2) To avoid the problem of ambiguous region of linear discriminant function for c categories, we can

○ (✓)

Define c linear function $g_i(x)$, one for each class for $i = 1, 2, \ldots, c$

○

Assign x to $w_j$ if $g_i(x) < g_j(x)$ for all $i \neq j$ (✓)

○ Take a linear machine classifier ✓

✓ All the above

$x \in w_j$
$g_j(x) > g_i(x);\ i \neq j$

You want to predict
$x \in c_1 :\ g_{c_1}(x) \uparrow \uparrow$

$g_{c_2}(x)$

$g_{c_3}(x)$

$g_{c_1}(x) > g_{c_2}(x) > g_{c_3}(x)$

linear function $g_{c_1}(x) \uparrow :\ x \in c_1$

**3) Which of the following statements is true about the learning rate in Gradient Descent?**

✓ A very high learning rate may lead to oscillation
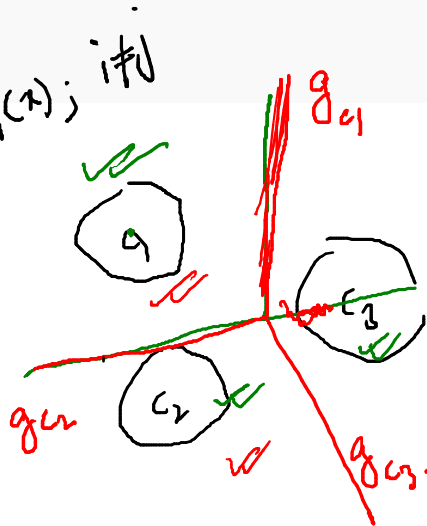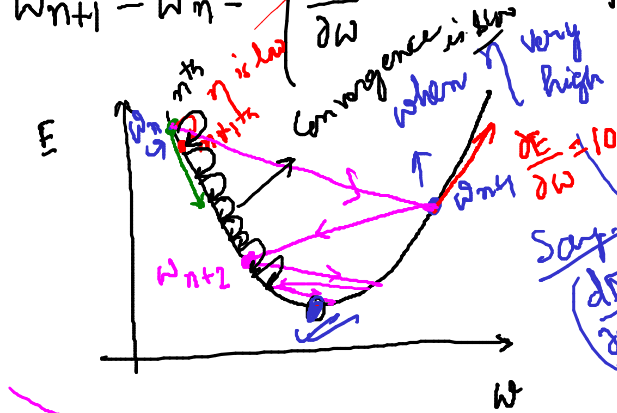
✗ A lower learning rate may lead to faster convergence ✗ ⟹ Convergence is slow -

✗ The learning rate doesn't determine the size of the steps taken towards the minimum

✗ The learning rate has no effect on the convergence of Gradient Descent

$$W_{n+1} = W_n - \eta \frac{\partial E}{\partial W}$$

$E = $ error function

Convergence is slow when $\eta$ is high. $\eta$ very high.

$\eta$ is very slow. $(\eta = 0.00001)$

→ The increment of the weight parameter or change in the value of $W$ will be very slow.

→ Convergence of the gradient optimization process is slow.

$$W_{n+1} = W_n - \eta \left(\frac{\partial E}{\partial W}\right) \to (-100)$$

$$= W_n - (0.00001 \times -100)$$

$$= W_n + 0.001 \to \text{Very small value.}$$

$\frac{\partial E}{\partial W} = -10$

$W_{n+1} \frac{\partial E}{\partial W}$

Saddle $\left(\frac{\partial E}{\partial W} = -100\right)$

$W_{n+2}$

E

W

$\eta \uparrow\uparrow$

$\eta = 100$
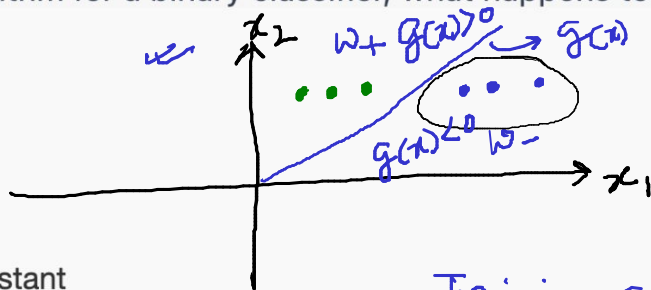
$W_{n+1} = W_n - (100 \times -100)$

$W_{n+1} = W_n + 10^4$

$W_{n+2} = W_{n+1} - (100 \times 10)$

$= W_{n+1} - 10^3$

$\eta$ is very high

→ There will be large no of oscillations

4) In the Perceptron algorithm for a binary classifier, what happens to the weights when a positive misclassified point is encountered?

- ○ It remains the same
- ○ It is increased
- ○ It is decreased
- ○ It is multiplied by a constant

$$W + g(x) > 0 \rightarrow g(x)$$

$g(x) < 0 \quad W_-$

Rule: 
$$\begin{cases} g(x) > 0 : x \in W_+ \\ g(x) < 0 : x \in W_- \end{cases}$$

$\rightarrow$ Decision Rule is not unified

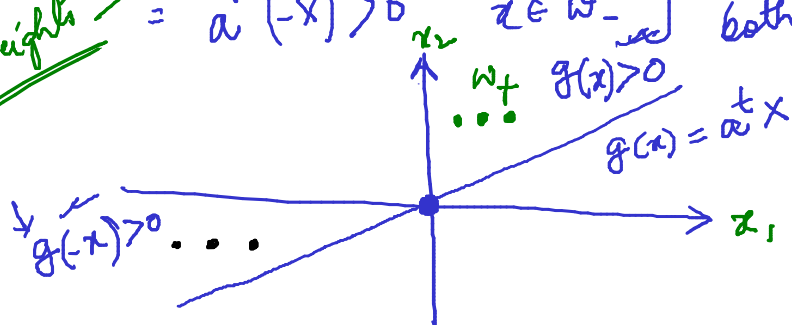To make a uniform decision Rule:-

Training set $\rightarrow$

$x \in W_+ : \quad g(x) > 0 \rightarrow$ That signifies proper classification

$x \in W_- \rightarrow$ You will negate the value of $x$ : $g(-x) > 0 \rightarrow$ Then again it signifies that you have done correct classification

Unified decision Rule

Learnable $g(x) = a^t x > 0 : x \in W_+$

Weights $\quad = a^t(-x) > 0 \quad x \in W_-$

$\left. \right\}$ Your decision Rule for both the class are now uniform.

$$W_+ \quad g(x) > 0$$

$$g(x) = a^t x$$

$$g(-x) > 0 \quad \cdots$$

For misclassification:-

when $x \in W_+ : - \quad g(x) = a^t x = (-ve) \rightarrow$ Misclassification

$$g(x) = a^t x < 0 \rightarrow \text{Loss}$$

$$\text{Loss} = \sum -a^t x \quad \text{for } \forall x \text{ misclassified}$$

4) In the Perceptron algorithm for a binary classifier, what happens to the weights when a positive misclassified point is encountered?
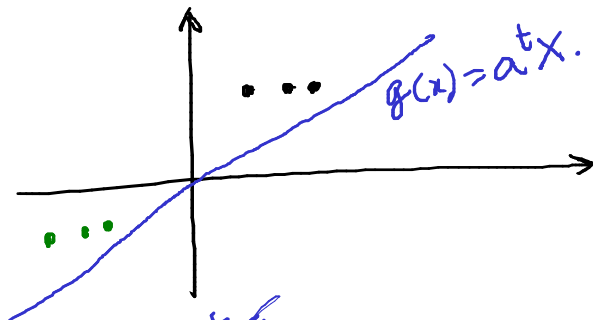
- ○ It remains the same
- ◉ It is increased
- ○ It is decreased
- ○ It is multiplied by a constant

Loss function $(L(a)) = \sum -a^t x$

$\forall x$ misclassified.

Training decision
Decision Rule $\Rightarrow$ $g(x) = a^t x > 0$ ; $x \in w_1$

$x \in w_2$

$= a^t(-x) > 0$ ;

$g(x) = a^t x$.

Misclassification :-
$x \in w_+ : g(x) = a^t x < 0$
$= -a^t x$

$a_{n+1} = a_n - \eta \nabla L(a)$

$a_{n+1} = a_n - \eta \frac{\partial L}{\partial a}$

$L(a) = \sum_{x \in w_1} -a^t x$

$= \sum (+ve) x (-ve)$

$= \sum -a^t x$.

$\frac{\partial L}{\partial a} = \sum -x$.

$a_{n+1} = a_n - \eta \left( \sum -x \right)$

$= a_n + \eta \sum x$

$a_{n+1}$

5) Let $w_{ij}$ represents weight between node i at layer k and node j at layer (k-1) of a given multilayer perceptron. The weight updation using gradient descent method is given by: ($\alpha$ and E represent learning rate and Error in the output respectively)

○ $W_{ij}(t+1) = W_{ij}(t) + \alpha \frac{\partial E}{\partial W_{ij}}, 0 \leq \alpha \leq 1$

○ $W_{ij}(t+1) = W_{ij}(t) - \alpha \frac{\partial E}{\partial W_{ij}}, 0 \leq \alpha \leq 1$

○ $W_{ij}(t+1) = \alpha \frac{\partial E}{\partial W_{ij}}, 0 \leq \alpha \leq 1$

○ $W_{ij}(t+1) = -\alpha \frac{\partial E}{\partial W_{ij}}, 0 \leq \alpha \leq 1$

$W_{ij}(t+1) = W_{ij}(t) - \alpha \frac{\partial E}{\partial W_{ij}} \; ; \; 0 \leq \alpha \leq 1$
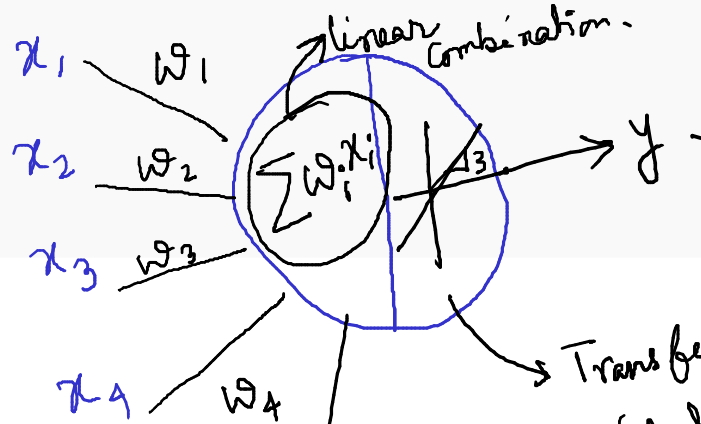
Gradient descent equation

6) A 4-input neuron has weights 3, 4, 5 and 6. The transfer function is <u>linear</u> with the constant of proportionality being equal to 3. The inputs are 6, 12, 10 and 20 respectively. What will be the output?

- ○ 238
- ○ 76
- ○ 708
- ○ 123

$x_1$

$w_1$

linear Combination.

$\sum w_i x_i$

$\frac{y}{x} = 3.$

$y = 3\hat{x}$

$= 3 \times 236$

$\boxed{y = 708.}$

$x_2$ $w_2$

$x_3$ $w_3$

$3$ $\to y$.

$\to$ Sgn $(x)$

$x_4$ $w_4$

$\to$ Transfer function.
(Activation function)

$\frac{y}{x} = 3.$

$y = 3x.$

$y = mx$

$\hat{x} = \sum w_i x_i$
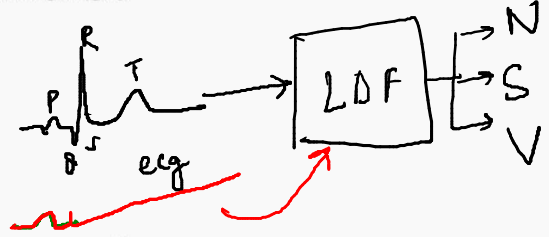
$= (w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4)$

$= (6 \times 3 + 4 \times 12 + 5 \times 10 + 20 \times 6)$

$\hat{x} = 18 + 48 + 50 + 120 = 236.$

$\boxed{\hat{x} = 236}$

7) Which of these is true about discriminant classifiers?

○ Assume conditional independence of features

○ Robust to outliers

○ Can perform classification if some missing data points are present

○ All the above

$$g(x) = P(\omega/x) = P(\omega) \boxed{P(x/\omega)}$$

→ Distribution
  → Gaussian Distribution

$x = [x_1, x_2 \cdots x_d]$.

$$P(x_i/\omega) = \frac{1}{\sqrt{2\pi} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^t \Sigma^{-1}(x-\mu)\right)$$
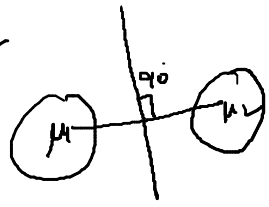
Covariance matrix

Multivariate gaussian dist.

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \\ & & \sigma_3^2 \cdots \sigma_d^2 \end{bmatrix}$$

$\Sigma_i = \Sigma_j$
$\forall i,j$

statistically.
→ independent

$\Sigma$ → arbitrary in nature

$(\mu_1 - \mu_2) \cdot W^t = 0$

LDF → N, S, V

ecg

7) Which of these is true about discriminant classifiers?
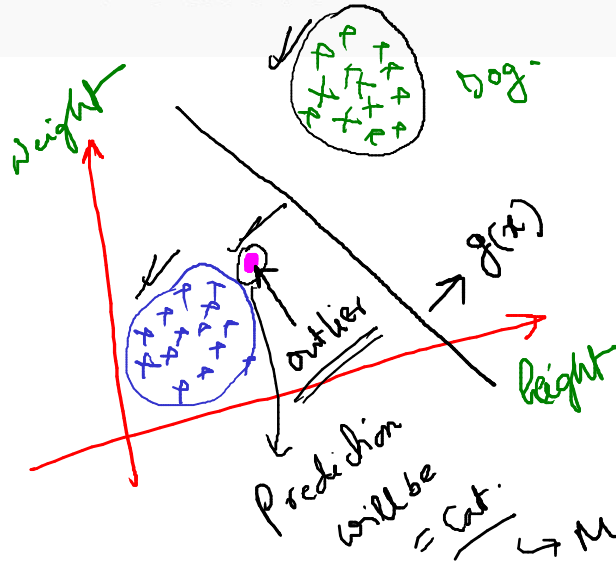
○ Assume conditional independence of features
○ Robust to outliers
○ Can perform classification if some missing data points are present
○ All the above

→ Imputation of data

Dog.

Cat → height,
Dog. → weight.

Weight

outlier

g(x)

height

Prediction will be = Cat.

→ Misclassification

→ Actually Dog.

→ But it's a data representation Seems to have similar characteristics on the Cat class.

8) A set of training samples are given below-

| $x_1$ | $x_2$ | y | $\lambda_i$ |
|-------|-------|---|-------------|
| 0.38 | 0.47 | + | 65.52 |
| 0.49 | 0.61 | - | 65.52 |
| 0.92 | 0.41 | - | 0 |
| 0.74 | 0.89 | - | 0 |
| 0.18 | 0.58 | + | 0 |
| 0.41 | 0.35 | + | 0 |
| 0.93 | 0.81 | - | 0 |
| 0.21 | 0.10 | + | 0 |

→ Support vectors from each of the classes.

$$W_i = \sum_i \lambda_i y_i x_i$$

Class specific equation would be :→

$$W^T x + b = 1; \quad y \in t$$

$b < D_1 + D_2$



$$W_1 x_1 + W_2 x_2 + b = 0$$

Generic equation of line that SVM predicts.

Using Support vector machine algorithm, the Marginal line for the classification can be calculated as-

○ $-5.32x_1 - 7.193x_2 + 9.09 = 0$

○ $-6.67x_1 + 8.134x_2 - 9.09 = 0$

○ $-7.21x_1 - 9.173x_2 + 9.09 = 0$

○ $8.21x_1 + 7.12x_2 - 9.09 = 0$

$$W_i = \sum_i \lambda_i y_i x_{1i}$$

$$W_1 = \lambda_1 y_1 x_{11} + \lambda_2 y_2 x_{12}$$

$$= 65.52 \times 1 \times (0.38) + 65.52 \times (-1) \times (0.49)$$

$$W_1 = -7.20$$

$$W_2 = \sum_i x_{2i} y_i \lambda_i = \lambda_1 y_1 x_{21} + \lambda_2 y_2 x_{22}$$

$$= 65.52 \times 1 \times 0.47 + 65.52 \times (-1) \times 0.61$$

$$W_2 = -9.17$$

8)  A set of training samples are given below-

| $x_1$ | $x_2$ | y | $\lambda_i$ |
|---|---|---|---|
| 0.38 | 0.47 | + | 65.52 |
| 0.49 | 0.61 | - | 65.52 |
| 0.92 | 0.41 | - | 0 |
| 0.74 | 0.89 | - | 0 |
| 0.18 | 0.58 | + | 0 |
| 0.41 | 0.35 | + | 0 |
| 0.93 | 0.81 | - | 0 |
| 0.21 | 0.10 | + | 0 |

Using Support vector machine algorithm, the Marginal line for the classification can be calculated as-

○ $-5.32x_1 - 7.193x_2 + 9.09 = 0$

○ $-6.67x_1 + 8.134x_2 - 9.09 = 0$

○ $-7.21x_1 - 9.173x_2 + 9.09 = 0$

○ $8.21x_1 + 7.12x_2 - 9.09 = 0$

*Handwritten work:*

$w_1 x_{11} + w_2 x_{21} + b_1 = 1$

$w_1 x_{12} + w_2 x_{22} + b_2 = -1$

$\Rightarrow b_1 = -(-7.2 \times 0.38 - 9.173 \times 0.47) + 1$
$= 8.05$

$b_2 = 10.12$

$b = \dfrac{b_1 + b_2}{2}$
$= 9.085$
$\approx 9.09$

The optimal eqn of line

$-7.2x_1 - 9.173x_2 + 9.09 = 0$

9) In refer to Q.8, A new test sample (0.5,0.5) is found. The class of the given sample is-
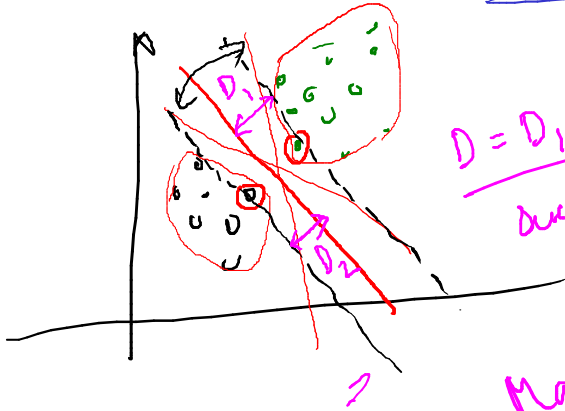
- ☑ Positive
- ⃝ Negative
- ⃝ Both class
- ⃝ Can't say

Q.8. Ans:- $\underline{-7.21x_1 - 9.173x_2 + 9.09 = 0}$ .

$$\underbrace{-7.21x_1 - 9.173x_2 + 9.09}_{g(x_1, x_2)}$$

$$g(0.5, 0.5) = -7.21 \times 0.5 - 9.173 \times 0.5 + 9.09$$

$$= 0.8985 = (+ve)$$

$$g(\tfrac{1}{2}, \tfrac{1}{2}) > 0$$

$$(+ve) \rightarrow \text{class identify} = +\text{ class}.$$

10) What is the main objective of a Support Vector Machine (SVM)?

~~To maximize the number of support vectors~~

~~To minimize the margin between classes~~

~~To maximize the training accuracy~~ *effect*

To find a hyperplane that separates classes with the maximum margin
*Cause*



$D = D_1 + D_2$

Such that $y_i (w^T x + b) > 0$

Maximize the margin Corresponding to a particular decision boundary.