# SUPERVISED LEARNING
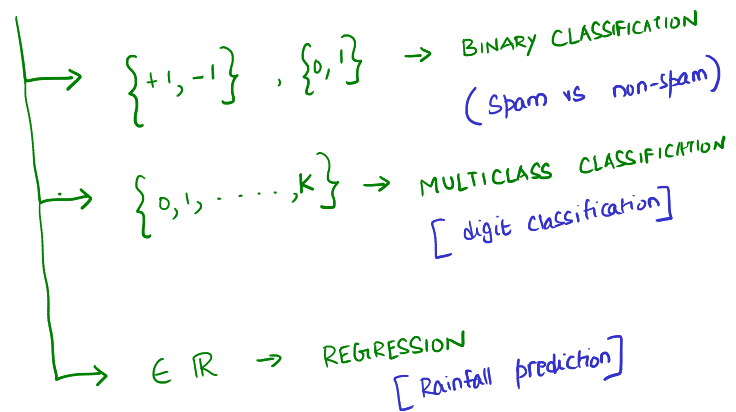
Input:  $\{x_1, \ldots, x_n\}$     $x_i \in \mathbb{R}^d$    ← Features / Attributes

        $\{y_1, \ldots, y_n\}$    ← LABELS (supervision)

→ $\{+1, -1\}$ , $\{0, 1\}$ → BINARY CLASSIFICATION (Spam vs non-spam)

→ $\{0, 1, \ldots, K\}$ → MULTICLASS CLASSIFICATION [digit classification]

→ $\in \mathbb{R}$ → REGRESSION [Rainfall prediction]

---

# REGRESSION

INPUT / TRAINING DATA    $\{x_1, \ldots, x_n\}$    $x_i \in \mathbb{R}^d$

                $\{y_1, \ldots, y_n\}$    $y_i \in \mathbb{R}$.

Goal:   Learn   $h : \mathbb{R}^d \to \mathbb{R}$

- How do we measure "goodness" of a function $h : \mathbb{R}^d \to \mathbb{R}$

- $\text{error}(h) = \sum_{i=1}^{n} (h(x_i) - y_i)^2$
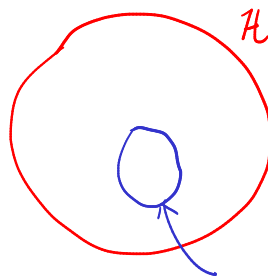
- How small can this error be? $\boxed{0}$

- Which $h$ achieves zero error?

$$h(x_i) = y_i \quad \forall i$$

---

What is the problem?

- By "memorizing", we can get zero error on training data

- What we care is about test performance.

- Impose "STRUCTURE" to reduce search space.

SIMPLEST STRUCTURE — LINEAR STRUCTURE.

$\mathcal{H}$

← Set of all functions $h: \mathbb{R}^d \to \mathbb{R}$.

$$\mathcal{H}_{linear} = \left\{ \begin{array}{l} h_w : \mathbb{R}^d \to \mathbb{R} \quad s.t \\ h_w(x) = w^T x \quad \text{for all } w \in \mathbb{R}^d \end{array} \right\}$$

GOAL:

$$\min_{h_w \in \mathcal{H}_{linear}} \sum_{i=1}^{n} \left( h_w(x_i) - y_i \right)^2$$

(or) equivalently.

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^{n} \left( w^T x_i - y_i \right)^2 \quad \to \quad \text{LINEAR REGRESSION}$$

---

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^{n} \left( w^T x_i - y_i \right)^2 \quad = \quad \| X^T w - y \|_2^2$$

$$X^T = \begin{bmatrix} - x_1 - \\ - x_2 - \\ \vdots \\ - x_n - \end{bmatrix}_{n \times d} \qquad \begin{bmatrix} | \\ w \\ | \end{bmatrix}_{d \times 1} \qquad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

$$\min_{w \in \mathbb{R}^d} \left( X^T w - y \right)^T \left( X^T w - y \right) \quad \leftarrow \quad \text{unconstrained Quadratic (in } w\text{) Optimisation problem}$$

Solution: Take derivative (gradient) and set to zero.

$$f(w) = \left( X^T w - y \right)^T \left( X^T w - y \right)$$

$$\nabla f(w) = 2 (X X^T) w - 2 (X y)$$

Solution satisfies $\boxed{(X X^T) w^* = X y}$
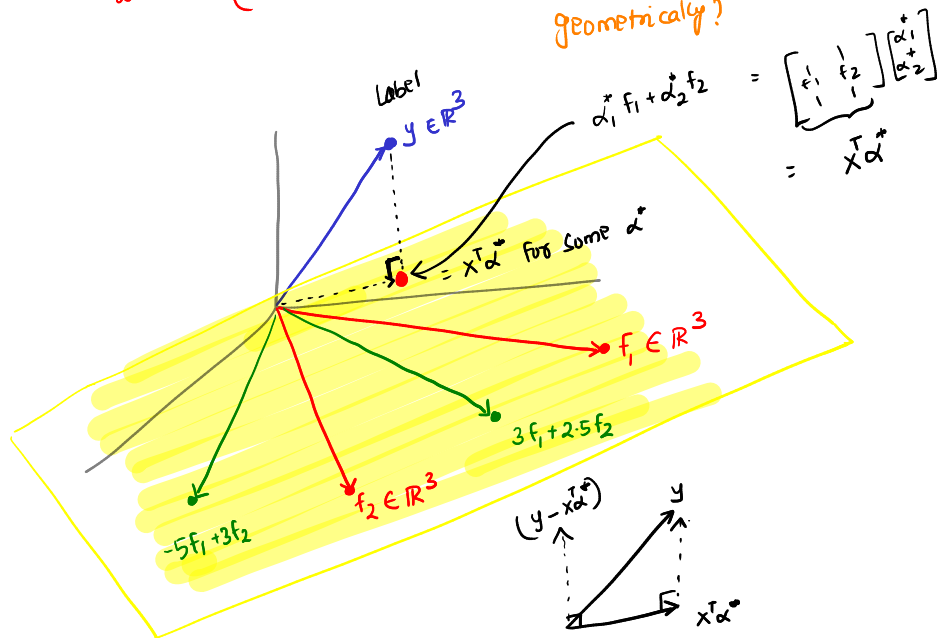
Like PCA, $w^*$ depends on a "covariance" like matrix. But it also involves $y$.

$$w^* = (xx^T)^\dagger xy$$

Pseudo-inverse.

→ Lin reg has closed form solution

→ GEOMETRIC VIEW ?

→ COMPUTATIONAL CONSIDERATIONS ?

→ NON-LINEAR FEATURE → LABEL RELATIONSHIP ?

→ PROBALISTIC VIEW ?

$$w^* = (xx^T)^\dagger xy \quad \leftarrow \quad \text{How can we interpret this geometrically ?}$$

$f_1$ height $\quad f_2$ weight

$$x_1 \begin{bmatrix} [] \end{bmatrix} \begin{bmatrix} [] \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \begin{matrix} d = 2 \\ n = 3 \end{matrix}$$

Label $y \in \mathbb{R}^3$

$$\alpha_1^* f_1 + \alpha_2^* f_2 = \begin{bmatrix} f_1 & f_2 \end{bmatrix} \begin{bmatrix} \alpha_1^* \\ \alpha_2^* \end{bmatrix} = x^T \alpha^*$$

$= x^T \alpha^*$ for some $\alpha^*$

$f_1 \in \mathbb{R}^3$

$3f_1 + 2.5 f_2$

$-5f_1 + 3f_2$

$f_2 \in \mathbb{R}^3$

$(y - x^T \alpha^*)$ $\quad y \quad x^T \alpha^*$

$$\left(y - x^T \alpha^*\right)^T \left(x^T \alpha^*\right) = 0$$

$$\underbrace{y^T x^T \alpha^* - \alpha^{*T} (xx^T) \alpha^* = 0}_{} \qquad — ①$$

Recall, $\quad w^* = (xx^T)^\dagger xy$

Substituting $w^* = \alpha^*$ on L.H.S, we get

$$y^T x^T \left((xx^T)^\dagger xy\right) - \left((xx^T)^\dagger xy\right)^T (xx^T)\left((xx^T)^\dagger xy\right)$$

$$= 0 \qquad \left[\begin{matrix} \text{use} \\ xx^T \text{ is symmetric} \end{matrix}\right]$$

CONCLUSION : $X^T \omega^+$ is The PROJECTION of the Labels onto the Subspace spanned by the features.