position to find the direction of steepest descent.
- Take a step in the direction of this gradient.
- After the step, you will be at the new position with a potentially lower error.
- Reevaluate the gradient at this new position and again move in the direction of steepest descent.

→ This iterative evaluating the gradient and taking steps in its direction, we

- Progressively move closer to the point of minimum Errors.
- This Iteration process is the main in gradient descent.

The Descent Rule & Learning Rule :-

- In neural networks, apart from weights, there are other parameters known as hyperparameters that help guide the training process. one such hyperparameter is "Learning Rate".
- The "Learning Rate" controls how much we adjust the weights of our Networks wrt to the gradient.
- When using the gradient descent the learning rate determines how big a step we take at each iteration.
- The steepness of the error surface can indicate how close we are to minimum error.
 - when Surface is flat → we are closer to minimum. Should take smaller steps.
 - Surface is steep → take larger steps.

- If the learning Rate is too small, the training process becomes very slow.
- If too high, there's a risk of overshooting the minimum.

Delta Rule:

The delta Rule is the method of updating the weights in a neural network.

It involves calculating the gradient of the Error function wrt. each weight.

By applying the gradient and the learning rate, the delta rule adjusts each weight to reduce the error.

We know that

$$\Delta w_k = -\eta \frac{\partial E}{\partial w_k} \quad \text{--- (1)}$$

We also know that

$$E = \frac{1}{2} \sum (t_k - \hat{o}_k)^2 \quad \text{--- (2)}$$

(2) in (1)

$$\Delta w_k = -\epsilon \frac{\partial}{\partial w_k} \left[\frac{1}{2} \sum (t_i - o_d)^2 \right]$$

$$= -\epsilon \frac{1}{2} \cancel{\frac{\partial}{\partial}} (t_i - o_d) \frac{\partial}{\partial w} (t_i - o_d)$$

$$= -\epsilon (t_i - o_d) \frac{\partial}{\partial w} \left(\frac{t_i - (w_i \cdot x)}{1} \right) \quad \because o_d = w_i \cdot x$$

$$= -\epsilon (t_i - o_d) (-x)$$

$$\boxed{\Delta w_k = \epsilon x (t_i - o_d)}$$

Applying this method of changing the weights at every iteration, we are finally able to utilize gradient descent.

Gradient Descent with Sigmoidal Neuron

We

Let's recall the mechanism by which logistic neurons compute their output value from their inputs:

$$z = \sum_k w_k x_k$$

$$y = \frac{1}{1 + e^{-z}}$$

The neuron computes the weighted sum of its inputs, the bias "z", It then feeds the bias into the input function to compute y , its final output.

To do so we start by taking the derivatives of the bias w.r.t input & the weight:

$$\nabla_{w_k} E = \frac{\partial E}{\partial w_k}$$

~~$\nabla_b E = \frac{\partial E}{\partial b}$~~

$$E = \frac{1}{2} (f(x) - y)^2$$

$$f(x) = \frac{1}{1 + e^{-(wx+b)}}$$

$$\frac{\partial E}{\partial w_k} = \frac{1}{2} \left[\frac{1}{1 + e^{-(wx+b)}} - y \right]^2$$

$$= \frac{1}{2} \left[2 \times (f(x) - y) * \frac{\partial}{\partial w} (f(x) - y) \right]$$

$$\frac{\partial E}{\partial w_k} = (f(x) - y) * \frac{\partial}{\partial w} (f(x) - y) \rightarrow D$$

$$\therefore \frac{\partial}{\partial w} (f(x) - y)$$

$$= \frac{\partial}{\partial w} \left[\frac{1}{1 + e^{-(wx+b)}} - y \right]$$

$$= \frac{-1}{(1 + e^{-(wx+b)})^2} \frac{\partial}{\partial w} (e^{-(wx+b)})$$

$$= \frac{-1}{(1 + e^{-(wx+b)})^2} * e^{-(wx+b)} \frac{\partial}{\partial w} (-\underbrace{(wx+b)}_{1 \downarrow 0})$$

$$= \frac{-1}{(1 + e^{-(wx+b)})^2} * e^{-(wx+b)} (-x)$$

$$= (+x) \frac{-1}{(1 + e^{-(wx+b)})^2} * e^{-\frac{(wx+b)}{2}}$$

$$= \frac{-1}{1 + e^{-z}} * \frac{e^{-z}}{1 + e^{-z}} (-x) \quad (z = wx + b)$$

$$\frac{\partial f}{\partial w_k} = f(x) * (1 - f(x)) (-x) \rightarrow ②$$

keep ② in ① $f(x) = y$

$$= [f(x) - y] [f(x) * (1 - f(x)) (-x)]$$

$$= -x f(x) (f(x) - y) (1 - f(x)) \frac{\partial y}{\partial w_k}, \Delta w_k = -e^{\frac{\partial E}{\partial w_k}}$$

$$\Delta w_k = -x f(x) (f(x) - y) (1 - f(x))$$

$$w_{t+1} = w_t - \eta \nabla_{w_t}$$

$$\leftarrow w_t - \eta \sum \cancel{e_n} y (1-y) (t-y).$$

$$w_{t+1} = w_t - \eta (f(a)-y) f'(a) (1-f(a)) x$$

As you notice, the new modification rule is just like the delta rule, except with extra multiplicative terms included to account for the logistic component of the Sigmoidal neuron.

Backpropagation algorithm: for reducing Err.

algorithm:

- Step 1: calculate the Error at the Output
- * Start with the output layer of the neural network
 - * for each neuron i in the output layer, calculate the Error (E) as the sum of the squares of the difference between the predicted output (y_i) and the true target (t_i)

$$E = \frac{1}{2} \sum_i (t_j - y_j)^2$$

* compute the derivative of the Error wrt each output neuron's activity:

$$\frac{\partial E}{\partial y_j} = \frac{1}{2} \times 2(t_j - y_j) \frac{\partial}{\partial y_j} (t_j - y_j)$$

$$\frac{\partial E}{\partial y_j} = \frac{1}{2} \times 2(t_j - y_j) \frac{\partial}{\partial y_j} (t_j - y_j)$$

$$\frac{\partial E}{\partial y_j} = -(t_j - y_j)$$

Step 2: Back propagate the Error.

- * Presume that the Error derivatives for layer j are known.
- * To calculate the Error derivatives for the preceding layer i ; accumulate info on how each neuron in layer i influences the neuron in layer j .
- * This can be done as follows using the Partial derivative of the bias wrt incoming output data from the layer beneath is merely the weight of the connection w_{ij} .

$$\frac{\partial E}{\partial y_i} = \frac{\partial E}{\partial z_j} \cdot \frac{\partial z_j}{\partial y_i}$$

$$= w_{ji} \underbrace{\frac{\partial E}{\partial z_j}}_{\textcircled{1}} \left(\frac{\partial z_j}{\partial y_i} = w_{ji} \right)$$

$$\frac{\partial E}{\partial z_j} = \frac{\partial E}{\partial y_i} \underbrace{\frac{\partial y_i}{\partial z_j}}_{\textcircled{2}}$$

$$\frac{\partial E}{\partial z_j} = \frac{\partial E}{\partial y_i} y_i (1-y_i) \rightarrow \textcircled{2}$$

$\textcircled{1} \oplus \textcircled{2}$

$$\frac{\partial E}{\partial y_i} = \sum_j w_{ji} y_i (1-y_i) \frac{\partial E}{\partial z_j}$$

Step 3: update the weights

* After calculating the partial derivatives for all units, update the weights to minimize the error, using the learning rate "n"

rate "n"

$$\Delta w_{ji} = -n \frac{\partial E}{\partial w_{ji}} \text{ --- } \textcircled{3}$$

We know

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial z_j}{\partial w_{ji}} \frac{\partial E}{\partial z_j}$$

From ②

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial E}{\partial y_j} y_j(1-y_j) \times \frac{\partial z_j}{\partial w_{ji}}$$

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial E}{\partial y_j} y_j(1-y_j) \times y_i \quad \left[\because \frac{\partial z_j}{\partial w_{ji}} = y_i \right]$$

~~Now keep the~~

Now for all the training examples, sum up the partial derivatives to adjust the weights:

$$\Delta w_{ji} = -\eta \cdot \frac{\partial E}{\partial y_j} y_j(1-y_j) \times y_i.$$

$$\boxed{\Delta w_{ji} = -\eta \sum_{\text{dataset}} y_i y_j (1-y_j) \frac{\partial E}{\partial y_j}}$$

Step 4: Iterate the process:

Repeat steps 1 to 3 for each layer in the network, moving from the output layer towards the input layer.

Iterate the process across all layers for multiple epochs or until the network converges to a satisfactory level of error.

Stochastic and Minibatch Gradient

after Evaluating each individual training Example.

It is a compromise b/w batch gradient descent & stochastic gradient descent.

It computes the gradient against a subset of the data known as mini batch.

$$\frac{1}{B} \sum_{i=1}^B \nabla J(\theta)$$

minibatch gradient descent

Different techniques to prevent overfitting.

1) Regularization

Overfitting :-

the model becomes very good at predicting or fitting the data it has been trained on but fails to generalize to data that hasn't been trained on.

1) Regularization :-

* Regularization is a technique to prevent overfitting by adding a penalty for larger weights in the objective function.

* The modified objective function becomes:

Error + $\lambda f(\theta)$, where $f(\theta)$ ↑ with larger weights
& λ is the regularization strength, a hyperparameter.

* choosing λ is crucial : $\lambda=0$ means no Regularization while a large λ might prioritize small weights over fitting the training data.

L2 Regularization (Weight Decay) :

The most common regularization in ML, L2 Reg.
Adds $\frac{1}{2} \lambda W^2$ to the Error function for
every weight "W".

It penalizes large weights, Preferring more evenly distributed weight values and Prevents Overfitting by encouraging the whole network to use all inputs rather than a few inputs a lot.

During gradient descent, L2 causes weights to decay linearly towards zero.

3) L1 Regularization :

- * Adds the term $|W|$ for each weight in the network.
- * Leads to Sparsity in the weight vector, making the model less sensitive to noise in the inputs.
- * useful for feature Selection as it indicates which features are most important for predictions.

4) Max Norm Constraints:-

- * Enforce an absolute upper bound on the magnitude of the weight vectors.
- * If the weight exceeds the bound during the gradient descent update, it is projected back onto a predefined bound.
- * Prevents weights from growing too large, ensuring stability in the learning process.

5) Dropout:-

- * A technique where each neuron is kept active with a certain probability P during training, & set to zero otherwise.
- * makes the network less reliant on any single neuron, effectively combining many different neural network architectures.
- * Dropout can prevent overfitting by ensuring the network remains accurate. Even when certain neurons are dropped.

6) Inverted Dropout

- * A variation of dropout where the neuron's output is scaled at training time, rather than test time.
- * Ensure that at test time, the neuron's output do not need to be scaled, maintaining the expected output level.