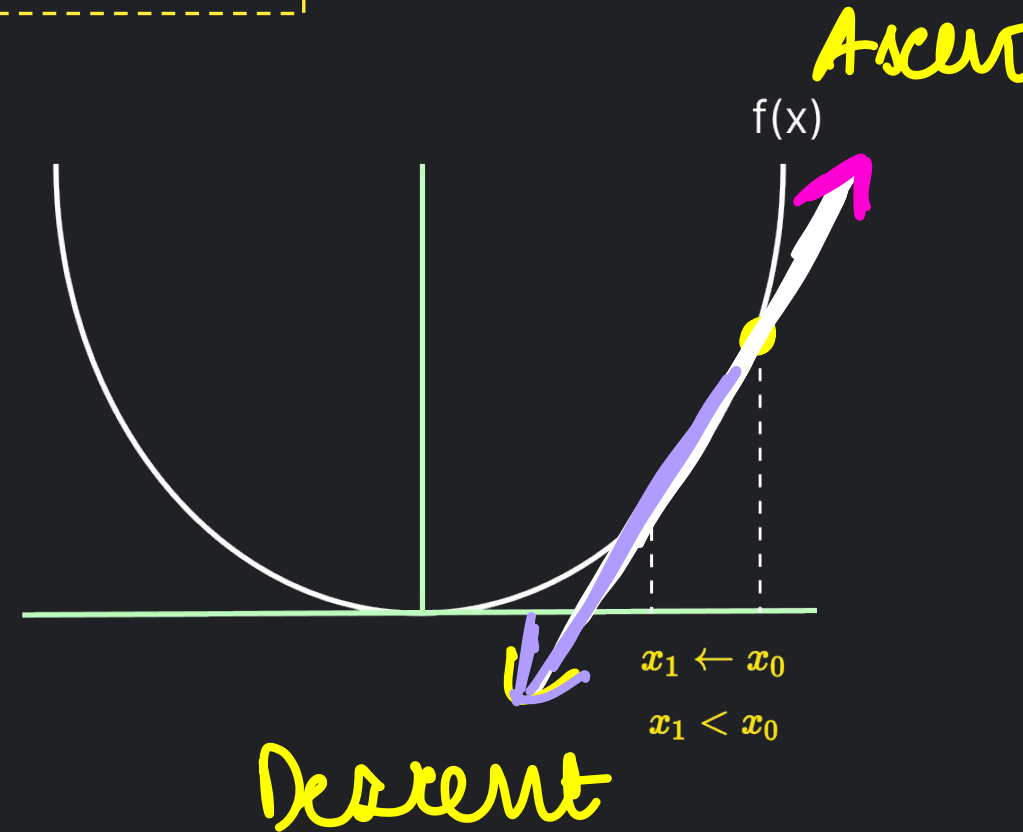


Gradient Descent Revision

https://colab.research.google.com/drive/1KxleE7hhOlx-Zp-hdAdhN_6uGJ3UKwdc?usp=sharing

$$y = x^2$$



① Standing at a random point

② $\frac{dy}{dx} = +ve$
 $x = x_0$

③ Go to towards minima, direction of steepest Descent

$$x_1 \rightarrow x_0 = \eta \frac{dy}{dx} x_0$$

learning rate

gradient

Why GD?

To optimize the loss fn and find its minima

To minimize MSE

$$\sum_{i=0}^m (y^i - \hat{y}^i)^2$$

Quadratic function

i-th sample

$$\sum_{i=0}^m (y^i - (w_0 + w_1 x^i))^2$$

constant

constant

variable

Using GD
find w_0, w_1
s.t that
MSE is
minimized



Linear Regression Helper Functions

https://colab.research.google.com/drive/1KxleE7hhOlx-Zp-hdAdhN_6uCJ3UKwdc?usp=sharing

—>>> We define our linear regression class and set LR & no. of iterations

```
import numpy as np
class LinearRegression():
    def __init__(self, learning_rate=0.01, iterations=5):
        self.learning_rate = learning_rate
        self.iterations = iterations
```

eta

5 steps



—>>> Next we define our predicted Fn

$$w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + \dots + w_dx_d$$

Remember: $\hat{y} = w_t x + w_0$ —> We will call this as bias

```
] def predict(self, X):
    return np.dot(X, self.W) + self.b
```

LinearRegression.predict = predict

$$\begin{matrix} \langle w_1, w_2, \dots, w_d \rangle \\ \langle x_1, \dots, x_d \rangle \end{matrix}$$

Broadcasted m times
→ $(m, 1)$



https://colab.research.google.com/drive/1KxleE7hhOlx-Zp-hdAdhN_6uGJ3UKwdc?usp=sharing

—>>> Finally we evaluate our evaluation metric $R^2 Score$

```
def r2_score(self, X, y):  
    y_ = predict(self,X)  
    ss_res = np.sum((y-y_)**2)  
    ss_tot = np.sum((y- y.mean())**2)  
    score = (1- ss_res/ss_tot)  
    return score
```

```
LinearRegression.score=r2_score
```



EXAMPLE

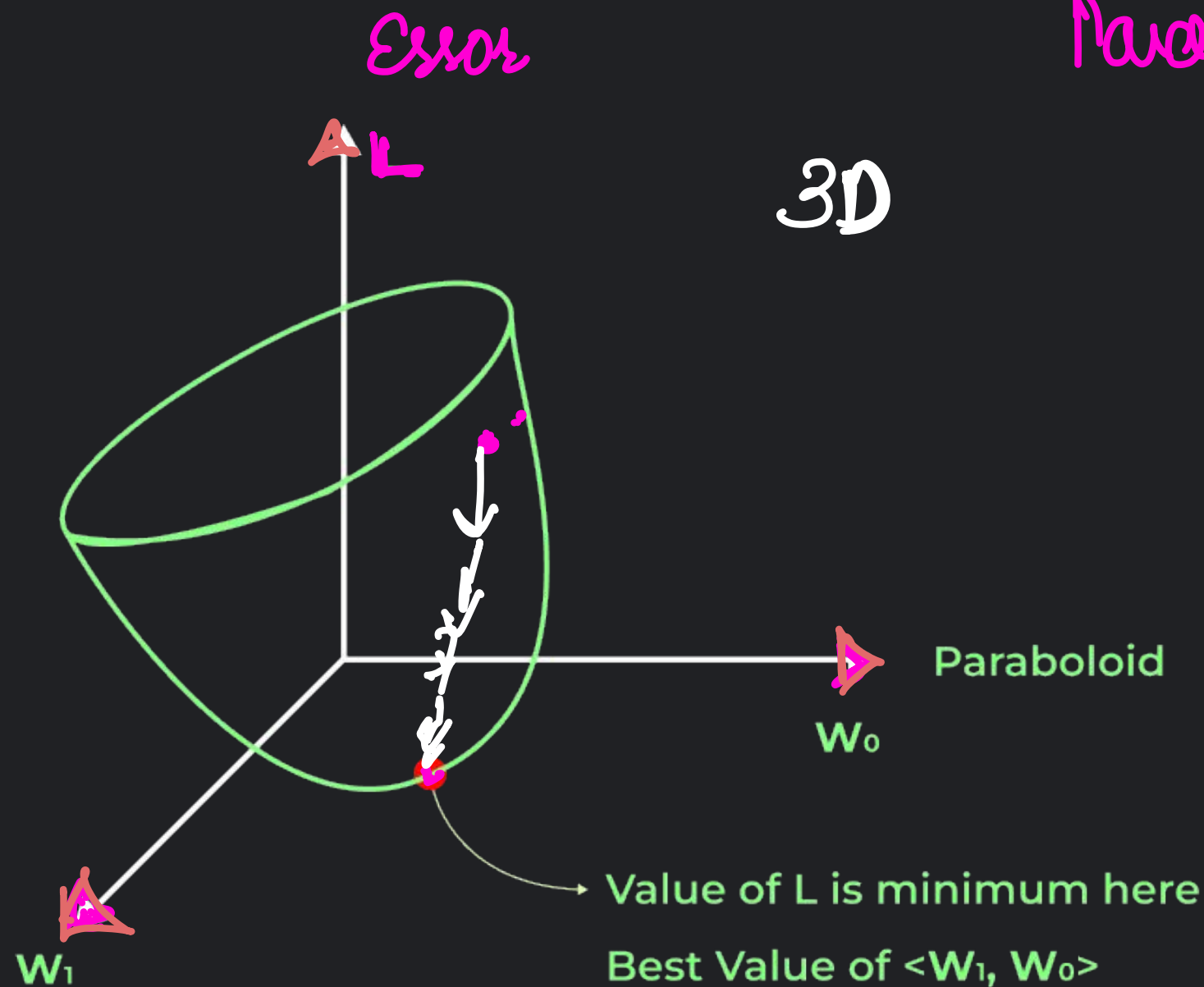
of a single predictor

$$L = \min_{w_0, w_1} \frac{1}{m} \sum_{i=1}^m [y^i - (w_0 + w_1 x^i)]^2$$

Quadratic

Paraboloid

2 parameters
 w_0, w_1



How will L change for multiple predictors ??

$$L = \min_{w_0, w_1, \dots, w_d} \frac{1}{m} \sum_{i=1}^m [y^{(i)} - (w_0 + w_1 x_1^{(i)} + w_2 x_2^{(i)} + \dots + w_d x_d^{(i)})]^2$$

d+1 parameters - $\langle w_0, w_1, w_2, \dots, w_d \rangle$

loss function - d+2 dimensional paraboloid

$$\frac{\partial L}{\partial w_0}$$

$$\frac{\partial L}{\partial w_d}$$

$$\frac{\partial L}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2}$$

How to find global minima?

$$\frac{\partial L}{\partial w_0} = 0$$

$$\frac{\partial L}{\partial w_1} = 0$$

$$\frac{\partial L}{\partial w_2} = 0$$

\vdots

$$\frac{\partial L}{\partial w_d} = 0$$

System of linear
eq.

$$w_1 \rightarrow w_1 - \alpha \frac{\partial L}{\partial w_1}$$

$$w_0 \rightarrow w_0 - \alpha \frac{\partial L}{\partial w_0}$$

$$w_d \rightarrow w_d - \frac{\partial L}{\partial w_d} \alpha$$



② features.

Optimization

Take a partial derivative of L w.r.t each variable

$$\frac{\partial L}{\partial w_0} = \frac{\partial \sum (y^i - \hat{y}^i)^2}{\partial w_0}$$

$$\frac{\partial x^2}{\partial x} = 2x$$

$$= \sum \frac{\partial (y^i - \hat{y}^i)^2}{\partial w_0}$$

$$= \sum 2(y^i - \hat{y}^i) \frac{\partial (y^i - \hat{y}^i)}{\partial w_0}$$

constant

$$= \sum -2(y^i - \hat{y}^i) \frac{\partial \hat{y}^i}{\partial w_0}$$

$$w_2 x_2 + w_1 x_1 + w_0$$

constant

$$= \sum -2(y^i - \hat{y}^i) \frac{\partial (w_2 x_2 + w_1 x_1 + w_0)}{\partial w_0}$$

$$= -\sum 2(y^i - \hat{y}^i) = -2 \cdot \sum \text{error}$$



$$L(w_2, w_1, w_0) = (y - (w_d x_2 + w_1 x_1 + w_0))^2$$

$$\frac{\partial L}{\partial w_1} = \sum -2(y^i - \hat{y}^i) \frac{\partial (w_2 x_2 + w_1 x_1 + w_0)}{\partial w_1}$$

$$= \sum -2(y^i - \hat{y}^i) x_1$$

$$= -2 \sum \text{error} \cdot x_1$$

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_d x_d$$

\downarrow
 w_1
 $-2.e$

\downarrow
 w_2
 $-2.e.x_1$

\downarrow
 w_3
 $-2.e.x_2$

\downarrow
 $-2.e.x_d$

Notice a pattern?

The partial derivative

-2 . Error . Input value of wt.



For m points

$$\frac{\partial L}{\partial w_0} = \frac{1}{m} \sum_{i=1}^m -2(y - \hat{y})$$

$$\frac{\partial L}{\partial w_d} = \frac{1}{m} \sum_{i=1}^m -2(y - \hat{y})x_d$$



How to update the weights ??

$$w_0 \rightarrow w_0 - \alpha \frac{\partial L}{\partial w_0}$$

$$w_1 \rightarrow w_1 - \alpha \frac{\partial L}{\partial w_1}$$

learning rate

$$w_d \rightarrow w_d - \alpha \frac{\partial L}{\partial w_d}$$

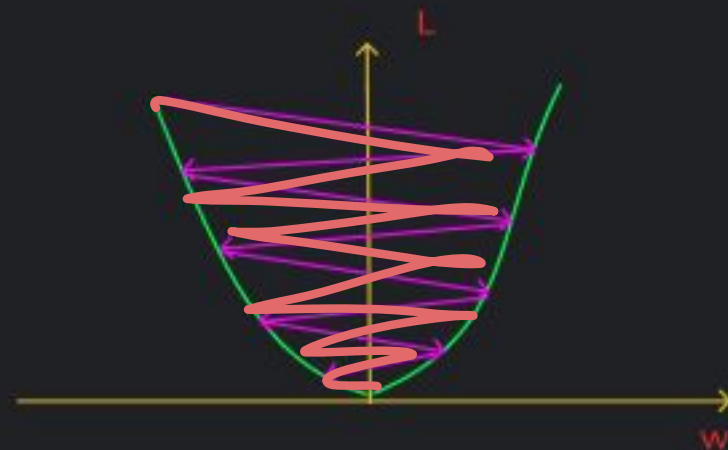
$\alpha \rightarrow$ v. large \rightarrow Big steps \rightarrow sprint (running)

$\alpha \rightarrow$ v. small \rightarrow Baby steps \rightarrow walk slowly

$\alpha \rightarrow$ mid way \rightarrow Right speed \rightarrow Jogging

How to decide the right Linear Regression ??

Run ($d = 1$)



v. large α

Jog ($d = 0.1$)



Just right

Walk ($d = 0.01$)



Takes too long to reach downhill

v. small

Implementation of Wt.s update and GD

Weights update

Figure out dimension matching

```
[ ] def update_weights(self):  
    Y_pred = self.predict( self.X )  
    # calculate gradients  
    dW = - (2*(self.X.T ).dot(self.Y - Y_pred))/self.m  
    db = - 2*np.sum(self.Y - Y_pred)/self.m  
    # update weights  
    self.W = self.W - self.learning_rate * dW  
    self.b = self.b - self.learning_rate * db  
    return self  
  
LinearRegression.update_weights=update_weights
```

Note:

here we: Hyperparameter Tuning

~~Using 1 point of updation is called stochastic gradient descent.~~

~~Whereas if we update the weights after m points, it's called batch gradient descent.~~

