Ans 1)

<a>

Gradient descent in an optimization algorithm used to minimise the loss function of an ML model. The gradient is a vector of partial derivatives that represents the rate of change of loss function w.r.t. weights. By updating the weights in the direction of negative gradient, the algorithm moves towards the minimum of the loss function.

<b>

| Batch gradient descent | Stochastic gradient descent |
|---|---|
| Deterministic in nature | (1) Stochastic in nature |
| Convergence is slow | (2) Convergence is fast |
| Slow & computationaly expensive | (3) fast & less expensive |
| Uses the whole training sample | (4) Uses a single training sample |

Gradient descent fails to converge when :—

① Chosing a very large learning rate which can cause gradient descent to overshoot the minimum & diverge.

② Chosing a very small learning rate, it will take a long time to reach the minimum.

③ Stochastic gradient descent may lead to slow convergence.

④ Encountering a non-convex function, which can have multiple local minima that trap the gradient descent.

⑤ A high dimensional parameter space which increases the complexity of finding the global minima.

Learning rate determines how much we adjust weights w.r.t. loss gradient descent. It affects how fast or slow we move towards the optimal weights.

If it's too small, we may never converge & if it's too big we may overshoot & increase our training error.

Ans 2)

| $x$ | $y$ | $xy$ | $x^2$ | $y^2$ |
|---|---|---|---|---|
| 0 | 72 | 0 | 0 | NA |
| 5 | 66 | 330 | 25 | NA |
| 10 | 70 | 700 | 100 | NA |
| 15 | 64 | 960 | 225 | NA |
| 20 | 60 | 1200 | 400 | NA |
| $\Sigma x = 50$ | $\Sigma y = 332$ | | $\Sigma x^2 = 750$ | |

$m = 5$

$\Sigma xy = 3190$

Equations for linear regression line :—

① $na + b \sum_{i=1}^{m} x_i = \sum_{i=1}^{m} y_i$

② $a \sum_{i=1}^{m} x_i + b \sum_{i=1}^{m} x_i^2 = \sum_{i=1}^{m} x_i y_i$

substituting values :—

$[5(a) + b(50) = 332] \times 10$

$- \quad 50(a) + b(750) = 3190$

$\underline{\hspace{8cm}}$

$250(b) = 3190 - 3320 = -130$

$b = -\dfrac{130}{250} = -\dfrac{13}{25}$

$\therefore a = \dfrac{332 - 50\left(-\dfrac{13}{25}\right)}{5} = \dfrac{332 + 26}{5}$

$\Rightarrow a = 71.6$

$\therefore$ Line $:= y = 71.6(x) + \left(-\dfrac{13}{25}\right)$

Ans 3)

| | x | y | $k = 2$ |
|---|---|---|---|
| A: | 5 | 7 | |
| B: | 11 | 45 | |
| C: | 10 | 6 | |
| D: | 18 | 29 | |
| E: | 10 | 25 | |
| F: | 9 | 3 | |

Assume B and F as initial clusters :-

| Point | x | y | $\Delta_B$ | $\Delta_F$ |
|---|---|---|---|---|
| A | 5 | 7 | $\sqrt{1444}$ | $\sqrt{17}$ |
| B | 11 | 45 | 0 | ~~(come.)~~ $\sqrt{1813}$ |
| C | 10 | 6 | $\sqrt{1522}$ | $\sqrt{45}$ |
| D | 18 | 29 | $\sqrt{305}$ | $\sqrt{872}$ |
| E | 10 | 25 | $\sqrt{401}$ | $\sqrt{520}$ |
| F | 9 | 3 | $\sqrt{1813}$ | 0 |

Now;

A → F
B → B
C → F
D → B
E → B
F → F

∴ New clusters = (ACF)   (BDE)

$\frac{19}{3}, \frac{16}{3}$   13, 33

For 2$^{nd}$ iteration :-

| Point | x | y | $\Delta_{ACF}$ | $\Delta_{BDE}$ |
|---|---|---|---|---|
| A | 5 | 7 | 2.134 | 27.202 |
| B | 17 | 45 | 39.94 | 12.16 |
| C | 10 | 6 | 3.73 | 27.166 |
| D | 18 | 29 | 26.39 | 6.403 |
| E | 10 | 25 | 20.009 | 8.344 |
| F | 4 | 3 | 3.295 | 31.320 |

After 2$^{nd}$ iteration :—

$A \rightarrow ACF$

$B \rightarrow BDE$

$C \rightarrow ACF$

$D \rightarrow BDE$

$E \rightarrow BDE$

$F \rightarrow ACF$

∴ Two clusters after 2 iterations are :—
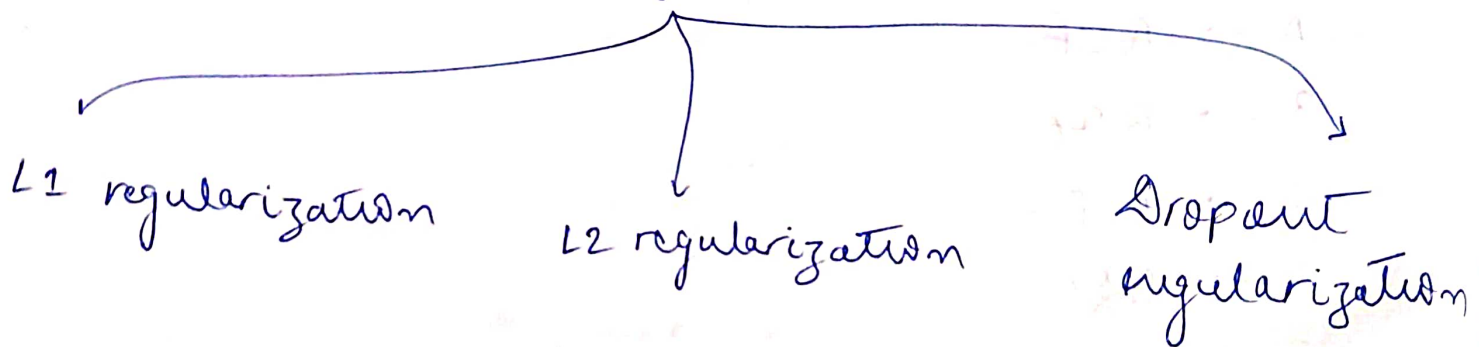
$\boxed{ACF}$ and $\boxed{BDE}$

Ans 4)    Target = (3, 7)

| | x | y | Label | $\Delta_{target}$ |
|---|---|---|---|---|
| A | 7 | 7 | 1 | 4 |
| B | 7 | 4 | 1 | 5 |
| C | 3 | 4 | 2 | 3 |
| D | 1 | 4 | 2 | $\sqrt{13} = 3.8$ |

Label of 3 nearest neighbors is:-

2, 2, 1

Target Label $= \dfrac{2+2+1}{3} = \textcircled{2}$

Ans 5) _Regularization_ : technique used to reduce the errors by fitting the function appropriately on the given training set & avoid overfitting.

L1 regularization          L2 regularization          Dropout regularization

L1 regularization (LASSO) : L1 adds absolute value of magnitude of coefficient as penalty term to the loss function.

L2 regularization (Ridge) : L2 adds squared magnitude of co-efficient as penalty term to the loss function.

Ans continued :-

$$RSS = \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{n} \beta_j x_{ij} \right)^2$$

Ridge regression adds the shrinkage quantity to modify the RSS :-

$$\sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{n} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{n} \beta_j^2 = RSS + \lambda \sum_{j=1}^{n} \beta_j^2$$

at $\lambda = 0$, the penalty term has no effect but as $\lambda \to \infty$ the penalty influence increases.

$\Rightarrow$    $\lambda \to 0$ == $\hat{\beta}_{ridge} = \hat{\beta}_{OLS}$

$\Rightarrow$    $\lambda \to \infty$ == $\hat{\beta}_{ridge} \to 0$

In case of Lasso :-
   • Penalty term is changed from $\beta^2$ to $|\beta|$ :-

$$\sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{n} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{n} |\beta_j| = RSS + \lambda \sum_{j=1}^{n} |\beta_j|$$

$\hookrightarrow$ for ridge : $RSS + \lambda \sum_{j=1}^{n} \beta_j^2$

vs

$\hookrightarrow$ for Lasso : $RSS + \lambda \sum_{j=1}^{n} |\beta_j|$

**Ans 6)** Derivation of Linear regression's loss function :—

Linear function b/w dependent & independent variables :—

$$\hat{y}^{(i)} = W_o + W_1 x_1^{(i)} + W_2 x_2^{(i)} + \ldots W_n x_n^{(i)}$$

Some notations :—

N : no. of datapoints

n : no. of features

$x_j^{(i)}$ : $j^{th}$ feature in $i^{th}$ observation / datapoint

Some terms :—

cost $f^n$ : calculates the error of the entire dataset

loss $f^n$ : calculates error on every observation of the dataset.

$$Error^{(i)} = e^{(i)} = y^{(i)} - \hat{y}^{(i)}$$

$$e^{(i)} = y^{(i)} - \sum_{j=0}^{n} W_j x_j^{(i)} \quad ; \quad x_o^{(i)} = 1$$

A.c.t. Least square method :—

$$\left(e^{(i)}\right)^2 = \left(y^{(i)} - \sum_{j=0}^{n} W_j x_j^{(i)}\right)^2 \quad ; \quad x_o^{(i)} = 1$$

$$RSS = E = \sum_{i=1}^{N} \left(e^{(i)}\right)^2 = \sum_{i=1}^{N} \left(y^{(i)} - \sum_{j=0}^{n} W_j x_j^{(i)}\right)^2 \quad ; \quad x_o^{(i)} = 1$$

$$X = \begin{bmatrix} x_0 & x_0 & & x_0 \\ x_1^{(1)} & x_1^{(2)} & \cdots & x_1^{(N)} \\ x_2^{(1)} & x_2^{(2)} & & x_2^{(N)} \\ \vdots & \vdots & & \vdots \\ x_m^{(1)} & x_n^{(2)} & & x_n^{(N)} \end{bmatrix} = X = [\bar{x}^{(1)} \ \bar{x}^{(2)} \cdots \bar{x}^{(N)}]$$

$$\forall \ \bar{x} = [1 \ x_1, \ x_2 \cdots x_n]^T$$

$$Y = [y^{(1)} \ y^{(2)} \cdots y^{(N)}]^T$$

$$\bar{W} = [W_0 \ W_1 \ W_2 \cdots W_n]^T$$

for optimum sol$^n$ :—

$$y^{(1)} = \bar{W}^T \bar{x}^{(1)}$$

$$\vdots$$

$$y^{(N)} = \bar{W}^T \bar{x}^{(N)}$$

$$[y^{(1)} \ y^{(2)} \cdots y^{(N)}] = \bar{W}^T [\bar{x}^{(1)} \ \bar{x}^{(2)} \cdots \bar{x}^{(N)}] = \bar{W}^T X$$

or $\quad Y = (\bar{W}^T x)^T$

2) $E(\bar{w}) = \bar{W}^T [x x^T] \bar{w} - 2\bar{w}^T xy + y^T y$

b) Find $\bar{w}^*$ which reduces the error to $0$ :—

$$\frac{\partial E(\bar{w})}{\partial \bar{w}} = 2(x x^T)\bar{w} - 2xy = 0$$

$$\Rightarrow \bar{w}^* = (x x^T)^{-1} xy$$

fitted output :—

$$\hat{y} = x^T \bar{w}^* = \underbrace{x^T (x x^T)^{-1} xy}_{\text{b} \ x^+ y}$$

$$\Rightarrow \bar{w}^* = x^+ y$$

**Ans7)** Weight updating method in gradient descent :—

$(\cdot)$ $\hat{y}^{(i)} = \sum\limits_{j=1}^{m} W_j x_j^{(i)} + W_0$

$(\cdot)$ Error$^{(i)} = y^{(i)} - \hat{y}^{(i)} = y^{(i)} - \left[\sum\limits_{j=1}^{n} W_j x_j^{(i)} + W_0\right] = e^{(i)}$

$E = \sum\limits_{i=1}^{N} e^{(i)} = \sum\limits_{i=1}^{N}\left[y^{(i)} - \sum\limits_{j=1}^{n} W_j x_j^{(i)} + W_0\right]$

$(\cdot)$ Applying least cost squared :—

$$E \rightarrow \frac{1}{2} E^2$$

2) $E = \frac{1}{2}\sum\limits_{i=1}^{N}\left[y^{(i)} - \left[\sum\limits_{j=1}^{n} W_j x_j^{(i)} + W_0\right]\right]^2$

$\hookrightarrow$ E is our cost function, we must minimise `E` :—

① Gradient descent updates every weight so we must find out how much every weight affects E :—

$$\frac{\partial E}{\partial W_j} = \frac{\partial\left[\frac{1}{2}\sum\limits_{i=1}^{N}\left[y^{(i)} - \left[\sum\limits_{j=1}^{n} W_j x_j^{(i)} + W_0\right]\right]^2\right]}{\partial W_j}$$

We know, $y^{(i)} - \left[\sum\limits_{j=1}^{n} W_j x_j^{(i)} + W_0\right] = e^{(i)}$

$\Rightarrow \dfrac{\partial E}{\partial W_j} = \dfrac{\partial\left[\frac{1}{2}\sum\limits_{i=1}^{N} (e^{(i)})^2\right]}{\partial W_j}$

$$\Rightarrow \quad \frac{\partial E}{\partial w_j} = \frac{1}{2} \sum_{i=1}^{N} \frac{\partial (e^{(i)})^2}{\partial w_j}$$

remember, chain rule :-

$$\frac{d(f(x))^2}{dx} = 2 f(x) \cdot f'(x)$$

$$\frac{\partial E}{\partial w_j} = 2 \cdot \frac{1}{2} \sum_{i=1}^{N} e^{(i)} \cdot \frac{\partial e^{(i)}}{\partial w_j}$$

$$\frac{\partial E}{\partial w_j} = \sum_{i=1}^{N} e^{(i)} \cdot \frac{\partial e^{(i)}}{\partial w_j}$$

for 
$$\frac{\partial e^{(i)}}{\partial w_j} = \frac{\partial \left[ y^{(i)} - \left[ \sum_{j=1}^{n} w_j x_j^{(i)} + w_o \right] \right]}{\partial w_j}$$

$$= 0 - x_j \qquad = - x_j^{(i)}$$

$$\therefore \quad \frac{\partial E}{\partial w_j} = \sum_{i=1}^{N} e^{(i)} \cdot (- x_j^{(i)})$$

$$= - \sum_{i=1}^{N} e^{(i)} x_j^{(i)}$$

$$\Rightarrow \quad \frac{\partial E}{\partial w_j} = - \sum_{i=1}^{N} \left[ y^{(i)} - \left[ \sum_{j=1}^{n} w_j x_j + w_o \right] \right] x_j^{(i)}$$

Now, $\dfrac{\partial E}{\partial w_o} = \dfrac{\partial \left[ \frac{1}{2} \sum\limits_{i=1}^{N} (e^{(i)})^2 \right]}{\partial w_j}$ ; $j = 0$

$$= 2 \cdot \frac{1}{2} \sum\limits_{i=1}^{N} e^{(i)} \cdot \dfrac{\partial e^{(i)}}{\partial w_j} \qquad ; \ j = 0$$

2) $\dfrac{\partial E}{\partial w_o} = \sum\limits_{i=1}^{N} e^{(i)} \cdot \dfrac{\partial e^{(i)}}{\partial w_o}$

for, $\dfrac{\partial e^{(i)}}{\partial w_o} = \dfrac{\partial \left[ y^{(i)} - \sum\limits_{j=1}^{n} w_j x_j + w_o \right]}{\partial w_o}$

$$= 0 - 0 + 1$$

2) $\dfrac{\partial E}{\partial w_o} = \sum\limits_{i=1}^{N} e^{(i)} \cdot 1$

Now, weight updates :—

$$w_j \leftarrow w_j + \eta \dfrac{\partial E}{\partial w_j}$$

2) $w_j = w_j + \eta \sum\limits_{i=1}^{N} \left[ y^{(i)} - \left[ \sum\limits_{j=1}^{n} w_j x_j + w_o \right] \right] x_j^{(i)}$

$$w_o \leftarrow w_o + \eta \dfrac{\partial E}{\partial w_o}$$

2) $w_o = w_o + \eta \sum\limits_{i=1}^{N} \left[ y^{(i)} - \left[ \sum\limits_{j=1}^{n} w_j x_j + w_o \right] \right]$