

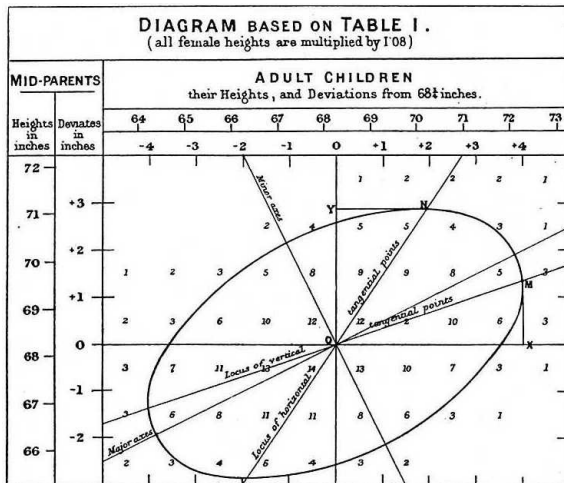
# Introduction to data science & artificial intelligence (INF7100)

Arthur Charpentier

#321 (Simple) Regression

été 2020

# Linear Regression

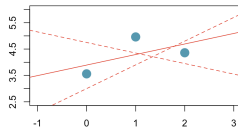
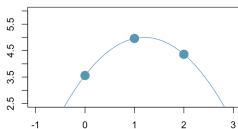
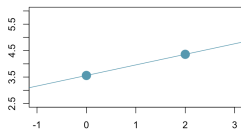


Galton regression towards mediocrity in hereditary stature, 1886.

# Regression

$\{(x_i, y_i)\}$  for  $i = 1, \dots, n$   $y_i = \alpha + \beta x_i + \varepsilon_i$

- ▶  $y$  is the variable of interest
- ▶  $x$  is the explanatory variable



For **Ordinary Least Squares**, solve

$$\min_{\alpha, \beta} \left\{ \sum_{i=1}^n \varepsilon_i^2 \right\} = \min_{\alpha, \beta} \left\{ \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 \right\}$$

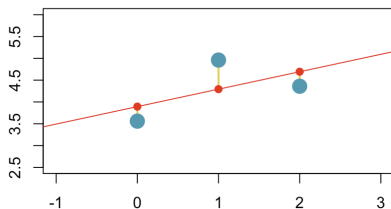
# Regression

$\{(x_i, y_i)\}$  for  $i = 1, \dots, n$

$$y_i = \alpha + \beta x_i + \varepsilon_i$$

- ▶  $y$  is the variable of interest
- ▶  $x$  is the explanatory variable

For **Ordinary Least Squares**, solve



$$\min_{\alpha, \beta} \left\{ \sum_{i=1}^n \varepsilon_i^2 \right\} = \min_{\alpha, \beta} \left\{ \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 \right\}$$

Solutions are (see #411)

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x},$$

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = r_{xy} \frac{s_y}{s_x}.$$

# Significant Explanatory Variable

With centered and scaled variables,  $\hat{\beta} = r_{xy}$ ,

$$\frac{\hat{y} - \bar{y}}{s_y} = r_{xy} \frac{x - \bar{x}}{s_x}.$$

$$t = \frac{\hat{\beta} - \beta}{s_{\hat{\beta}}} \simeq \mathcal{N}(0, 1),$$

where

$$s_{\hat{\beta}} = \sqrt{\frac{\frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

Significant:  $H_0 : \beta = 0$  ? use  $t = \frac{\hat{\beta}}{s_{\hat{\beta}}} \simeq \mathcal{N}(0, 1)$ ,

# Regression

```
1 > import numpy as np
2 > import statsmodels.api as sm
3 > x = np.array([5, 15, 25, 35, 45, 55])
4 > x = x.reshape((-1, 1))
5 > x = sm.add_constant(x)
6 > y = np.array([5, 20, 14, 32, 22, 38])
7 > model = sm.OLS(y, x)
8 > results = model.fit()
9 > print(results.summary())
```

```
10 =====
11                coef.  std err          t      P>|t|      [0.025   0.975]
12 -----
13 const          5.6333    5.872     0.959    0.392    -10.670    21.936
14 x1              0.5400    0.170     3.175    0.034     0.068     1.012
15 -----
16 Dep. Variable:    y                R-squared:            0.716
17 Model:            OLS              Adj. R-squared:       0.645
18                                F-statistic:            10.08
```

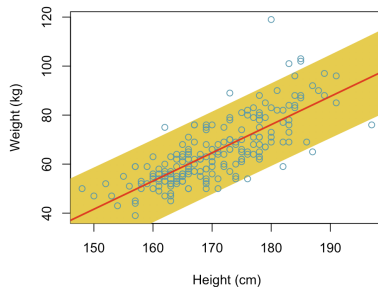
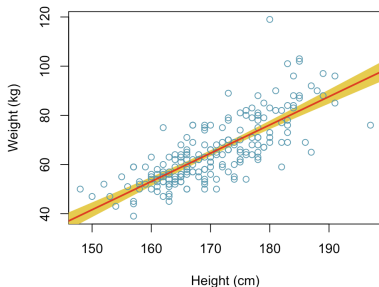
# Regression

```
1 > df = data.frame(x=c(5, 15, 25, 35, 45, 55),
2                   y=c(5, 20, 14, 32, 22, 38))
3 > model = lm(y~x, data=df)
4 > summary(model)
5
6 Coefficients:
7             Estimate Std. Error t value Pr(>|t|)
8 (Intercept)   5.6333      5.8719   0.959   0.3917
9 x             0.5400      0.1701   3.175   0.0337 *
10 ---
11 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.'
12
13 Residual standard error: 7.116 on 4 degrees of freedom
14 Multiple R-squared:  0.7159, Adjusted R-squared:  0.6448
15 F-statistic: 10.08 on 1 and 4 DF,  p-value: 0.03371
```

# Confidence ?

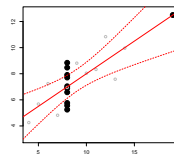
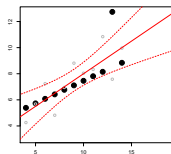
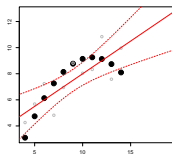
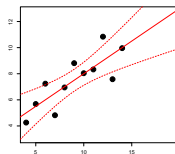
$$\hat{y} = \hat{\alpha} + \hat{\beta}x \quad \text{and} \quad y = \hat{\alpha} + \hat{\beta}x + \hat{\varepsilon}$$

We can get the 95% confidence band for  $\hat{y}$  and  $y$





# Anscombe's Quartet



1 Coefficients:

```
2           Estimate Std. Error t value Pr(>|t|)
3 (Intercept)   3.0000     1.1247   2.667  0.02573 *
4 x1            0.5000     0.1179   4.241  0.00217 **
```

5 ---

6 Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.'

7

8 Residual standard error: 1.237 on 9 degrees of freedom

9 Multiple R-squared: 0.6665, Adjusted R-squared: 0.6295

10 F-statistic: 17.99 on 1 and 9 DF, p-value: 0.00217