⑨ Linear Classification    (Linear Transformation)

Goal: Is to classify data points into classes using a linear boundary. We predict a class label.

$$f(x) = w^T x + b \;; \quad \hat{y} = \begin{cases} 1, & w^T x + b > 0 \\ 0, & \text{otherwise} \end{cases}$$

(Decision boundary → $w^T x + b = 0$) ↳ it seperates classes.

Linear classifier gives hard class outputs [0|1] not probability.
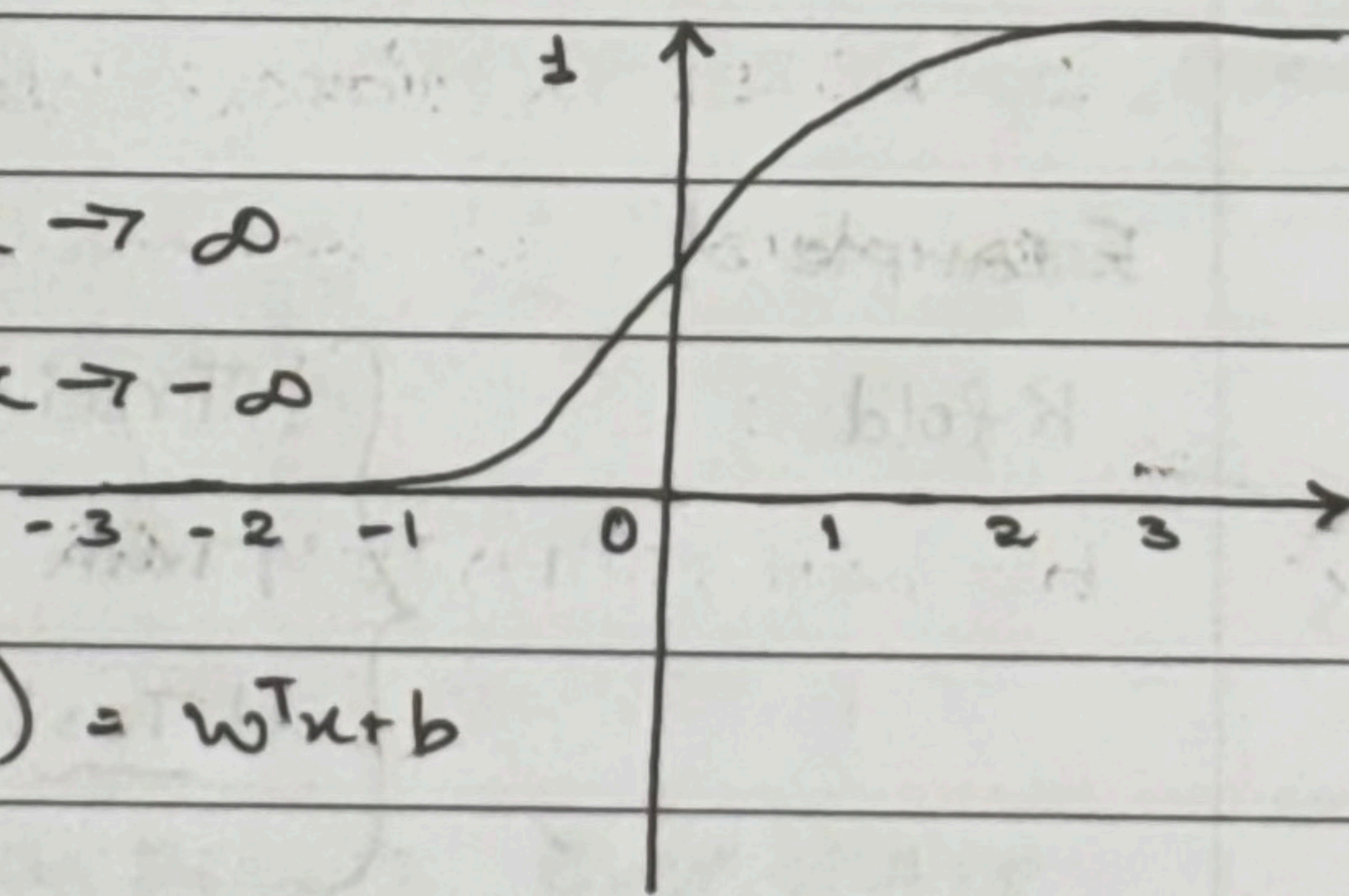
FIX

Predicts boundary directly

→⑧ Logistic Regression    (Discriminative Model)    $P(y|x)$

↳ we map linear output → probability range [0,1] using sigmoid function

(large Data)

(Convex)

$$\begin{cases} P(y=1 \mid x) = \sigma(z) = \dfrac{1}{1+e^{-z}}, \quad z = w^T x + b \\[2mm] P(y=0 \mid x) = 1 - P(y=1 \mid x) \end{cases}$$

$\sigma(z)$ tends → 1 as $z \to \infty$

$\sigma(z)$ tends → 0 as $z \to -\infty$

→ Log odds (Logit): $\log(P / 1-P) = w^T x + b$

1. Each prediction is a Bernoulli random variable with success probability $\hat{p}_i$

2. We want to find 'w' and 'b' that maximize probability of observing data.

Likelihood:    $L(w) = \prod_{i=1}^{n} P(y_i \mid x_i \,; w)$

Minimizing loss $\propto$ Maximizing Likelihood

_/_/_

$$\therefore L(w) = \prod_{i=1}^{n} (\hat{P_i})^{y_i} (1 - \hat{P_i})^{(1-y_i)}$$

3. For numerical stability, we take log:

$$l(w) = \log L(w) = \sum_{i=1}^{n} \left[ y_i \log(\hat{P_i}) + (1-y_i) \log(1-\hat{P_i}) \right]$$

This is the log likelihood function.

4. We want to maximize likelihood $\rightarrow$ maximize log-likelihood

In optimization, we minimize a cost, we take negative log likelihood

$$J(w) = -l(w) = -\sum_{i=1}^{n} \left[ y_i \log(\hat{P_i}) + (1-y_i) \log(1-\hat{P_i}) \right]$$

This is log loss / Binary Cross Entropy Loss

5. Gradient Descent: $w = w - \alpha \, \partial J / \partial w$, $b = b - \alpha \, \partial J / \partial b$

predicts how each class generates data

⑨ Linear Discriminant Analysis (Generative Model)

$P(x|y)$ & $P(y)$

(small Data)

LDA: Classification, Dimensionality Reduction.

It finds a linear combination of features that best seperates two or more classes.

★

Assumptions of LDA: 1. Gaussian / Normal Distribution
2. Same Covariance Matrices
3. Linear Separability.

Goal is to find a projection vector $w$,

such that $\quad$ seperation $= \dfrac{\text{Between class variance}}{\text{Within class variance}}$ is maximised

Fisher Criterion, $J(w)$: $S_b / S_w \quad \approx \uparrow$

measures seperability of classes along projection $w$.

$M_0, M_1$ : mean vectors for class 0 and 1

$\Sigma_0, \Sigma_1$ : covariance matrices

$S_W$ : $\Sigma_0 + \Sigma_1$ : within class scatter matrix

$S_B$ : $(M_1 - M_0)(M_1 - M_0)^T$ : between class scatter matrix

we find $\quad w^* = S_W^{-1}(M_1 - M_0) \Rightarrow \boxed{\dfrac{M_1 - M_0^{\,2}}{S_W}}$

The discriminant function for class k,

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \tfrac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log(\pi_k)$$

$\pi_k$ : prior probability class k.

Predict class = $\hat{y} = \underset{k}{\arg\max}\ \delta_k(x)$

→ Finds $(k-1)$ new axes for k classes.

↳ Linear combinations of original features.

---

⑩ Seperating Hyperplane ~ Perceptron (Historical - Base of NN)

(Discriminative Model)

Goal: Finding a linear decision boundary (hyperplane)
that divides data points belonging to different classes.

1. Initialise weights and bias

2. Iterate through each training sample:

   - Predict $\hat{y}$

   - If $\hat{y} \neq y$ : update weights

      $w : w - \alpha\, n y_i x_i$

      $b : b - \alpha\, n y_i$

3. Repeat until convergance.

→ Outputs a hard classification

→ Finds some hyperplane that seperates classes.