# Linear Regression

- Goal of regression is to estimate the impact of the variation in X (independant variables) on central tendency of Y (outcome variable).
- There are many measures of central tendency, for now let us consider the mean:

$$\text{Population: } \mu = E(Y)$$

$$\text{Sample: } X_1, X_2, \ldots, X_n \implies \hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} X_i$$

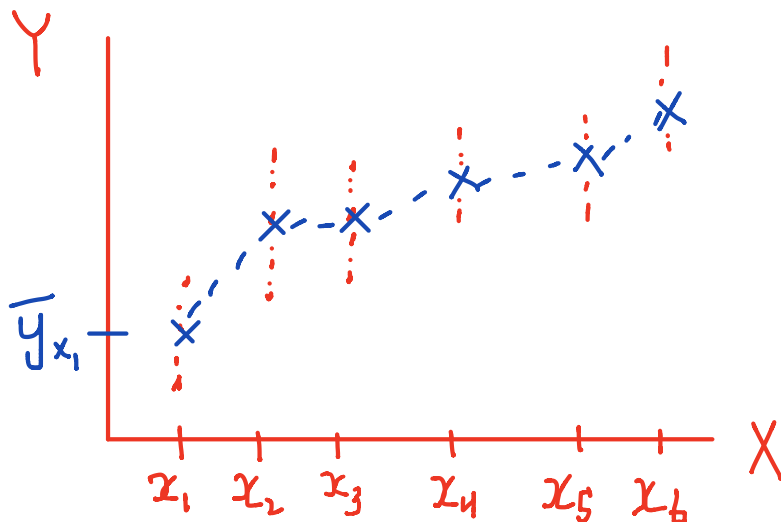- Although the above definitions are intuitive, an equivalent definition of the mean is as follows:

$$\text{Population: } \min_{\mu} E\left[(Y-\mu)^2\right] \implies \mu = E(Y)$$

$$\text{Sample: } \min_{\mu} \frac{1}{n} \sum_{i=1}^{n} (X_i - \mu)^2 \implies \hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} X_i$$

- The relationship between X and Y can be summarized by the conditional expectation function E(Y|X). The CEF is interesting because it measures the mean tendency of Y|X (Y given X).

$$\min_{\mu(x)} E\left[(y-\mu(x))^2 \mid X\right] \implies \mu(x) = E(Y|X)$$

- The CEF is a function of X and represents the average of Y for a given value of X.



$$\underline{CEF: } E[Y \mid X=x]$$

$$\bar{y}_{x_1} = \frac{1}{n_{x_1}} \sum_{j:X=X_1} y_j$$

$$\underbrace{\phantom{\frac{1}{n_{x_1}} \sum_{j:X=X_1} y_j}}_{\text{Avg. of Y at } X=X_1}$$
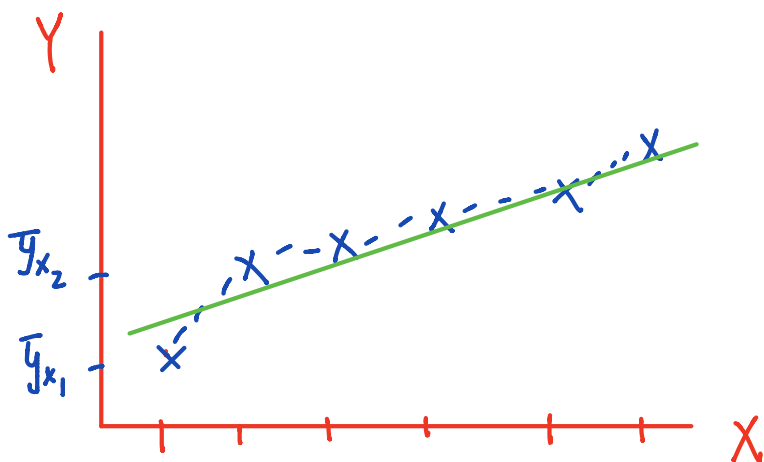
- It can be shown that Y can be decomposed into the CEF plus an orthogonal (unrelated to X) error term.

$$Y_i = E[Y_i \mid X_i] + \varepsilon_i, \quad E[\varepsilon_i \mid X_i] = 0$$

$$\underbrace{\phantom{Y_i}}_{\text{Outcome}} \quad \underbrace{\phantom{E[Y_i|X_i]}}_{\text{CEF}} \quad \underbrace{\phantom{\varepsilon_i}}_{\text{Error}} \quad \underbrace{\phantom{E[\varepsilon_i|X_i]=0}}_{\text{Mean Independance}}$$

- Linear regression imposes parametric assumption on the CEF.

$$E[Y_i|X_i] = \beta X_i \implies Y_i = \beta X_i + \varepsilon_i$$

- We are essentially approximating the CEF with the standard linear regression model.



$$\underline{OLS:} \quad \hat{Y}_i = \hat{\beta} X_i$$

$$\underline{CEF:} \quad \{\bar{Y}_{x_1}, \bar{Y}_{x_2}, \dots\}$$

- We can estimate the linear regression parameter as follows:

$$\text{Since } \mu(x) = \beta X \implies \min_{\beta} E\left[(y - X\beta)^2\right]$$

$$\implies \min_{\beta} \frac{1}{n} \sum_{i=1}^{n} (y_i - \beta X_i)^2 \implies \hat{\beta} = \frac{S_{xy}}{S_x S_y}$$

$$\text{where } S_{xy} = \text{sample } COV(X_i, Y), \quad S_x \, \& \, S_y \text{ are stdev.}$$

- We estimate the linear regression parameter above by minimizing the sum of square errors. This uses the square loss function (also known as L2 loss).
- As discussed above, using the square loss will result in a model that approximates the CEF. Using a different loss function will result in a different interpretation of the model.
- If we have multiple regressors x1, x2, ..., xk, all the math above can be extended using matrices. Now the parameter vector can be estimated as follows:

$$\text{Since } \mu(X_1, \dots, X_R) = \beta_0 + \beta_1 X_1 + \dots + \beta_R X_R = X\beta$$

$$\implies \min_{\beta} \frac{1}{n} (Y - X'\beta)'(Y - X'\beta) \implies \hat{\beta} = (X'X)^{-1} X'y$$

- Note above that above dimensions are: dim(beta) = (k+1)x1, dim(X) = n x (k+1), and dim(y) = n x 1