

Stat 134 Sec 3.9

Sec 6.4 Correlation

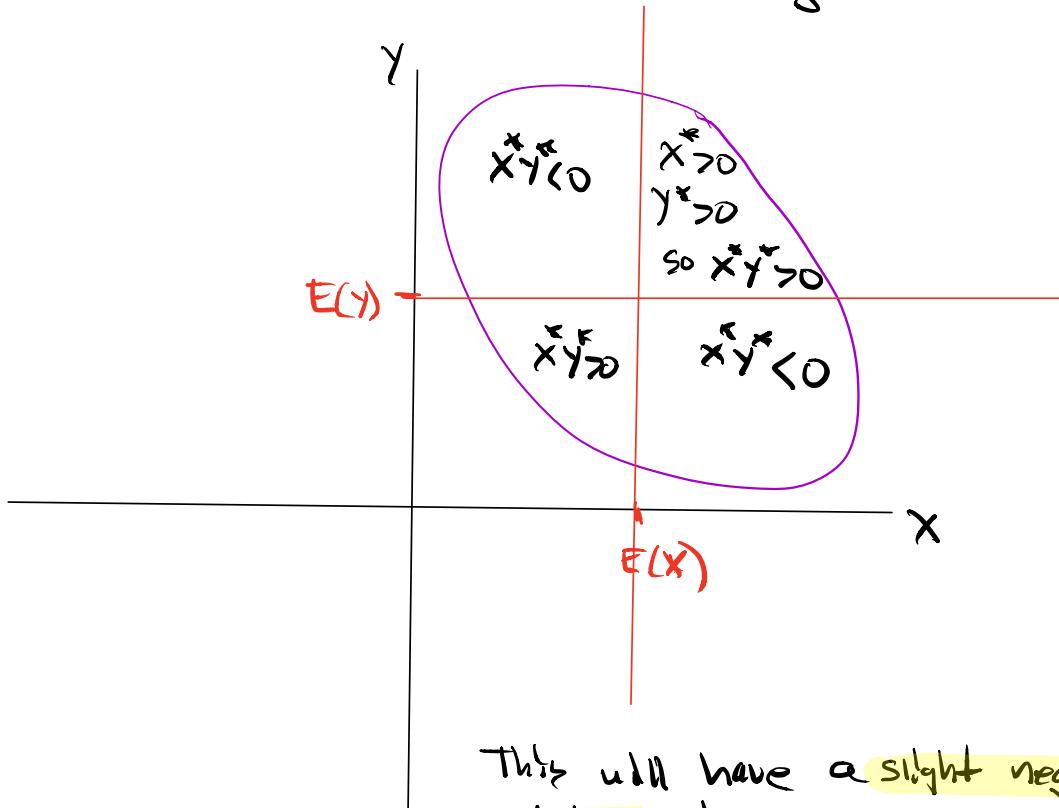
last time

$$\text{Cov}(x, y) = E((x - \mu_x)(y - \mu_y))$$

$$\text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\text{SD}(x)\text{SD}(y)} = E\left(\left(\frac{x - \mu_x}{\text{SD}_x}\right)\left(\frac{y - \mu_y}{\text{SD}_y}\right)\right) = E(x^* y^*)$$

x, y in S.U.

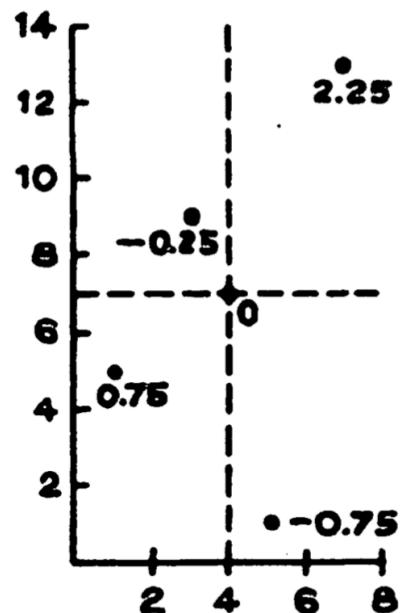
so suppose you have a scatter diagram



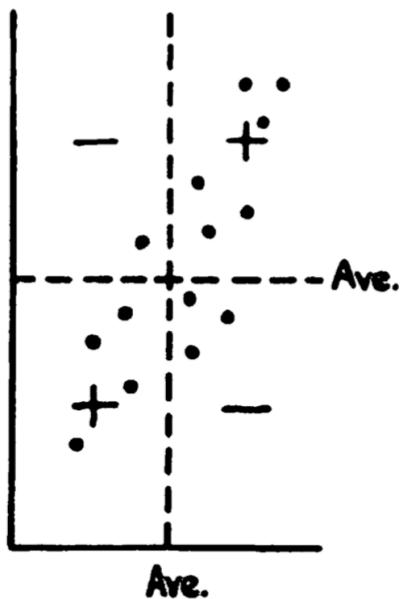
This will have a slight negative correlation since more of the points are in the 2nd and 4th quadrants (of the red axes).

Figure 9. How the correlation coefficient works.

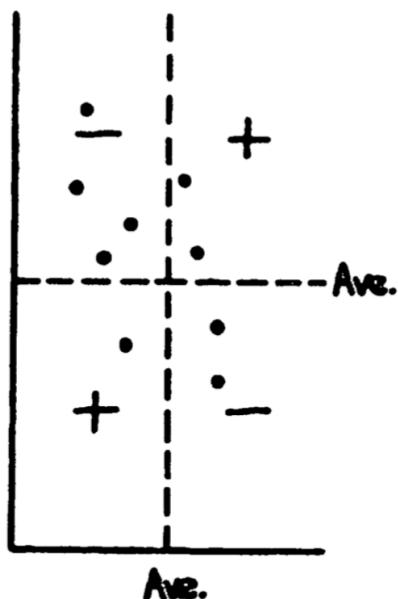
(a) Scatter diagram
from Table 1



(b) Positive r



(c) Negative r



Properties of correlation

- 1) Pure numbers ✓
- 2) Correlation is invariant to linear changes at scale except possibly by a sign.

In other words:

$$\text{corr}(x, y) = |\text{corr}(ax+b, cy+d)| \quad \begin{matrix} \text{for constants } a, b, c, d \\ \text{with } a \neq 0, c \neq 0. \end{matrix}$$

Proof

$$|\text{corr}(ax+b, cy+d)| = \frac{|\text{cov}(ax+b, cy+d)|}{\text{SD}(ax+b)\text{SD}(cy+d)} \\ = \frac{|ac\text{cov}(x, y)|}{|a||c|\text{SD}(x)\text{SD}(y)} = \text{corr}(x, y).$$

3) $\text{corr}(x, y) = \text{corr}(y, x)$ ✓

4) $-1 \leq \text{corr}(x, y) \leq 1$

Proof

$$\left. \begin{array}{l} E(x^*) = 0 = E(y^*) \\ \text{SD}(x^*) = 1 = \text{SD}(y^*) \\ E(x^{*2}) = 1 = E(y^{*2}) \end{array} \right\} \text{since } x^*, y^* \text{ are standard units}$$

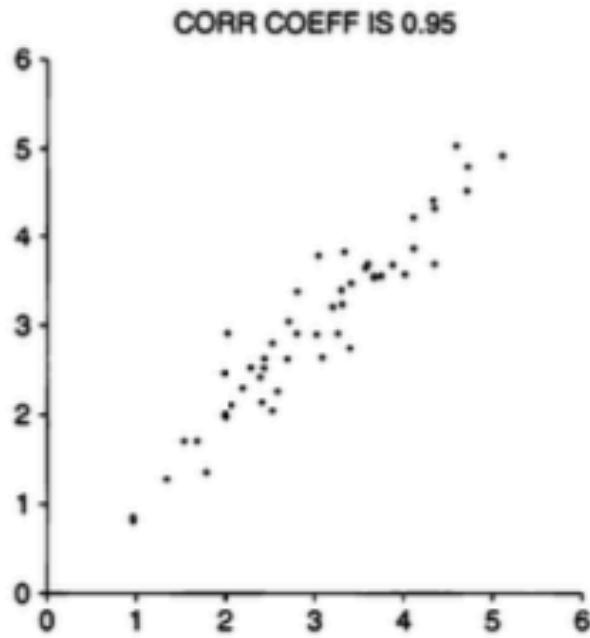
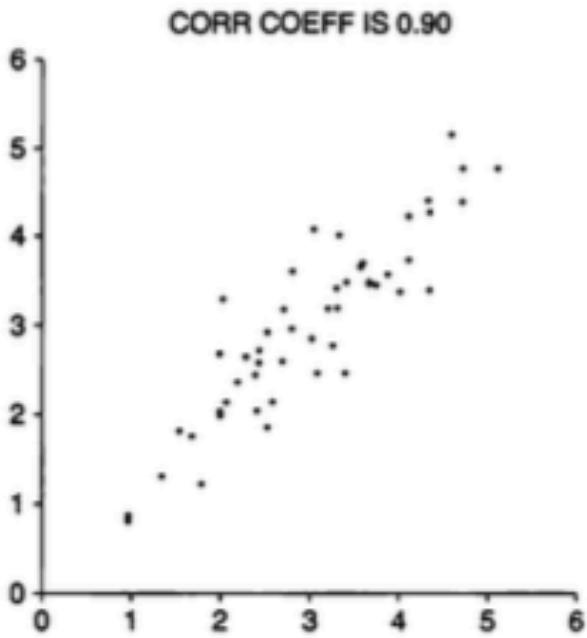
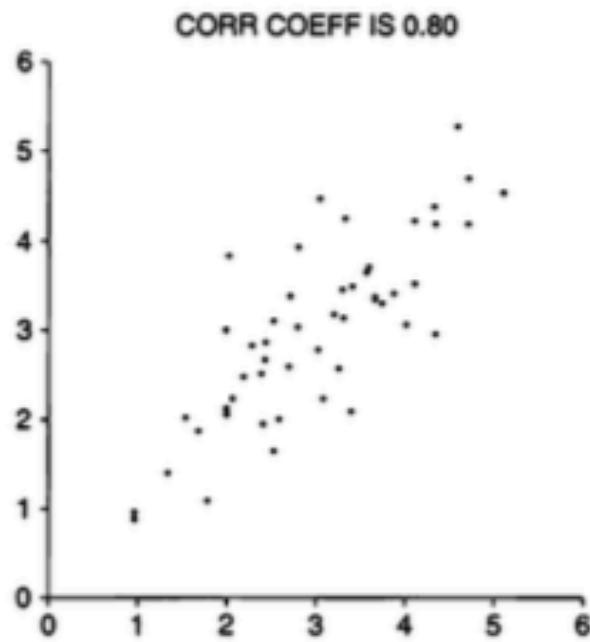
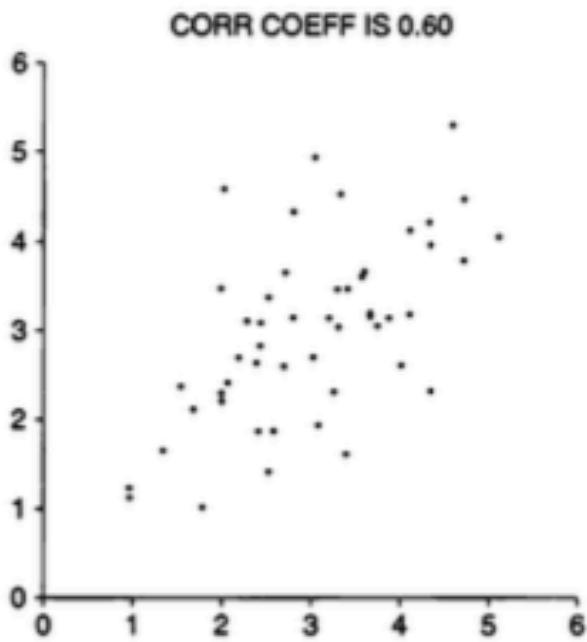
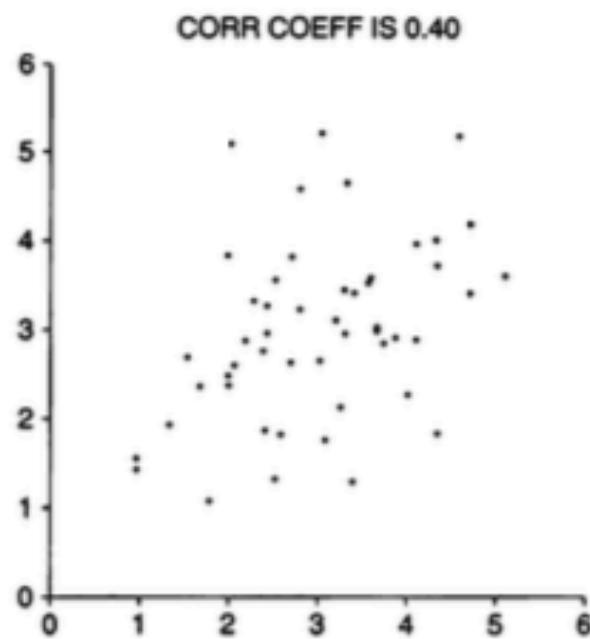
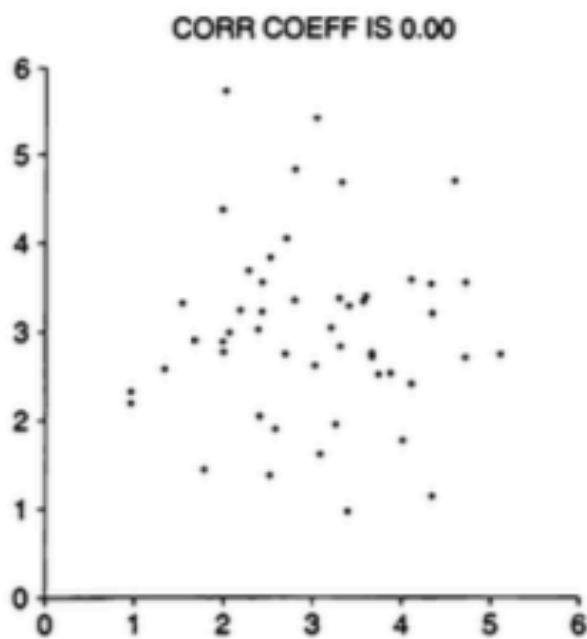
$$\begin{array}{ll} (x^* + y^*)^2 \geq 0 & | \quad (x^* + y^*)^2 \geq 0 \\ \text{so } E((x^* + y^*)^2) \geq 0 & | \quad \text{so } E((x^* + y^*)^2) \geq 0 \\ E(x^{*2} + y^{*2} + 2x^*y^*) \geq 0 & | \quad E(x^{*2} + y^{*2} + 2x^*y^*) \geq 0 \\ 1 + 1 + 2E(x^*y^*) \geq 0 & | \quad 1 + 1 + 2E(x^*y^*) \geq 0 \\ E(x^*y^*) \geq -1 & | \quad E(x^*y^*) \geq -1 \end{array}$$

$$\Rightarrow -1 \leq E(x^*y^*) \leq 1$$

$-1 \leq \text{corr}(x, y) \leq 1$

- 5) Correlation measures linear association, how tightly clustered scattered diagram is around a straight line.

see next page



\Leftarrow Roll a die n times.

$$N_1 = \# \text{ of ones}$$

$$N_6 = \# \text{ of sixes}.$$

Find $\text{Corr}(N_1, N_6)$

method 1 (Indicators, bilinearity of Cov)

$$N_1 = \sum_{j=1}^n I_{1,j} \leftarrow \text{get 1 on } j\text{th roll}$$

$$N_6 = \sum_{j=1}^n I_{6,j} \leftarrow \text{get 6 on } j\text{th roll.}$$

$$\text{Cov}(N_1, N_6) = \text{Cov}\left(\sum_{j=1}^n I_{1,j}, \sum_{j=1}^n I_{6,j}\right)$$

$$= n \cdot \text{Cov}(I_{1,1}, I_{6,1}) + n(n-1)\text{Cov}(I_{1,1}, I_{6,2}) \\ \text{etc}$$

method 2 (Var of $N_1 + N_6$)

$$N_1 \sim \text{Bin}(n, \frac{1}{6}) \quad \text{Var}(N_1) = n\left(\frac{1}{6}\right)\left(\frac{5}{6}\right)$$

$$N_6 \sim \text{Bin}(n, \frac{1}{6})$$

$$N_1 + N_6 \sim \text{Bin}\left(n, \frac{2}{6}\right) \quad \text{Var}(N_1 + N_6) = n\left(\frac{2}{6}\right)\left(\frac{4}{6}\right)$$

$$\text{Var}(N_1 + N_6) = \text{Var}(N_1) + \text{Var}(N_6) + 2\text{Cov}(N_1, N_6)$$

$$\text{Cov}(N_1, N_6) = n\left(\frac{1}{6}\right)\left(\frac{4}{6}\right) - 2n\left(\frac{1}{6}\right)\left(\frac{5}{6}\right) = -\frac{n}{36}$$

$$\text{Corr}(N_1, N_6) = \frac{-\frac{n}{36}}{\sqrt{n\left(\frac{1}{6}\right)\left(\frac{5}{6}\right)}}^2 = \frac{-\frac{n}{36}}{\sqrt{n\left(\frac{1}{6}\right)\left(\frac{5}{6}\right)}} = \frac{-\frac{1}{5}}{\frac{1}{5}}$$

More generally; (sum of exchangeable RVs)
 See p238

Suppose the sum of K exchangeable (i.e. identically distributed) RVs is a constant

$$N_1 + N_2 + \dots + N_K = C$$

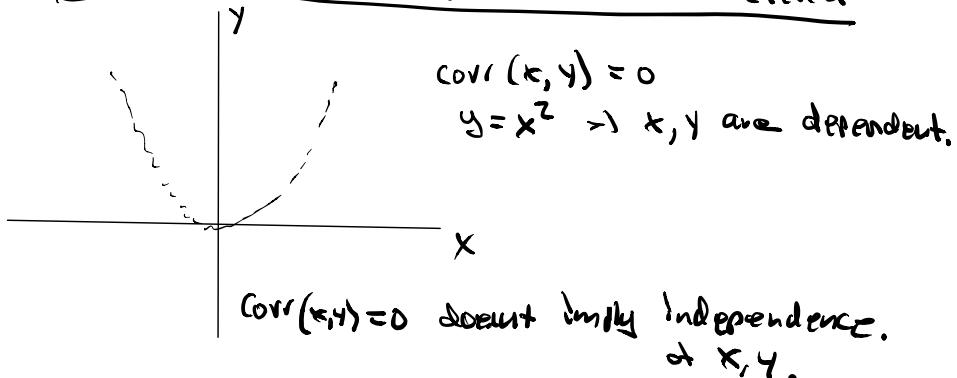
then $\boxed{\text{Corr}(N_1, N_2) = -\frac{1}{K-1}}$

Ex You roll a die n times.
 Let N_i = the number of rolls where you get i

$$N_1 + N_2 + N_3 + N_4 + N_5 + N_6 = n$$

then $\text{Corr}(N_1, N_2) = -\frac{1}{5} \quad \checkmark$

Conditions for X and Y to be uncorrelated



$$\left. \begin{aligned} \text{Corr}(x, y) &= 0 \\ \text{Cov}(x, y) &= 0 \\ E(xy) &= E(x)E(y) \end{aligned} \right\} \text{equivalent.}$$

$$x, y \text{ indep} \Rightarrow E(xy) = E(x)E(y) \Leftrightarrow \text{Corr}(x, y) = 0.$$

Stat 134

Friday April 23 2018

1. Consider a Poisson(λ) process. Let $T_r \sim \text{gamma}(r, \lambda)$ be the rth arrival time. $\text{Corr}(T_1, T_3)$ equals:

$$\underline{\text{review}} \quad \text{if } T_r = \text{gamma}(r, \lambda)$$

$$E(T_r) = \frac{r}{\lambda}$$

$$SD(T_r) = \sqrt{r}/\lambda.$$



$$T_3 = T_1 + w$$

d none of the above

$$\text{Cov}(T_1, T_2) = \text{Cov}(T_1, T_1 + w)$$

$$= \text{Cov}(T_1, T_1) + \text{Cov}(T_1, w)$$

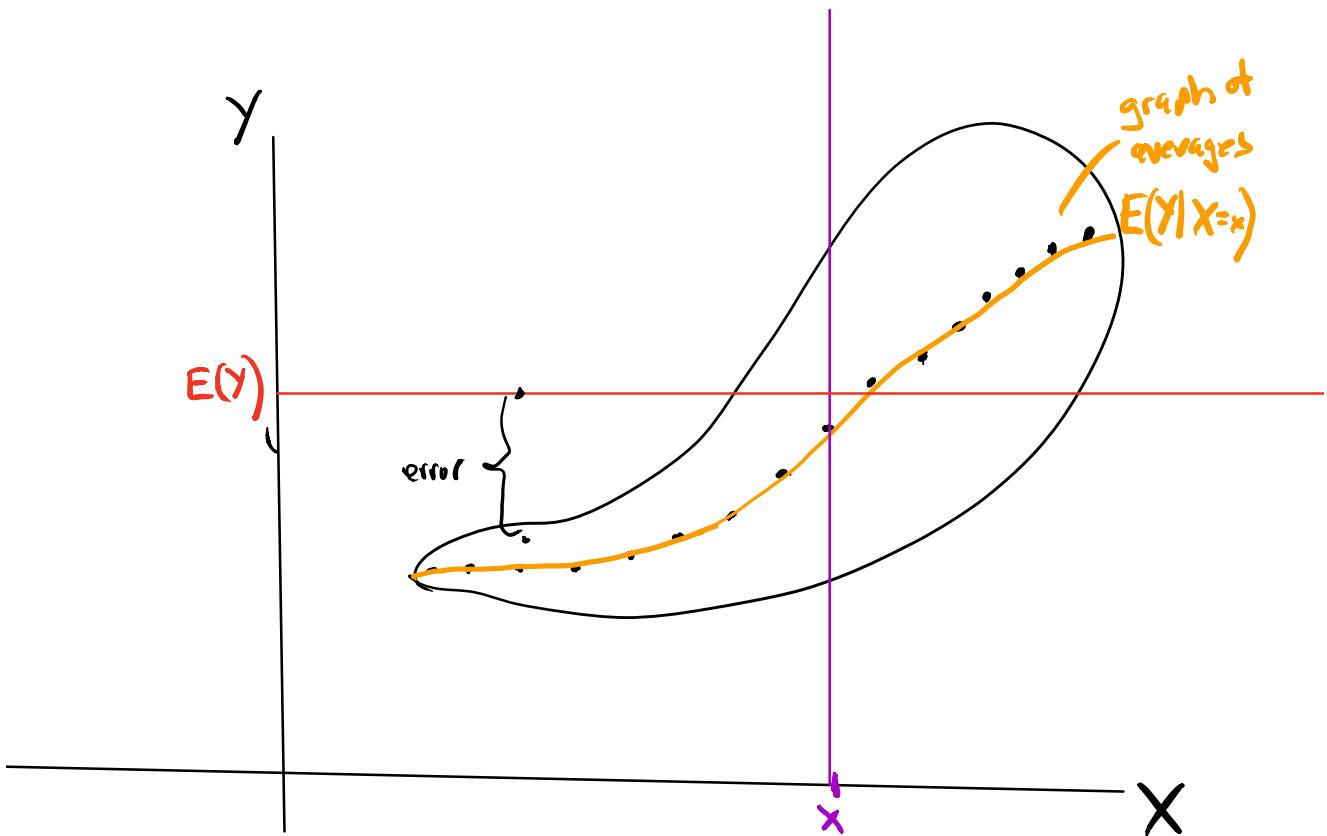
$$\text{Var}(T_1)$$

$$0$$

$$\frac{1}{\lambda^2}$$

$$\text{Cov}(T_1, T_3) = \frac{\text{Cov}(T_1, T_1 + w)}{SD(T_1)SD(T_3)} = \frac{\frac{1}{\lambda^2}}{\frac{1}{\lambda} \cdot \frac{1}{\sqrt{3}}} = \frac{1}{\sqrt{3}}$$

Predicting Y from X



We wish to predict Y from X .
What is the best predictor? — ^{smallest} mean square error.
(mse)

Predictor #1 $E(Y)$ straight line.
for every x same prediction.

$$\text{error} = Y - E(Y)$$

$$E(\text{error}) = E(Y) - E(Y) = 0$$

$$E(\text{error}^2) = \text{Var}(\text{error}) = \text{Var}(Y - E(Y)) \\ = \boxed{\text{Var}(Y)}$$

Predictor #2

graph of averages , $E(Y|X)$ → see picture above.

error = $Y - E(Y|X)$

Predictor
2

$$mse = E(\text{error}^2) = E((Y - E(Y|X))^2)$$

We will show that not only is $E(Y|X)$ a better predictor of Y than $E(Y)$ but it is a better predictor of Y than any function $b(X)$.

That is cool !!

Stay tuned...