# Supervised Learning Setup

<u>Note:</u> I've covered some derivations but this will help in the long run.

\* If this book seems hand to approach, start from this part. I tried to simplify as much as I could. \*

Supervised Learning: Think of this as a child who is walking by his father in a garden. His father is feeding him about properties about the flowers in the garden and after enough knowledge, the kid will be able to identify flowers all by himself. This is supervised learning.

Now, there is a thing called training set.

$$\{(x^1, y^1), (x^2, y^2), \ldots (x^n, y^n)\}$$

\* Let's assume the most commonly used stat data which is house data.

(Suppose we have a dataset of house size vs house prize)

\* Now, our training set is really important, it is a part of a large population. We care about new $x$s so that we can predict new $y$s.

* If 'y' is discrete, it is a classification problem.
  ( Ex: Is this a Rose? Is this a tulip)
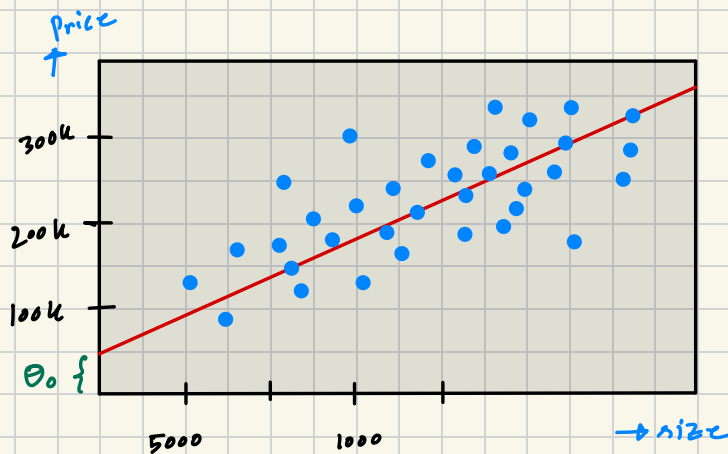  If 'y' is continuous, it is a regression problem.

⊛ Now. in the housing dataset, we have price vs area and our target is to find the price of a given lot area. We can call this a hypothesis (function)
& we can plot this as

$$h(x) = \theta_0 + \theta_1 (x)$$

{ Following this notation is a bit hard for me; so, later I used $f(x) = \beta_0 + \beta_1 (x)$ }
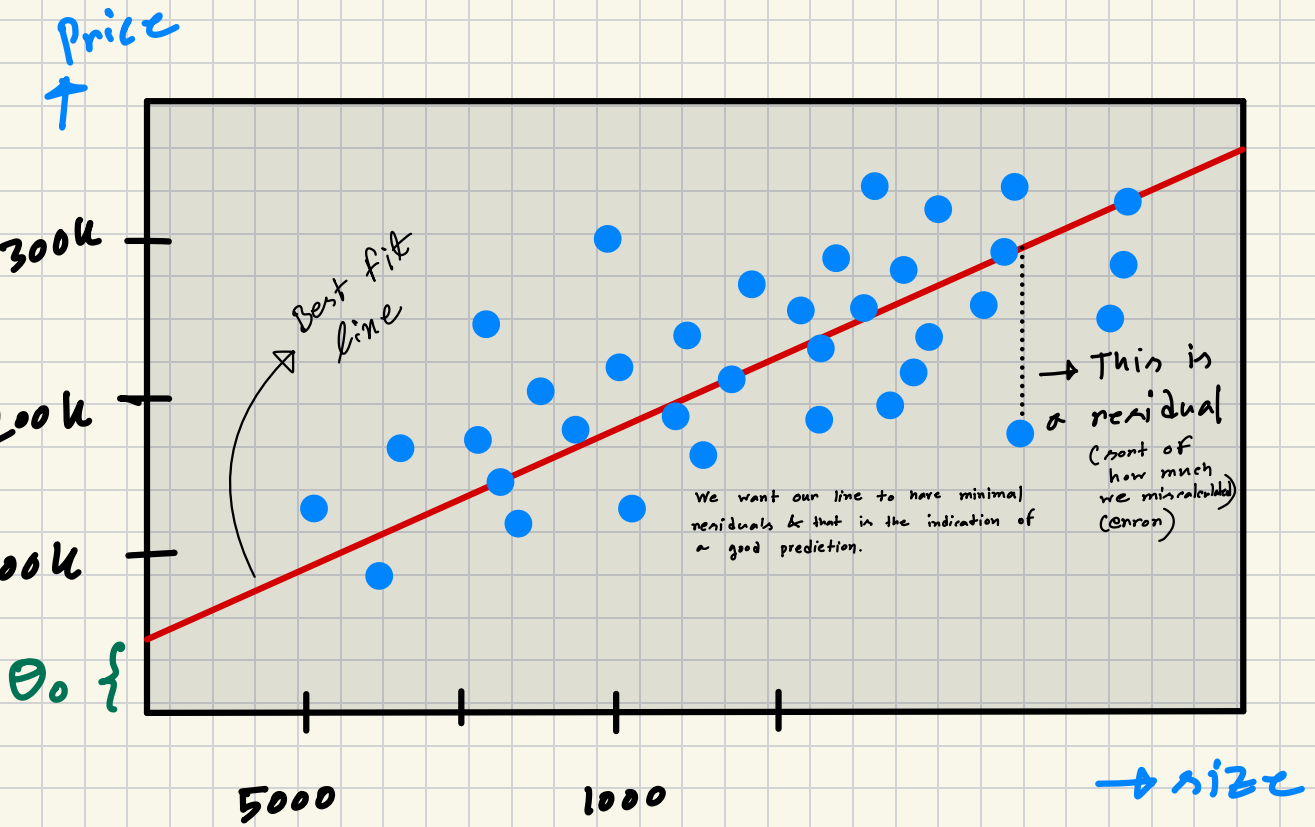


These are the Xs →                                          This is y

| Size | Bedroom Count | lot size | | | Price |
|------|------|------|------|------|------|
| 2200 | 4 | 45k | | | 900k |
| 2500 | 3 | 30k | | | 900k |

# How do we find an example of something good?



Best fit line

We want our line to have minimal residuals & that is the indication of a good prediction.

→ This is a residual (sort of how much we miscalculated) (error)

Price

300k

200k

100k

$\theta_0$ {

5000    1000    → size

We'll use sum of the squares of the residuals & this is for historical purposes only. (You can use absolutes)

$$h_\theta(x) = \sum_{j=0}^{d} \theta_i x_j \qquad h_\theta(x) \approx y$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^{n} \left( h_\theta(x^i) - y^i \right)^2$$

# Linear Regression

1. **Simple Linear Regression :** One input column vs one output column.

2. **Multiple LR :** Multiple input columns.

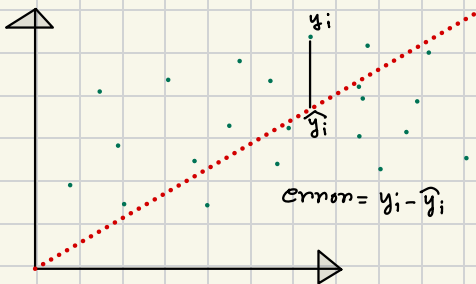3. **Polynomial LR :** For non linear data.

## * Simple Linear Regression :

The main idea was finding the slope & intercept. There are two ways to do it

1) Closed Form Solution { We use direct formula } [OLS]

2) Non closed form { Kind of derivation } [Gradient Descent]
   ↳ Better for higher dimensions.

## OLS

$$b = \overline{y} - m\overline{x} \qquad m = \frac{\sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})}{\sum_{i=1}^{n} (x_i - \overline{x})^2}$$



error = $y_i - \hat{y_i}$

1) We have sort of linear points.
2) We want to find the best fit line.
3) We find the sum of the square of the residuals.

$$RSS = e_1^2 + e_2^2 + e_3^2 + \cdots + \cdots e_n^2$$

But why square & why not modulus ?

↳ Mod's graph is not differentiable.

$$\therefore E = e_1^2 + e_2^2 + \cdots e_n^2$$

$$E = \sum_{i=1}^{n} d_i^2$$

But, the errore we're finding it is in terms of $y$.

$\therefore$ $d_i = y_i - \widehat{y}_i$

$\widehat{y}_i = mx_i + b$

Total error
$$E = \sum_{i=1}^{n} (y_i - \widehat{y}_i)^2$$
$$= \sum_{i=1}^{n} (y_i - mx_i - b)^2$$

Average error
$$E = \frac{1}{2} \sum_{i=1}^{n} (y_i - \widehat{y}_i)^2$$

Now, our target is to minimize the error function & for that we need derivation.

$$\therefore E_{(m,b)} = \sum_{i=1}^{n} (y_i - mx_i - b)^2$$

Now,

$$\frac{\delta}{\delta b} E = 0$$

$$\Rightarrow \frac{\delta}{\delta b} E \equiv \sum_{i=1}^{n} \frac{\delta}{\delta b} (y_i - mx_i - b)^2 = 0$$

$$\equiv \sum_{i=1}^{n} (2(y_i - mx_i - b) \times -1) = 0$$

$$\equiv \sum_{i=1}^{n} (y_i - mx_i - b) = 0$$

$$\equiv \sum y_i - \sum mx_i - \sum b = 0$$

$$\equiv \frac{\sum y_i}{n} - \frac{\sum mx_i}{n} - \frac{\sum b}{n} = 0$$

$$\equiv \bar{y} - m\bar{x} - b = 0$$

$$\frac{\delta}{\delta m} E = \frac{\delta}{\delta m} \sum (y_i - mx_i - \bar{y} + m\bar{x})^2$$

$$= \sum (2(y_i - mx_i - \bar{y} + m\bar{x}) \times (-x_i + \bar{x}))$$

$$\equiv \sum [(y_i - \bar{y}) - m(x_i - \bar{x})](x_i - \bar{x}) = 0$$

$$\equiv \sum [(y_i - \bar{y})(x_i - \bar{x}) = m\sum(x_i - \bar{x})^2$$

$$\therefore m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})}$$

# Multiple Linear Regression

→ When multiple input columns exist.

Let's assume a dataset

| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| cgpa | iq | salary |

$$y = B_0 + B_1 X_1 + B_2 X_2 + \cdots + B_n X_n$$

$$y = B_0 + B_1 X_1 + B_2 X_3$$

$$y = B_0 + \sum_{i=1}^{n} B_i x_i$$

## * Mathematical Proof:

* I'll work with a dataset that has 9 input columns & one output column.

$$\hat{y} = B_0 + B_1 X_1 + B_2 X_2 + B_3 X_3$$

$$\hat{Y} = \begin{bmatrix} \hat{y_1} \\ \hat{y_2} \\ \vdots \\ \hat{y_n} \end{bmatrix} = \begin{bmatrix} B_0 & B_1 X_{11} & B_2 X_{12} & B_3 X_{13} & B_m X_{1m} \\ B_0 & B_1 X_{21} & B_2 X_{22} & B_3 X_{23} & B_m X_{2m} \\ \vdots & & & & \\ B_0 & B_1 X_{n1} & B_2 X_{n2} & B_3 X_{n3} & B_n X_{nm} \end{bmatrix}$$

$$= \begin{bmatrix} 1 & X_{11} & X_{12} & X_{13} & X_{1m} \\ 1 & X_{21} & X_{22} & X_{23} & X_{2m} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{n1} & X_{n2} & X_{n3} & X_{nm} \end{bmatrix} \begin{bmatrix} B_1 \\ B_2 \\ \vdots \\ B_m \end{bmatrix}$$

$\therefore$ Prediction Matrix = Unknown matrix $\times$ Co-efficient matrix

Error Matrix = Actual Y matrix − predicted Y matrix

$$Mat(e) = \begin{bmatrix} y_1 - \widehat{y_1} \\ y_2 - \widehat{y_2} \\ \vdots \\ y_3 - \widehat{y_3} \end{bmatrix}$$

Now, $E = e^T e$

$$\Rightarrow \sum_{i=1}^{n} (y_i - \widehat{y}_i)^2 = \begin{pmatrix} y_1 - \widehat{y_1} & y_2 - \widehat{y_2} & \cdots & y_n - \widehat{y_n} \end{pmatrix} \begin{pmatrix} y_1 - \widehat{y_1} \\ y_2 - \widehat{y_2} \\ \vdots \\ y_n - \widehat{y_n} \end{pmatrix}$$

$$(1 \times n)(n \times 1) = 1 \times 1$$

Now,

**Formulae**

$$E = e^T e = (Y - \widehat{Y})^T (Y - \widehat{Y})$$

$$= (Y^T - \widehat{Y}^T)(Y - \widehat{Y})$$

$(A \pm B)^T = A^T \pm B^T$

$$= (Y^T - (XB)^T)(Y - XB)$$

$(AB)^T = B^T A^T$

matrix diff formulae

$$= Y^T Y - \underbrace{Y^T X B - (XB)^T Y}_{Equal} + (XB)^T (XB)$$

$$y = A^T X A$$
$$\frac{dy}{dA} = 2XA^T$$

$$\boxed{E = Y^T Y - 2Y^T X B + (XB)^T (XB)}$$

$\hookrightarrow$ Loss function

Now, we need to find its lowest point using differentiation

$$\frac{dE}{dB} = \frac{d}{dB}\left[ Y^T Y - 2Y^T X B + B^T X^T X B \right] = 0$$
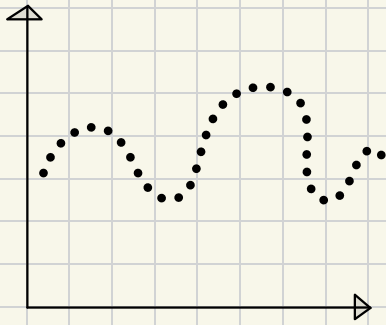
$$= 0 - 2Y^T X + 2X^T X B^T$$

$$((A^T A)^{-1})^T = A^T A^{-1}$$

Now, $2X^T X B^T = 2Y^T X$

$$B^T = Y^T X [X^T X]^{-1} \Rightarrow B = \left[ (X^T X)^{-1} \right]^T X^T Y$$

$$\therefore B = (X^T X)^{-1} X^T Y$$

# Polynomial Linear Regression



$$y = ax^2 + bx + c$$

$$\hat{y} = \hat{a}x^2 + \hat{b}x + c$$

$$e_i = (y_i - \hat{y}_i)$$
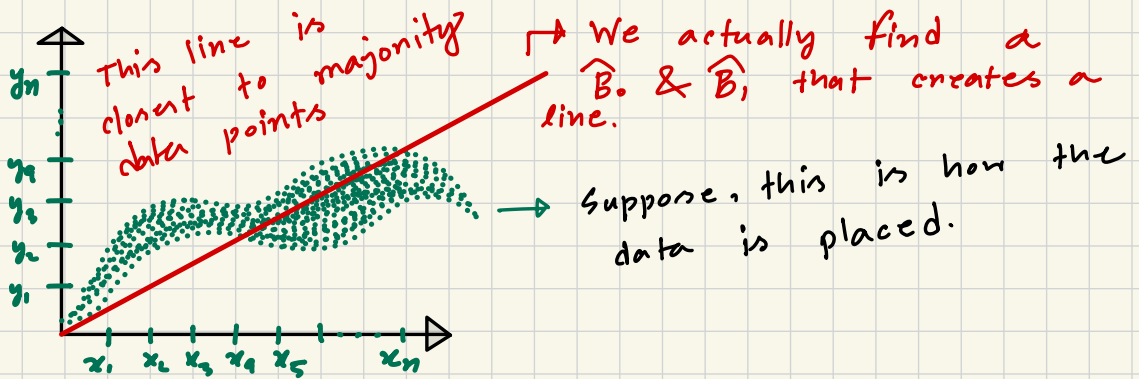
$$RSS = e_1^2 + e_2^2 + e_3^2 + \ldots\ldots + e_n^2$$

$$= \sum_{i=1}^{n} e_i^2$$

$$= \sum (ax^2 + bx + c - \hat{a}x^2 - bx - c)$$

$$y = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \theta_3 x^3 + \ldots\ldots \theta_n x^n$$

Q: Find the best fit second degree polynomial for the given data $\{(1,3), (2,4), (3,8)\}$

→ we had 200 datapoints. (of two dimennional data)

→ We anumed some nelationnhip $y \approx B_0 + B_1 X$

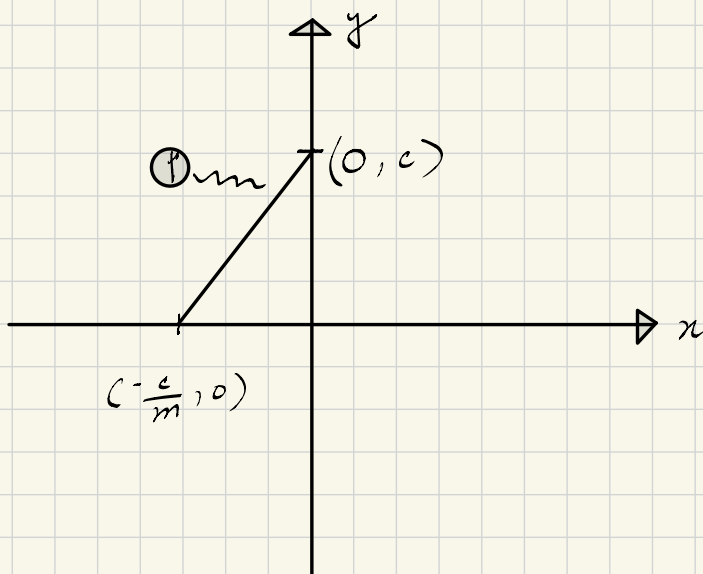→ Data points were $(x_1, y_1), (x_2, y_2), (x_3, y_3) \dots (x_n, y_n)$

This line is clonent to majonity data points

We actually find a $\hat{B}_0$ & $\hat{B}_1$ that creates a line.

Suppose, this is how the data is placed.



- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Note:** The nimple linear equation is $y = mx + c$ where $m$ is the slope and $c$ is the y inteneept.

If we want to plot this line, we'll have to

$$y = mx + c$$

$$\Rightarrow mx - y = -c \Rightarrow \frac{mx}{-c} + \frac{y}{c} = 1 \Rightarrow \frac{x}{-c/m} + \frac{y}{c} = 1$$

$$\longmapsto ①$$

$$① \quad (0, c)$$

$$\left(-\frac{c}{m}, 0\right)$$

$$\sum_{i=0}^{n} 3i - 1 = (3 \times 0 - 1) + (3 \times 1 - 1) + (3 \times 2 - 1)$$
$$+ \cdots + (3 \times n - 1)$$

$n \rightarrow$ End value

$i = 0 \quad \llcorner\rightarrow$ Start value

Now, $\hat{y}_i = \hat{B}_0 + \hat{B}_1 x_i$ predicted value of $y$ based on the $i$th value of $x$.

$y_i =$ Observed $i$th response

Residue $e_i = y_i - \hat{y}_i$

$\therefore$ RSS (Residual Sum of Squares) $= e_1^2 + e_2^2 + \ldots + e_n^2$

We try to minimize RSS through $\hat{B}_0$ & $\hat{B}_1$

$$\hat{B}_1 = \frac{\sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})}{\sum_{i=0}^{n} (x_i - \overline{x})^2}$$

$$\hat{B}_0 = \overline{y} - \hat{B}_1 x$$

$\Big\}$ 3.4

$\llcorner\rightarrow$ We can use this as formulas for the slope & Intercept. This $\hat{B}_1$ & $\hat{B}_0$ give the smallest RSS.

# ✳ Assessing Precision

Standard error $\Rightarrow$

$$SE\left(\hat{B_1}\right) = \frac{\sigma^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \quad \rightarrow \text{Standard error for slope}$$

$$SE\left(\hat{B_0}\right) = \sigma^2\left[\frac{1}{n} + \frac{\bar{x}^2}{n \sum_{i=1}^{} (x_i - \bar{x})^2}\right]$$

$\llcorner\!\rightarrow$ Standard error for intercept

$$\sigma^2 = Var(E)$$