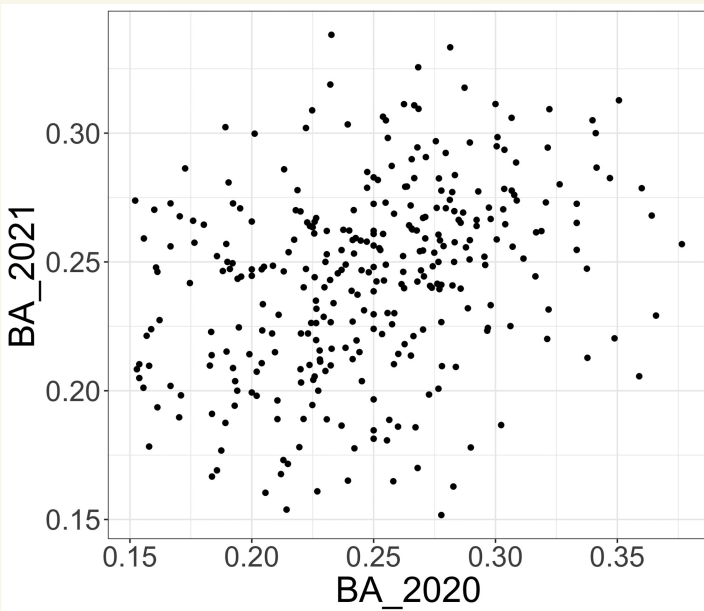


Simple Linear Regression

Q Suppose we have access to each MLB player's 2020 batting average and 2021 batting average, and no other info, Predict BA_{2021} from BA_{2020} .

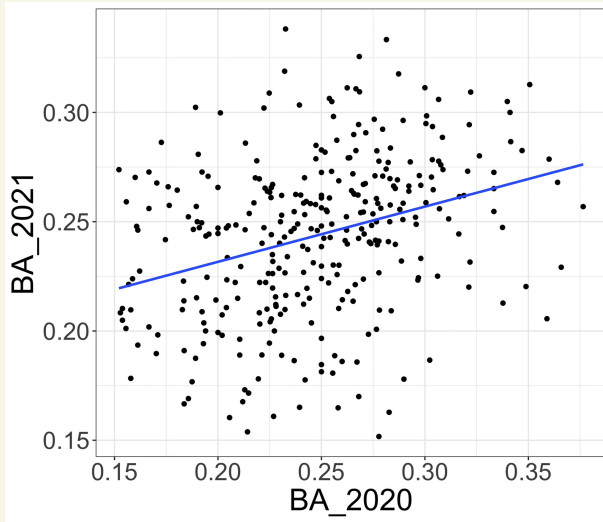
Always good to begin with plotting visualization,



Looks linear-ish

Positive slope

can imagine drawing a best fit line



How do we get a best fit line?

Model index each baseball player
in the dataset by i

let $X_i = BA_i^{(2020)}$ predictor variable

$Y_i = BA_i^{(2021)}$ response variable
outcome

Assume a Linear Relationship

$$Y_i = \beta_0 + \beta_1 \cdot X_i + \varepsilon_i$$

b + mx

noise

ε_i is a random variable (noise)
with $\mathbb{E}\varepsilon_i = 0$

β_0, β_1 are unknown constants

goal is to estimate β_0, β_1
to get the best fit line

$$\begin{aligned}\hat{Y}_i &:= \mathbb{E}(Y_i | X_i) = \mathbb{E}(\beta_0 + \beta_1 X_i + \varepsilon_i) \\ &= \beta_0 + \beta_1 X_i + \cancel{\mathbb{E}(\varepsilon_i)} \\ &\approx \hat{\beta}_0 + \hat{\beta}_1 X_i\end{aligned}$$

How to estimate β_0 and β_1 ?

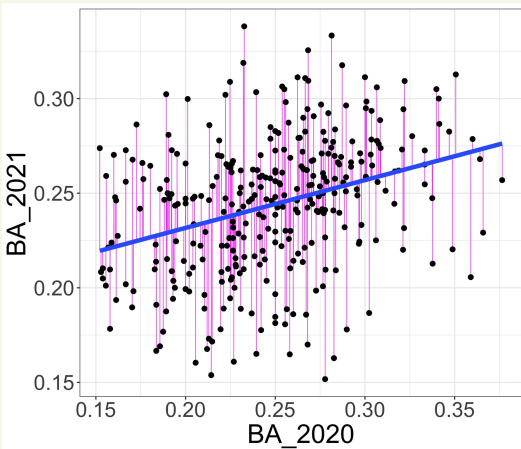
Ordinary least squares (OLS) — choose the values of β_0, β_1 which minimize the Residual Sum of Squares

i.e. minimize mean square error (MSE)

$$\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i))^2$$

$$RSS(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

$$\hat{\beta}_0, \hat{\beta}_1 = \arg \min_{(\beta_0, \beta_1)} RSS(\beta_0, \beta_1)$$



$$\hat{\beta}_0, \hat{\beta}_1 = \arg \min_{(\beta_0, \beta_1)} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Calculus: to minimize a function, set the derivative equal to 0 and solve.

$$\begin{cases} \frac{d}{d\beta_0} \text{RSS} = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i) \cdot (-1) = 0 \\ \frac{d}{d\beta_1} \text{RSS} = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i) \cdot (-x_i) = 0 \end{cases}$$

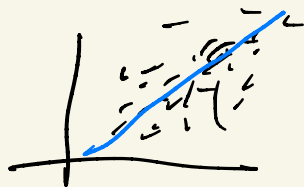
$$\frac{1}{n} \sum_{i=1}^n (y_i - \beta_1 x_i) = \frac{1}{n} \sum_{i=1}^n \beta_0 \implies \bar{y} - \beta_1 \bar{x} = \beta_0$$

$$-\frac{1}{n} \sum x_i y_i + \underbrace{\beta_0 \frac{1}{n} \sum x_i}_{\beta_0 \bar{x}} + \beta_1 \frac{1}{n} \sum x_i^2 = 0$$

$$\underbrace{(\bar{y} - \beta_1 \bar{x}) \bar{x}}_{(\bar{y} - \beta_1 \bar{x}) \bar{x}}$$

$$\beta_1 \left(\frac{1}{n} \sum x_i^2 - (\bar{x})^2 \right) = \frac{1}{n} \sum x_i y_i - \bar{x} \bar{y}$$

$$\left\{ \begin{aligned} \beta_1 &= \frac{\frac{1}{n} \sum X_i Y_i - \bar{X} \bar{Y}}{\frac{1}{n} \sum X_i^2 - (\bar{X})^2} \approx \frac{\text{COV}(X, Y)}{\text{VAR}(X)} \\ \beta_0 &= \bar{Y} - \beta_1 \bar{X} \end{aligned} \right.$$



$$\text{VAR}(X) = \mathbb{E}(X - \mathbb{E}X)^2 = \mathbb{E}X^2 - \mathbb{E}X)^2$$

$$\text{COV}(X, Y) = \mathbb{E}[(X - \mathbb{E}X) \cdot (Y - \mathbb{E}Y)]$$

$$= \mathbb{E}[X \cdot Y - X \cdot \mathbb{E}Y - \mathbb{E}X \cdot Y + \mathbb{E}X \cdot \mathbb{E}Y]$$

$$= \mathbb{E}[X \cdot Y] - \mathbb{E}[X \cdot \mathbb{E}Y] - \mathbb{E}[\mathbb{E}X \cdot Y] + \mathbb{E}[\mathbb{E}X \cdot \mathbb{E}Y]$$

$$= \mathbb{E}[X \cdot Y] - \mathbb{E}Y \mathbb{E}X - \mathbb{E}X \mathbb{E}Y + \mathbb{E}X \mathbb{E}Y$$

$$\approx \sum_{i=1}^n (X_i - \bar{X}) \cdot (Y_i - \bar{Y})$$

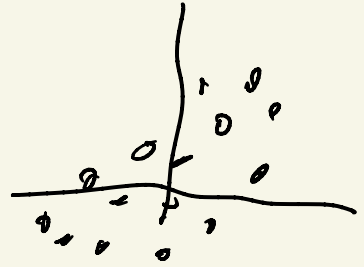
$$\text{COV}(X, Y) = \mathbb{E}(XY) - \mathbb{E}X \mathbb{E}Y$$

$$\approx \frac{1}{n} \sum X_i Y_i - \bar{X} \bar{Y}$$

pretend $\mathbb{E}X=0$, $\mathbb{E}Y=0$

$$\text{COV}(X, Y) = \mathbb{E}(X \cdot Y)$$

X, Y positive covariance:



$$\mathbb{E}(X \cdot Y) > 0$$

$X \cdot Y > 0$ on average

then on average,

when $X > 0$ $Y > 0 \rightarrow X \cdot Y > 0$

when $X < 0$ $Y < 0 \rightarrow X \cdot Y > 0$

X, Y negative covariance:

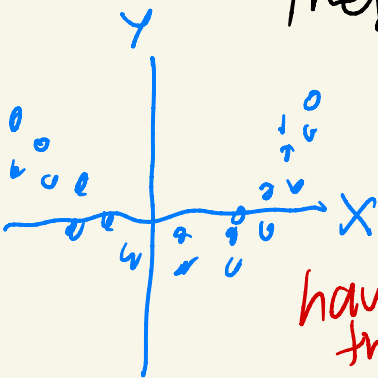
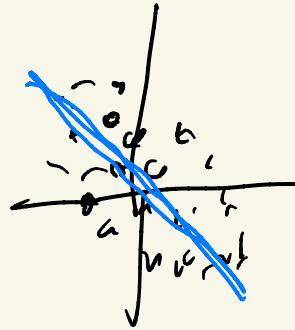
$$\mathbb{E}(X \cdot Y) < 0$$

$X \cdot Y < 0$ on average

then on average,

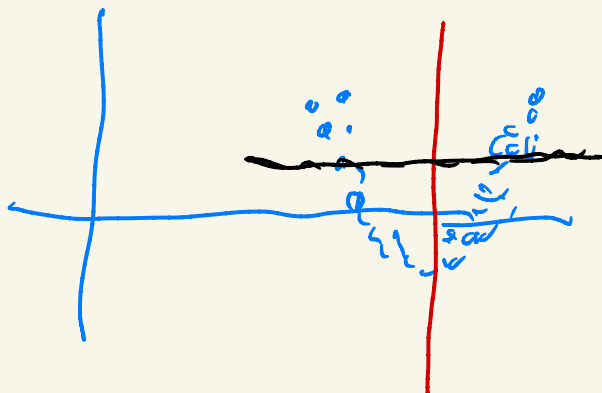
when $X < 0$ $Y > 0$

when $X > 0$ $Y < 0$



having 0 covariance does not mean there's no relationship.

means can't glean much info
abt the relationship of the signs
(assuming mean 0)



Covariance is a measure of
linear relationship

slope

$$\hat{\beta}_1 \approx \frac{\text{cov}(X, Y)}{\text{var}(X)}$$

$$\left\{ \begin{aligned} \beta_1 &= \frac{\frac{1}{n} \sum X_i Y_i - \bar{X} \bar{Y}}{\frac{1}{n} \sum X_i^2 - (\bar{X})^2} = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2} \\ \beta_0 &= \bar{Y} - \beta_1 \bar{X} \end{aligned} \right.$$

Showed $\text{VAR}(X) = \mathbb{E}(X - \mathbb{E}X)^2 = \mathbb{E}X^2 - (\mathbb{E}X)^2$

$$\text{COV}(X, Y) = \mathbb{E}[(X - \mathbb{E}X) \cdot (Y - \mathbb{E}Y)]$$

$$= \mathbb{E}(X \cdot Y) - \mathbb{E}X \mathbb{E}Y$$

Sample covariance

$$\sigma_{XY} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

Sample variance

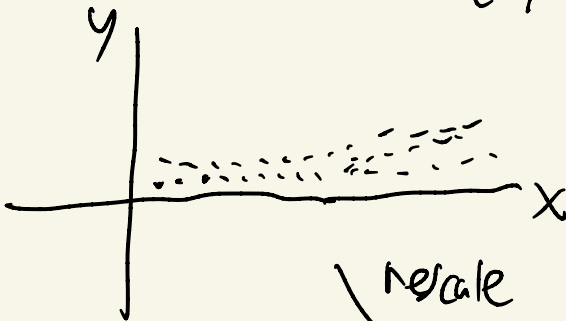
$$\sigma_X^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

Correlation

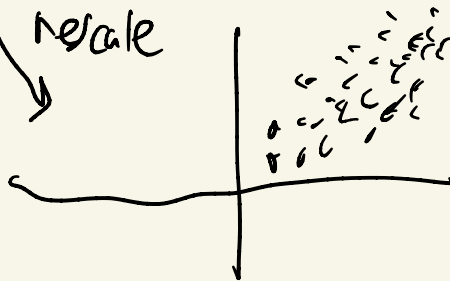
$$\text{COR}(X, Y) = \frac{\text{COV}(X, Y)}{\text{SD}(X) \cdot \text{SD}(Y)}$$

$$\text{SD}(X) \cdot \text{SD}(Y)$$

$$-1 \leq \text{COR}(X, Y) \leq 1$$



↙ relate



Sample correlation $r_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

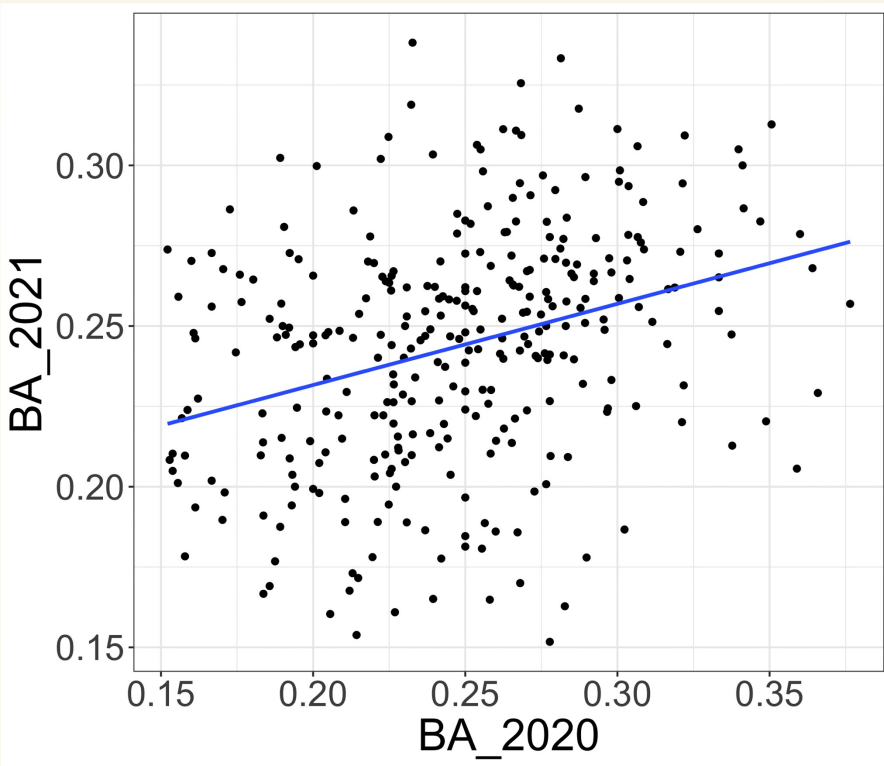
$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \cdot \frac{\sqrt{\sum_i (y_i - \bar{y})^2}}{\sqrt{\sum_i (y_i - \bar{y})^2}}$$

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} \cdot \frac{\sqrt{\sum (y_i - \bar{y})^2}}{\sqrt{\sum (y_i - \bar{y})^2}}$$

$$= r_{xy} \cdot \frac{\sigma_y}{\sigma_x} = \hat{\beta}_1$$

$$r_{xy} = \hat{\beta}_1 \frac{\sigma_x}{\sigma_y}$$

If x, y are normalized (have same) variance
then linear regression slope is the sample correlation.



$$\hat{\beta}_1 = \frac{1}{4}$$

$$E(x, y)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}, \quad \hat{\beta}_1 = \frac{1}{4}$$

$$\bar{x} \approx 0.25 = \frac{1}{4}, \quad \bar{y} \approx \frac{1}{4}$$

$$\beta_1 = \frac{1}{4}, \quad \beta_0 = \frac{3}{16}$$

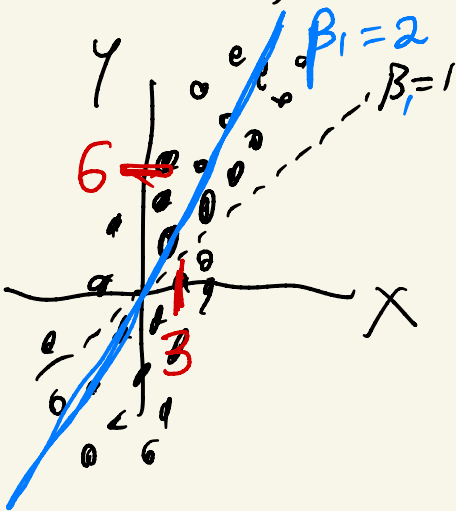
$$\hat{y} = \beta_0 + \beta_1 \cdot x \approx \frac{3}{16} + \frac{1}{4}x$$

Regression to the mean

if $x = \bar{x} = \frac{1}{4}$, then $\hat{y} = \frac{1}{4}$

if $x > \bar{x}$, $x = 0.3$, $\hat{y} = \frac{13}{48} > \frac{12}{48}$

if $x < \bar{x}$, $x = 0.2$, $\hat{y} = \boxed{\frac{19}{80}} < \frac{20}{80}$



- Regression code
- 2nd example → Pythagorean win Percentage

Pythagorean Win Percentage

RA = runs allowed by a (baseball) team in one season

RS = runs scored in season

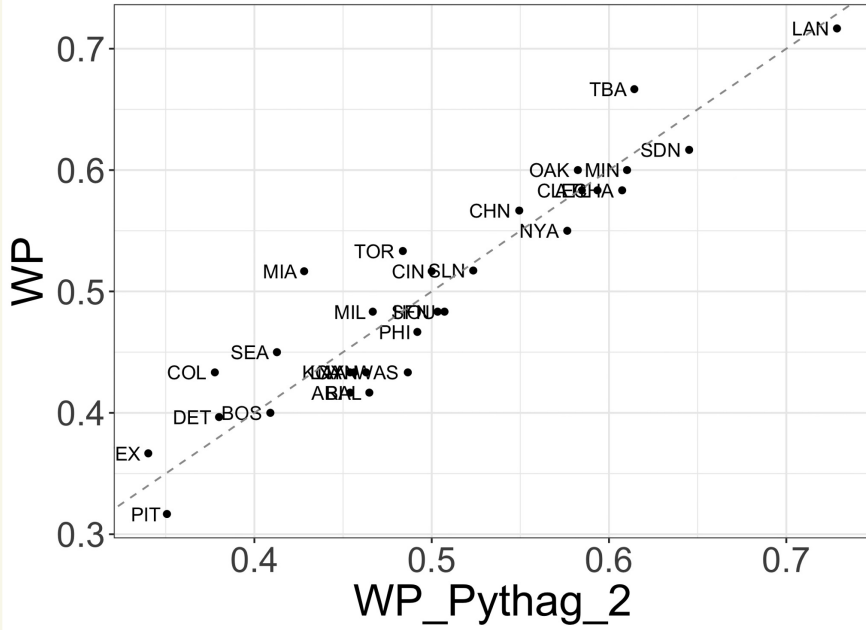
WP = team's win percentage that season

Bill James

$$\widehat{WP} = \frac{RS^2}{RS^2 + RA^2}$$

Pythagorean

2020 win percentage vs. pythagorean win percentage



$$\widehat{WP} = \frac{RS^\alpha}{RS^\alpha + RA^\alpha}$$

find α ?

Model $E WP = \frac{RS^\alpha}{RS^\alpha + RA^\alpha} = \frac{1}{1 + \left(\frac{RA}{RS}\right)^\alpha}$

$$\Rightarrow 1 + \left(\frac{RA}{RS}\right)^\alpha = \frac{1}{EWP}$$

$$\Rightarrow \left(\frac{RA}{RS}\right)^\alpha = \frac{1}{EWP} - 1$$

$$= \frac{1 - EWP}{EWP}$$

$$\Rightarrow \log \left(\frac{RA}{RS}\right)^\alpha = \log(\cdot)$$



$$\beta_0 + \alpha \cdot \log\left(\frac{RA_i}{RS_i}\right) = \log\left(\frac{1 - WP_i}{WP_i}\right)$$

Diagram illustrating the mapping of variables in the final equation:

- β_0 (intercept) is indicated by a red arrow pointing to the constant term.
- β_1 (slope) is indicated by a red arrow pointing to the coefficient α .
- X_i (independent variable) is indicated by a red arrow pointing to the log term $\log\left(\frac{RA_i}{RS_i}\right)$.
- Y_i (dependent variable) is indicated by a red arrow pointing to the log term $\log\left(\frac{1 - WP_i}{WP_i}\right)$.

Simple linear regression!

close to 2

$$\hat{\alpha} = 1.867$$

using 2020 data
pretty close to 2!

