

Gradient Descent for Neural Networks

This is the heart of supervised learning.

x	y
2	4
3	6
5	10
7	14
9	18

The relationship we may see is that
 $y = 2 * x$.

Similarly, computers use Technique called Gradient Descent to find, $y = x * w_1 - w_2$

$w_1 = \text{weight}$
 $w_2 = \text{bias}$.

* Gradient descent is used in the training of a Neural Network.

So understanding Gradient Descent better with a Binary Classification:-

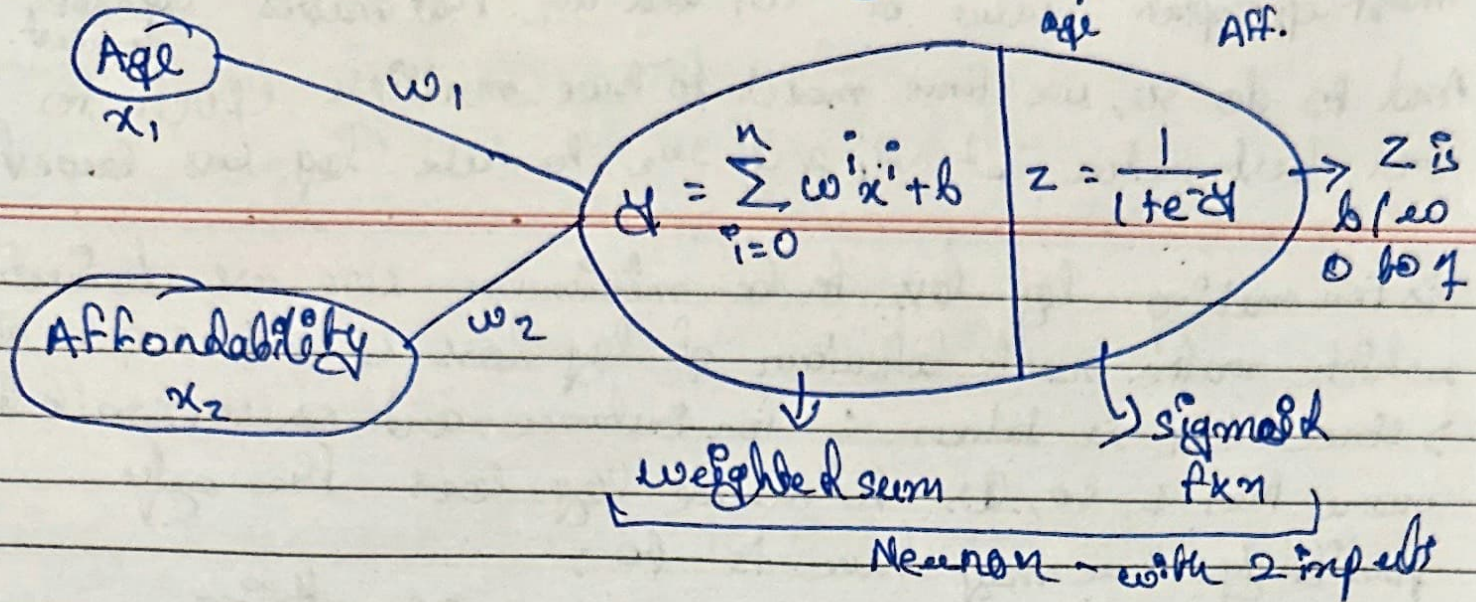
age	affordability	have-insurance
23	1	0
25	0	0
47	1	0
52	0	1
46	1	1
56	1	1
65	0	0
⋮	⋮	⋮
x		y

Based on age and affordability, come up with a fn that can predict if person will buy insurance or not.

Neural Diagram:-

$$y = w_1 * x_1 + w_2 * x_2 + \text{bias}$$

\uparrow
Age
 \uparrow
Aff.



→ Since we are terminally/finally going to make a prediction model hence Gradient Descent is a technique to find patterns withing sample training data to then make a prediction model.

Step-1:- input the training sample data with random w_1 and w_2 value let here be $w_1 = 1, w_2 = 1$

hence, $\&$ weighted sum $= y = w_1 * x_1 + w_2 * x_2 + \text{bias}$ (bias = 0).

Taking age = 22, Aff. = .1, $y = 1 * 22 + 1 * 1 + 0 = 23$

$z \text{ or } \hat{y} = 0.99$ But have insurance = 0

predicted value
Truth value.

Thus there is error:-

In ~~Gradient Descent~~ ^{Logistic Regression} we consider log-loss, to measure loss

Here, error 1 = $-(y \log(\hat{y}) + (1-y) \log(1-\hat{y})) = 4.6$

↳ we find individual errors then total error

Repeating for further samples:- error 2 = 4.6

Total error = error 1 + error 2 + ... + error 13

error 13 = 0.01

Log loss on binary cross entropy $= -\frac{1}{n} \sum_{i=0}^n y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$

Here we get log loss as: 4.31

Now \rightarrow our main goal is to reduce log-loss and to find the most appropriate value of w_1 and w_2 that makes log-loss lowest. And to do so, we force model to have multiple epoch to find best value of w_1 and w_2 to make log-loss lowest.

\rightarrow For making log-loss to be minimum we use derivatives which make mark behaviour of log loss w.r.t. w_1 or w_2 .

\rightarrow Since w_1 is taken 1 for instance and can't be more than 1 so, As to reduce log loss the only possibility we may have is to :-

$$w_1 = w_1 - \text{something}$$

$$\rightarrow w_1 = w_1 - \text{learning rate} * \frac{d}{d w_1} \quad \text{or} \quad \frac{d(\text{learning rate})}{d(w_1)}$$

$$\rightarrow w_1 = w_1 - \frac{d(\text{learning rate})}{d(w_1)}$$

$$w_1 = w_1 - \text{something}$$

$$= w_1 - \text{learning rate} * \frac{d(\text{loss})}{d(w_1)}$$

(usually = 0.01)

$$\text{similarly } w_2 = w_2 - \text{learning rate} * \frac{d(\text{loss})}{d(w_2)}$$

$$b = b - \text{learning rate} * \frac{d(\text{loss})}{d(b)}$$

$$\text{Now!} - \frac{d(\text{loss})}{d(w_1)} = \frac{1}{n} \sum_{i=1}^n x_i (\hat{y}_i - y_i) \quad \left. \vphantom{\frac{1}{n} \sum_{i=1}^n x_i (\hat{y}_i - y_i)} \right\} \text{ same for } w_2$$

$$\text{But for Bias!} - \frac{d(\text{loss})}{d(\text{bias})} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)$$

Through this we get values $\Rightarrow w_1 = 0.8$

Initially values were: $w_1 = 1$
 $w_2 = 1$
 $\text{bias} = 0$ \rightarrow $w_2 = 0.7$
 $\text{bias} = -0.2$

→ Here we are using Batch Gradient Descent.

Now, we will find individual errors then find loss

Remember our first epoch has given loss = 4.31

Now, using new w_1 and w_2 and bias

we calculate new loss = that may be something } second epoch

→ Our main goal is to minimize loss and find appropriate values of w_1, w_2 and bias.

Using TensorFlow keras to find best values the code goes like this:-

```
model = keras.Sequential([  
    keras.layers.Dense(1, input_shape=(2,), activation='sigmoid',  
        kernel_initializer='ones', bias_initializer='zeros')  
])
```

```
model.compile(optimizer='adam',  
              loss='binary_crossentropy',  
              metrics=['accuracy'])
```

```
model.fit(X_train_scaled, y_train, epochs=5000)
```

we initialized w_1, w_2 as kernel_initializer → 'ones'
and " bias as bias_initializer → 'zeros'.
log loss = binary_crossentropy.

→ After training the model with TF keras we may want to know values of w_1, w_2 and bias for the
so, →

coef, intercept = model.get_weights()

we got:-
 $w_1 = 5.060867$
 $w_2 = 1.4086$
bias = -2.9137

making Neural Network without Tensor Flow:-

```
def sigmoid(x):  
    import math  
    return 1 / (1 + math.exp(-x))
```

```
def prediction_fxn (age, affordability):  
    weighted_sum = coef[0] * age + coef[1] * affordability  
    # intercept  
    return sigmoid (weighted_sum)
```

→ we successfully made a single Neuron.

Gradient Descent Function in python from scratch:-

```
def log_loss (y_true, y_pred):  
    epsilon = 1e-15  
    y_pred_new = [max(i, epsilon) for i in y_pred]  
    y_pred_new = [min(i, 1-epsilon) for i in y_pred_new]  
    y_pred_new = np.array(y_pred_new)  
    return -np.mean (y_true * np.log (y_pred_new)  
                    + (1-y_true) * np.log (1-y_pred_new))
```

```
def sigmoid_numpy (x):  
    return 1 / (1 + x)  
    return 1 / (1 + np.exp(-x))
```

sigmoid_numpy (np.array ([12, 0, 1])) → This gives
the computation
of sigmoid in
array form.

def gradient-descent (age, affordability, y-true, epochs
loss-threshold):

here we will be getting $w_1, w_2, bias$.

$$w_1 = w_2 = 1$$

$$bias = 0$$

$$rate = 0.5$$

$$n = \text{len}(\text{age})$$

for i in range (epochs):

$$w_1 * \text{age} + w_2 * \text{aff.} + \text{bias} = \text{weighted_sum}$$

weighted sum

$$y\text{-pred} = \text{sigmoid-numpy}(\text{weighted_sum})$$

sigmoid fn

$$\text{loss} = \text{log-loss}(y\text{-true}, y\text{-pred})$$

$$w_1\text{-derivative} = \left(\frac{1}{n}\right) * \text{np.dot}(\text{np.transpose}(\text{age}), (y\text{-pred} - y\text{-true}))$$

This is same as:- $\frac{d(\text{loss})}{d(w_1)} = \frac{1}{n} \sum_{i=1}^n x_i^1 (y_i^1 - y_i)$

$$w_2\text{-derivative} = \left(\frac{1}{n}\right) * \text{np.dot}(\text{np.transpose}(\text{affordability}), (y\text{-pred} - y\text{-true}))$$

$$\text{bias-derivative} = \text{np.mean}(y\text{-pred} - y\text{-true})$$

This is same as:- $\frac{d(\text{loss})}{d(\text{bias})} = \frac{1}{n} \sum_{i=1}^n (y_i^1 - y_i)$

$$w_1 = w_1 - (\text{rate} * w_1\text{-derivative})$$

$$w_2 = w_2 - (\text{rate} * w_2\text{-derivative})$$

$$\text{bias} = \text{bias} - \text{rate} * \text{bias-derivative}$$

if (loss < loss-threshold)

break

return w_1, w_2, bias

gradient-descent(X-train-scaled['age'], X-train-scaled['affordability'],
y-train, 1000, 0.4631)