

M1 INTERMEDIATE ECONOMETRICS

Linear regression: Sampling

Koen Jochmans François Poinas

2025 — 2026

The regression slope as a random variable

With data $(Y_1, X_1), \dots, (Y_n, X_n)$ the fitted regression line has slope coefficient

$$\hat{\beta} = \arg \min_b \sum_{i=1}^n (Y_i - X_i' b)^2 = \left(\sum_{i=1}^n X_i X_i' \right)^{-1} \left(\sum_{i=1}^n X_i Y_i \right).$$

When the data change the regression line changes because $\hat{\beta}$ changes.

We see the data as a random sample from the distribution of (Y, X) .

This means that the (Y_i, X_i) are all

mutually independent, and

distributed according to the same distribution as (Y, X) .

The population regression line

Given the distribution of (Y, X) we can look for the best linear predictor of Y given $X = x$.

With mean squared error loss this is $x'\beta$ for

$$\beta = \arg \min_b \mathbb{E}((Y - X'b)^2) = \mathbb{E}(XX')^{-1}\mathbb{E}(XY).$$

Indeed, the first-order condition for a minimum here is

$$\mathbb{E}(X(Y - X'b)) = 0.$$

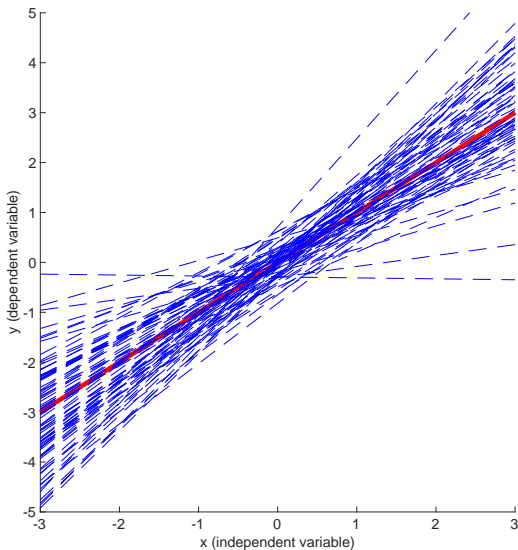
We can define the errors

$$e = Y - X'\beta$$

and write

$$Y = X'\beta + e, \quad \mathbb{E}(Xe) = 0.$$

Then $\hat{\beta}$ is an estimator of β obtained on solving the sample version of this minimization problem.



Regression lines from many different samples (blue) and the population regression line (red).

If we would know the sampling distribution of $\hat{\beta}$ we could calculate measures of spread, such as the variance, or of extreme events via tail probabilities.

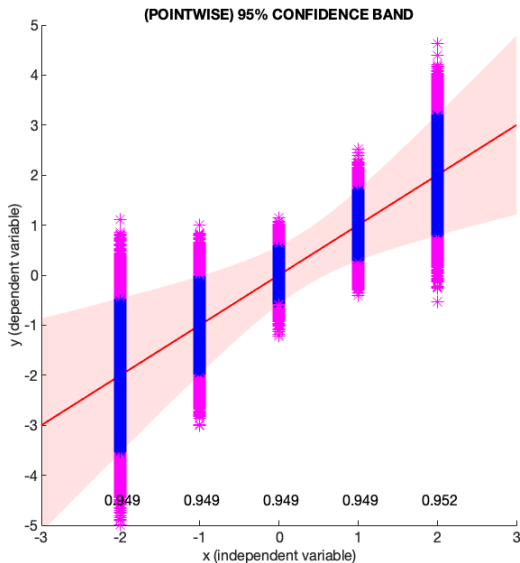
This is not possible, in general.

An exception is the classical linear regression model, where we have that

$$Y|X = x \sim N(x'\beta, \sigma^2).$$

More generally, we will need to work with an approximation to this distribution that becomes better as the sample size grows large (an asymptotic argument as $n \rightarrow \infty$).

Classical linear regression model



Here

$$\mathbf{Y} = \mathbf{X}\beta + \mathbf{e}, \quad \mathbf{e}|\mathbf{X} \sim N(0, \sigma^2 \mathbf{I}_n)$$

and so

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Y}) = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{e} = \beta + \mathbf{A}\mathbf{e}.$$

The source of randomness is the second term.

We condition on \mathbf{X} , so that \mathbf{A} is no longer random.

Then

$$\hat{\beta} - \beta|\mathbf{X} \sim \mathbf{A}\mathbf{e}|\mathbf{A} \sim N(0, \sigma^2 \mathbf{A}\mathbf{A}').$$

Because

$$\mathbf{A}\mathbf{A}' = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = (\mathbf{X}'\mathbf{X})^{-1}$$

the variance is simply $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$.

The independent normal errors and the conditioning on the regressors are both important for this result to go through. The key is that, then, $\hat{\beta} - \beta$ is a linear combination of zero-mean i.i.d. normals, which is again normal.

It follows that the regression line at x , $x'\hat{\beta}$, is (conditionally) unbiased for the population regression line $x'\beta$, with normal deviations around it.

Let $\hat{\theta} = x'\hat{\beta}$ and $\theta = x'\beta$ for a chosen value x .

Then

$$\hat{\theta} - \theta | \mathbf{X} \sim x'(\hat{\beta} - \beta) | \mathbf{X} \sim x' \mathbf{A} \mathbf{e} | \mathbf{X} \sim N(0, \sigma^2 x' (\mathbf{X}' \mathbf{X})^{-1} x).$$

This can be used to compute the confidence bands in the previous figure (where we have taken σ^2 to be known).

The standard deviation of this distribution,

$$\sigma \sqrt{x' (\mathbf{X}' \mathbf{X})^{-1} x},$$

is a measure of spread of $\hat{\theta}$ around θ .

Note how the spread depends on the point x at which the regression line is evaluated.

