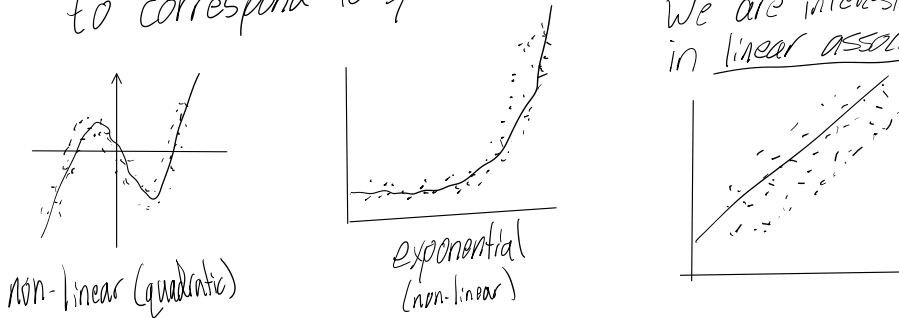


Association

Two variables, (X, Y) where each value of one variable is associated with another (due to being from the same individual) are associated if values of one variable tend to correspond to specific values of the other variable.

We are interested in linear association



We're interested in Linear association which we can measure using the linear correlation coefficient

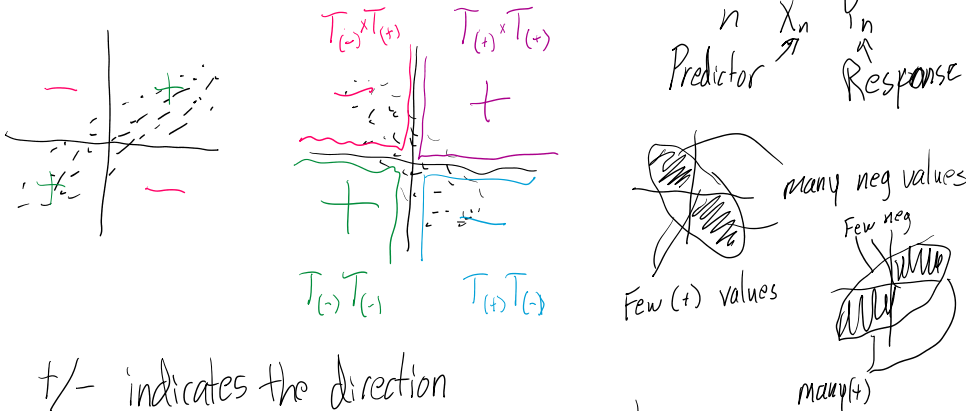
Ref: [1]

$$R = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{X_i - \bar{X}}{S_x} \right) \left(\frac{Y_i - \bar{Y}}{S_y} \right)$$

$$R = \frac{1}{n-1} \sum_{i=1}^n T_{X_i} \times T_{Y_i}$$

ID	X	Y
1	X_1	Y_1
2	X_2	Y_2
3	X_3	Y_3
4	X_4	Y_4
...
n	X_n	Y_n

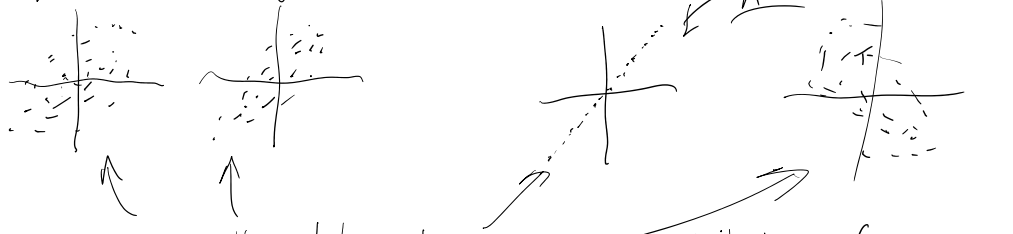
Predictor \nearrow Response



+/- indicates the direction

|R| indicates the strength

Lower |R| Higher |R|



A scatter plot with a vertical Y-axis and a horizontal X-axis. A diagonal line representing a linear regression model passes through the origin. Several data points, represented by blue circles, are scattered around this line. The text "We try to predict" is written in the bottom right corner of the plot area.

$\beta_0, \beta_1, \beta_2$ Regression Coefficients

$$y = b_0 + b_1 x$$

= Regression Coefficients

We try
to predict
the average
value of y_i values
for specific values
of x_i

Intercept term

Slope term

$$\hat{y} = b_0 + b_1 x$$

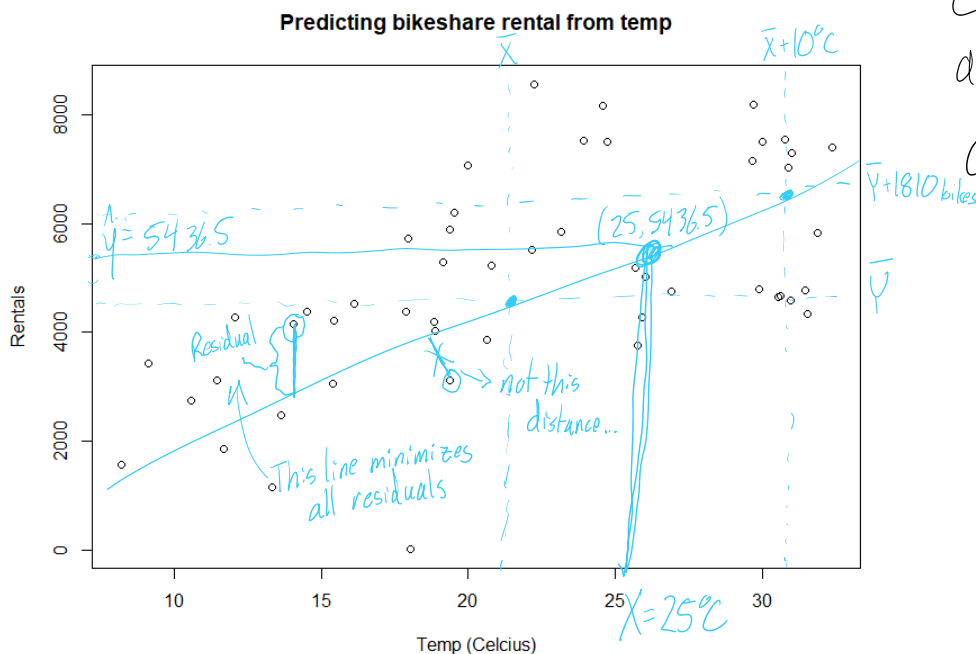
Consider

Bikeshare rentals → how do rental rates (# of bikes rented) relate to temperature?

Y - # of bikes rented

X - temperature avg ($^{\circ}\text{C}$)

→ Describes a day



Overall the trend appears fairly linear, can we find an equation to describe it?

Variables needed for finding the regression Eq.

$$\begin{aligned}\bar{X} &= 20.91^{\circ}\text{C} & \bar{Y} &= 4696.0 \text{ bikes} \\ S_x &= 7.429^{\circ}\text{C} & S_y &= 2032.4 \text{ bikes} \\ R &= 0.6616\end{aligned}$$

The line will always pass through the point of averages (\bar{X}, \bar{Y}) , $(20.91^{\circ}\text{C}, 4696.0 \text{ bikes})$ Slope

$$\Rightarrow \bar{Y} = \hat{b}_0 + \hat{b}_1 \bar{X} \quad \hat{b}_0 = \bar{Y} - \hat{b}_1 \bar{X} \quad \hat{b}_1 = \frac{\text{Rise}}{\text{Run}} = \frac{S_y}{S_x} \cdot R$$

↖ intercept

$$\begin{aligned}\hat{b}_0 &= 4696 \text{ bikes} - \frac{181.0 \text{ bikes}}{^{\circ}\text{C}} \times 20.91^{\circ}\text{C} \\ &= 911.5 \text{ bikes}\end{aligned}$$

$$\begin{aligned}\hat{b}_1 &= \frac{2032.4 \text{ bikes}}{7.429^{\circ}\text{C}} \cdot 0.6616 \\ &= 181.0 \text{ bikes}/^{\circ}\text{C}\end{aligned}$$

Interpret: For a day that is 1 degree hotter, we expect 181 more bike rentals on average.

Interpret (if it makes sense...) ↗ not here
this is the expected # of rentals for a day with temp 0°C

This is outside the range of data

$$\hat{Y} = 911.5 \text{ bikes} + 181.0 \frac{\text{bikes}}{^{\circ}\text{C}} \times X$$

Regression equation describing bike rentals for different temps.

How many bike rentals

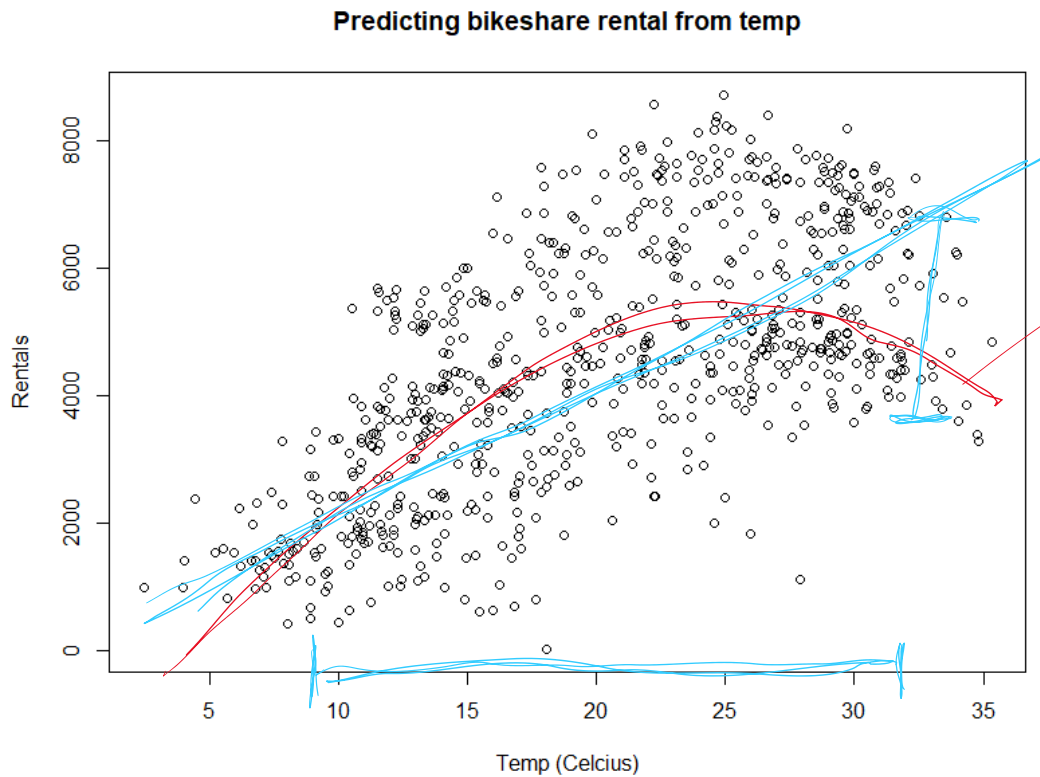
should we expect if $X = 25^{\circ}\text{C}$

$$\begin{aligned}\hat{Y} &= 911.5 \text{ bikes} + 181.0 \frac{\text{bikes}}{^{\circ}\text{C}} \times 25^{\circ}\text{C} \\ \hat{Y} &= 5436.5 \text{ bikes}\end{aligned}$$

This is only valid

within the range of

within the range of
data used to fit the model.



This is
a non-linear
relationship

Using a model outside the range of data used to fit it
is extrapolation and should be avoided.