Recall:  Supervised learning $\begin{cases} \text{discrete } y_n: & \text{classification} \\ \text{continuous } y_n: & \text{regression} \end{cases}$

## 1. Linear Regression setup:

↳ __Training set__:

$$D = \{ (\underline{x}_1, y_1), (\underline{x}_2, y_2) \cdots (x_d, y_d) \}$$

$$\underline{x}_n \in \mathbb{R}^d, \quad y_n \in \mathbb{R}$$

2. __Decision Rule__: (aka "Hypothesis set")

$$\hat{y} = h(\underline{x}) = w_0 + w_1 x_1 + \cdots + w_d x_d$$

Redefined augumented form: $\quad \underline{w} = (w_0, w_1, \cdots w_d)$

$$\underline{x} = (x_0 = 1, x_1, \cdots x_d)$$

$$\therefore \hat{y} = h(\underline{x}) = w^T \underline{x}$$

<span style="color:blue">↑
h for hypothesis</span>

↳ __criterion for learning__:

$$E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^{N} (y_n - \hat{y})^2 = \frac{1}{N} \sum_{n=1}^{N} \underbrace{(y_n - w^T \underline{x})^2}_{\color{red} e_n(\underline{w})}$$

<span style="color:red">↑ "averaged squared error"</span>

// $e_n(\underline{w})$: squared error on the $n^{th}$ example

↳ __Goal__:

Given $D$, find $\underline{w} \in \mathbb{R}^{d+1}$ to minimize $E_{in}(\underline{w})$
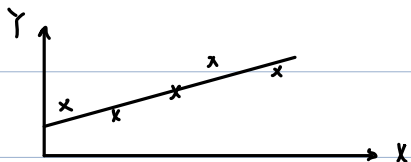
## 2. Example:

↳ **Q:** aim is to find the impact of advertisement on sales.

Let $X$ = adv. cost in one week $\quad (d = 1)$

$y$ = sales in one week

$\mathcal{D} = \{(x_1, y_1), (x_2, y_2) \cdots (x_N, y_N)\}$

Find a linear model $\quad y = w_0 + w_1 x$



↳ **Refined Model:**

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} \text{adv on } N_a \text{ of TV} \\ \text{adv on } N_a \text{ of Radio} \\ \text{adv on } N_a \text{ of Newspaper} \end{bmatrix} \qquad (d = 3)$$

$y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3$

larger $w_i \Rightarrow$ More profitable $x_i$

## 3. Algebra on how to solve:

↳ **Data Matrix:**

$$\underline{\overline{X}} = \begin{bmatrix} \underline{x_1}^T \\ \underline{x_2}^T \\ \vdots \\ \underline{x_N}^T \end{bmatrix} \in \mathbb{R}^{N \times (d+1)} \qquad \text{i.e.} \begin{bmatrix} X_{1,1} & \cdots & X_{1,d} \\ X_{2,1} & & X_{2,d} \\ \vdots & & \\ X_{N,1} & & X_{N,d} \end{bmatrix}$$

↳ **Target vector:**

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^N$$

↳ **Weight vector:**

$$\underline{w} \in \mathbb{R}^{d+1}$$

↳ __Linear Regression Model:__

$$\hat{\underline{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \begin{bmatrix} \underline{w}^T \underline{x}_1{}^T \\ \underline{w}^T \underline{x}_2{}^T \\ \vdots \\ \underline{w}^T \underline{x}_N{}^T \end{bmatrix} = \begin{bmatrix} \underline{x}_1{}^T \\ \underline{x}_2{}^T \\ \vdots \\ \underline{x}_N{}^T \end{bmatrix} \underline{w}$$

∴ $\boxed{\hat{\underline{y}} = \overline{X}\,\underline{w}}$

↳ __Error:__

$\boxed{E_{in}(\underline{w}) = \dfrac{1}{N} \| \underline{y} - \hat{\underline{y}} \|^2}$

- when is $E_{in}(\underline{w}) = 0$, i.e. $\underline{y} = \hat{\underline{y}}$ ?

  ∵ $y_n = \underline{w}^T \underline{x}_n$ , for $n = 1, 2, \cdots N$

  #. of linear equations: $N$

  #. of variables : $d+1$

  ∴ In practice, $N \gg d$ ⟹ No exact solution

  ∴ Instead, we try to minimize $E_{in}(\underline{w})$

  ⟹ least square solution.

# Lec 5.2 Least squares solution

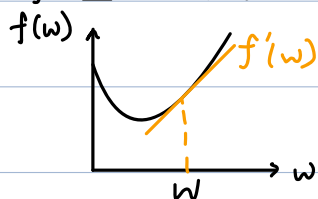↳ To minimize $E_{in}(\underline{w}) = \frac{1}{N} \|\underline{y} - \hat{\underline{y}}\|^2$

## 1. Least square sol:

↳ def: $f(\underline{w}) = \|\underline{y} - \hat{\underline{y}}\|^2$

$$= \|\overline{X} \underline{w} - \underline{y}\|^2$$

$$= \sum_{n=1}^{N} \underbrace{(\underline{w}^T \underline{x}_n - y_n)^2}_{e_n(\underline{w})}$$

↳ def: gradient of $f(\underline{w})$:

- $\nabla f(\underline{w}) = \begin{bmatrix} \partial f / \partial w_1 \\ \partial f / \partial w_2 \\ \vdots \\ \partial f / \partial w_d \end{bmatrix}$ ,  It points in the direction of the steepest increase.

- If $\underline{w}$ is 1-dimension, only two direction: left / right



↳ Claims:

① $\boxed{\nabla f(\underline{w}) = 2 \overline{X}^T (\overline{X} \underline{w} - \underline{y})}$

② Least squares solution $\underline{w}_{LS}$ is such that $\nabla f(\underline{w}_{LS}) = 0$

$$\therefore 2 \overline{X}^T (\overline{X} \underline{w} - \underline{y}) = 0$$

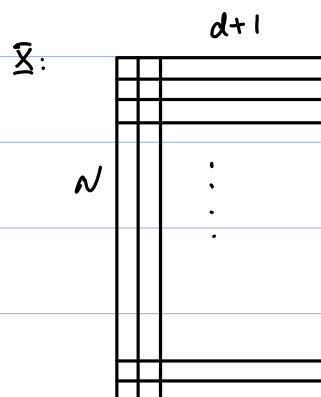$$\overline{X}^T \overline{X} \underline{w} = \overline{X}^T \underline{y}$$

Assume: The $d+1$ col of $\overline{X}$ are lin. indept.
↓ data samples
$\equiv \exists$ at least $d+1$ rows of $\overline{X}$ that are lin. indept.

$\therefore$ rank$(\overline{X}) = d+1 \iff \overline{X}^T \overline{X}$ is invertible

$\therefore \boxed{\underline{w}_{LS} = (\overline{X}^T \overline{X})^{-1} \overline{X}^T \underline{y}}$  // "Win in textbook"

$\overline{X}$:

↳ __Pseudo - inverse of $\bar{X}$__ :

$$\underline{\bar{X}^+ \triangleq (\bar{X}^T \bar{X})^{-1} \bar{X}^T}$$

why we called "pseudo" inverse ?

- $\underline{\bar{X}}^+ \underline{\bar{X}} = 1$

- But $\underline{\bar{X}} \underline{\bar{X}}^+ = \underline{\bar{X}} (\bar{X}^T \underline{\bar{X}})^{-1} \underline{\bar{X}}^T \neq 1$

__Note:__

① We have $\underline{y} = \bar{X} \underline{w}$ , and want to find $\underline{w}$

     $\underline{w} = \boxed{?} \underline{y}$ , where $\boxed{?}$ is kind of $\bar{X}^{-1}$

     But $\bar{X}$ is not a square matrix $\Rightarrow$ Not invertible

     ∴ We use pseudo inverse to do instead

② $\underline{\hat{y}}_{LS} = \bar{X} \underline{w}_{LS} = \underline{\bar{X} \bar{X}^+ y}$

                 ↑

      projection matrix from $\underline{y}$ to $\underline{\hat{y}}$