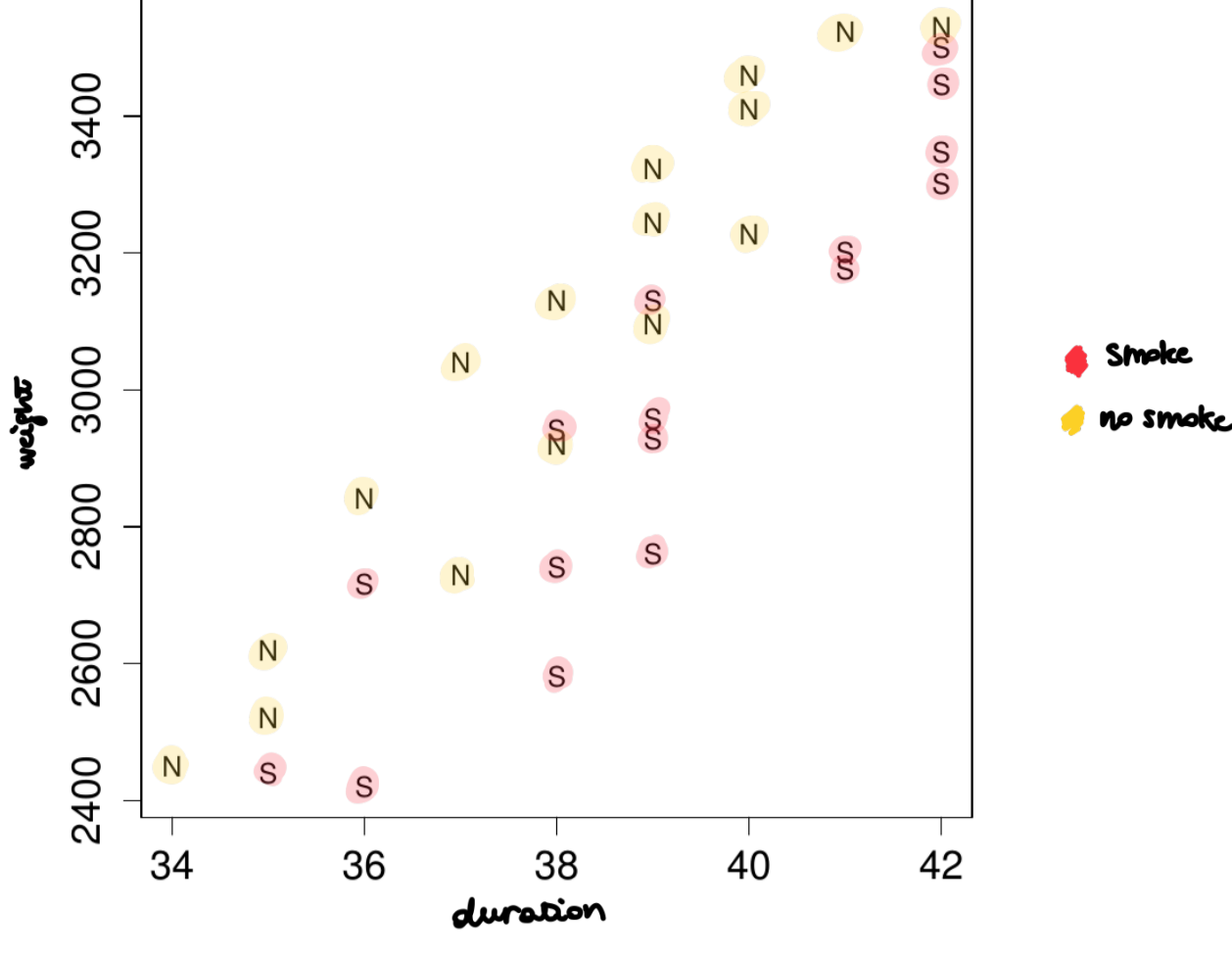


ANALYSIS OF COVARIANCE (ANCOVA)

• Assume that now we have **two groups** and a **continuous covariate**  $x$

- Example:
- $y_i$  = weight of a baby at birth
  - $x_i$  = duration of the pregnancy
  - group = smoke / no smoke (of the mother)



The interest is understanding whether smoking affects the weight at birth, while controlling for the pregnancy duration.

Indeed the weight is clearly influenced by the duration: premature babies have a lower weight compared to babies born later. This is true in general, regardless of the smoking factor.

Hence, it does not make sense to compare the weight of a child whose mother smokes with the weight of a child whose mother does not smoke, if the duration of the pregnancy is different. In that case it would not be clear if an observed difference in the weight is due to smoke or to the duration.

The individual effect of smoke is obtained only if we consider pregnancies with a similar duration.

We have two groups and a covariate.

Consider first modelling the two groups separately

SMOKE GROUP "S":  $Y_i^S = \beta_1^S + \beta_2^S x_i + \epsilon_i \quad i=1, \dots, n_S$

NO-SMOKE GROUP "N":  $Y_i^N = \beta_1^N + \beta_2^N x_i + \epsilon_i \quad i=1, \dots, n_N$

We have noticed how the weight depends on the smoking habit, given the duration  $x$

If we fix a duration  $x_0$

$\mu_0^S = E[Y^S] = \beta_1^S + \beta_2^S x_0$

$\mu_0^N = E[Y^N] = \beta_1^N + \beta_2^N x_0$

Given the specific duration  $x_0$ , is there an effect of smoking?

This question corresponds to a test

$$\begin{cases} H_0: \mu_0^S = \mu_0^N \\ H_1: \mu_0^S \neq \mu_0^N \end{cases}$$

However, we are not interested in the effect of smoking only on babies born at the duration  $x_0$ . We want to study the effect of smoking, given the duration, for all durations.

We can do it using a unique linear model for the two groups.

We define, for  $i=1, \dots, n_S+n_N$

$Y_i = \beta_1 + \beta_2 x_i + \beta_3 s_i + \beta_4 x_i \cdot s_i + \epsilon_i \quad \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$

with

- $x_i$  = duration  $i$
- $s_i$  : indicator of smoke  $s_i = \begin{cases} 1 & \text{if woman } i \text{ smokes} \\ 0 & \text{if woman } i \text{ does not smoke} \end{cases}$
- $x_i \cdot s_i$  : interaction  $x_i \cdot s_i = \begin{cases} x_i & \text{if } s_i = 1 \\ 0 & \text{if } s_i = 0 \end{cases}$

In matrix form we get:

$X = \begin{bmatrix} 1 & x_1 & s_1 & x_1 \cdot s_1 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n_S} & s_{n_S} & x_{n_S} \cdot s_{n_S} \\ 1 & x_{n_S+1} & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n_S+n_N} & 0 & 0 \end{bmatrix}$

Labels: duration (under 1), dummy smoke (under  $s_i$ ), interaction (under  $x_i \cdot s_i$ ).  
 Groups: Smoke group ( $i=1, \dots, n_S$ ), no-smoke group ( $i=n_S+1, \dots, n_S+n_N$ )

Let's look at the mean of  $Y_i$  for different combinations of  $x_i$  and  $s_i$

- if individual  $i$  smokes:  $\mu_i = \beta_1 + \beta_2 x_i + \beta_3 \cdot 1 + \beta_4 \cdot x_i \cdot 1 = (\beta_1 + \beta_3) + (\beta_2 + \beta_4) x_i$
- if individual  $i$  doesn't smoke:  $\mu_i = \beta_1 + \beta_2 x_i + \beta_3 \cdot 0 + \beta_4 \cdot x_i \cdot 0 = \beta_1 + \beta_2 x_i$

INTERPRETATION OF THE PARAMETERS

- $\beta_1$  is the intercept in the "no-smoke" group
- $\beta_1 + \beta_3$  is the intercept in the "smoke" group
- $\beta_2$  is the effect on  $E[Y_i]$  of increasing the duration  $x_i$  by 1 unit in the "no-smoke" group
- $\beta_2 + \beta_4$  is the effect on  $E[Y_i]$  of increasing the duration  $x_i$  by 1 unit in the "smoke" group

We are interested in whether smoking has an effect on the weight, while controlling for the pregnancy duration.

If there is no effect, the two groups will have the same estimated regression line.

ie, equal intercept and slope:  $\beta_1^S = \beta_1^N$  and  $\beta_2^S = \beta_2^N$

With the model for both groups it means:

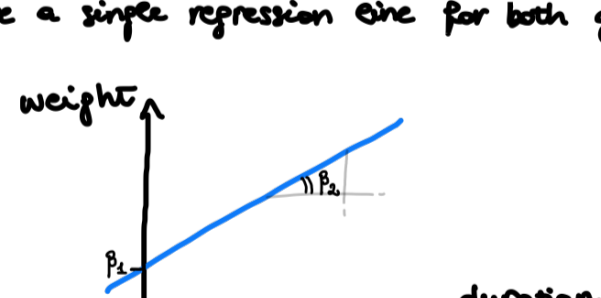
$\beta_1^N = \beta_1^S \Rightarrow \beta_3 = \beta_1^S - \beta_1^N \Rightarrow \beta_3 = 0$   
 $\beta_2^N = \beta_2^S \Rightarrow \beta_4 = \beta_2^S - \beta_2^N \Rightarrow \beta_4 = 0$

Hence, to test the absence of an effect of smoking on the weight, we test

$\begin{cases} H_0: \beta_3 = \beta_4 = 0 \\ H_1: \beta_3 \neq 0 \text{ or } \beta_4 \neq 0 \end{cases}$

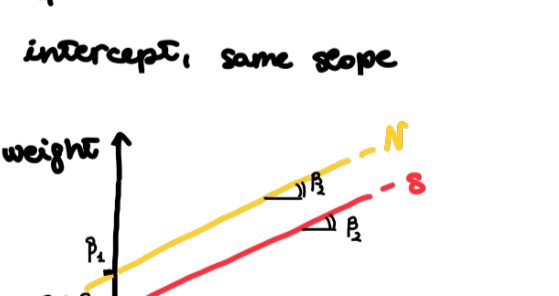
Test on whether smoking affects the weight at birth, controlling for the duration

- under  $H_0$ : no effect of smoking we have a single regression line for both groups



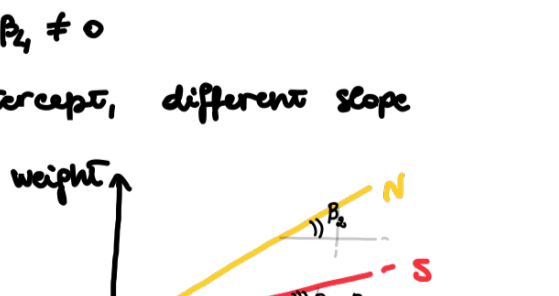
- Under  $H_1$ , if I reject  $H_0$ , I can have different scenarios

- $\beta_3 \neq 0, \beta_4 = 0$   
different intercept, same slope



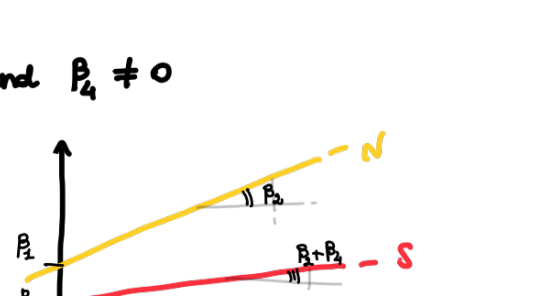
In this example,  $\beta_3 < 0$ . The effect of smoking is constant, regardless of the duration.  $\rightarrow$  SMOKING REDUCES THE EXPECTED WEIGHT BY  $\beta_3$ , FOR ALL DURATIONS.

- $\beta_3 = 0, \beta_4 \neq 0$   
same intercept, different slope



Here,  $\beta_4 < 0$ . At a duration  $x=0$  (not meaningful here) smoking has no effect. The effect increases for increasing duration.

- $\beta_3 \neq 0$  and  $\beta_4 \neq 0$



In the example,  $\beta_3 < 0$  and  $\beta_4 < 0$ . At a duration  $x=0$  the two groups have different means. Moreover, the effect of smoking changes for different durations.