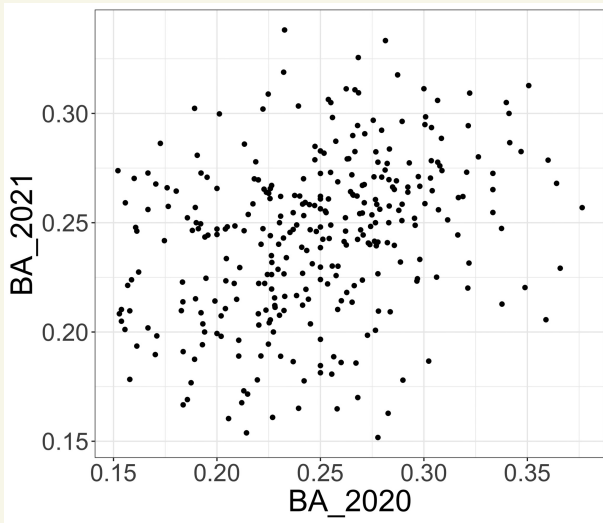


Regression Part 1: Simple Linear Regression

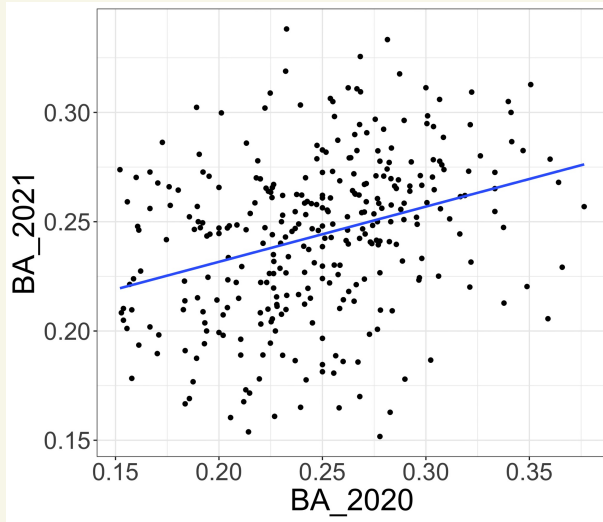
Q Suppose we have access to each MLB player's 2020 batting average and 2021 batting average and no other information.
Predict BA_{2021} from BA_{2020} .

Generally a good idea to begin with exploratory data analysis :



What does the relationship look like?

- Looks linear with a positive slope
 - can imagine drawing a best fit line through the points
 - positive slope (relationship) because, on average, you'd expect that a higher BA_{2020} is associated with a higher BA_{2021}



- not the most perfect Relationship there is a lot of noise but, still some correlation

- So, how do we get this best fit line?

Model

Index each baseball player by i

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

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

Assume a linear relationship

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

β_0 is an unknown constant intercept

β_1 is an unknown constant slope

ε_i is Random independent and identically distributed noise

$$\begin{aligned} \mathbb{E}[\varepsilon_i] &= 0 && \text{mean zero} \\ \varepsilon_i &&& \text{iid} \end{aligned}$$

with unknown constant variance σ^2

- We are interested in the conditional expectation

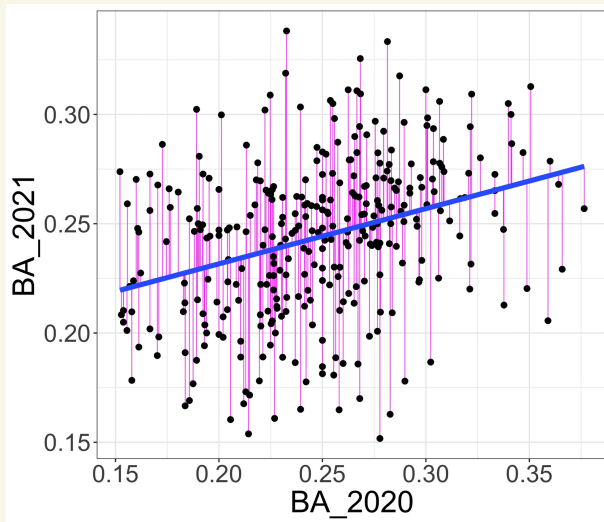
$$\mathbb{E}[Y_i | X_i] = \beta_0 + \beta_1 X_i$$

which is the "true" underlying line

- How do we estimate this best fit line?
How to obtain estimates $\hat{\beta}_0, \hat{\beta}_1$ of β_0, β_1 ?

Ordinary Least Squares — find the values β_0, β_1 which minimize the Residual Sum of Squares (RSS) i.e. minimize the mean squared error,

$$\begin{aligned} \text{RSS}(\beta_0, \beta_1) &= \sum_{i=1}^n [Y_i - \mathbb{E}(Y_i | X_i)]^2 \\ &= \sum_{i=1}^n [Y_i - (\beta_0 + \beta_1 X_i)]^2 \end{aligned}$$



Find the intercept β_0 and slope β_1 (i.e., the blue line) which minimizes the sum of the squares of the lengths of the pink line segments.

$$(\hat{\beta}_0, \hat{\beta}_1) = \underset{(\beta_0, \beta_1)}{\operatorname{argmin}} \operatorname{RSS}(\beta_0, \beta_1)$$

$$= \underset{(\beta_0, \beta_1)}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Calculus: set the derivative equal to zero

$$\frac{\partial}{\partial \beta_0} \operatorname{RSS}(\beta_0, \beta_1) = \sum_{i=1}^n (-2)(y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n \beta_0 = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_1 x_i)$$

$$\Rightarrow \beta_0 = \bar{y} - \beta_1 \bar{x}$$

$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
 $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$

$$\frac{\partial}{\partial \beta_1} \operatorname{RSS}(\beta_0, \beta_1) = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i)(-x_i) = 0$$

$$\Rightarrow -\frac{1}{n} \sum y_i x_i + \beta_0 \frac{1}{n} \sum x_i + \beta_1 \frac{1}{n} \sum x_i^2 = 0$$

$$= \beta_0 \bar{x}$$

$$= (\bar{y} - \beta_1 \bar{x}) \bar{x}$$

$$= \bar{x} \bar{y} - \beta_1 \bar{x}^2$$

$$\Rightarrow \beta_1 \left(\frac{1}{n} \sum X_i^2 - \bar{X} \right)^2 = \frac{1}{n} \sum X_i Y_i - \bar{X} \bar{Y}$$

\Rightarrow by some algebra...

$$\beta_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

Coefficient estimates

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}, \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

What do these formulas mean?

Covariance of two random variables X, Y is

$$\sigma_{XY} = \text{cov}(X, Y) := \mathbb{E}[(X - \mathbb{E}X) \cdot (Y - \mathbb{E}Y)]$$

EX Suppose $\mathbb{E}X = 0$ and $\mathbb{E}Y = 0$

• positive covariance:

if when X is positive Y is positive
and when X is negative Y is negative
then $\mathbb{E}(XY) > 0$

• negative covariance:

if when X is negative Y is positive
and when X is positive Y is negative
then $E(XY) < 0$

Thm X, Y independent $\Rightarrow \text{Cov}(X, Y) = 0$

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

$$= E[X - E(X)] \cdot E[Y - E(Y)] \text{ by independence}$$

$$= (E[X - E(X)]) \cdot (E[Y - E(Y)]) \text{ by linearity of } E$$

$$= 0 \cdot 0 = 0$$

Sample Covariance $\hat{\sigma}_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$

need $n-1$ since we use observed data $\{x_i\}, \{y_i\}$
and $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ and $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$

to get an unbiased estimate $E(\hat{\sigma}_{xy}) = \sigma_{xy}$

Correlation is normalized covariance, ranging from -1 to 1:

$$\rho_{xy} = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$$

Where $\sigma_x = \sqrt{\text{var}(x)} = \sqrt{E X^2 - (E X)^2}$
is the standard deviation of X

sample standard deviation $S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

sample correlation

$$r_{xy} = \frac{\hat{\sigma}_{xy}}{S_x S_y} = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}}$$

$$\text{Then } \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

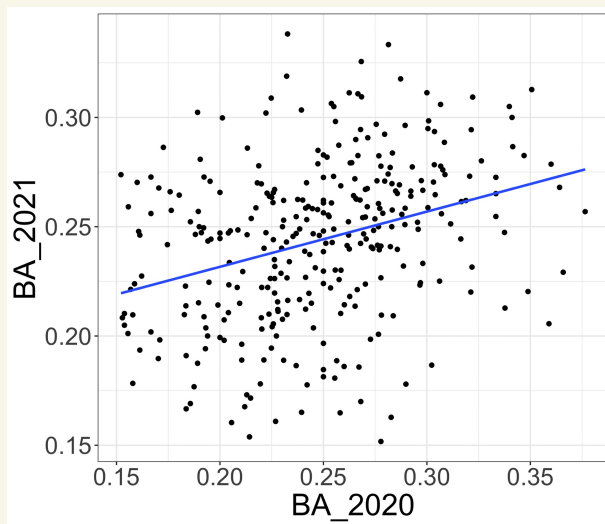
$$= \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}} \cdot \frac{\sqrt{\sum_i (y_i - \bar{y})^2}}{\sqrt{\sum_i (x_i - \bar{x})^2}}$$

$$= r_{xy} \frac{S_y}{S_x}$$

$$\hat{\beta}_1 = r_{xy} \frac{S_y}{S_x}$$

$$r_{xy} = \hat{\beta}_1 \frac{S_x}{S_y}$$

Insight: If X, Y are normalized (same variance) then the linear regression slope is simply the sample correlation!



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

- So, $\hat{\beta}_1$ is a measure of how correlated BA_{2020} is with BA_{2021}
- $\hat{\beta}_1 = \frac{1}{4}$ reflects that a batting average increase of $\bullet 020$ in 2020 is associated with a predicted batting average increase of $\bullet 005$ in 2021.

Regression Line $Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + \varepsilon_i$

$$\begin{aligned}\hat{Y}_i = E(Y_i | X_i) &= \hat{\beta}_0 + \hat{\beta}_1 X_i \\ &= (\bar{Y} - \hat{\beta}_1 \bar{X}) + \hat{\beta}_1 X_i \\ &= \bar{Y} + \hat{\beta}_1 (X_i - \bar{X})\end{aligned}$$

If $X_i = \bar{X}$ then $\hat{Y}_i = \bar{Y}$

→ (\bar{X}, \bar{Y}) lies on the Regression Line!!

● Now, suppose $X_i > \bar{X}$

i.e. $X_i = \bar{X} + \delta$ and $\delta > 0$

$$\begin{aligned}\text{Then } \hat{Y}_i &= \bar{Y} + \hat{\beta}_1 (X_i - \bar{X}) \\ &= \bar{Y} + \hat{\beta}_1 \delta\end{aligned}$$

$$= \bar{Y} + \frac{\delta}{4} \quad \text{since } \hat{\beta}_1 = \frac{1}{4} \text{ in our example}$$

This is Regression to the Mean :

BA_i^{2020} is δ greater than \overline{BA}^{2020}

but \widehat{BA}_i^{2021} is only $\frac{1}{4} \delta$ greater than \overline{BA}^{2021}

In other words, our prediction of player i 's batting average in 2021 is somewhere in between the 2020 BA and the mean BA.

Regression Part 1a: Transformations

Pythagorean Win Percentage

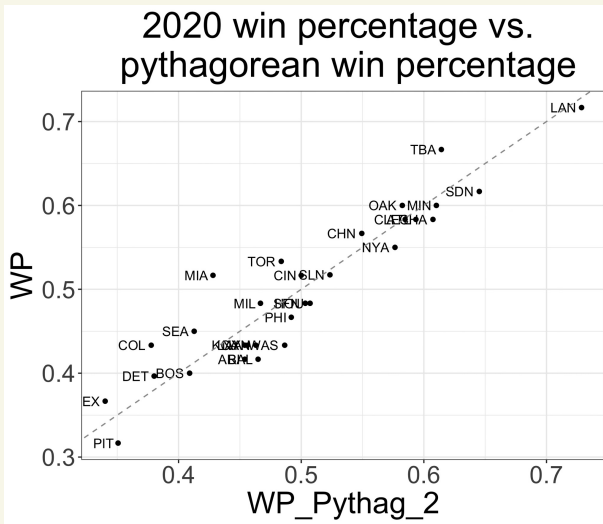
RA = runs allowed by a baseball team in one season.

RS = runs scored.

WP = a team's win percentage.

Bill James' Pythagorean Identity says

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



Pretty damn good!!

Q The pythagorean exponent is quite good, but is arbitrary. Can we find a better exponent α so that

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

Using $\{RS, RA, WP\}$ data from the 2020 MLB season, estimate α .

Idea Transform $E(WP) = \frac{RS^\alpha}{RS^\alpha + RA^\alpha}$ into something that looks linear ($Y = \beta_0 + \beta_1 X$) and then try Linear Regression.

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

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

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

$$\Rightarrow \alpha \log\left(\frac{RA}{RS}\right) = \log\left(\frac{1 - E(WP)}{E(WP)}\right)$$

So, let $EY = \log\left(\frac{1 - E(WP)}{E(WP)}\right)$, $X = \log\left(\frac{RA}{RS}\right)$,

$$\beta_1 = \alpha, \text{ and } \beta_0 = 0.$$

Then we have $EY = \beta_1 X$
Use linear regression!

```
> m2 = lm(data=data_train, log((1-WP)/WP) ~ log(RA/RS) + 0 )  
> m2
```

Call:

```
lm(formula = log((1 - WP)/WP) ~ log(RA/RS) + 0, data = data_train)
```

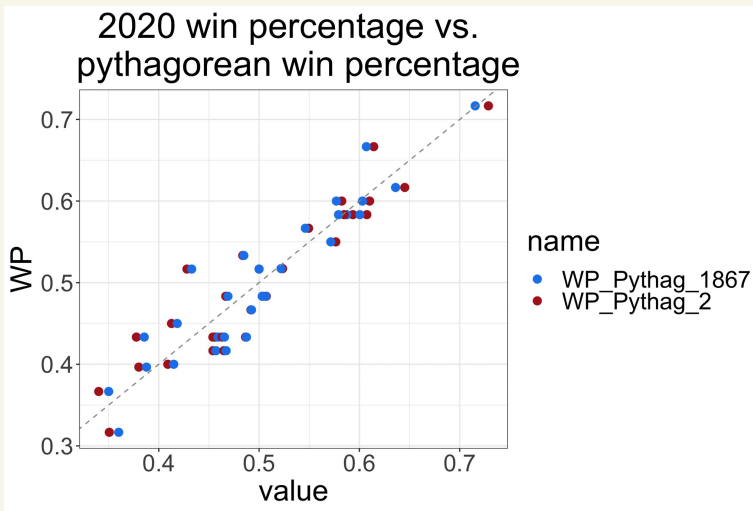
Coefficients:

```
log(RA/RS)
```

```
1.867
```

$$\alpha = 1.867$$

using 2020 data
pretty close to 2!



pretty similar, but blue is slightly better since it is closer to the line $y=x$.

Lesson Even with a nonlinear relationship like
$$WP = \frac{RS^\alpha}{RS^\alpha + RA^\alpha}$$
, we can still use linear
Regression by transforming our data!