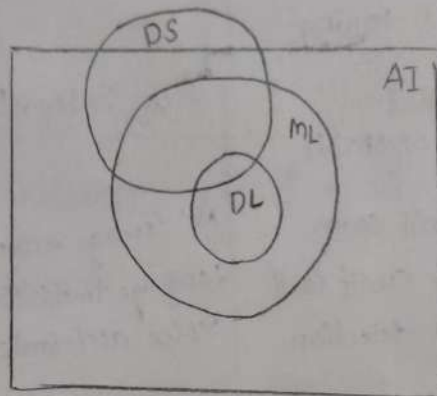


Deep Learning Assignment-1

22ATIABISI

CAI-A

1. Differentiate AI vs ML vs DL vs DS



VEN DIAGRAM.

Artificial Intelligence	Machine Learning	Deep Learning	Data science.
→ Creating machines that can perform tasks requiring human intelligence	A subset of AI where machines learn from data without explicit programming	A subset of ML that uses deep neural networks to process data.	A multidisciplinary field that focuses on extracting insights from data.
→ Broad (includes ML, DL, rule based system & more)	Narrower than AI, focus on learning from data	Very narrow specifically uses neural networks	Broad (includes AI, ML, statistics & big data)
→ Rule based Systems, Expert Systems, Search algorithms	Supervised and unsupervised & reinforcement learning.	Convolutional Neural Networks, Recurrent Neural Networks, Transformers.	Statistical analysis, data mining, visualization, ML and DL
→ Data dependency is Low-Medium	→ Data dependency is high (requires lot of data to learn)	Data dependency is very high (requires large labelled dataset)	Data dependency is High (works with structured & unstructured data)
→ Low to High computational power	→ Medium to high computational power	very high computation power	medium to high computation power.

May required predefined rules.	Requires human intervention for feature selection & tuning	Minimal Minimal intervention, learns features automatically	Requires human expertise in data preprocessing & analysis.
Interpretability can be explainable	Somewhat interpretable	Low Interpretability	Medium interpretability
Ex: chatbots, self driving cars, recommendation systems	Ex: Email spam filters, credit card fraud detection	Ex: Image recognition, language translation, voice assistants	Predictive analysis, customer segmentation, market trend analysis.

2. Explain principal component analysis with an example using Decomposition method.

Principal Component Analysis (PCA)

PCA technique was introduced by mathematician Karl Pearson in 1901. It works on the condition that why the data in a bipolar dimensional space is mapped to data in a lower dimensional space the variance of data in the lower dimensional space should be maximum.

- PCA is a statistical procedure that uses an orthogonal transformation that converts a set of correlated variables, to a set of uncorrelated variables.
- PCA is an unsupervised ML algorithm technique used to examine the inter relations among a set of variables.
- Main goal of PCA is to reduce the dimensionality of a dataset by preserving the most important patterns or relationship. For example given the following dataset use PCA to reduce dimensions from 2 to 1.

2
Data set:

feature	x_1	x_2	x_3	x_4
x	4	8	13	7
y	11	4	5	14

Step 1: The Given Dataset, No. of features $\Rightarrow n = 2$
No. of samples $N = 4$

Step 2: computation of mean of variable, $\bar{x} = 8$
 $\bar{y} = 8.5$

Step 3: computation of covariance matrix,
 $(x, x) (x, y) (y, x) (y, y)$

$$\begin{aligned} \text{(i) covariance of } (x, x) &= \frac{1}{N-1} \sum_{k=1}^N (x_k - \bar{x})^2 \\ &= \frac{1}{4-1} [4-8]^2 + [8-8]^2 + [13-8]^2 + [7-8]^2 \\ &= \frac{1}{3} [16 + 0 + 25 + 1] \\ &= \frac{1}{3} [42] \\ &= \frac{14}{1} \end{aligned}$$

$$\begin{aligned} \text{(ii) covariance of } (x, y) &= \frac{1}{N-1} \sum_{k=1}^N (x_k - \bar{x})(y_k - \bar{y}) \\ &= \frac{1}{4-1} [(4-8)(11-8.5) + (8-8)(4-8.5) + (13-8)(5-8.5) \\ &\quad + (7-8)(14-8.5)] \\ &= \frac{1}{3} [-10 + 0 + (-17.5) + (-5.5)] \\ &= \frac{1}{3} [-10 - 17.5 - 5.5] = -11 \end{aligned}$$

$$(iii) \text{ covariance } (y, x) = -11$$

$$(iv) \text{ covariance } (y, y) = \frac{1}{N-1} \sum_{k=1}^N (y_i - \bar{y})^2$$

$$= \frac{1}{4-1} [11-8.5]^2 + [4-8.5]^2 + [5-8.5]^2 + [14-8.5]^2$$

$$= \frac{1}{3} [69]$$

$$= 23$$

Covariance Matrix :

$$S = \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix} = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

Step 4: To FIND EIGEN values and EIGEN vector and Normalized Eigen Value.

$$\begin{bmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{bmatrix} = \begin{aligned} & (14-\lambda)(23-\lambda) - 121 = 0 \\ & 322 - 14\lambda - 23\lambda + \lambda^2 - 121 = 0 \\ & \lambda^2 - 37\lambda + 201 = 0 \end{aligned}$$

$$\lambda_1 = 30.3849 \quad \lambda_2 = 6.6151$$

$$\lambda_1 = 30.3849$$

$$\begin{bmatrix} 14-30.3849 & -11 \\ -11 & 23-30.3849 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0 \quad \times$$

$$\begin{bmatrix} -16.3849 & -11 \\ -11 & -7.3849 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0 \quad \times$$

$$-16.3849x - 11y = 0$$

$$\lambda_1 = 30.3849$$

$$\begin{bmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = 0$$

$$(14-\lambda)u_1 - 11u_2 = 0$$

$$-11u_1 + (23-\lambda)u_2 = 0$$

$$(14-\lambda)u_1 = 11u_2$$

$$\frac{u_1}{11} = \frac{u_2}{14-\lambda}$$

$$\frac{u_1}{11} = \frac{u_2}{14-30.3849}$$

$$\begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} 11 \\ -16.3849 \end{bmatrix}$$

Normalize vector,

$$e_1 = \begin{bmatrix} \frac{11}{\sqrt{(11)^2 + (-16.3849)^2}} \\ \frac{-16.3849}{\sqrt{(11)^2 + (-16.3849)^2}} \end{bmatrix}$$

$$e_1 = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$$

$$\lambda_2 = 6.6151$$

$$\frac{u_1}{11} = \frac{u_2}{14-6.6151}$$

$$\frac{u_1}{11} = \frac{u_2}{8.6151}$$

$$\begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} 11 \\ 8.6151 \end{bmatrix}$$

$$e_2 = \begin{bmatrix} \frac{11}{\sqrt{(11)^2 + (8.6151)^2}} \\ \frac{8.6151}{\sqrt{(11)^2 + (8.6151)^2}} \end{bmatrix}$$

$$e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$$

Step 5: The New dataset.

First principle component	P_{11}	P_{12}	P_{13}	P_{14}
	-4.3052	3.7361	5.6923	-5.1235

we find P_{11}

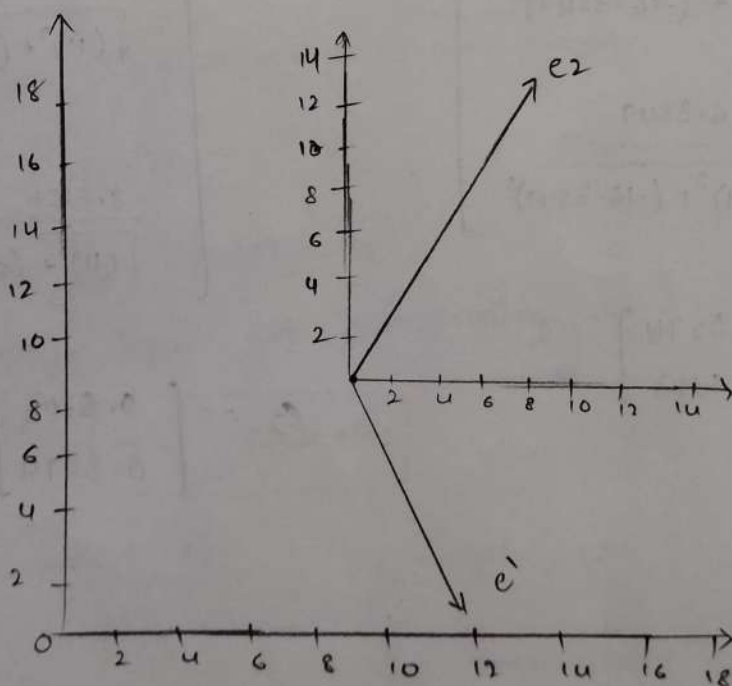
$$P_{11} = e_1^T \begin{bmatrix} 4 - 8 \\ 11 - 8.5 \end{bmatrix} = \begin{bmatrix} 0.5574 & -0.8303 \end{bmatrix} \times \begin{bmatrix} -4 \\ 2.5 \end{bmatrix} = -4.3052$$

$$P_{12} = e_1^T \begin{bmatrix} 8 - 8 \\ 4 - 8.5 \end{bmatrix} = \begin{bmatrix} 0.5574 & -0.8303 \end{bmatrix} \times \begin{bmatrix} 0 \\ -4.5 \end{bmatrix} = 3.7361$$

$$P_{13} = e_1^T \begin{bmatrix} 13 - 8 \\ 5 - 8.5 \end{bmatrix} = \begin{bmatrix} 0.5574 & -0.8303 \end{bmatrix} \times \begin{bmatrix} 5 \\ -3.5 \end{bmatrix} = 5.6928$$

$$P_{14} = e_1^T \begin{bmatrix} 7 - 8 \\ 14 - 8.5 \end{bmatrix} = \begin{bmatrix} 0.5574 & -0.8303 \end{bmatrix} \times \begin{bmatrix} -1 \\ +5.5 \end{bmatrix} = -5.1235$$

Step 6: Graphical Representation.

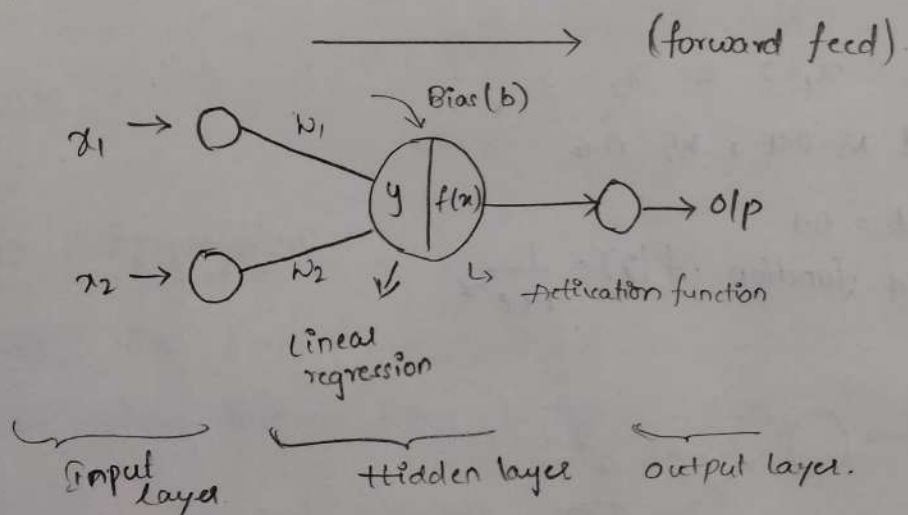


30 Explain forward and backward propagation in Artificial Neural Network (ANN)?

1. Forward propagation:

Forward Propagation is the process where input data is passed through the network layer by layer to generate output. It involves.

- Multiplication of Input with weights.
- Adding Bias.
- Applying an activation function.
- passing the result to the next layer.



Steps of forward propagation:

Input layer: Take the input values

Weighted sum: Multiple inputs with corresponding weights and add bias

$$y = \sum_{i=1}^n x_i w_i + b$$

$$y = (w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n) + b$$

3. Activation function: Apply an activation function like RELU, Sigmoid, Tanh etc. to introduce non-linearity for sigmoid Activation function,

$$Z = f(y) = \frac{1}{1+e^{-y}}$$

4. Output layer:

The final layer computes the output.

Example for forward propagation:

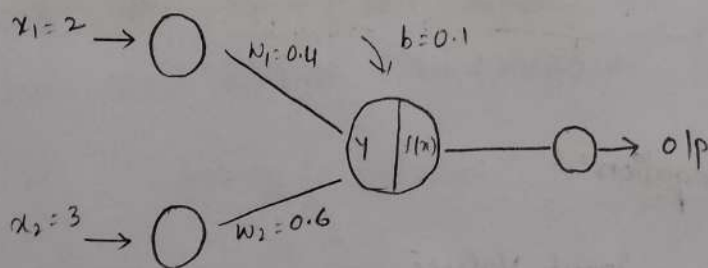
Consider,

input $x_1 = 2$ & $x_2 = 3$

weight $w_1 = 0.4$, $w_2 = 0.6$

Bias $b = 0.1$

Sigmoid function $f(z) = \frac{1}{1+e^{-z}}$



$$\begin{aligned} \text{ie } y &= w_1 x_1 + w_2 x_2 + b \\ &= (0.4)(2) + (0.6)(3) + 0.1 \end{aligned}$$

$$y = 2.7$$

Sigmoid function,

$$f(y) = \frac{1}{1+e^{-y}} \Rightarrow \frac{1}{1+e^{-2.7}}$$

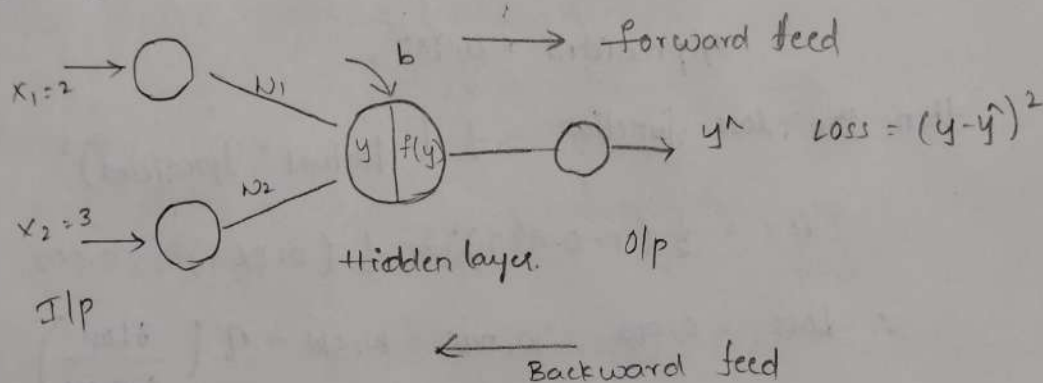
$$f(y) = 0.937$$

Where $\geq 0.5 \approx 1$
 $< 0.5 \approx 0$

\therefore Output of forward propagation is $0.937 \approx 1$

Backward propagation:

Backward propagation is the process of adjusting the weights to minimize the error using Gradient Descent



Steps for Back propagation:

1. calculate the Error (Loss function):

Use Mean Squared Error (MSE) or Cross Entropy Loss

$$\text{Loss} = \frac{1}{2} (y_{\text{actual}} - y_{\text{predicted}})^2$$

2. compute gradient:

- Differentiate loss function with respect to weights using chain rule and update the weights using Gradient Descent

$$w_n = w_0 - \eta \frac{\partial \text{Loss}}{\partial w_0}$$

$$\text{i.e., } w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial \text{Loss}}{\partial w_{\text{old}}}$$

where $\eta \rightarrow$ Learning rate, $\frac{\partial L}{\partial w} \rightarrow$ slope (or) Gradient Descent

3. Adjust Weights and Biases!

Repeat until loss is minimized

Example:

from the forward propagation, the actual output is

$$\therefore y_{\text{actual}} = 1$$

then the predicted value = 0.937

$$\therefore y_{\text{predicted}} = 0.937,$$

$$\text{then the loss function} = \frac{1}{2} (y_{\text{actual}} - y_{\text{predicted}})^2$$

$$\text{ie } = \frac{1}{2} (1 - 0.937)^2 \Rightarrow \frac{1}{2} (0.063)^2 \Rightarrow 0.002$$

$$\therefore \text{Loss} = 0.002, \quad w_{1, \text{new}} = w_{1, \text{old}} - \eta \left(\frac{\partial \text{Loss}}{\partial w_{1, \text{old}}} \right)$$

$$\text{then, } \frac{\partial \text{Loss}}{\partial w_{1, \text{old}}} = \frac{\partial \text{Loss}}{\partial y_{\text{predicted}}} \times \frac{\partial y_{\text{predicted}}}{\partial z} \times \frac{\partial z}{\partial w_1}$$

$$\text{ie } \frac{\partial \text{Loss}}{\partial y_{\text{predicted}}} = y_{\text{predicted}} - y_{\text{actual}} \\ = 0.937 - 1 \Rightarrow -0.063$$

Derivative of output wrt Activation function (Sigmoid)

$$f(z) = \frac{1}{1 + e^{-z}}$$

$$f'(z) = f(z) \times (1 - f(z))$$

$$= 0.937 \times (1 - 0.937)$$

$$= 0.059$$

(6)

wrt weight w_1 ,

$$\frac{\partial z}{\partial w_1} = x_1 = 2, \quad \frac{\partial z}{\partial w_2} = x_2 = 3$$

then,

$$\frac{\partial \text{Loss}}{\partial w_1} = \frac{\partial \text{Loss}}{\partial y_{\text{predicted}}} \times \frac{\partial y_{\text{predicted}}}{\partial z} \times \frac{\partial z}{\partial w_1}$$

$$(-0.063) \times (0.059) \times (2)$$

$$= -0.007434 //$$

Similarly,

$$\frac{\partial \text{Loss}}{\partial w_2} = \frac{\partial \text{Loss}}{\partial y_{\text{predicted}}} \times \frac{\partial y_{\text{predicted}}}{\partial z} \times \frac{\partial z}{\partial w_2}$$

$$= (-0.063) \times (0.059) \times (3)$$

$$= -0.011151$$

then, Here $\eta = 0.1$

$$w_{1 \text{ new}} = w_{1 \text{ old}} - \eta \times (-0.007434)$$

$$= 0.4 - 0.1(-0.007434)$$

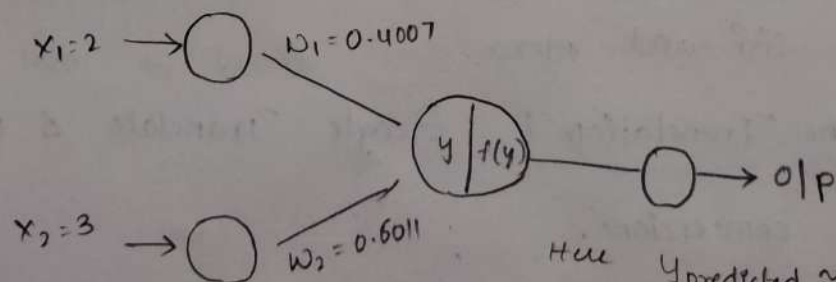
$$\boxed{w_{1 \text{ new}} = 0.40074}$$

$$w_{2 \text{ new}} = w_{2 \text{ old}} - 0.1 \times (-0.011151)$$

$$= 0.4 - 0.1(-0.011151)$$

$$\boxed{w_{2 \text{ new}} = 0.6011}$$

then,

Here $y_{\text{predicted}} \approx y_{\text{actual}}$

Here y predicted and y actual values would be similar (or) same after the updation of weights, where the error rate (or) the loss function value minimizes (or) reduces.

4Q Write a short note on,

- A) Applications of Deep Learning.
- B) Activation function
- C) Optimizers.
- D) Loss functions.

1. Applications of Deep Learning:

Deep Learning (DL) is widely used in various fields due to its ability to process large amount of data and complex patterns.

1. Computer Vision:

- A) Image classification: Identifies objects in images (eg cat or dog)
- B) Object detection: Detects & labels multiple objects in an image (eg self-driving cars, detecting pedestrians)
- C) Facial Recognition: Used in security systems and social media. (eg: Facebook auto-tagging)

2. Natural Language processing (NLP):

- chatbots & virtual Assistants: AI-powered bots like chat GPT, Siri and Alexa.
- Machine Translation: Google Translate & DeepL for language conversions.

Example:

- Sentiment analysis: Analyzing customer feedback, reviews and Social Media Sentiment.

3. Healthcare:

- Disease Diagnosis: AI models detect diseases like cancer from X-rays and MRI's.
- Drug Discovery: predicts chemical interactions for new medicine.
- Personalized Treatment: AI recommends treatments based on patient history.

4. Finance:

- Fraud Detection: identifies fraudulent transactions in banking.
- Stock market prediction: AI-powered trading bots analyze stock trends.

5. Speech Recognition:

- Voice Assistants: Alexa, Siri, Google Assistant convert speech to text.
- Automated Transcription: AI tools like Otter.ai transcribe meetings.

6. Recommendation systems:

Netflix, Youtube, Amazon: suggests movies, videos & products based on user behaviour.

(B) Activation Functions:

Activation function introduces non-linearity in neural networks, allowing them to learn complex relationships.

Types of Activation functions:

1. Sigmoid function: (Good for binary classification)

$$f(z) = \frac{1}{1+e^{-z}}$$

Output range = $(0, 1)$

Example

if $z=2$ then

$$f(2) = \frac{1}{1+e^{-2}} \approx 0.88$$

problem: causes Vanishing gradient in deep networks

2. RELU (Rectified Linear Unit)

→ It mostly used activation function

$$f(z) = \max(0, z)$$

Output range = $(0, \infty)$

Example

$$f(3) = 3$$

$$f(-2) = 0$$

Advantage: faster learning & avoids vanishing gradients

problem: Dying RELU (neurons stop updating for negative inputs)

3. Tanh (Hyperbolic Tangent):

It is better than sigmoid

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

Output range = $(-1, 1)$

Example:

$$f(2) = 0.96 \quad f(-2) = 0.96$$

Advantage: Helps in zero-centered learning

4. Softmax function: (used in multi-class classification)

$$f(z_i) = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

→ Converts outputs into probabilities (Sum = 1)

Example for 3 classes (A, B, C):

$$i) z_A = 2, \quad z_B = 1, \quad z_C = 0.5$$

Softmax assigns: A = 57%, B = 29%, C = 14%

(c) Optimizers:

Optimizers adjust the neural networks weights to minimize the loss function.

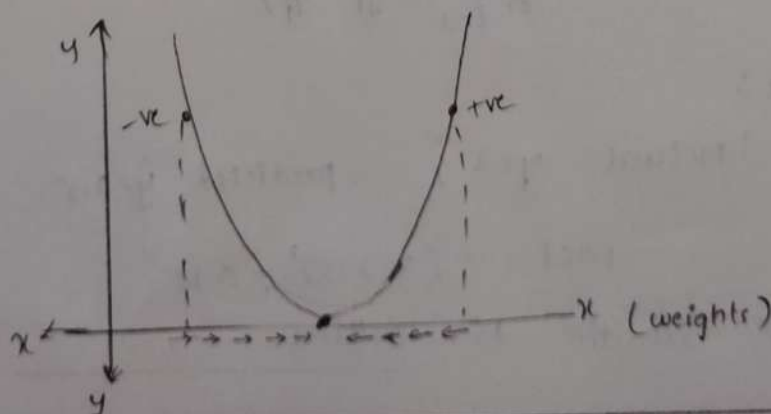
Types of Optimizers:

1. Gradient Descent (Basic Method)

$$w = w - \eta \frac{\partial \text{loss}}{\partial w}$$

η = learning rate

problem: slow convergence.



2. Stochastic Gradient Descent (SGD)

- Updates weights using a random sample instead of the whole dataset.
- Faster than regular gradient descent but has high variance

3) Momentum-Based Gradient Descent:

- Uses momentum to accelerate learning and avoid getting stuck in local minima.

4) Adam (Adaptive moment estimation)

- most widely used

$$w = w - \frac{\eta}{\sqrt{v_1} + \epsilon} m_1$$

- combines momentum & adaptive learning rate
- works well for most deep learning models.

①) Loss functions: Loss function measures the error between actual and predicted values.

Types of loss functions:

1) Mean-Squared Error (MSE)

used for Regression

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Example:

→ Actual $y = 3$, predicted $\hat{y} = 2.5$

$$MSE = (3 - 2.5)^2 = 0.25$$

problem: Sensitive to outliers.

2) Mean Absolute Error (MAE):

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

It is less sensitive to outliers than MSE

3) Cross-Entropy Loss (used in classification):

•) for binary classification:

$$loss = -(y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$$

•) for multi-class classification

$$loss = - \sum_{i=1}^n y_i \log(\hat{y}_i)$$

Example: If actual class = 1, predicted probability = 0.9:

$$loss = -\log(0.9) = 0.10$$

4) Huber loss (Hybrid of MSE and MAE)

•) used when outliers are present

•) used MSE for small error and MAE for large errors.