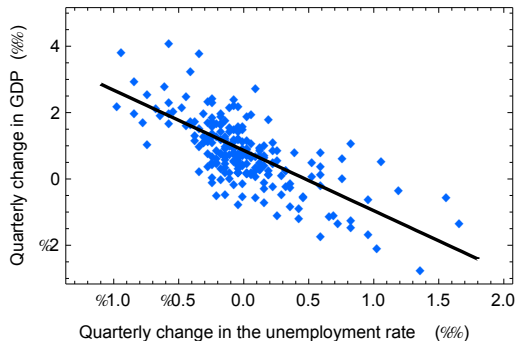


Linear regression

Eléments de Statistiques pour les Data Sciences

Simple linear regression



In economics, Okun's law is an empirical relationship between the increase in unemployment rate x and the increase in GDP y .

In statistics x and y are called

y the *response* (or *dependent variable*)

x the *explanatory* (or *independent*) *variable*

- What is the "best" linear function of x , so of the form $ax + b$, to approximate y ?
- We will define "the best" as the one which minimizes the *mean squared error* (MSE).

This is the problem of linear regression. We talk about *simple* linear regression when there is a single explanatory variable.

Simple linear regression from a sample : statement

We consider a collection of observations (x_i, y_i) and consider the question: What is the linear (or more precisely affine) transformation of x that best approximates y in the least square sense?

$$\min_{a,b \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n (y_i - a x_i - b)^2$$

$$\min_{a,b \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n (y_i - a x_i)^2 - 2b \frac{1}{n} \sum_{i=1}^n (y_i - a x_i) + b^2$$

so that $\hat{b} = \bar{y} - a\bar{x}$ for a given value of a . Replacing b by its optimal value we get:

$$\min_{a \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n (\check{y}_i - a \check{x}_i)^2 \quad \text{with} \quad \check{x}_i := x_i - \bar{x}, \check{y}_i := y_i - \bar{y}.$$

$$\min_{a \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \check{y}_i^2 - 2a \frac{1}{n} \sum_{i=1}^n \check{x}_i \check{y}_i + a^2 \frac{1}{n} \sum_{i=1}^n \check{x}_i^2.$$

Simple linear regression from a sample : statement

$$\min_{a \in \mathbb{R}} \frac{1}{n-1} \sum_{i=1}^n \check{y}_i^2 - 2a \frac{1}{n-1} \sum_{i=1}^n \check{x}_i \check{y}_i + a^2 \frac{1}{n-1} \sum_{i=1}^n \check{x}_i^2$$

We consider the sample variances and covariance

$$s_x^2 := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad s_y^2 := \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2, \quad s_{xy} := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

and the correlation coefficient $r := \frac{s_{xy}}{s_x s_y}$.

We can rewrite the optimization problem as

$$\min_{a \in \mathbb{R}} s_y^2 - 2a s_{xy} + a^2 s_x^2, \quad \text{so that} \quad \hat{a} = \frac{s_{xy}}{s_x^2} = \frac{r s_x s_y}{s_x^2} = r \frac{s_y}{s_x}.$$

Simple linear regression in the sample case: solution

We found that the best affine function of X to approximate Y in the least square sense is of the form

$$\hat{y} := \hat{f}(x) := \hat{a}x + \hat{b} \quad \text{with} \quad \hat{a} = r \frac{s_y}{s_x}, \quad \hat{b} = \bar{y} - \hat{a}\bar{x}.$$

So that for any value x we have the *estimated response* \hat{y} .

$$\hat{y} = \hat{f}(x) = \bar{y} + r \frac{s_y}{s_x} (x - \bar{x}).$$

If we consider the residual $e_i := y_i - f^*(x_i)$, then we can write

$$y_i = \hat{y}_i + e_i = \bar{y} + r s_y \frac{x_i - \bar{x}}{s_x} + e_i.$$

Properties of the residuals $e_i = y_i - \hat{y}_i$

Given that $\hat{b} = \bar{y} - \hat{a}\bar{x}$, we have $\hat{y}_i - \bar{y} = \hat{a}x_i + \hat{b} - \bar{y} = \hat{a}(x_i - \bar{x})$, we have

$$\sum_{i=1}^n e_i = \sum_{i=1}^n (\hat{y}_i - y_i) = \sum_{i=1}^n (\hat{y}_i - \bar{y}) = \sum_{i=1}^n \hat{a}(x_i - \bar{x}) = 0.$$

$$\begin{aligned} \sum_{i=1}^n e_i x_i &= \sum_{i=1}^n e_i (x_i - \bar{x}) = \sum_{i=1}^n (\hat{y}_i - y_i)(x_i - \bar{x}) = \sum_{i=1}^n (\hat{y}_i - \bar{y})(x_i - \bar{x}) - \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) \\ &= \sum_{i=1}^n \hat{a}(x_i - \bar{x})(x_i - \bar{x}) - (n-1)rs_x s_y = (n-1)\hat{a}^2 s_x^2 - (n-1)rs_x s_y = 0. \end{aligned}$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y} + e_i)^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n e_i (\hat{y}_i - \bar{y}) + \sum_{i=1}^n e_i^2,$$

$$\text{but } \sum_{i=1}^n e_i (\hat{y}_i - \bar{y}) = \hat{a} \sum_{i=1}^n e_i (x_i - \bar{x}) = 0$$

Summary: properties of the residuals $e_i = y_i - \hat{y}_i$

Let $(x_1, y_1), \dots, (x_n, y_n)$ be a collection of datapoints for linear regression.

Let

- s_x^2 and s_y^2 be the *sample variances*.
- $s_{xy} := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$ the *sample covariance*.
- $r := \frac{s_{xy}}{s_x s_y}$ the *sample correlation*.

Let $\hat{a} := r \frac{s_y}{s_x}$, $\hat{b} := \bar{y} - \hat{a}\bar{x}$, $\hat{y}_i := \hat{a}x_i + \hat{b}$.

Then the *residuals* $e_i := y_i - \hat{y}_i$ satisfy the following properties:

- They are *centered*: $\bar{e} = \frac{1}{n} \sum_{i=1}^n e_i = 0$.
- They are *empirically decorrelated* from the x_i : $\sum_{i=1}^n e_i(x_i - \bar{x}) = 0$

Empirical variances of the estimated \hat{y}_i and of the residuals e_i

We have $\bar{\hat{y}} = \bar{y} - \bar{e} = \bar{y}$. So the sample variance of \hat{y}_i is

$$\frac{1}{n-1} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \frac{1}{n-1} \sum_{i=1}^n \hat{a}^2 (x_i - \bar{x})^2 = r^2 \frac{s_y^2}{s_x^2} s_x^2 = r^2 s_y^2.$$

But we have proven that

$$s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n-1} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \frac{1}{n-1} \sum_{i=1}^n e_i^2,$$

so that

$$\frac{1}{n-1} \sum_{i=1}^n e_i^2 = s_y^2 - r^2 s_y^2 = s_y^2 (1 - r^2).$$

Pythagoras and a decomposition between explained and residual variance

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

$\overbrace{\sum_{i=1}^n (y_i - \bar{y})^2}^{(n-1)s_y^2} = \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{(n-1)r^2 s_y^2} + \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{(n-1)(1-r^2) s_y^2}$

$$\text{TSS} = \text{ESS} + \text{RSS}$$

$$\text{Total SoS} = \text{Explained SoS} + \text{Residual SoS}$$

with SoS=sum of squares.

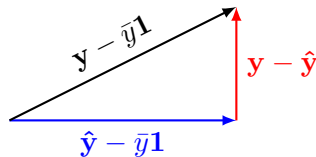
Note that

$$\hat{\mathbf{y}} - \bar{y}\mathbf{1} = \hat{a}(\mathbf{x} - \bar{x}\mathbf{1}) = r \frac{s_y}{s_x}(\mathbf{x} - \bar{x}\mathbf{1}).$$

Let

- $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$
- $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)^\top \in \mathbb{R}^n$,
- $\mathbf{1} = (1, \dots, 1)^\top \in \mathbb{R}^n$.

The decomposition of the TSS corresponds to the Pythagorean triangle:



Coefficient of determination

The *coefficient of determination* noted R^2 is defined as the fraction of the variance explained by the *explanatory variable*

$$R^2 := \frac{\text{ESS}}{\text{TSS}} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{r^2 s_y^2}{s_y^2} = r^2.$$

So the coefficient of determination is the *square of the correlation coefficient* between x and y in the data.

Linear regression with a vector of explanatory variables

Given a dataset

$$\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\},$$

with $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i \in \mathbb{R}$, we consider

- the *design matrix* \mathbf{X} and
- the vector of *responses* \mathbf{y}

defined as

$$\mathbf{X} = \begin{bmatrix} \text{---} & \mathbf{x}_1^\top & \text{---} \\ \text{---} & \mathbf{x}_2^\top & \text{---} \\ \text{---} & \vdots & \text{---} \\ \text{---} & \mathbf{x}_n^\top & \text{---} \end{bmatrix} \quad \text{and} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Remark: most of the time it is relevant to

- center the data: $\mathbf{x}_i^c = \mathbf{x}_i - \bar{\mathbf{x}}$
- normalize via e.g. $x_{ij}^s = x_{ij}^c / \hat{\sigma}_j$ or mapping \mathbf{x}_{ij}^c to $[0, 1]$, etc

Linear regression *aka* ordinary least square regression (OLS)

Given a dataset $\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ we have

$$\mathbf{y} - \mathbf{X}\boldsymbol{\beta} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \text{---} & \mathbf{x}_1^\top & \text{---} \\ \text{---} & \mathbf{x}_2^\top & \text{---} \\ \text{---} & \vdots & \text{---} \\ \text{---} & \mathbf{x}_n^\top & \text{---} \end{bmatrix} \boldsymbol{\beta} = \begin{bmatrix} y_1 - \mathbf{x}_1^\top \boldsymbol{\beta} \\ y_2 - \mathbf{x}_2^\top \boldsymbol{\beta} \\ \vdots \\ y_n - \mathbf{x}_n^\top \boldsymbol{\beta} \end{bmatrix}$$

So that we have

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \boldsymbol{\beta}^\top \mathbf{x}_i)^2 = \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$$

with

- the vector of responses $\mathbf{y}^\top = (y_1, \dots, y_n) \in \mathbb{R}^n$
- the design matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ whose i th row is equal to \mathbf{x}_i^\top .

Solving linear regression

We can rewrite the MSE as $\frac{1}{n}Q(\beta)$ with

$$Q(\beta) = \beta^\top \mathbf{X}^\top \mathbf{X} \beta - 2 \beta^\top \mathbf{X}^\top \mathbf{y} + \|\mathbf{y}\|^2.$$

A minimum has to be stationary point, i.e., such that $\nabla Q(\beta) = 0$. To compute the gradient, we can use the property that for Q differentiable, we have

$$Q(\beta + \mathbf{h}) = Q(\beta) + \nabla Q(\beta)^\top \mathbf{h} + o(\|\mathbf{h}\|),$$

where $o(\|\mathbf{h}\|)$ is a higher order term in \mathbf{h} . In our case we have

$$Q(\beta + \mathbf{h}) = Q(\beta) + \mathbf{h}^\top \mathbf{X}^\top \mathbf{X} \beta + \beta^\top \mathbf{X}^\top \mathbf{X} \mathbf{h} + \mathbf{h}^\top \mathbf{X}^\top \mathbf{X} \mathbf{h} - 2 \mathbf{h}^\top \mathbf{X}^\top \mathbf{y}$$

from which we deduce that $\nabla Q(\beta) = 2\mathbf{X}^\top \mathbf{X} \beta - 2\mathbf{X}^\top \mathbf{y}$.

Normal equations

We have thus established that the stationary points of Q satisfy the

Normal equations:

$$\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} - \mathbf{X}^\top \mathbf{y} = \mathbf{0}$$

Given that $\mathbf{X}^\top \mathbf{X}$ is a positive semi-definite matrix, the curvature of the function is non-negative everywhere and so all stationary points are global minima. The normal equation thus characterizes exactly the vectors $\boldsymbol{\beta}$ which are solutions of the linear regression problem.

If $\mathbf{X}^\top \mathbf{X}$ is invertible, then there is a unique solution to the normal equations and $\hat{\boldsymbol{\beta}}$ is given by:

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

Remarks:

- $\mathbf{X}^\top \mathbf{X}$ is invertible iff the columns of \mathbf{X} are linearly independent
- they are linearly dependent that one of them is a linear combination of the others.
- $\mathbf{X}^\top \mathbf{X}$ is never invertible for $p > n$.

Linear or affine regression?

Compare the linear vs affine functions of \mathbf{x}

$$f_{\boldsymbol{\beta}}(\mathbf{x}) = \boldsymbol{\beta}^{\top} \mathbf{x} \quad \text{vs} \quad f_{\boldsymbol{\beta},b}(\mathbf{x}) = \boldsymbol{\beta}^{\top} \mathbf{x} + b = \tilde{\boldsymbol{\beta}}^{\top} \tilde{\mathbf{x}}$$

With a new definition of the variables

$$\tilde{\boldsymbol{\beta}} = \begin{bmatrix} \boldsymbol{\beta} \\ b \end{bmatrix} \quad \text{and} \quad \tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

we can rewrite an affine model in dimension p as a linear model in dimension $p + 1$, in which the last column of the design matrix is $\mathbf{1} = (1, \dots, 1)^{\top} \in \mathbb{R}^n$.

Exercise: What is the value of \hat{b} if the data is centered?

Gaussian conditional model and linear regression

We decide to model the conditional distribution of Y given X by

$$Y \mid X \sim \mathcal{N}(\beta^\top X + b, \sigma^2)$$

or equivalently $Y = \beta^\top X + \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, \sigma^2)$.

Given a dataset $\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ we can consider the likelihood of β in the conditional model of Y given X and estimate β using the maximum likelihood principle.

Likelihood for one pair

$$p(y_i \mid \mathbf{x}_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(y_i - \beta^\top \mathbf{x}_i)^2}{\sigma^2}\right)$$

Negative log-likelihood

$$-\ell(\beta, \sigma^2) = -\sum_{i=1}^n \log p(y_i \mid \mathbf{x}_i) = \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2} \sum_{i=1}^n \frac{(y_i - \beta^\top \mathbf{x}_i)^2}{\sigma^2}.$$

Gaussian conditional model and linear regression

$$\min_{\sigma^2, \beta} \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2} \sum_{i=1}^n \frac{(y_i - \beta^\top \mathbf{x}_i)^2}{\sigma^2}$$

The minimization problem in β

$$\min_{\beta} \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$$

that we recognize as the usual linear regression.

Optimizing over σ^2 , we find:

$$\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\beta}_{\text{MLE}}^\top \mathbf{x}_i)^2$$

Properties if the data is actually Gaussian

Assume that $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\varepsilon}$ with

Full column rank *fixed* design: $\text{rank}(\mathbf{X}) = p$ (which implies $n \geq p$).

I.i.d. centered **Gaussian** noise: $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$

then

- $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \sim \mathcal{N}(\boldsymbol{\beta}^*, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$
- $S^2 = \frac{1}{n-p} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2 \sim \frac{\sigma^2}{n-p} \chi_{n-p}^2$
- $\hat{\boldsymbol{\beta}}$ and S^2 are independent

All of these are used for

- t-test and to construct confidence intervals
- **Only valid if the data is Gaussian** (= model is well-specified)