#### Machine learning and statistical learning

#### Chapter 2. Regression and statistical learning

#### Ewen Gallic ewen.gallic@gmail.com







MASTER in Economics - Track EBDS - 2nd Year

2018-2019

This chapter presents some concepts of statistical learning, through the prism of regression.

## 1. Some context

## Model specification

In a regression problem, the aim is to understand how a **response variable** y varies, **conditionally** on the available information on some **predictors**  $\mathbf{x}$ .

Let us take an example, that of the salaries for Professors in the US in 2008-09.

The salary of a professor may be linked, among other things, to the number of years since he or she obtained their Ph.D.



Machine learning and statistical learning 5/160

Here, the linear regression suggests that on the average, the salary increases with the number of years since Ph.D:

• the slope of 985.3 indicates that each additional year since Ph.D leads to an increase of 985 dollars of 9-month salary.

But the relationship does not seem to be linear...

It should be noted here that:

- the regression analysis does not depend on a **generative model** here (a model explaining how the are generated)
- there is no **causal** claims regarding the way mean salary would change if the number of years since Ph.D is altered
- there is no statistical inference

We could **add some predictors** to the model to get a better story on what is going on with salary :

• some ommitted variables may play an important role in explaining the variations.

We can also perform some regression analysis if the **response variable is categorical**.

Let us look at the salary in a different way: let us split it into two categories, either <\$100k or  $\ge\$100k$ .

For each decile of years since Ph.D, we can plot the **conditional propor-tions**.



#### Levels of regression analysis

Berk (2008) mentions three levels of regression analysis:

- Level I regression analysis:
  - aiming at describing the data
  - assumption free
  - should not be neglected
- Level II regression analysis:
  - based on statistical inference
  - uses results from level I regression analysis
  - use with real data may be challenging
  - allows to make predictions
- Level III regression analysis:
  - based on causal inference
  - uses level I analysis, sometimes coupled with level II
  - rely more on algorithmic methods rather than model-based methods.

# 2. The linear regression

## Some references

- Berk (2008). Statistical learning from a regression perspective, volume 14. Springer.
- Cornillon and Matzner-Løber (2007). Régression: théorie et applications. Springer.
- James et al. (2013). An introduction to statistical learning, volume 112. Springer.

#### The linear regression

Linear regression combines level I and level II perspectives.

It is useful when one wants to predict a quantitative response.

A lot of newer statistical learning approaches can be seen as generalizations or extensions of linear regression, as reminded in James et al. (2013).

## 2.1 Simple linear regression

## Principle

Let us consider first the case of simple linear regression.

We aim at predicting a quantitative response variable y using a single predictor  $\mathbf{x}$  (or regressor).

- y is a  $n\times 1$  numerical response variable, where n represents the numbr of observations
- $\mathbf{x}$  is a  $n \times 1$  predictor.

We assume there exists a linear relationship between y and  $\mathbf{x}$  such that:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n, \tag{1}$$

where  $\varepsilon_i$  is an error term normally distributed with 0 mean and variance  $\sigma^2$ , *i.e*,  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ .

# Principle

In Eq. 1, the **coefficients** (or **parameters**)  $\beta_0$  (*i.e.*, the constant) and  $\beta_1$  (*i.e.*, the slope) are unknown parameters to be estimated.

These coefficients are **estimated** using a **training sample**.

The estimates of  $\beta_0$  and  $\beta_1$  are, respectively,  $\hat{\beta_0}$  and  $\hat{\beta_1}$ .

Once they are estimated using a learning procedure (in this case using linear regression), they can used to **predict** values for y for some value  $x_0$ :

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0 \tag{2}$$

2. The linear regression

2.1. Simple linear regression

 $\square_{2.1.1.}$  Estimating the coefficients

## 2.1.1 Estimating the coefficients

## Estimating the coefficients

To estimate  $\beta_0$  and  $\beta_1$ , we rely on a set of training examples  $\{(x_1, y_1), \ldots, (x_n, y_n)\}.$ 

For example, let us go back to our data describing the 9 month salary of professors (the response variable) and look at the relationship between the salary and years since Ph.D (x).

└─2. The linear regression └─2.1. Simple linear regression └─2.1.1. Estimating the coefficients

#### Estimating the coefficients

Figure 1: Varying the intercept.

Figure 2: Varying the slope.

There is an infinity of possibles values that one can pick for  $\hat{\beta}_0$  and  $\hat{\beta}_1$ .

However, we want to find an estimation that leads to a line being as close as possible to the points: but what does "close" mean?

Ewen Gallic

## Estimating the coefficients

The most common metric we want to minimize is known as the **least** square criterion.

The predictions  $\hat{y}_i$  for each of the  $x_i$ , i = 1, ..., n are given by  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ .

Let  $e_i = y_i - \hat{y}_i$  the *i*th residual, *i.e.*, the difference between the osberves value and its prediction by the linear model.

The residual sum of square is defined as:

$$\mathsf{RSS} = \sum_{i=1}^{n} \varepsilon_i^2 = \sum_{i=1}^{n} (y_i - \hat{\beta}_1 x_i - \hat{\beta}_0)^2.$$
(3)

We aim at minimizing this metric.

#### Least squares coefficient estimates

It can easily be shown that the minimization of the RSS leads to:

$$\begin{cases} \hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} \\ \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \end{cases}$$
(4)

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

2. The linear regression
 2.1. Simple linear regression
 2.1.1. Estimating the coefficients

#### Least squares coefficient estimates

In our example, the least squares coefficient estimates  $be\hat{t}a_0$  and  $\beta_1$  are, respectively, 9.1719 and 0.0985.





2. The linear regression
2.1. Simple linear regression
2.1.1. Estimating the coefficients

#### Residual sum of squares

We can have a look at the RSS when we vary the values of  $\hat{\beta_0}$  and  $\hat{\beta_1}$ :



Figure 4: Surface plot of the RSS depending on the values of  $\hat{\beta_0}$  and  $\hat{\beta_1}$ .

Figure 5: Contour plot of the RSS depending on the values of  $\hat{\beta_0}$  and  $\hat{\beta_1}$ .

2. The linear regression

2.1. Simple linear regression

 $\_$  2.1.2. Accuracy of the coefficient estimates

## 2.1.2 Accuracy of the coefficient estimates

## Accuracy of the coefficient estimates

The estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are point estimates.

When they are estimated by least squares, they are:

- unbiased
  - $\mathbb{E}(\hat{\beta}_0) = \beta_0 \text{ and } \mathbb{E}(\hat{\beta}_1) = \beta_1$
- efficient
  - $\mathbb{V}(\hat{\beta}_0)$  and  $\mathbb{V}(\hat{\beta}_1)$  are minimal
- convergent
  - $\lim_{n \to +\infty} \mathbb{V}(\hat{\beta}_0) = 0$  and  $\lim_{n \to +\infty} \mathbb{V}(\hat{\beta}_1) = 0$

They are called **BLUE** (Best Linear Unbiased Estimator).

2. The linear regression
2.1. Simple linear regression
2.1.2. Accuracy of the coefficient estimates

#### Accuracy of the coefficient estimates

#### It is easy to show that:

$$\begin{cases} \mathbb{V}(\hat{\beta}_0) &= \sigma^2 \left[ \frac{1}{n} + \frac{\overline{x}}{\sum_{i=1}^n (x_i - \overline{x})^2} \right] \\ \mathbb{V}(\hat{\beta}_1) &= \frac{\sigma^2}{\sum_{i=1}^n (x_i - \overline{x})^2} \end{cases}$$

where  $\sigma^2$  can be estimated:

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2} = \frac{\sum_{i=1}^n e_i^2}{n-2}.$$

(5)

L The linear regression
2.1. Simple linear regression
L 2.1.2. Accuracy of the coefficient estimates

## Accuracy of the coefficient estimates

Figure 6: A: True relationship (in red), Observed values of y (points) and Least Squares line (in blue). B: True relationship (in red), Current Least Squares line (in blue), Previous Least Squares lines (in gray).

2. The linear regression
 2.1. Simple linear regression
 2.1.2. Accuracy of the coefficient estimates

## Accuracy of the coefficient estimates



Figure 7: Mean of estimates of  $\beta_0$  and  $\beta_1$  depending on the number of resampling.

## Hypothesis tests

We wish to test if a coefficient  $\theta$ ,  $\theta \in \{\beta_0, \beta_1\}$  is equal to a specific value  $\theta_0$ :

$$\begin{cases} H_0: \theta = \theta_0 \\ H_1: \theta \neq \theta_0 \end{cases}$$

We know that 
$$\hat{\theta} \sim \mathcal{N}\left(\theta, \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)$$
, so:
$$\frac{\hat{\theta} - \theta}{\sigma/\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \sim \mathcal{N}(0, 1).$$

2. The linear regression
2.1. Simple linear regression
2.1.2. Accuracy of the coefficient estimates

#### Hypothesis tests

As 
$$\frac{\sum_{i=1}^{n} \epsilon_i^2}{\sigma^2} \sim \chi^2_{n-2}$$
, we can define a variable  $T$  as:

$$T = \frac{\frac{\bar{\theta} - \theta}{\sigma/\sqrt{\sum_{i=1}^{n} (x_i - \bar{x})^2}}}{\sqrt{\frac{\sum_{i=1}^{n} \varepsilon_i^2}{\sigma_u^2}}/\sqrt{n-2}} \sim \mathcal{S}t(n-2)$$

We can show that the expression of T can be simplified to:

$$T = \frac{\hat{\theta} - \theta}{\hat{\sigma}_{\hat{\theta}}}$$

2. The linear regression
 2.1. Simple linear regression
 2.1.2. Accuracy of the coefficient estimates

## Hypothesis tests

It is thus possible to perform the following test:

$$\begin{cases} H_0: \theta = \theta_0 \\ H_1: \theta \neq \theta_0 \end{cases}$$

knowing that  $\frac{\hat{\theta}-\theta}{\sigma^{\theta}} \sim \mathcal{S}t(n-2)$ 

2. The linear regression
 ↓ 2.1. Simple linear regression
 ↓ 2.1.2. Accuracy of the coefficient estimates

## Hypothesis tests

And we need to find the following probability:

$$\mathbb{P}\left(-t_{\alpha/2} < \frac{\hat{\theta} - \theta}{\hat{\sigma}_{\hat{\theta}}} < t_{\alpha/2}\right)$$

We therefore need to compute a t-statistic, that measures the number of standard deviations that  $\hat{\theta}$  is away from  $\theta_0$ :

$$t_{\rm obs.} = \frac{\hat{\theta} - \theta_0}{\hat{\sigma}_{\hat{\theta}}}$$

if t<sub>obs.</sub> ∈ [-t<sub>α/2</sub>, t<sub>α/2</sub>]:
 we do not reject the null hypothesis (H<sub>0</sub>) with a first-order risk of α%

• if 
$$t_{\text{obs.}} \notin \left[-t_{\alpha/2}, t_{\alpha/2}\right]$$
:

• we reject the null hypothesis  $(H_0)$  with a first-order risk of  $\alpha\%$ 

## Hypothesis tests

Most of the time, we are interested in a specific case:

$$\begin{cases} H_0 : \alpha = 0\\ H_1 : \alpha \neq 0, \end{cases}$$

In such a case, the t-statistic becomes:

$$T = \frac{\hat{\theta} - 0}{\hat{\sigma}_{\theta}} = \frac{\hat{\theta}}{\hat{\sigma}_{\hat{\theta}}}$$

The observed value is  $t_{\rm obs.} = \frac{\hat{\alpha}}{\hat{\sigma}_{\hat{\alpha}}}$ .

2. The linear regression
 ↓ 2.1. Simple linear regression
 ↓ 2.1.2. Accuracy of the coefficient estimates

#### Hypothesis tests: confidence interval

We can also use the standard error of the coefficient estimates to construct a **confidence interval**:

$$I.\widehat{C._{\theta}(1-\alpha)} = \left[\hat{\theta} \pm t_{\alpha/2} \times \hat{\sigma}_{\hat{\theta}}\right].$$
 (6)

If the intervals contain 0, then we can conclude that the coefficient  $\theta$  is not statistically different from zero (at the  $\alpha\%$  level of significance).

We can also compute the probability of observing any number equal to |t| or larger while assuming  $\theta = 0$  (this probability is known as **the p-value**).

2. The linear regression
2.1. Simple linear regression
2.1.2. Accuracy of the coefficient estimates

#### Hypothesis tests

	Least squares
(Intercept)	$9.17^{***}$
	(0.28)
yrs.since.phd	$0.10^{***}$
	(0.01)
$R^2$	0.18
Adj. R $^2$	0.17
Num. obs.	397
RMSE	2.75

\*\*\* p < 0.001, \*\* p < 0.01, \*p < 0.05

Table 1: Statistical models

## 2.1.3 Accuracy of the model
Recall that the linear regression is a supervised learning method. Hence, we can compare the predictions we obtain with the observed values of the output variable.

We want to have an idea of the quality of the estimation, to know how well the model fits the data.

To that end, we usually use several metrics, among which:

- the root mean squared error (RMSE)
- the residual standard error (RSE)
- the  $R^2$  statistic.

The mean squared error (MSE) is an estimate of the average of the squares of the errors:

$$\mathsf{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 \tag{7}$$

The root mean squared error is the square root of the MSE:

$$\mathsf{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2} = \sqrt{\frac{\mathsf{RSS}}{n}},$$
(8)

where  $\mathsf{RSS} = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$ 

The value of the RMSE is always non-negative.

A value of 0 indicates a perfect fit to the data.

Recall that the linear model contains an error term ( $\varepsilon$ ). Hence, we will not be able to perfectly predict the response variable.

The **Residual Standard Error** is the average amount that the response will deviate from the true regression line. It is an estimate of the standard deviation of  $\varepsilon$ :

$$\mathsf{RSE} = \sqrt{\frac{1}{n-2}RSS} = \sqrt{\frac{1}{n-2}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}.$$
 (9)

In our example of the regression of salaries onto years since Ph.D, the value of the RSE is 2.7534.

This means that the actual salary can deviate from the true regression line by approximately 2.7534 thousand dollars, on average.

The mean salary in the data is \$11.37065 thousand dollars. Hence, the percentage error for any prediction, using our estimation would be  $2.7534/11.37065 \approx 25\%$ .

Now, let us turn to the  ${\cal R}^2$  statistic, which provides another method to assess the quality of fit.

The  $R^2$  measures the **proportion of variance explained**. It takes a value between 0 and 1.

Let us illustrate this.

The variations of y are only partially explained by those of  $\mathbf{x}$ 



Figure 8: Variation from  $y_2$  to  $y_1$ 

As shown in Figure 8, the variation from  $y_1$  to  $y_2$  is partially explained by the variation from  $x_1$  to  $x_2$ .

The **quality** of fit at each point, as measured by the total variation, can therefore be broken down into two parts:

- the explained variation
- the residual variation

using the average point  $(\overline{x}, \overline{y})$  as reference, *i.e.*:



The closer  $\hat{A}$  is to A, the stronger the explained variation is, relatively.



Figure 9: Decomposition of the variation.

Thus, one way to assess the quality of the adjustment is to measure the following ratio:

#### explained variance

total variance

Or, for all observations:

$$R^{2} = \frac{\sum_{i=1}^{n} (\hat{y}_{i} - \bar{y})^{2}}{\sum_{i=1}^{n} (y_{i} - \bar{y})^{2}} = \frac{\text{ESS}}{\text{TSS}} = \frac{\text{explained sum of squares}}{\text{total sum of squares}}$$
(10)

We can write the  $R^2$  differently, as we know that:

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

Thus:

$$R^{2} = \frac{\sum_{i=1}^{n} (y_{i} - \bar{y})^{2} - \sum_{i=1}^{n} (y_{i} - \hat{y}_{i})^{2}}{\sum_{i=1}^{n} (y_{i} - \bar{y})^{2}}$$
$$= 1 - \frac{\sum_{i=1}^{n} (y_{i} - \hat{y}_{i})^{2}}{\sum_{i=1}^{n} (y_{i} - \bar{y})^{2}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$
(11)

The value of the  $R^2$  lies between 0 and 1:

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} \qquad \Rightarrow \qquad 0 \le R^2 \le 1.$$

- When the economic **theory** suggests that the relationship between the response and its predictor should be **linear**, we expect the value of the  $R^2$  to be really close to zero, otherwise, it suggests there might be something wrong with the generation of the data.
- In other situations, when the **linear relationship** can be at best a rough approximation of the real form, we expect to find low values of the  $R^2$ .

# $\mathbb{R}^2$ and correlation

It can be noted that in the case of simple linear regression, the  $R^2$  is equal to the squared correlation coefficient.

Indeed:

$$y_{i} - \hat{y}_{i} = y_{i} - \bar{y} + \bar{y} - \hat{y}_{i}$$
  
=  $(y_{i} - \bar{y}) - (\hat{y}_{i} - \bar{y})$   
=  $(y_{i} - \bar{y}) - (\hat{\beta}_{1}x_{i} + \hat{\beta}_{0} - \hat{\beta}_{1}\bar{x} - \hat{\beta}_{0})$   
=  $(y_{i} - \bar{y}) - \hat{\beta}_{1}(x_{i} - \bar{x}).$ 

Taking the squared value:

$$(y_i - \hat{y}_i)^2 = (y_i - \bar{y})^2 + \hat{\beta}_1^2 (x_i - \bar{x})^2 - 2\hat{\beta}_1 (y_i - \bar{y})(x_i - \bar{x})$$

# $\mathbb{R}^2$ and correlation

Which leads to:

$$(y_i - \hat{y}_i)^2 = (y_i - \bar{y})^2 + \hat{\beta}_1^2 (x_i - \bar{x})^2 - 2\hat{\beta}_1 (y_i - \bar{y})(x_i - \bar{x})$$

Summing on all individuals:

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

└─2. The linear regression └─2.1. Simple linear regression └─2.1.3. Accuracy of the model

# $\mathbb{R}^2$ and correlation

In can indeed be shown that

$$2\hat{\beta}_{1}^{2} \sum_{i=1}^{n} (x_{i} - \bar{x})^{2} = 2\hat{\beta}_{1} \sum_{i=1}^{n} (x_{i} - \bar{x})^{2} \frac{\sum_{i=1}^{n} (x_{i} - \bar{x})(y_{i} - \bar{x})}{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}}$$
$$= 2\hat{\beta}_{1} \sum_{i=1}^{n} (y_{i} - \bar{y})(x_{i} - \bar{x}).$$

We also have:

$$(\hat{y}_i - \bar{y}) = \hat{\beta}_1 x_i + \hat{\beta}_0 - \hat{\beta}_1 \bar{x} - \hat{\beta}_0 = \hat{\beta}_1 (x_i - \bar{x}).$$

By taking the squared value and summing for all individuals:

$$\sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2 = \hat{\beta}_1^2 \sum_{i=1}^{n} (x_i - \bar{x})^2.$$
(12)

└─2. The linear regression └─2.1. Simple linear regression └─2.1.3. Accuracy of the model

# $\mathbb{R}^2$ and correlation

Then, introducing (12) in (10), we get:

$$R^{2} = \frac{\hat{\alpha}^{2} \sum_{i=1}^{n} (x_{i} - \bar{x})^{2}}{\sum_{i=1}^{n} (y_{i} - \bar{y})^{2}}$$

$$= \underbrace{\left( \underbrace{\sum_{i=1}^{n} (x_{i} - \bar{x})(y_{i} - \bar{y})}_{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}} \right)^{2}}_{\hat{\alpha}^{2}} \times \underbrace{\sum_{i=1}^{n} (x_{i} - \bar{x})^{2}}_{\sum_{i=1}^{n} (y_{i} - \bar{y})^{2}}$$

$$R^{2} = \frac{\left( \sum_{i=1}^{n} (x_{i} - \bar{x})(y_{i} - \bar{y}) \right)^{2}}{\left( \sum_{i=1}^{n} (x_{i} - \bar{x})^{2} \right) \left( \sum_{i=1}^{n} (y_{i} - \bar{y})^{2} \right)}$$

$$= \frac{\left( \mathbb{C}ov(x, y) \right)^{2}}{\mathbb{V}(x) \times \mathbb{V}(y)}$$
(13)

## 2.2 Multiple linear regression

## Principle

We have considered so far only one predictor in the design matrix  $\mathbf{x}$ . Let us now look at the case where we want to use **multiple predictors**:  $\mathbf{x}$ becomes a  $n \times p$  matrix, with n observations and p predictors.

We assume there exists a relationship between the response y and the predictors  $\mathbf{x}$  such that:

$$y_i = \beta_0 + \beta_1 \mathbf{x}_{1i} + \ldots + \beta_p \mathbf{x}_{pi} + \varepsilon_i, \quad i = 1, \ldots, n,$$
(15)

where  $\varepsilon_i$  is an error term normally distributed with 0 mean and variance  $\sigma^2$ , *i.e.*,  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ , and where  $\mathbf{x}_{ji}$  represents the *i*th observation for the *j*th predictor,  $j = 1, \ldots, p$ .

## Principle

In Eq. 22 the **coefficients**  $\beta_0$  (*i.e.*, the constant) and  $\beta_j$  are unknown parameters to be estimated.

We interpret the coefficients  $\beta_j$  as the average effect on y of a one unit increase in  $\mathbf{x}_j$ , *ceteris paribus* (*i.e.*, holding all other predictors fixed).

2. The linear regression

2.2. Multiple linear regression

 $\square_{2.2.1.}$  Estimating the coefficients

### 2.2.1 Estimating the coefficients

The coefficients of the multiple linear regression, can once again be estimated so that they minimize the sum of squared residuals:

$$RSS = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$
(16)  
=  $\sum_{i=1}^{n} (y_i - \hat{\beta}_0 - \hat{\beta}_1 \mathbf{x}_{i1} - \dots - \hat{\beta}_p \mathbf{x}_{ip}),$ (17)

where  $\hat{y}_i = \hat{eta}_0 + \hat{eta}_1 \mathbf{x}_{i1} + \ldots + \hat{eta}_p \mathbf{x}_{ip}$ 

m

Using matrix algebra, it is easy to estimate the coefficients. First, we can write:

$$\hat{\boldsymbol{y}} = \boldsymbol{X}\hat{\boldsymbol{\beta}},$$
 (18)

where 
$$\hat{\boldsymbol{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$
,  $\boldsymbol{X} = \begin{bmatrix} \mathbf{x}_{11} & \mathbf{x}_{12} & \cdots & \mathbf{x}_{1p} & 1 \\ \mathbf{x}_{21} & \mathbf{x}_{22} & \cdots & \mathbf{x}_{2p} & 1 \\ \vdots & \vdots & \ddots & \vdots & 1 \\ \mathbf{x}_{n1} & \mathbf{x}_{n2} & \cdots & \mathbf{x}_{np} & 1 \end{bmatrix}$ , and  $\hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_p \\ \hat{\beta}_0 \end{bmatrix}$ 

 $\sim$ 

Let 
$$\boldsymbol{y} - \hat{\boldsymbol{y}}$$
 denote the column vector  $\boldsymbol{y} - \hat{\boldsymbol{y}} = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$ ,

and 
$$\boldsymbol{y} - \hat{\boldsymbol{y}}^{\top}$$
 the vector column  $\boldsymbol{y} - \hat{\boldsymbol{y}}^{\top} = \begin{bmatrix} y_1 - \hat{y}_1 & y_2 - \hat{y}_2 & \cdots & y_n - \hat{y}_n \end{bmatrix}^{\top}$ 

By definition:

$$(\boldsymbol{y} - \hat{\boldsymbol{y}}) \top (\boldsymbol{y} - \hat{\boldsymbol{y}}) = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{n} e_i^2 = \text{RSS},$$

2. The linear regression
2.2. Multiple linear regression
2.2.1. Estimating the coefficients

#### Estimating the coefficients

By replacing y by its expression given in Eq. 18:

$$\sum_{i=1}^{n} e_i^2 = (\boldsymbol{y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^\top (\boldsymbol{y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})$$
$$= (\boldsymbol{y}^\top - \hat{\boldsymbol{\beta}}^\top \boldsymbol{X}^\top) (\boldsymbol{y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})$$
$$= \boldsymbol{y}^\top \boldsymbol{y} - \boldsymbol{y}^\top \boldsymbol{X}\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^\top \boldsymbol{X}^\top \boldsymbol{y} + \hat{\boldsymbol{\beta}}^\top \boldsymbol{X}^\top \boldsymbol{X}\hat{\boldsymbol{\beta}}$$

2. The linear regression
2.2. Multiple linear regression
2.2.1. Estimating the coefficients

#### Estimating the coefficients

Figure 10: Regression of salaries on years since Ph.D and years of service. The red dots represent the observed values. The plane minimizes the sum of squared distances represented by the red (overstimated values) and blue segments (underestimted values).

Ewen Gallic

	Model 1	Model 2
(Intercept)	$9.17^{***}$	8.99***
	(0.28)	(0.28)
yrs.since.phd	$0.10^{***}$	$0.16^{***}$
	(0.01)	(0.03)
yrs.service		$-0.06^{*}$
		(0.03)
$R^2$	0.18	0.19
Adj. R $^2$	0.17	0.18
Num. obs.	397	397
RMSE	2.75	2.74

\*\*\* p < 0.001, \*\* p < 0.01, \*p < 0.05

Table 2: Statistical models

#### Bias of the coefficients

Let us look at the **bias of the coefficients**. First, we can write the estimated vector of coefficients  $\hat{\beta}$  as:

$$\begin{split} \hat{\boldsymbol{\beta}} &= (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}\boldsymbol{y} \\ &= (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}(\boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) \\ &= (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}(\boldsymbol{X}^{\top}\boldsymbol{X})\boldsymbol{\beta} + (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}\boldsymbol{\varepsilon} \\ &= \boldsymbol{\beta} + (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}\boldsymbol{\varepsilon} \end{split}$$

─ 2. The linear regression
└─ 2.2. Multiple linear regression
└─ 2.2.1. Estimating the coefficients

#### Bias of the coefficients

Hence the expected value of  $\hat{\beta}$  is given by:

$$\mathbb{E}(\hat{\boldsymbol{\beta}}) = \mathbb{E}(\boldsymbol{\beta}) + \mathbb{E}\left[ (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}\boldsymbol{\varepsilon} \right]$$
$$= \boldsymbol{\beta} + (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}\mathbb{E}(\boldsymbol{\varepsilon})$$
$$= \boldsymbol{\beta}$$
(19)

since we assumed  $\mathbb{E}(\varepsilon) = 0$ .

As a consequence,  $\hat{\beta}$  is an unbiased estimator of  $\beta$ , *i.e.*:

$$\mathcal{B}\left(\hat{\boldsymbol{\beta}};\boldsymbol{\beta}\right) = \mathbb{E}\left(\hat{\boldsymbol{\beta}}\right) - \boldsymbol{\beta} = 0$$

2. The linear regression
2.2. Multiple linear regression
2.2.1. Estimating the coefficients

#### Variance of the coefficients

The variance of the coefficients writes:

$$\begin{split} \mathbb{V}(\hat{\boldsymbol{\beta}}) &= \mathbb{E}\left[\hat{\boldsymbol{\beta}} - \mathbb{E}(\hat{\boldsymbol{\beta}})\right]^2 \\ &= \mathbb{E}\left[(\hat{\boldsymbol{\beta}} - \mathbb{E}(\boldsymbol{\beta}))(\hat{\boldsymbol{\beta}} - \mathbb{E}(\boldsymbol{\beta}))^\top\right] \\ &= \mathbb{E}\left[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top\right]. \end{split}$$

Since:

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}\boldsymbol{\varepsilon}$$

$$(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^{\top} = \boldsymbol{\varepsilon}^{\top} \boldsymbol{X} (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1}.$$

Hence:

$$(\hat{oldsymbol{eta}}-oldsymbol{eta})(\hat{oldsymbol{eta}}-oldsymbol{eta})^{ op}=(oldsymbol{X}^{ op}oldsymbol{X})^{-1}oldsymbol{X}^{ op}oldsymbol{arepsilon}^{ op}oldsymbol{X}(oldsymbol{X}^{ op}oldsymbol{X})^{-1}$$

#### Variance of the coefficients

So in the end:

$$\mathbb{V}(\hat{\boldsymbol{\beta}}) = (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}\mathbb{E}(\boldsymbol{e}\boldsymbol{e}^{\top})\boldsymbol{X}(\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}$$
  
$$= (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}\sigma_{u}^{2}\boldsymbol{X}(\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}$$
  
$$= (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}(\boldsymbol{X}^{\top}\boldsymbol{X})(\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\sigma_{\varepsilon}^{2}$$
  
$$= (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\sigma_{\varepsilon}^{2}.$$
 (20)

So the variance of  $\hat{\beta}$  is equal to the variance  $\varepsilon$  multiplied by the *i*th term of the diagonal of  $(\mathbf{X}^{\top}\mathbf{X})^{-1}$ .

2. The linear regression

2.2. Multiple linear regression

 $\square_{2.2.2.}$  Accuracy of the estimation

#### 2.2.2 Accuracy of the estimation

## Accuracy of the estimation

As in the simple linear regression case, we are interested in measuring the **accuracy of the estimation**.

In particular, we will look at the following aspects:

- is there a significant relationship between the response and the predictors?
- which predictors should be kept in the model, and which should be discarded?
- how well does the model fit to the data?

## Relationship between y and ${\bf x}$

To infer whether there is a relationship between the response y and the predictors x, a statitiscal test can be performes. The **null hypothesis** of this test writes:

$$\mathsf{H}_0:\beta_1=\beta_2=\ldots=\beta_p=0$$

The alternative writes:

 $H_1$ : at least one  $\beta_j$  is non-zero,  $j = 1, \dots, p$ 

This test is based on the following F-statistic:

$$F = \frac{(\mathsf{TSS} - \mathsf{RSS})/p}{\mathsf{RSS}/(n-p-1)} \sim \mathcal{F}(p, n-p-1). \tag{21}$$

## Relationship between y and $\mathbf{x}$

- If the linear model assumptions are correct:
  - $\blacktriangleright \ \mathbb{E}(\mathsf{RSS}/(n-p-1)) = \sigma^2$
- and if  $H_0$  is true:
  - $\blacktriangleright \mathbb{E}\left[(\mathsf{TSS} \mathsf{RSS})/p\right] = \sigma^2$

As a consequence:

- when there is **no relationship** between the response and its predictors:
  - the value of the F-statistic should be close to zero
- when **H**<sub>1</sub> is true:
  - ▶  $\mathbb{E}\left[(\mathsf{TSS} \mathsf{RSS})/p\right] > \sigma^2$ , hence **F** should be greater than one.

## Relationship between y and ${\bf x}$

In the example of the salary regressed on years since Ph.D and years of services, the value of the F statistic is 45.71 (2 and 394 degrees of freedom). The p-value associated to the test is lower than  $2.2 \times 10^{-16}$ .

Hence, we reject the null hypothesis in favor of the alternative at the 1% level: at least  $\beta_1$  or  $\beta_2$  is different from zero.

### Variable selection

It is a thing to test whether at least one of the variables is related to the repsonse variable, it is another to find which of these is.

A first idea would be to test for each variable if its associated coefficient is statistically different from zero:

$$\begin{cases} \mathsf{H}_0: \beta = 0\\ \mathsf{H}_1: \beta_j \neq 0 \end{cases}, j = 1, 2, \dots, p.$$

The T statistic associated with this test writes:

$$T = \frac{\hat{\beta}_j - \beta_{j,H_0}}{\hat{\sigma}_{\hat{\beta}_j}} \sim St(n - p - 1,)$$

where  $\beta_{j,H_0}$  is the value of  $\beta_j$  under the null hypothesis.

Ewen Gallic
### Tests on the coefficients

To perform this bilateral test at an  $\alpha$  level, we can get the quantiles  $-t_{\alpha/2}$  and  $t_{\alpha/2}$  from a Student distribution such as:

$$\mathbb{P}\left(-t_{\alpha/2} < \frac{\hat{\beta}_j - \beta_{j,H_0}}{\hat{\sigma}_{\hat{\beta}_j}} < t_{\alpha/2}\right) = 1 - \alpha.$$

From the estimates, we can compute the observed value of the T statistic as:

$$t_{j,\text{obs.}} = \frac{\hat{\beta}_j}{\hat{\sigma}_{\hat{\beta}_j}}.$$

### Tests on the coefficients

The decision rule is:

- if  $t_{j,\text{obs.}} \in [-t_{\alpha/2}, t_{\alpha/2}]$ :
  - non-rejection region: we do not reject  $H_0$  at the  $\alpha$  level
  - $\beta_j$  is not statistically different from zero
- if  $t_{j,\text{obs.}} \notin [-t_{\alpha/2}, t_{\alpha/2}]$ :
  - rejection region: we reject  $H_0$  at the  $\alpha$  level
  - $\beta_j$  is statistically different from zero

### Tests on the coefficients

When the number p of predictors is larger than the number of observed values n, it is not even possible to fit the multiple linear regression model using least squares:

- testing the coefficients one by one is therefore not possible
- performing the F-test is not possible either.

In that case, choosing which variables to keep in the model requires a different approach, such as:

- forward/backward/bi-directional selection
- reducing the dimension.

# Selecting variables

Most of the time, not all predictors are associated with the response.

It is then possible to select a model with a subset of predictors. But then, which one should we choose?

The basic idea is to use a metric to compare models with each other, e.g.:

- the Akaike Information Criterion (AIC)
- the Bayesian Information Criterion (BIC)
- the adjusted  $R^2$
- Mallow's  $C_p$
- . . .

But, with p variables, there is a total of  $2^p$  different models that can be estimated using a subset of p :

• fitting all the possible subset is not to be considered.

# Selecting variables

Some **recursive algorithms** can be used to perform variable selection, without screening all the possible models:

#### • forward selection:

- starting with a model with an intercept but not predictor
- choosing the first variable to be included by fitting p regressions and selecting the one with the lowest RSS
- Finding another variable to be added by fitting p-1 regressions and selecting the one with the lowest RSS
- and so on, util a stopping rule is satisfied

#### • backward selection:

- starting with a model with an intercept and all predictors
- removing the variable with the largest p-value
- estimating the new model without that variable and remove the variable with the largest value
- and so on, until a stoppig rule is satisfied

#### • bidirectional elimination:

a combination of the forward and backward selection methods.

# Measuring the quality of fit

The quality of fit can be assessed using some metrics (RSE,  $R^2, \ldots$ ), as in the simple linear case.

- While in the simple linear regression case, the  $R^2$  is equal to the squared value of the correlation of the response and the variable.
- In the multiple linear regression, it can be shown that it is equal to the square of the correlation between the response and the fitted linear model

In the multiple linear regression case, the value of the  $R^2$  increases with the number of predictor introduced in the model:

• adding another variable allows fitting better the training data

# Prediction error

The prediction given by the linear model carries multiple sources of errors.

One is related to the reducible error.

Recall that the least squares plane is given by

$$\hat{y} = X\hat{\beta},$$

which is an estimation for the true population regression plane:

$$f(\mathbf{X}) = \mathbf{X}\boldsymbol{\beta}$$

Providing a confidence interval to the prediction allows us to determine how close the prediction is to the true value.

# Prediction error

The model itself carries a reducible error: we assume a linear model, which is usually an approximation for the true form of the relationship between the response and the predictors.

Finally, the prediction contains an irreductible error coming from the error term  $\varepsilon$  of the model. By using prediction intervals, we can account for this error term.

# 2.3 Qualitative Predictors

# Qualitative predictors

We have so far used two predictors in our example: years since Ph.D and years of service. These two variabls were considered as real-valued.

Now, we will consider another type of predictor: the qualitative predictors.

In the example of the salaries, some information regarding the gender of the professor is provided (Female/Male), the rank (Professor, Associate Professor, Assistant Professor), and the discipline (Theoretical/Applied departmens).



# 2.3.1 Two levels

### Two levels

Let us first focus on qualitative predictors with only two levels. This is the case of the gender variable in the data.

It is a factor variable, created as a dummy variable:

$$x_i = \begin{cases} 1 & \text{if the } i \text{th person is female} \\ 0 & \text{if the } i \text{th person is male} \end{cases}, \quad i = 1, \dots, n$$

The model thus writes:

.

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i = \begin{cases} \beta_0 + \beta_1 + \varepsilon_i & \text{if the } i \text{th person is female} \\ \beta_0 + \varepsilon_i & \text{if the } i \text{th person is male} \end{cases}$$

### Two levels

The interpretation of the constant  $\beta_0$  therefore changes. It should now be viewed as the average salary for male professors. The average salary among female professors is equal to  $\beta_0 + \beta_1$ .

	Least squares
(Intercept)	$11.51^{***}$
	(0.16)
genderFemale	$-1.41^{**}$
	(0.51)
$R^2$	0.02
Adj. $R^2$	0.02
Num. obs.	397
RMSE	3.00

\*\*\*p < 0.001, \*\*p < 0.01, \*p < 0.05

Table 3: Statistical models



#### Two levels

The average salary for male professor is therefore \$11.509 thousand dollars for 9 months, while it is only \$11.509 - 1.409 = 10.1 for women. This difference is significative at the 5% level.

Coding "Female" as 0 and "Male" as 1 does not change the regression fit, but it changes the interpretation:

	Least squares
(Intercept)	$10.10^{***}$
	(0.48)
gender $Male$	$1.41^{**}$
	(0.51)
$R^2$	0.02
Adj. $R^2$	0.02
Num. obs.	397
RMSE	3.00

\*\*\*p < 0.001, \*\*p < 0.01, \*p < 0.05

Table 4: Statistical models

# 2.3.2 More than two levels

## More than two levels

Now, let us consider an example with a qualitative predictor with more than two levels: the rank (Professor, Assistant Professor, Associate Professor).

In this situation, we can create an additional dummy variable:

The first one would be, let us say:

$$x_{1i} = \begin{cases} 1 & \text{if the } i\text{th person is professor} \\ 0 & \text{if the } i\text{th person is not professor} \end{cases}, \quad i = 1, \dots, n$$

And the second:

$$x_{2i} = \begin{cases} 1 & \text{if the } i\text{th person is associate professor} \\ 0 & \text{if the } i\text{th person is not associate professor} \end{cases}, \quad i = 1, \dots, n$$

# More than two levels

The model then writes:

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i \\ &= \begin{cases} \beta_0 + \beta_1 + \varepsilon_i & \text{if the } i\text{th person is professor} \\ \beta_0 + \beta_2 + \varepsilon_i & \text{if the } i\text{th person is associate professor} \\ \beta_0 + \varepsilon_i & \text{if the } i\text{th person is assistant professor} \end{cases} \end{aligned}$$

- $\beta_0$  is the average 9 months salary for assistant professor
- $\beta_1$  is the difference in the average 0 months salary between the assistant professor and professor categories
- $\beta_0 + \beta_1$  is the average 9 months salary for professor
- β<sub>2</sub> is the difference in the average 0 months salary between the assistant professor and associate professor categories
- $\beta_0 + \beta_2$  is the average 9 months salary for associate professor

└─2. The linear regression └─2.3. Qualitative Predictors └─2.3.2. More than two levels

## More than two levels

	Least squares	
(Intercept)	8.08***	
	(0.29)	
rankProf	$4.60^{***}$	
	(0.32)	
rankAssocProf	$1.31^{**}$	
	(0.41)	
$R^2$	0.39	
Adj. $R^2$	0.39	
Num. obs.	397	
RMSE	2.36	
***p < 0.001, **p < 0.01, *p < 0.05		

Table 5: Statistical models

Let us consider adding some **interaction terms** to the model. In previous estimations, we have assumed that the effect on the response of changing the value of one predictor is independent of the values of the other predictors. This assumption is known as the **additive assumption**.

Let us suppose that the salary of professors depends on the number of years since Ph.D and on the gender of the individual. The model writes:

$$\begin{split} \mathsf{salary}_i &= \beta_0 + \beta_1 \mathsf{Years\ since\ Ph}.\mathsf{D}_i + \beta_2 \mathsf{Gender}_i + \varepsilon_i \\ &= \begin{cases} (\beta_0 + \beta_2) + \beta_1 \mathsf{Years\ since\ Ph}.\mathsf{D}_i + \varepsilon_i & \text{if\ the\ }i\text{th\ person\ is\ female} \\ \beta_0 + \beta_1 \mathsf{Years\ since\ Ph}.\mathsf{D}_i + \varepsilon_i & \text{if\ the\ }i\text{th\ person\ is\ male} \end{cases} \end{split}$$

This corresponds to fitting two slopes: one for the females and another for males.



Now, let us consider that the the effect of a unit increase in the number of years since Ph.D may be different depending on the gender of the professor. The model now writes:

 $\begin{aligned} \mathsf{salary}_{i} &= \beta_{0} + \beta_{1} \mathsf{Years} \text{ since } \mathsf{Ph}.\mathsf{D}_{i} + \beta_{2} \mathsf{Gender}_{i} + \beta_{3} \mathsf{Years} \text{ since } \mathsf{Ph}.\mathsf{D}_{i} \times \mathsf{Gender}_{i} + \varepsilon_{i} \\ &= \begin{cases} (\beta_{0} + \beta_{2}) + (\beta_{1} + \beta_{3}) \mathsf{Years} \text{ since } \mathsf{Ph}.\mathsf{D}_{i} + \varepsilon_{i} & (\mathsf{female}) \\ \beta_{0} + \beta_{1} \mathsf{Years} \text{ since } \mathsf{Ph}.\mathsf{D}_{i} + \varepsilon_{i} & (\mathsf{male}) \end{cases} \end{aligned}$ 





As the slope of the line for women is larger than that of men, this suggests that the effect on salary of an additional year since Ph.D is larger for women than it is for men.

This result may sound odd, as we would have (unfortunately) expected the contrary.

Why do we observe such a result?

Two reasons may explain that:

- we can look at the coefficient of the interaction term between the number of years since Ph.D and gender (next slide): it is not significant;
- 2. contrary to men, there is no observations for women who got their Ph.D more than 39 years ago. We saw that the relationship between salary and the number of years since Ph.D does not seem linear and shows a hill-shaped effect. Hence, the non-linearity not well accounted for in the estimation lowers the slope for men due to values corresponding to large number of years since Ph.D. This is not observed for women, since such values are not in the data.

	Without interaction	With interaction
(Intercept)	9.31***	$9.41^{***}$
	(0.29)	(0.29)
yrs.since.phd	$0.10^{***}$	$0.09^{***}$
	(0.01)	(0.01)
genderFemale	-0.79	$-2.02^{*}$
-	(0.47)	(0.92)
yrs.since.phd:genderFemale		0.07
		(0.05)
$R^2$	0.18	0.19
Adj. $R^2$	0.18	0.18
Num. obs.	397	397
RMSE	2.75	2.74

 $p^{***}p < 0.001, p^{**}p < 0.01, p^{*}p < 0.05$ 

Table 6: Statistical models

# Accounting for non-linear effects

The relationship between salary and years since Ph.D. does not seem to be linear. It may be a good idead to try to look at a quadratic effect instead, by introducing the squared value of number of years since Ph.D:

salary<sub>i</sub> =  $\beta_0 + \beta_1$ Years since Ph.D<sub>i</sub> +  $\beta_2$ Years since Ph.D<sub>i</sub><sup>2</sup> +  $\varepsilon_i$ 

	Polynomial regression
(Intercept)	$6.51^{***}$
	(0.39)
yrs.since.phd	$0.41^{***}$
	(0.04)
yrs.since.phd_squared	$-0.01^{***}$
	(0.00)
$R^2$	0.31
Adj. R <sup>2</sup>	0.31
Num. obs.	397
RMSE	2.52

\*\*\*p < 0.001, \*\*p < 0.01, \*p < 0.05

Table 7: Statistical models

# 2.5 Working with wrong models

# Working with wrong models

Among the problems that may occur when we fit a linear regression model (James et al., 2013):

- non-linearity of the relationship between y and  ${\bf x}$
- correlation of error terms
- non-constant variance of error terms
- outliers
- high-leverage points
- collinearity

When facing models that are wrong, Berk (2008) recalls that two approaches can be used:

- 1. Patching up models that are mispecified
- 2. Working with misspecified models

Let us begin by talking about the last point, using some illustration.

#### Illustration: linear regression

- red lines: true conditional means (nature's response surface)
- vertical black dotted lines: distribution of y values around each conditional mean (also from nature), assuming the same variance for each conditional distribution
- The relationship between y and x seems approximatively quadratic
- Red circle: an observed value, realization of *y*



Figure 11: Estimation of a nonlinear response surface under the true linear model perspective. (Source: Berk 2008).

#### Illustration: linear regression

- Let us assume a linear model:  $y_i = \beta_0 + \beta_1 \mathbf{x}_i + \varepsilon_i$  for which we obtain the estimates  $\hat{\beta}_0$ ,  $\hat{\beta}_1$ and  $\hat{\sigma}^2$
- Dashed blue line: estimated mean function
- Solid blue line: expectation of the mean function

Estimation Using a Linear Function



Figure 12: Estimation of a nonlinear response surface under the true linear model perspective. (Source: Berk 2008).

### Illustration: linear regression

- Blue arrow: bias at a value  $\mathbf{x}_i$ (Bias  $(\hat{f}(\mathbf{x}_0))$ )
- Magenta arrow: random variation  $\left( \mathbb{V}ar\left( \hat{f}(\mathbf{x}_{0}) \right) \right)$
- Green arrow: irreducible error  $(\mathbb{V}ar(\varepsilon))$

Estimation Using a Linear Function



Figure 13: Estimation of a nonlinear response surface under the true linear model perspective. (Source: Berk 2008).

## Illustration: non-linear function

- The three sources of error remains when using a nonlinear function
- Still not possible to know the bias...



Figure 14: Estimation a nonlinear response surface under the true nonlinear model perspective. (Source: Berk 2008).

2. The linear regression

2.5. Working with wrong models

-2.5.1. Correlation of the error terms

## 2.5.1 Correlation of the error terms

### Correlation of the error terms

Recall the linear model:

$$y_i = \beta_0 + \beta_1 \mathbf{x}_{1i} + \ldots + \beta_p \mathbf{x}_{pi} + \varepsilon_i, \quad i = 1, \ldots, n,$$
(22)

where  $\varepsilon_i$  is an error term normally distributed with 0 mean and variance  $\sigma^2$ , *i.e.*,  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ , and where  $\mathbf{x}_{ji}$  represents the *i*th observation for the *j*th predictor,  $j = 1, \ldots, p$ .

We assume that the error terms are uncorrelated, *i.e.*,  $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ .

## Correlation of the error terms

If there is correlation between the  $\varepsilon_i$ :

- the estimation of the **standard errors** underestimate the true standard errors
- **confidence/prediction intervals** are therefore narrower than they should be
- p-values associated with the model will be lower than they should be
- -2. The linear regression
  - -2.5. Working with wrong models
    - 2.5.2. Non-constant variance of error terms

#### 2.5.2 Non-constant variance of error terms

#### Non-constant variance of error terms

In the linear model, we also assume that the variance of the error term is constant:  $\mathbb{V}ar(\varepsilon_i) = \sigma^2$  for all i = 1, ..., n.

Once again, if it is not the case (if there is **heteroscedasticity**), this has consequences on the estimation of the standard errors, on the confidence and prediction intervals and also on the p-values associated with the model.

In presence of **heteroscedasticity**, a way to tackle the issue is to transform the data using a concave function (such as the log function).

Another way of getting around the problem is to estimate the model by **weighted least squares**, where the weights are proportional to the inverse variance.



# 2.5.3 Outliers

The prediction of some points may be relatively far from the observed value. These points are called **outliers**.

They can be the result of an incorrect recording, or the observation can come from a sub-population.

To detect such points, Cornillon and Matzner-Løber (2007) suggests using **standardized residuals**.

Normalized residuals are given by:

$$r_i = \frac{e_i}{\sigma\sqrt{1 - h_{ii}}},\tag{23}$$

where  $h_{ij}$  is the (i, j)th element of the matrix  $\boldsymbol{X} \left( \boldsymbol{X}^{\top} \boldsymbol{X} \right)^{-1} \boldsymbol{X}^{\top}$ Replacing  $\sigma$  by its estimate  $\hat{\sigma}$  gives the standardized residuals:

$$t_i = rac{e_i}{\hat{\sigma}\sqrt{1-h_{ii}}},$$
 (24)  
Machine learning and statistical learning 112/160

Ewen Gallic

The standardized residuals are not independent by construction (the residual variance  $\hat{\sigma}^2$  was estimated with all data):

- they cannot be representative of an absence or a presence of autocorrelation
- but they have the same variance unit and can therefore be used to detect residuals with high variance.

However, Cornillon and Matzner-Løber (2007) suggest that we should use **studentized residual** (obtained by cross validation) instead of the standardized residuals:

$$t_i^{\star} = \frac{e_i}{\hat{\sigma}_{(i)}\sqrt{1 - h_{ii}}},\tag{25}$$

where  $\hat{\sigma}_{(i)}$  is the estimation of  $\sigma$  by least squares based on all observed valued except for the  $i{\rm th}.$ 

It can be shown, assuming the residuals are normally distributed, that  $t_i^\star\sim \mathcal{S}\sqcup (n-p-1).$ 

Using these studentized residual, we can define an **outlier** as a point  $(\mathbf{x}_i, y_i)$  for which the value associated with  $t_i^*$  is high, compared to the threshold given by a Student distribution, *i.e.*:

$$|t_i^{\star}| > t_{n-p-1} (1 - \alpha/2)$$



As an illustration, let us consider the cas in which we regress the salary of professors on the number of years since their Ph.D, the same value squared, the gender and the discipline (*i.e.*, a total of 5 regressors).



2. The linear regression

2.5. Working with wrong models

2.5.4. High-leverage points

# 2.5.4 High-leverage points

# High-leverage points

While **outliers** are observations for which the response  $y_i$  is unusual given the predictors, **high leverage points** are observations which have unusual value for  $x_i$ .

Let us recall that:

$$\hat{oldsymbol{y}} = oldsymbol{X} \left(oldsymbol{X}^ op oldsymbol{X}
ight)^{-1}oldsymbol{X}^ op oldsymbol{y}$$

For the *i*th observation:

$$\hat{y}_i = \sum_{j=1}^n h_{ij} y_j = h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j,$$

where  $h_{ij}$  is the (i, j)th element of the matrix  $oldsymbol{X} \left( oldsymbol{X}^{ op} oldsymbol{X} 
ight)^{-1} oldsymbol{X}^{ op}.$ 

This allows us to know the weight of the observation on its prediction, through  $h_{ii}$ .

Ewen Gallic

# High-leverage points

This leads to the definition of a **leverage point**, provided by Cornillon and Matzner-Løber (2007):

- A point is a leverage point if the values  $h_{ii}$  of the projection matrix  $\boldsymbol{X} \left( \boldsymbol{X}^{\top} \boldsymbol{X} \right)^{-1} \boldsymbol{X}^{\top}$  are greater than:
  - $h_{ii} > 2p/n$  according to Hoaglin and Welsch (1978)
  - ▶  $h_{ii} > 3p/n$  for p > 6 and n p > 12 according to Velleman and Welsch (1981)
  - $h_{ii} > 0.5$  according to Huber (1981)

# 2.5.5 Collinearity

# Collinearity

When two or more predictors are closely related, we face a phenomenon known as **collinearity**.

This is the case of the predictors "Years since Ph.D" and "Years of service":





#### Collinearity

The presence of collinearity may be the source of problems when estimating a linear model:

- it can then become difficult to disentangle the individual effects of collinear variables on the response
- the variance of at leat one of the estimated coefficients  $\hat{\beta}_j$  tends to be inflated

As a consequence, since the t-statistic for each predictor uses the estimated variance of the coefficient, it can lead to a p-value lower tan it should be.

Looking at the correlation matrix of the predictors may help identifying possible problems of collinearity.

But collinearity can exist between three or more variables. In that case, known as **multicollinearity**, looking at the correlation matrix does not help.

## Multicollinearity

There are multiple ways of detecting the presence of multicollinearity. One of those consists in computing the variance inflation factor (VIF):

$$\mathsf{VIF}\left(\hat{\beta}_{j}\right) = \frac{1}{1 - R_{\mathbf{x}_{j}|\mathbf{x}_{-j}}^{2}},\tag{26}$$

where  $R_{\mathbf{x}_j|\mathbf{x}_{-j}}^2$  is the  $R^2$  obtained from a regression of  $\mathbf{x}_j$  onto all the other predictors  $\mathbf{x}_{-j}$ .

- The smallest value for VIF is 1: complete absence of collinearity
- When the value is high (> 5 or > 10): we can suspect the presence of multicollinearity, due to the predictor  $\mathbf{x}_i$

When facing multicollinearity, a simple solution consists in dropping one of the problematic variables.

# 3. Quantile Regression

### Some references

- Arellano, M (2009). Quantile methods, Class notes
- Charpentier, A. (2018). Big Data for Economics, Lecture 3
- Givord, P., D'Haultfoeuillle, X. (2013). La régression quantile en pratique, INSEE
- He, X., Wang, H. J. (2015). A Short Course on Quantile Regression.
- Koenker and Bassett Jr (1978). *Regression quantiles*. Econometrica: journal of the Econometric Society (46), 33–50
- Koenker (2005). Quantile regression. 38. Cambridge university press.

# 3.1 Introduction

### Introduction

In the linear regression context, we have focused on the conditional distribution of y, but only paid attention to the **mean effect**.

In many situations, we only look at the effects of a predictor on the conditional mean of the response variable. But there might be some asymetry in the effects across the quantiles of the response variable:

• the effect of a variable could not be the same for all observations (*e.g.*, if we increase the minimum wage, the effect on wages may affect low wages differently than high wages).

Quantile regression offers a way to account for these possible asymmetries.



## Quantiles

Let us consider a random variable Y, with cumulative distribution function F:

$$F_Y(y) = \mathbb{P}(Y \le y).$$

For any  $0 < \tau < 1$ , the  $\tau$ th quantile of Y is defined as:

$$Q\tau(Y) = F^{-1}(\tau) = \inf \left\{ x \in \mathbb{R} : F_Y(x) \ge \tau \right\}$$

The most used quantiles are:

- $\tau = 0.5$ : the median
- $\tau = \{0.1, 0.9\}$ : the first and last deciles
- $\tau = \{0.25, 0.75\}$ : the first and last quartiles



## Quantiles

Figure 15: Quantiles of the  $\mathcal{N}(0, 1)$  distribution.

# 3.2 Principles

# Principles

Let Y be the response variable we want to predict using p+1 predictors X (including the constant).

In the linear model, using least squares, we write:

$$Y = \boldsymbol{X}^{\top} \boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

with  $\varepsilon$  a zero mean error term with variance  $\sigma^2 I_n$ .

Thus, the conditional distribution writes:

$$\mathbb{E}\left(Y \mid X = \mathbf{x}\right) = \mathbf{x}^{\top} \boldsymbol{\beta}$$

Here,  $\beta$  represents the marginal change in the mean of the response Y to a marginal change in x.

### Principles of quantile regression

Instead of looking at the mean effect, we can look at the effect at a given quantile  $\tau$ . The conditional quantile is defined as:

$$Q_{\tau}(Y \mid X = \mathbf{x}) = \inf\{y : F(y \mid \mathbf{x}) \ge \tau\}$$

The linear quantile regression model assumes:

$$Q_{\tau}(Y \mid X = \mathbf{x}) = \mathbf{x}^{\top} \boldsymbol{\beta}_{\tau}$$
(27)

where  $\boldsymbol{\beta}_{\tau} = \begin{bmatrix} \beta_{0,\tau} & \beta_{1,\tau} & \beta_{2,\tau} & \dots & \beta_{p,\tau} \end{bmatrix}^{\top}$  is the quantile coefficient:

it corresponds to the marginal change in the *τ*th quantile following a marginal change in x.

If we assume that  $Q_{\tau}(\varepsilon \mid X) = 0$ , then Eq. (27) is equivalent to:

$$Y = \boldsymbol{X}^{\top} \boldsymbol{\beta}_{\tau} + \varepsilon_{\tau}$$
 (28)



#### Location shift model

Let us now consider a simple model:

$$Y = \boldsymbol{X}^{\top} \boldsymbol{\gamma} + \boldsymbol{\varepsilon}$$
 (29)

We have:

$$Q_{\tau}(Y \mid X = x) = \mathbf{X}^{\top} \gamma + Q_{\tau}(\varepsilon)$$

Hence, in this model known as the location model, the only coefficient that varies accordingly with  $\tau$  is the coefficient associated with the constant:  $\beta_0, \tau = \gamma_1 + Q_\tau(\varepsilon)$ 

- The conditional distribution  $F_{Y|X=x}$  are parallel when x varies.
- As a consequence, the conditional quantiles are linearly dependant on X, and the only coefficient that varies with the quantile is  $\beta_{0,\tau}$ , *i.e.*, the coefficient associated with the constant :

• 
$$\beta_{0,\tau} = \gamma_1 + Q_{\tau}(\varepsilon)$$
  
•  $\beta_{j,\tau} = \gamma_j$  for all the coefficients except the constant.

#### Location shif model: exemple

Consider the following true process:  $Y = \beta_1 + \beta_2 \mathbf{x}_2 + \varepsilon$ , with  $\beta_0 = 3$  and  $\beta_1 = -.1$ , and where  $\varepsilon \sim \mathcal{N}(0, 4)$ .

Let us generate 500 observations from this process.



#### Location-scale model

Now, let us consider a **location-scale model**. In this model, we assume that the predictors have an impact both on the mean and on the variance of the response:

$$Y = X^{\top} \boldsymbol{\beta} + (X^{\top} \boldsymbol{\gamma}) \varepsilon,$$

where  $\varepsilon$  is independent of X.

As  $Q_{\tau}(aY+b) = aQ_{\tau}(Y) + b$ , we can write:

$$Q_{\tau}(Y \mid X = \mathbf{x}) = \mathbf{x}^{\top} (\boldsymbol{\beta} + \gamma Q_{\tau}(\varepsilon))$$

- By posing  $\beta_{\tau} = \beta + \gamma Q_{\tau}(\varepsilon)$ , the assumption (27) still holds.
- The impact of the predictors will vary accross quantiles
- The slopes of the lines corresponding to the quantile regressions are not parallel

#### Location-scale model: example

Consider the following true process:  $Y = \beta_1 + \beta_2 \mathbf{x}_2 + \varepsilon$ , with  $\beta_0 = 3$  and  $\beta_1 = -.1$ , and where  $\varepsilon$  is a normally distributed error with zero mean and non-constant variance.

Let us generate 500 observations from this process, and then estimate a quantile regression on different quantiles.



## 3.3 Estimation



# 3.3.1 Definitions

## Asymetric absolute loss

Let us define the asymetric absolute loss function, also called the **check function**, as follows:

$$\rho_{\tau}(u) = (\tau - \mathbb{1}(u < 0)) \times u,$$
(30)

for  $0 < \tau < 1$ . This loss function is

- a continuous piecewise linear function
- non differentiable at u = 0

Figure 16: Quantile regression loss function  $\rho_{\tau}$ .

### Asymetric absolute loss

With  $\rho_{\tau}(u)$  used as a loss function, it can be shown that  $Q_{\tau}$  minimizes the expected loss, *i.e.*:

$$Q_{\tau}(Y) \in \operatorname*{arg\,min}_{m} \left\{ \mathbb{E} \left[ \rho_{\tau}(Y-m) \right] \right\}$$
(31)

We can note that in the case in which  $\tau=1/2,$  this corresponds to the median.

Quantiles may not be unique:

- any element of  $\{x \in \mathbb{R} : F_Y(x) = \tau\}$  minimizes the expected loss
- if the solution is not unique, we have an interval of  $\tau$ th quantiles
  - the smalest element is chosen (this way, the quantile function remains left-continuous).

# Empirically

Now, let us turn to the estimation. Let us consider a random sample  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ .

In the case of the least square estimation, the expectation  $\mathbb{E}$  minimizes the risk that corresponds to the quadratic loss function, *i.e.*:

$$\mathbb{E}(Y) = \operatorname*{arg\,min}_{m} \left\{ \mathbb{E}\left[ (Y - m)^2 \right] \right\}$$

- The sample mean solves  $\min_m \sum_{i=1}^n (y_i m)^2$
- The least squares estimates of the parameters are obtained by minimizing  $\sum_{i=1}^n (y_i \mathbf{x}_i^\top \beta)^2$

− 3. Quantile Regression ↓ 3.3. Estimation ↓ 3.3.1. Definitions

#### Empirically

We have seen that the  $\tau$ th quantile minimizes the risk associated with the asymetric absolute loss function, *i.e.*:

$$Q_{\tau}(Y) \in \operatorname*{arg\,min}_{m} \left\{ \mathbb{E} \left[ \rho_{\tau}(Y-m) \right] \right\}$$

The  $\tau$ th sample quantile of Y solves:

$$\min_{m} \sum_{i=1}^{n} \rho_{\tau}(y_i - m)$$

If we assume that  $Q_{\tau}(Y \mid X) = X^{\top} \beta_{\tau}$ , then, the quantile estimator of the parameters is given by

$$\hat{\boldsymbol{\beta}}_{\tau} \in \operatorname*{arg\,min}_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^{n} \rho_{\tau} (y_i - \mathbf{x}_i^{\top} \boldsymbol{\beta}) \right\}$$
(32)



# 3.3.2 Computation

# Computation

The Eq. (32) does not admit an explicit solution, which leads us to use numerical optimization.

Recall the expression of the loss function: it is non differentiable at u = 0, thus making standard optimization algorithms useless here.

However, we will see that the minimization program given by Eq. (32) can be rewritten as a linear optimization problem.

−3. Quantile Regression → 3.3. Estimation → 3.3.2. Computation

## Linear programming

Some reminders about linear programming.

Let us consider the following minimization problem:

$$egin{aligned} \min_{m{y}\in\mathbb{R}^n}ig\{m{y}^{ op}m{b}ig\} \ ext{s.t.}\ m{y}^{ op}m{A}\geqm{c}^{ op}, \end{aligned}$$

where  $y_1 \ge 0, \ldots, y_n \ge 0$ , A is an  $n \times p$  matrix,  $b \in \mathbb{R}^n$  and  $c \in \mathbb{R}^p$ The dual maximization problem writes:

$$\max_{\mathbf{x}\in\mathbb{R}^n} \left\{ \boldsymbol{c}^{\top} \mathbf{x} \right\}$$
  
s.t.  $A\mathbf{x} \leq \mathbf{x} \geq 0$ .
└─3. Quantile Regression └─3.3. Estimation └─3.3.2. Computation

#### Computation

The linear quantile regression model can be rewritten as:

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_\tau + \varepsilon_i = \mathbf{x}_i^\top \boldsymbol{\beta}_\tau + (u_i - v_i),$$
(33)

where  $u_i = \varepsilon_i \mathbb{1}(\varepsilon_i > 0)$  and  $v_i = |\varepsilon_i| \mathbb{1}(\varepsilon_i < 0)$ .

In that case, we have:

$$\begin{split} \rho_{\tau}(u_{i} - v_{i}) &= (\tau - \mathbb{1}(u_{i} > v_{i})) \cdot (u_{i} - v_{i}) \\ &= \tau u_{i} - \tau v_{i} - \mathbb{1}(u_{i} > v_{i}) \cdot u_{i} + \mathbb{1}(u_{i} > v_{i}) \cdot v_{i} \\ &= \begin{cases} \tau u_{i} - 1 & \text{if } u_{i} > v_{i} \\ v_{i} - \tau v_{i} & \text{if } u_{i} \leq v_{i} \end{cases} \end{split}$$

Indeed, if  $u_i > v_1$ , then  $v_i = 0$  and if  $u_i \le v_i$  then  $u_i = 0$ .

### Computation

Hence, the minimization program given by Eq. (32) can be rewritten as:

$$\min_{\boldsymbol{\beta}, \boldsymbol{u}, \boldsymbol{v}} \left\{ \tau \mathbf{1}^{\top} \boldsymbol{u} + (1 - \tau) \mathbf{1}^{\top} \boldsymbol{v} \right\}$$
s.t.  $\boldsymbol{y} = \boldsymbol{X} \boldsymbol{\beta} + \boldsymbol{u} - \boldsymbol{v},$ 
(34)

where  $oldsymbol{u},oldsymbol{v}\in\mathbb{R}^n_+$ 

The dual version of this program is:

$$\max_{d} \left\{ \boldsymbol{y}^{\top} \boldsymbol{d} \right\}$$
(35)  
s.t.  $\boldsymbol{y}^{\top} \boldsymbol{d} = (1 - \tau) \boldsymbol{X}^{\top} \boldsymbol{1},$ 

where  $\boldsymbol{d} \in [0,1]^n$ 

### Computation

Several methods exist to estimate this linear optimization problem.

- ? suggest using the simplex method
- ? suggest using the Frisch-Newton interior point method, which is more efficient when n become larger

# 3.4 Inference

### Asymptotic distribution

As there is no explicit form for the estimate  $\hat{\beta}_{\tau}$ , achieving its **consistency** is not as easy as it is with the least squares estimate.

First, we need some assumptions to achieve consistency:

- the observations  $\{(\mathbf{x}_1,y_1),\ldots,(\mathbf{x}_n,y_n)\}$  must be conditionnaly i.i.d.
- the predictors must have a bounded second moment, *i.e.*,  $\mathbb{E}\left[||X_i||^2\right] < \infty$
- the error terms  $\varepsilon_i$  must be continuously distributed given  $X_i$ , and centered, *i.e.*,  $\int_{-\infty}^0 f_{\varepsilon}(\epsilon) d\epsilon = 0.5$
- $\left[f\varepsilon(0)\boldsymbol{X}\boldsymbol{X}^{\top}\right]$  must be positive definite (*local identification* property).

### Asymptotic distribution

Under those weak conditions,  $\hat{\boldsymbol{\beta}}_{\tau}$  is asymptotically normal:

$$\sqrt{n} \left( \hat{\boldsymbol{\beta}}_{\tau} - \boldsymbol{\beta}_{\tau} \right) \xrightarrow{\mathrm{d}} \mathcal{N} \left( 0, \tau (1 - \tau) D_{\tau}^{-1} \Omega_x D_{\tau}^{-1} \right)$$
(36)

where

$$D_{\tau} = \mathbb{E}\left[f\varepsilon(0)\boldsymbol{X}\boldsymbol{X}^{\top}\right] \text{ and } \Omega_{x} = \mathbb{E}\left[\boldsymbol{X}\boldsymbol{X}^{\top}\right]$$

The asymptotic variance of  $\hat{m{eta}}$  writes, for the location shift model (Eq. 29):

$$\hat{\mathbb{V}}ar\left[\hat{\boldsymbol{\beta}}_{\tau}\right] = \frac{\tau(1-\tau)}{\left[\hat{f}_{\varepsilon}(0)\right]^2} \left(\frac{1}{n}\sum_{i=1}^n \mathbf{x}_i^{\top}\mathbf{x}_i\right)^{-1}$$
(37)

where  $\hat{f}_{\varepsilon}(0)$  can be estimated using an histogram (see ?).

Ewen Gallic

### Confidence intervals

If we have obtained an estimator of  $\mathbb{V}ar\left[\hat{\beta}_{\tau}\right]$ , the confidence interval of level  $1-\alpha$  is given by:

$$\left[\hat{\beta}_{\tau} \pm z_{1-\alpha/2} \sqrt{\hat{\mathbb{V}}ar\left[\hat{\boldsymbol{\beta}}_{\tau}\right]}\right]$$
(38)

Otherwise, we can rely on bootstrap or resampling methods (see Emmanuel Flachaire's course):

- Generate a sample  $\{(\mathbf{x}_1^{(b)}, y_1^{(b)}), \dots, (\mathbf{x}_n^{(b)}, y_n^{(b)})\}$  from  $\{(\mathbf{x}_1, y_1, \dots, (\mathbf{x}_n, y_n)\}$
- Then estimate  $\beta_{\tau}^{(b)}$  using  $\hat{\boldsymbol{\beta}}_{\tau}^{(b)} = \arg\min\left\{\rho_{\tau}\left(y_{i}^{(b)} \mathbf{x}_{i}^{(b)\top}\boldsymbol{\beta}\right)\right\}$
- Do the two previous staps B times
- The bootstrap estimate of the variance of  $\hat{eta}_{ au}$  is then computed as:

$$\hat{\mathbb{V}}ar^{\star}\left(\hat{eta}_{ au}
ight)=rac{1}{B}\sum_{b=1}^{B}\left(\hat{eta}_{ au}^{\ (b)}-\hat{eta}_{ au}
ight)^{2}$$
Machine learning and statistical learning 151/160

Ewen Gallic

# 3.5 Example

Let us consider an example provided in Koenker (2005), about the impact of demographic characteristics and maternal behaviour on the birth weight of infants botn in the US.

The data concerns 198,377 babies born in 1997 from american mothers aged bewteen 18 and 45.

The predictors used are:

- education (less than high school (ref), high school, some college, college graduate)
- prenatal medical care (no prenatal visit, first prenatal visit in the first trimester of pregnancy (ref), first visit in the second trimester, first visit in the last trimester)
- the sex of the infant
- the marital status of the mother
- the ethnicity of the mother
- the age of the mother (quadratic effect)
- whether the mother smokes
- the number of cigarette per day
- the gain of weight for the mother (quadratic effect)



Figure 17: Quantile regression for birth weight (Source: Koenker 2005).



Figure 18: Quantile regression for birth weight (Source: Koenker 2005).

- Intercept: the estimated conditional quantile function of the birth-weight distribution of a girl born to an unmarried, white mother with less than a high school education who is 27 years (sample mean) old and had a weight gain of 30 pounds (sample mean), did not smoke and had her first prenatal visit in the first trimester of the pregnancy.
- Boy:
  - OLS: boys are about 100 grams bigger than girls according to the OLS estimates of the mean effect
  - Quantile: the disparity is much smaller in the lower quantiles of the distribution and larger than 100 grams in the upper tail of the distribution

Example

Ethnicity: the difference in birth weight between a baby born to a black mother and a white mother at the 5th percentile of the conditional distribution is roughly 330 grams

Smoking: smoking during pregnancy is associated with a decrease of roughly 175 grams in birth weight

In the lower tail of the conditional distribution, mothers who are roughly 30 years of age have the largest children, but in the upper tail it is mothers who are 35–40 who have the largest children.



Figure 19: Mother's age effect on Birth weight (Source: Koenker 2005).

Berk, R. A. (2008). Statistical learning from a regression perspective, volume 14. Springer.

- Cornillon, P.-A. and Matzner-Løber, É. (2007). Régression: théorie et applications. Springer.
- Hoaglin, D. C. and Welsch, R. E. (1978). The hat matrix in regression and anova. The American Statistician, 32(1):17–22.
- Huber, P. J. (1981). Robust Statistics. John Wiley & Sons, Inc.
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). An introduction to statistical learning, volume 112. Springer.
- Koenker, R. (2005). Quantile regression. Number 38. Cambridge university press.
- Koenker, R. and Bassett Jr, G. (1978). Regression quantiles. Econometrica: journal of the Econometric Society, pages 33–50.
- Velleman, P. F. and Welsch, R. E. (1981). Efficient computing of regression diagnostics. The American Statistician, 35(4):234–242.