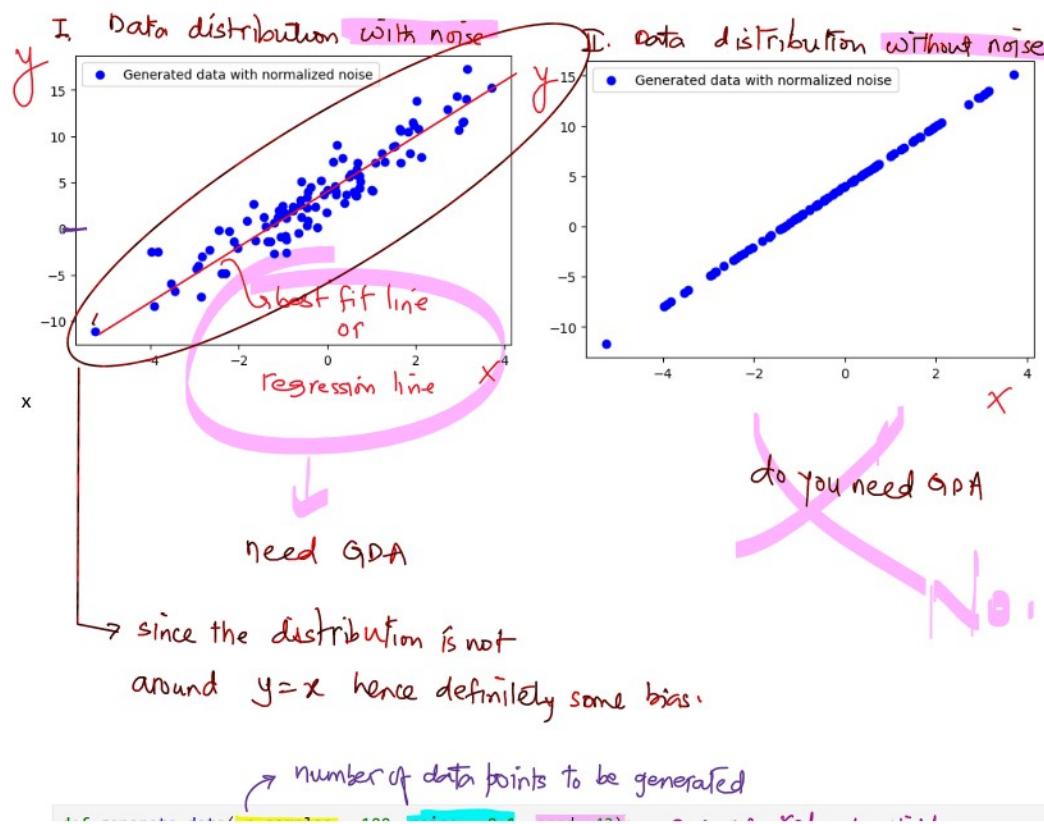
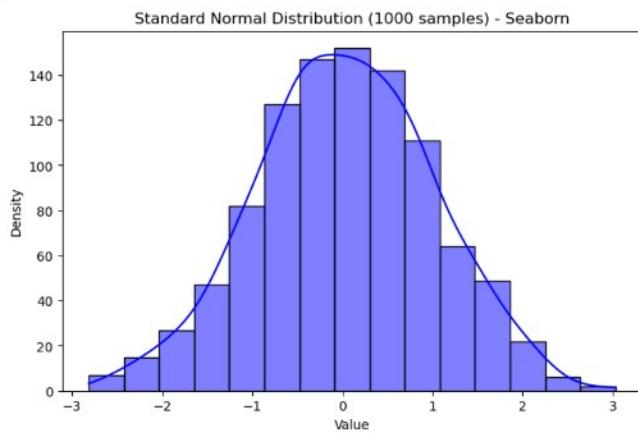
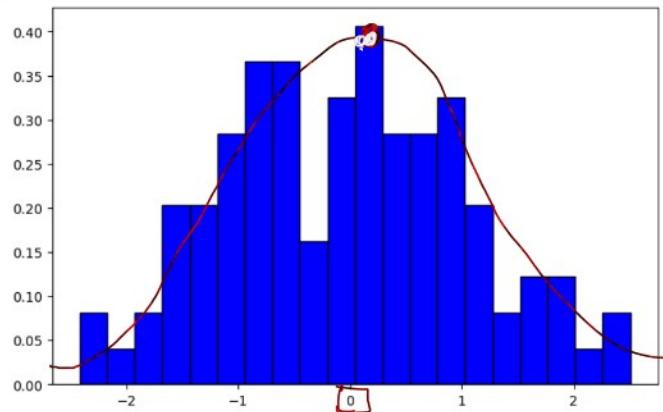


GDA Code Explanation

27 September 2025 11:55



Number of data points to be generated

```

def generate_data( n_samples = 100, noise = 0.1, seed =42):
    """
    Generate random samples (synthetic) linear data: y=4 + 3*X + noise
    """
    np.random.seed(seed) #to ensure the reproducibility of the same random data --> to fix the random numbers generated
    X = 2 * np.random.randn(n_samples, 1) #generates random numbers from `standard normal distribution`
    y = 4 + 3*X + noise*np.random.randn(n_samples, 1) # Creates a random distribution around a Line having some random noise added to it as well
    return X, y

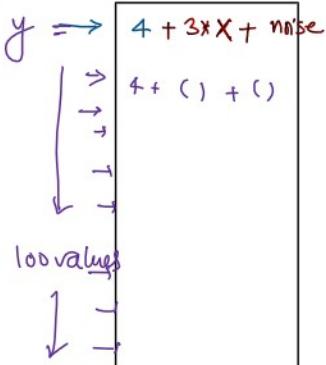
```

ensures Reproducibility
constant noise → 0.1 (default)

$0.1 \times (100 \text{ random})$

2* random std. normal values

values coming from std. normal distribution)



$$y = 4 + 3x$$

↓

intercept coefficient/slope
(bias) (weight)

Random values same

```

X = 2 * np.random.randn(n_samples, 1) #generates random numbers from `standard normal distribution`
y = 4 + 3*X + noise*np.random.randn(n_samples, 1) # Creates a random distribution around a Line having some random noise added to it as well
return X, y

```

generate_data(n_samples = 10, noise = 0, seed =42)

Random (same)

$X = \begin{bmatrix} 0.99342831 \\ -0.2765286 \\ 1.29537708 \\ 3.04605971 \\ -0.46830675 \\ -0.46827391 \\ 3.15842563 \\ 1.53486946 \\ -0.93894877 \\ 1.08512009 \\ 6.98028492 \end{bmatrix}$

$y = \begin{bmatrix} y_1 = 4 + 3 \cdot 0.99342831 + 0 \\ \dots \\ y_{10} = 4 + 3 \cdot 6.98028492 + 0 \end{bmatrix}$

$4 + 3 \cdot 0.99342831 = 6.98028493$

Mean Squared Error - Cost Function for regression problems

For m training examples:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

↑ weight ↑ bias
cost function

error ① Mean Squared Error
③ ← ② ← ①

Predicted actual
 $(P-A)^2$

where:

- $y^{(i)}$ = actual output for sample i .
- $\hat{y}^{(i)}$ = predicted output.
- w, b = parameters (weights, bias).
- m = number of samples.
- $i \rightarrow i$ th row / sample

why square ??

→ to prevent positive and negative errors cancelling / nullifying each other

→ squaring makes the error +ve

→ penalizes large errors more strongly than small errors

Division by m gives the mean avg. making it independent of dataset size
↓
to get the avg. model error.

Factor $\frac{1}{2}$ is to simplify the derivative output (2 cancels when differentiating)

$$\begin{aligned} y &= xe^2 & y &= \frac{1}{2}x^2 \\ \frac{dy}{dx} &= 2x & \frac{dy}{dx} &= \frac{1}{2}(2x) = x. \end{aligned}$$

Task

Do the below derivations:

$$\frac{\partial J}{\partial w} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \cdot x^{(i)} \quad \checkmark$$

$$\frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) \quad \checkmark$$

gradient

[Cost Function] → MSE

Linear Function → Model

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \rightarrow$$

(Predicted)

$$\hat{y} = \hat{\beta}_0 x^0 + \hat{\beta}_1 x$$

$$\hat{y} = [\hat{\beta}_0 \quad \hat{\beta}_1]_{1 \times 2} \times [1 \quad x]_{2 \times 1}$$

$$\hat{y} = X \cdot \theta$$

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

$$[\hat{\beta}_0 \quad \hat{\beta}_1 \quad \hat{\beta}_2 \dots] \rightarrow \theta \rightarrow \text{parameters}$$

```

def compute_cost(X, y, theta):
    """
    Compute the mean squared error cost function
    """
    m = len(y) #no. of rows in the data
    return np.sum((X.dot(theta) - y)**2)/(2*m)

```

Linear Algebra Videos:

<https://www.khanacademy.org/math/linear-algebra>

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \underbrace{(h_\theta(x^{(i)}) - y^{(i)})^2}_{\text{A}}$$

$$h_\theta(x) = \theta^T x = \hat{\beta}_0 + \hat{\beta}_1 x_1$$

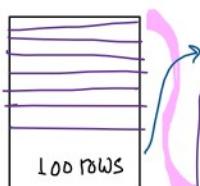
$$[\hat{\beta}_0 \quad \hat{\beta}_1]$$

$y = f(x)$

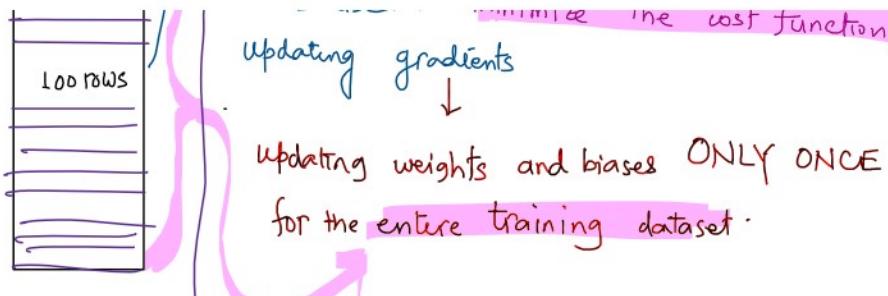
BATCH GRADIENT DESCENT ALGORITHM (BGD)

↳ (Vanilla Gradient Descent)

↳ by default → BGD



Batch Gradient Descent is an optimization algorithm that is used to minimize the cost function, updating gradients.



→ refers to the fact that the gradient is computed using the entire training dataset at each iteration (ONE EPOCH)

What is an Epoch?

An epoch → one complete pass through the entire training dataset by the model.

During one epoch, every training sample has been used once to

Update the model's parameters

Weights Biases

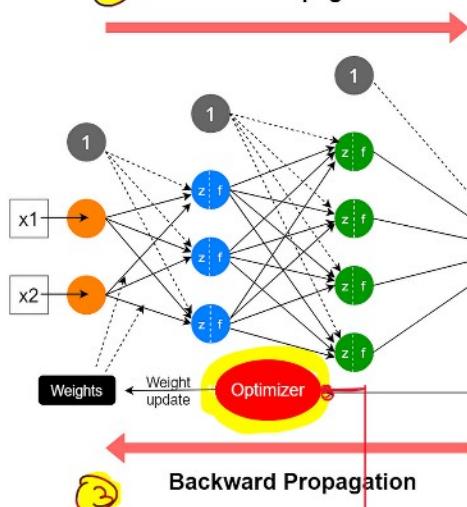


Pro-tip

Batch Gradient Descent: → uses all training samples / rows → the entire training dataset to compute the gradient of the cost function (loss)

①

Forward Propagation



Iterative process until loss function is minimized

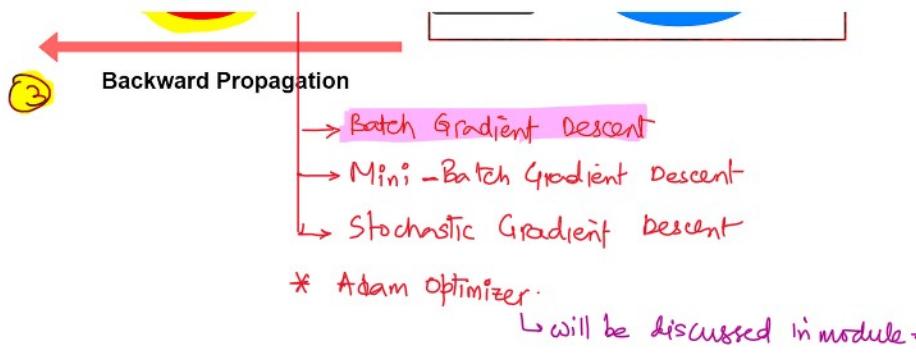
Predictions (y')

True Values (y)

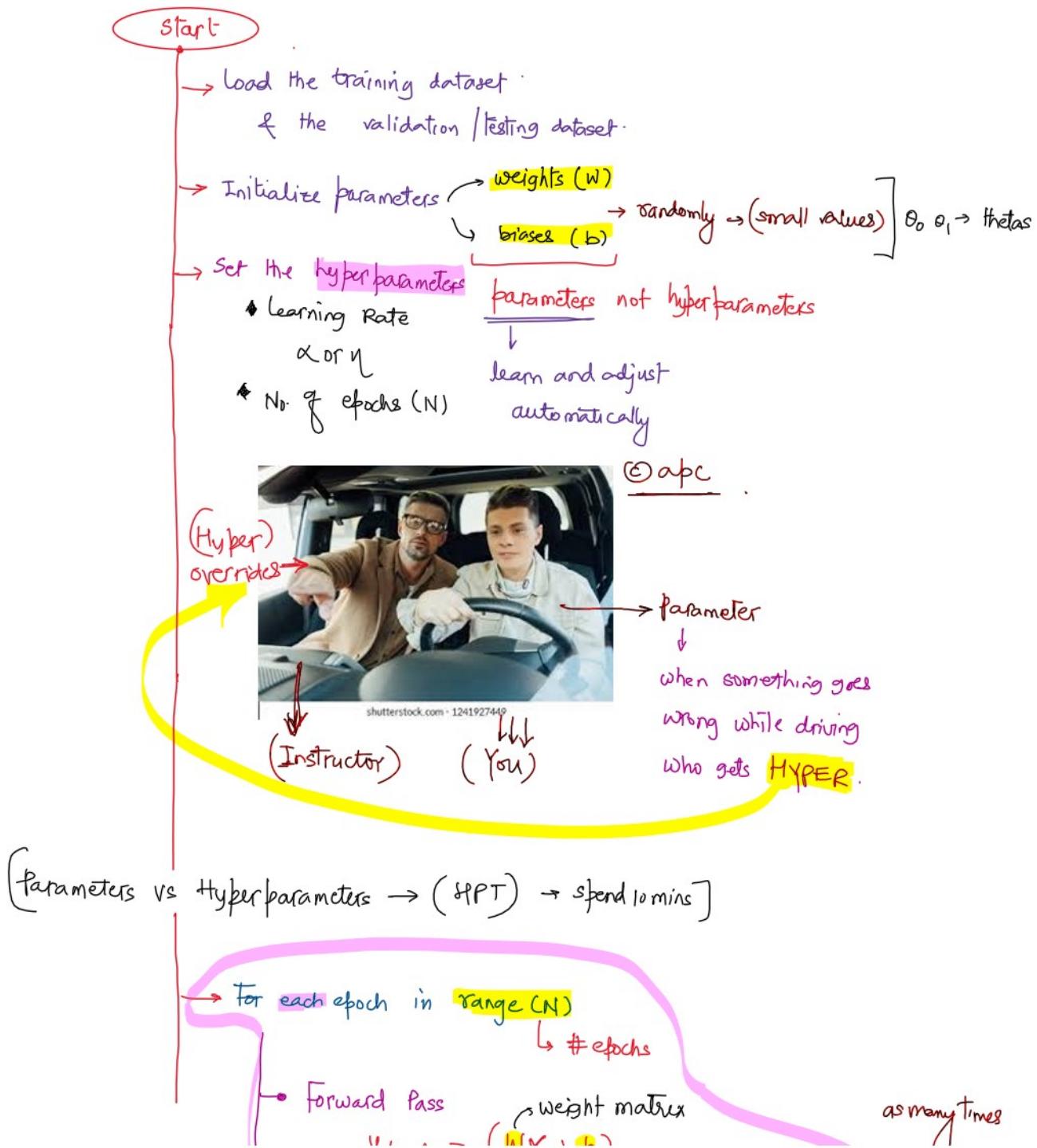
Model Evaluation

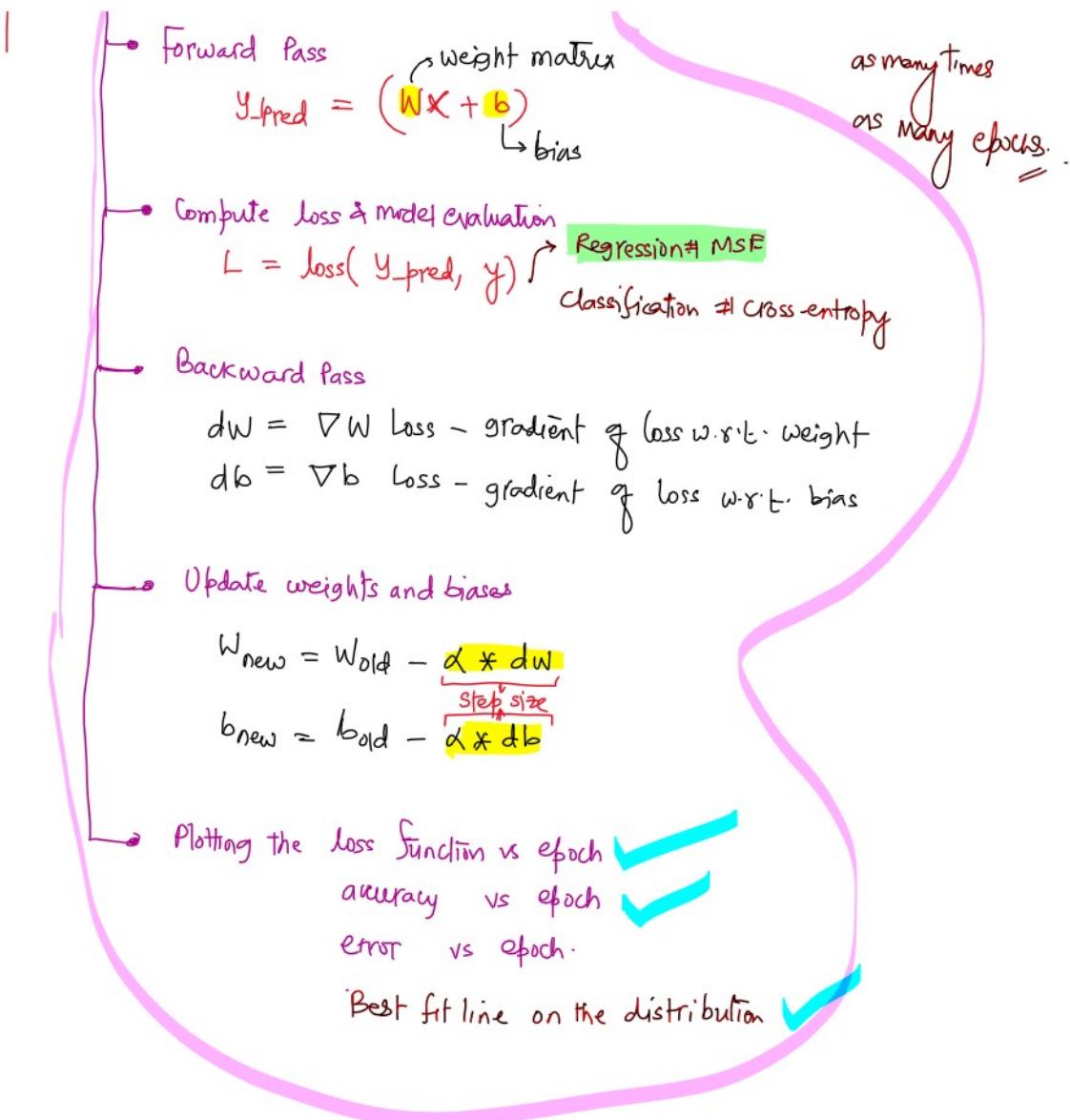
②

Loss Score ← Loss Function



Flowchart : Batch Gradient Descent (from scratch)



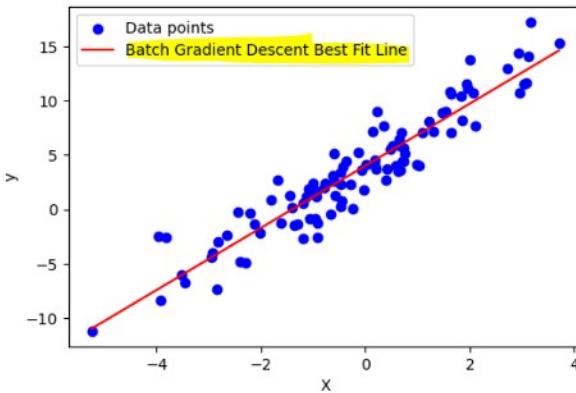
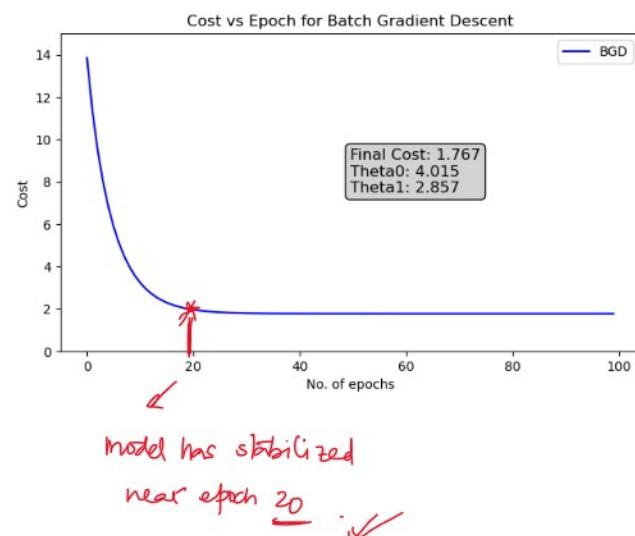
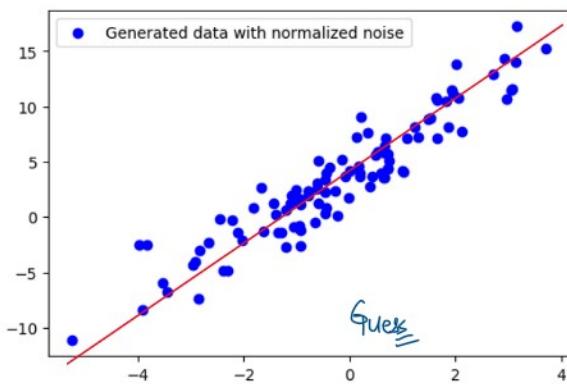


$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

```

: def compute_cost(X, y, theta):
    """
    Compute the mean squared error cost function
    """
    m = len(y) #no. of rows in the data
    return np.sum((X.dot(theta) - y)**2)/(2*m)

```



<https://github.com/scikit-learn/scikit-learn/tree/main/sklearn>

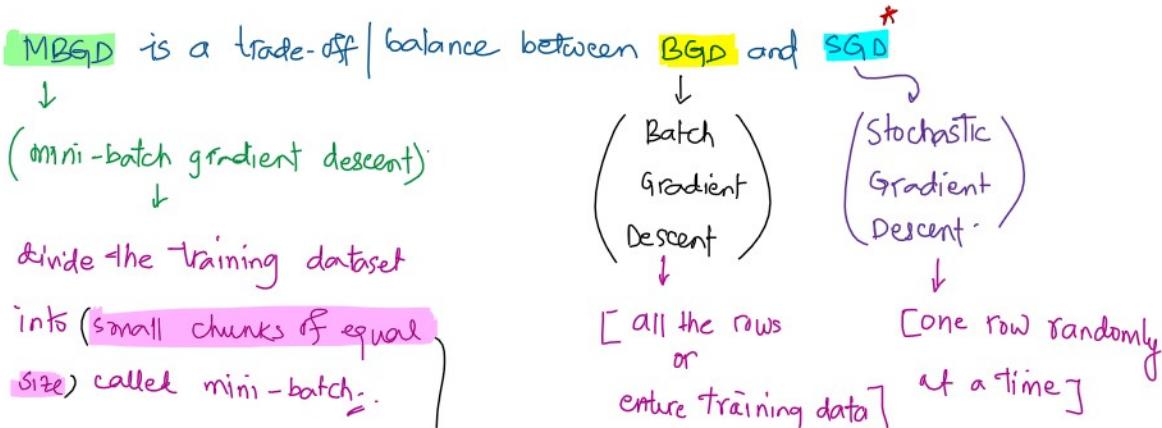
https://github.com/scikit-learn/scikit-learn/blob/main/sklearn/neural_network/_stochastic_optimizers.py

TASK: Create the documentation for the BGD concept

COOKBOOK

MODELING BOOK - *shared from my end*

MINI-BATCH GRADIENT DESCENT (MBGD)



size) called mini-batch.]
 or
 entire training data] at a time]
 (mini-batches)

Standard size: $n = 32$ rows | training samples or examples.
 (one mini-batch)

For ex: Training dataset size = 100 rows
 std. batch size: $n = 32$

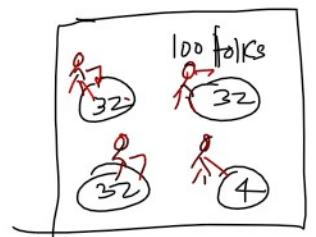
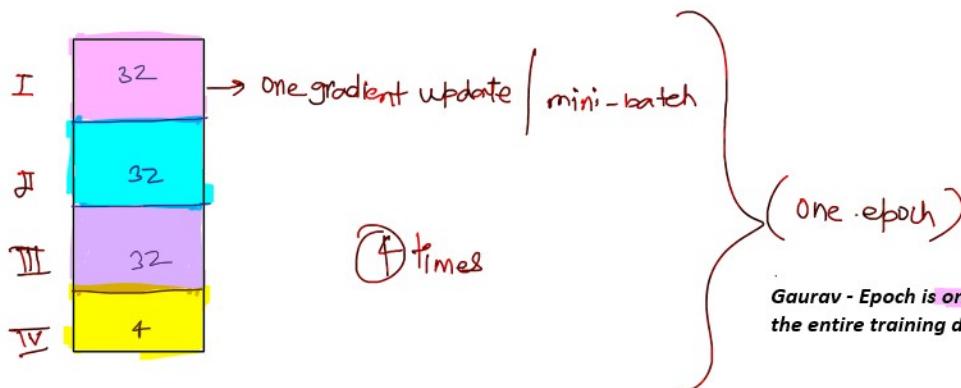
How many mini-batches??

$$\frac{100}{32} = \frac{100}{32} = 3.125$$

↓
take ceiling function

$$\left[\frac{100}{32} \right] = \left[3.125 \right] = 4 \rightarrow 4 \text{ mini-batches}$$

In one epoch, how many gradient updates would happen for the above training dataset.



Gaurav - Epoch is one full traverse of whole data or one full pass thru the entire training data

STOCHASTIC GRADIENT DESCENT (SGD)

refers to systems or processes that are random or probabilistic in nature

random or probabilistic in nature

↳ (uncertainty)

To predict: 1) In Bangalore, what's the chance of rain today evening?

(58% - probability of rainfall)

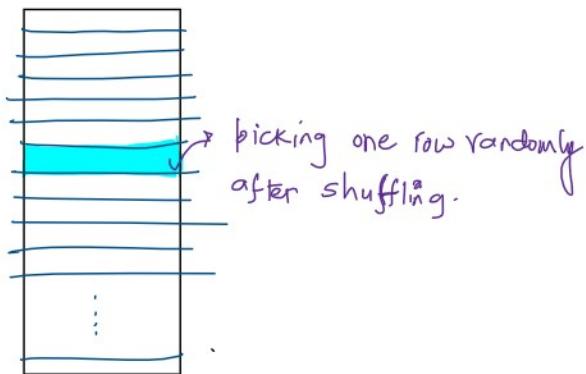
IMD: Indian Meteorological Deptt.

2) What's going to be traffic on ORR route around 9AM tomorrow?

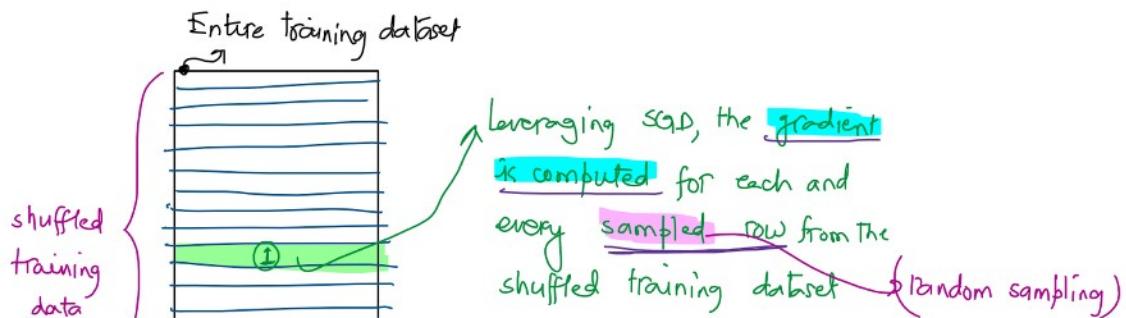
3) opening price of a company's stock everyday?

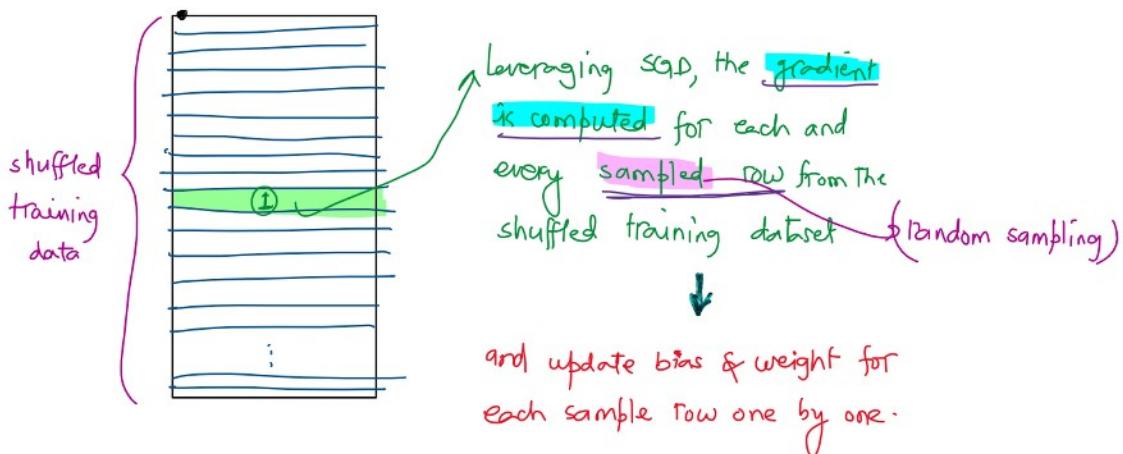
Uncertain

In ML/DL, stochastic describes algorithms or models that incorporate randomness in their operations.



SGD: It is an optimization technique (algorithm) used to minimize the loss function. To do so, SGD computes the gradient for a single row, randomly chosen from the shuffled training dataset at each iteration





Training dataset: $m = 100 \rightarrow$ 100 rows in the training dataset

Epochs = 100 ✓

$n = 32$ (batch size)

How many gradient updates or iteration would happen in:

BGD : 100 updates $\xrightarrow{\times 100} 1$ / epoch

MBGD: 400 updates $\xrightarrow{\times 100} 4$ / epoch

SGD : 10,000 updates $\xrightarrow{\times 100} 100$ / epoch

Conclusion

Training Rows = 100		
BGD	1 gradient update / epoch	$1 \times 100 = 100$ gradient updates
MBGD *	4 gradient updates / epoch	$4 \times 100 = 400$ gradient updates
SGD	100 gradient updates / epoch	$100 \times 100 = 10,000$ gradient updates

*: No. of training rows = 100

std. Batch size = 32: [default standard batch size]
decided by top DL researchers

- For linear regression problem, we use mean-squared error (MSE)

```
[1]: def compute_cost(X, y, theta):  
    ....  
    Compute the mean squared error cost function
```

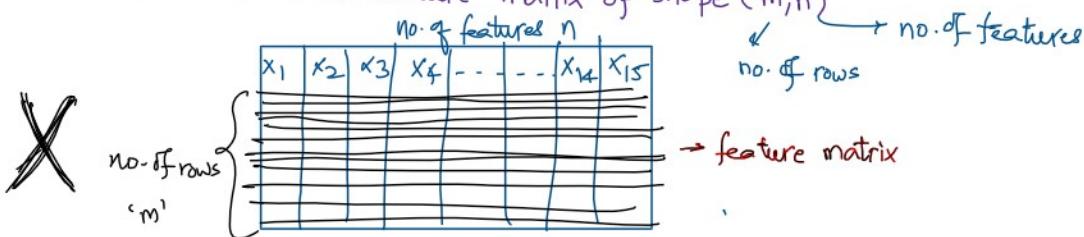
$$\frac{\sum_i^m (\text{Prediction} - \text{Actual})^2}{2^m}$$

In BGD

$$\hat{y} = x \cdot \theta$$

where 'X' is the feature matrix of shape (m, n)

(m, n) → no. of features
no. of rows



Θ : is the parameter array
(weights & biases)

[Cost Function] → MSE

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \rightarrow$$

(Predicted)

$$\hat{y} = \hat{\beta}_0 \cdot x^0 + \hat{\beta}_1 x$$

$$\hat{y} = [\hat{\beta}_0 \quad \hat{\beta}_1]_{1 \times 2} x \begin{bmatrix} 1 \\ x \end{bmatrix}_{2 \times 1}$$

$$\hat{y} = X \cdot \Theta$$

Gradient of cost Function

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (x \cdot \theta - y)^2$$

Let us find the gradient of cost function

$$\frac{\partial}{\partial \theta} J(\theta) = \frac{\partial}{\partial \theta} \left[\frac{1}{2m} \sum_{i=1}^m (x \cdot \theta - y)^2 \right]$$

$$= \frac{1}{2m} \cancel{x} \sum (x \cdot \theta - y) \cdot \frac{\partial}{\partial \theta} (x \cdot \theta - y)$$

$$= \frac{1}{m} \sum (x \cdot \theta - y) \cdot x$$

$$\frac{\partial}{\partial \theta} J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m x \cdot (x \cdot \theta - y) \right]$$

$$\frac{\partial}{\partial x} \frac{(x \cdot \theta - y)^2}{x}$$

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

$$= 2(x \cdot \theta - y) \cdot \frac{\partial}{\partial x} (x \cdot \theta - y)$$

$$= 2(x \cdot \theta - y) \cdot \theta$$

$$= 2\theta(x \cdot \theta - y)$$

$$= \frac{d}{dx} (k \cdot x - y)$$

$$= \frac{d}{dx} kx - \frac{d}{dx}(y)$$

= Q

AI generated

Derivation of Gradient of Cost Function (Linear Regression)

We start with the Mean Squared Error (MSE) cost function:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (x^{(i)}\theta - y^{(i)})^2$$

ith row | sample

Step 1: Differentiate with respect to θ

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{\partial}{\partial \theta} \left[\frac{1}{2m} \sum_{i=1}^m (x^{(i)}\theta - y^{(i)})^2 \right] \quad (1)$$

Step 2: Apply the power rule

$$= \frac{1}{2m} \sum_{i=1}^m 2(x^{(i)}\theta - y^{(i)}) \cdot \frac{\partial(x^{(i)}\theta - y^{(i)})}{\partial \theta} \quad (2)$$

Step 3: Simplify (the 2 cancels with 1/2)

$$= \frac{1}{m} \sum_{i=1}^m (x^{(i)}\theta - y^{(i)}) \cdot x^{(i)} \quad (3)$$

Step 4: Write in vectorized form

$$\nabla_{\theta} J(\theta) = \frac{1}{m} X^T (X\theta - y) \quad (4)$$

code snippet

```
# Compute gradient
gradients = (X.T.dot(X.dot(theta) - y)) / m
# Transpose
```

```
def bgd(X, y, theta, learning_rate=0.01, epochs=100):
    """
    Batch Gradient Descent (Vanilla) using the entire training dataset
    X = Array of X with the added bias (X_b)
    y = Vector of y
    theta: Array of weight & bias parameters randomly assigned
    learning_rate: alpha value set to default 0.01
    epochs: number of times model will run through the entire training dataset
    """
    m = len(y) # number of training rows

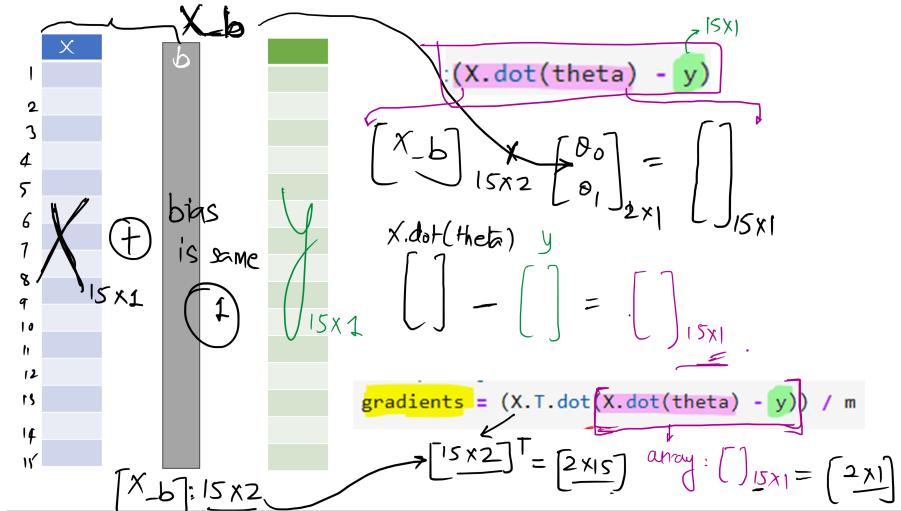
    # Arrays to track cost and theta history
    cost_history = np.zeros(epochs) # 1D array to store cost after each epoch
    theta_history = np.zeros((epochs, theta.shape[0])) # 2D array to store parameter values

    for epoch in range(epochs):
        # Compute gradient
        gradients = (X.T.dot(X.dot(theta) - y)) / m
```

(Prediction - Actual) = Error

$$\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = J(\theta)$$

$$\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



```

100
for epoch in range(epochs):
    # Compute gradient
    gradients = (X.T.dot(X.dot(theta) - y)) / m → calculate gradient for each epoch
    # Update parameters
    theta = theta - (learning_rate * gradients) → update θ value → weight
                                                → bias
    # Compute cost
    cost = compute_cost(X, y, theta) → compute cost

    # Store history
    cost_history[epoch] = cost
    theta_history[epoch, :] = theta.T

return theta, cost_history, theta_history

```