

# FinMath Academic Review

## Regression

Mark Hendricks  
Financial Mathematics, University of Chicago

September 2018

## Contents

<b>1</b>	<b>Projection</b>	<b>4</b>
1.1	Environment . . . . .	4
1.2	Conditional expectation . . . . .	4
1.3	Decomposition . . . . .	4
1.4	Linear approximation . . . . .	5
1.5	Example: Linear dependence . . . . .	6
1.6	Including a constant . . . . .	6
1.7	Projection . . . . .	7
1.8	Nonlinearity . . . . .	8
1.9	Partial derivative . . . . .	9
1.10	Confounding variables . . . . .	10
<b>2</b>	<b>Stochastic processes</b>	<b>12</b>
2.1	Stationarity . . . . .	12
2.2	Ergodicity . . . . .	12
2.3	Consistency . . . . .	13
2.4	Law of Large Numbers . . . . .	13
2.5	Variance of a sum . . . . .	15

2.6	Asymptotic Distribution . . . . .	16
2.7	Special case: i.i.d. . . . .	17
<b>3</b>	<b>The Projection Estimate</b>	<b>18</b>
3.1	Estimators . . . . .	18
3.2	Sample notation . . . . .	18
3.3	Projection estimate . . . . .	20
3.4	$\mathbf{b}^*$ as a random variable . . . . .	20
3.5	Consistent projection . . . . .	21
3.6	Variability of the estimation error . . . . .	22
3.7	Asymptotic distribution of $\mathbf{b}^*$ . . . . .	22
<b>4</b>	<b>Measuring variability of <math>\mathbf{b}^*</math></b>	<b>24</b>
4.1	Errors that are i.i.d. . . . .	24
4.2	Simplifying to conditional homoskedasticity . . . . .	24
4.3	Estimation with conditional heteroskedasticity . . . . .	26
4.3.1	Restricting for large $h$ . . . . .	27
4.3.2	Ensuring the covariance is positive definite . . . . .	27
4.3.3	Assuming no serial correlation . . . . .	28
4.3.4	Putting it together to get $\Sigma_{\mathbf{b}}$ . . . . .	28
4.4	Estimating the asymptotic variance . . . . .	28
<b>5</b>	<b>Inference</b>	<b>29</b>
5.1	Notation . . . . .	29
5.2	z-statistic . . . . .	29
5.3	t-statistic . . . . .	30
5.4	Hypothesis Testing . . . . .	30
5.4.1	Right-tailed test . . . . .	30
5.4.2	Left-tailed test . . . . .	31

5.4.3	Two-sided test . . . . .	31
5.5	Errors and Power . . . . .	31
<b>6</b>	<b>Descriptive statistics</b>	<b>32</b>
6.1	R-squared . . . . .	32
6.1.1	Caveat: Regressing on a constant . . . . .	32
6.2	Problems with multicollinearity . . . . .	32
6.3	Variation in regressors . . . . .	33

# 1 Projection

## 1.1 Environment

- We are interested in a random scalar variable,  $y$ .
- We observe a random  $k \times 1$  column vector,  $\mathbf{x}$ .
- We do not put any restriction on the joint distribution of  $(y, \mathbf{x})$ .

## 1.2 Conditional expectation

Let's consider the conditional expectation in the context of an optimal approximation.

- Suppose we want to approximate  $y$  having observed  $\mathbf{x}$ .
- Let  $f(\mathbf{x})$  denote an arbitrary function, resulting in the error  $\varepsilon = y - f(\mathbf{x})$ .
- If we specify a loss function that we wish to minimize,  $\ell(\varepsilon)$ , then we can consider the optimal approximating function,  $f^*(\mathbf{x})$ , that minimizes  $\ell(\varepsilon)$ .

For the **mean squared error** loss function,

$$\ell(\varepsilon) = \mathbb{E} \left[ (y - f(\mathbf{x}))^2 \right]$$

the **conditional expectation** is the optimal approximation,

$$f^*(\mathbf{x}) = \mathbb{E} [y | \mathbf{x}]$$

## 1.3 Decomposition

In general, the conditional expectation may,

- depend on the joint distribution of  $(y, \mathbf{x})$
- be a nonlinear function of  $\mathbf{x}$ .

We can decompose  $y$  into two parts, the conditional expectation and the approximation error.

$$\begin{aligned} y &= \mathbb{E} [y | \mathbf{x}] + \varepsilon \\ 0 &= \mathbb{E} [\varepsilon | \mathbf{x}] \end{aligned}$$

Note that the second statement holds by construction. It is a significant restriction on  $\varepsilon$ , implying the following:

$$\mathbb{E}[\varepsilon | \mathbf{x}] = 0 \iff 0 = \mathbb{E}[f(\mathbf{x}) \varepsilon], \forall f(\mathbf{x})$$

which implies,

$$\mathbb{E}[\mathbf{x} \varepsilon] = \mathbf{0}, \quad 0 = \mathbb{E}[\varepsilon]$$

Furthermore, the last condition restricts the covariance and correlation:

$$\begin{aligned} \text{cov}[f(\mathbf{x}), \varepsilon] &= \mathbf{0} = \text{corr}[f(\mathbf{x}), \varepsilon], \forall f \\ \text{cov}[\mathbf{x}, \varepsilon] &= \mathbf{0} = \text{corr}[\mathbf{x} \varepsilon] \end{aligned}$$

The approximation error cannot be improved through knowledge of  $\mathbf{x}$ : it is uncorrelated to any (nonlinear) function of  $\mathbf{x}$ , and certainly is uncorrelated to  $\mathbf{x}$  itself.

## 1.4 Linear approximation

Suppose we restrict ourselves to linear approximation. If we again use the MSE loss function, we require a vector  $\beta^*$  that minimizes,

$$\ell(\varepsilon) = \mathbb{E}[(y - \mathbf{x}'\beta)^2]$$

This is optimized for the choice,  $\beta^*$  which satisfies,<sup>1</sup>

$$\mathbb{E}[\mathbf{x}(y - \mathbf{x}'\beta^*)] = \mathbf{0}$$

That, the optimal  $\beta^*$  ensures the approximation error is orthogonal to  $\mathbf{x}$ . Existence requires the following assumption,

**Assumption 1** (Identified). The  $k \times k$  matrix,  $\mathbb{E}[\mathbf{x}\mathbf{x}']$ , is nonsingular, (and thus finite.)

Note,

- If this assumption does not hold, then the linear factors,  $\mathbf{x}$  are not **identified**. That is, one of them is linearly a function of the other and thus its impact on  $y$  is not well defined in relation to the other factors.

---

<sup>1</sup>Consider the MSE for an arbitrary  $\beta$ ,

$$\begin{aligned} \ell(\epsilon) &= \mathbb{E}[(y - \mathbf{x}'\beta)^2] \\ &= \mathbb{E}[(y - \mathbf{x}'\beta^*) + \mathbf{x}'(\beta^* - \beta)]^2 \\ &= \mathbb{E}[(y - \mathbf{x}'\beta^*)^2] + 2(\beta^* - \beta)' \underbrace{\mathbb{E}[\mathbf{x}(y - \mathbf{x}'\beta^*)]}_{0 \text{ by orthogonality of } \beta^*} + \mathbb{E}[(\mathbf{x}'(\beta^* - \beta))^2] \\ &= \mathbb{E}[(y - \mathbf{x}'\beta^*)^2] + \mathbb{E}[(\mathbf{x}'(\beta^* - \beta))^2] \\ &> \mathbb{E}[(y - \mathbf{x}'\beta^*)^2] \end{aligned}$$

Thus, any choice,  $\beta$  has MSE at least as large as  $\beta^*$ .

- This assumption is not too restrictive; if one of the  $k$  elements of  $\mathbf{x}$  is an exact linear function of the others, simply drop it from consideration. It additionally requires that  $\mathbf{x}$  have a well defined (finite) mean.

Rearranging, we have

$$\boldsymbol{\beta}^* = (\mathbb{E} [\mathbf{x}\mathbf{x}'])^{-1} \mathbb{E} [\mathbf{x}y] \quad (1)$$

Thus,  $\mathbb{L}(y | \mathbf{x})$  is the optimal linear approximation of  $y$  given  $\mathbf{x}$ .

$$\mathbb{L}(y | \mathbf{x}) = \mathbf{x}'\boldsymbol{\beta}^*$$

## 1.5 Example: Linear dependence

Consider an example where Assumption 1 is violated, and  $\mathbf{x}\mathbf{x}'$  is not invertible.

- Suppose in analyzing mean excess returns for S&P 500 stocks, we wanted to also condition on whether the stock is in the Information Technology (IT) sector. We set two indicator variables in  $\mathbf{x} = \begin{bmatrix} x^{[1]} \\ x^{[2]} \end{bmatrix}$ .

$$x^{[1]} = \mathbb{I}_{\text{IT sector}}, \quad x^{[2]} = \mathbb{I}_{\text{not IT sector}}$$

- But clearly,  $x^{[1]}$  and  $x^{[2]}$  are not linearly independent:

$$x_i^{[2]} = 1 - x_i^{[1]}, \quad \forall i$$

- We need to drop either linearly dependent variable.

This example is contrived: it should have been clear we did not need to set separate complementary indicator variables. However, this issue does arise if we are setting indicator variables for many categories, and forget to exclude one category to make it the baseline.

For instance, we might look at how membership in each of the 11 GICS-defined sectors impacts a stock's mean excess returns. Then  $y$  would be mean excess return, and  $\mathbf{x}$  should have indicator variables for just 10 of the 11 sectors. Whichever sector we exclude is seen as the baseline, for which each of the elements of  $\boldsymbol{\beta}^*$  is measuring relative impact.

## 1.6 Including a constant

Suppose we allow an affine approximation by including a constant as one of the elements of  $\mathbf{x}$ .

$$\mathbf{x} = \begin{bmatrix} 1 \\ \tilde{\mathbf{x}} \end{bmatrix}$$

where  $\tilde{\mathbf{x}}$  are the non-constant elements of  $\mathbf{x}$ . The second moments are then equivalent to centered second moments:

$$\begin{aligned}\mathbb{E}[\mathbf{x}\varepsilon] &= \text{cov}[\mathbf{x}\varepsilon] \\ \mathbb{E}[\mathbf{x}\mathbf{x}'] &= \text{var}[\mathbf{x}]\end{aligned}$$

It can be shown that the linear projection is simply,

$$\mathbb{L}(y | \mathbf{x}) = \mathbb{E}[y] + (\tilde{\mathbf{x}} - \mathbb{E}[\tilde{\mathbf{x}}])' (\text{var}[\tilde{\mathbf{x}}])^{-1} \text{cov}[\tilde{\mathbf{x}}y]$$

which in the case of scalar  $\tilde{\mathbf{x}} = x$  reduces to

$$\mathbb{L}\left[y \mid \begin{bmatrix} 1 \\ x \end{bmatrix}\right] = \mathbb{E}[y] + (x - \mathbb{E}[x]) \frac{\text{cov}[x, y]}{\text{var}[x]}$$

This makes clear that including a constant is equivalent to de-meaning all the variables and then approximating.

## 1.7 Projection

$\mathbb{L}[y | \mathbf{x}]$  is decomposing  $y$  into a portion spanned by  $\mathbf{x}$  and a portion orthogonal to  $\mathbf{x}$ . Thus, it is formally a projection of  $y$  onto  $\mathbf{x}$ .<sup>2</sup>

Note that

- It is a linear projection of  $y$  in that the operator is a linear function of  $\mathbf{x}$ .
- The coefficient vector of this linear function,  $\beta^*$ , depends only on the second moments of  $(y, \mathbf{x})$ , not the entire joint distribution.

Thus we can decompose  $y$  into two parts, the linear projection and the orthogonal piece:

$$\begin{aligned}y &= \mathbb{L}[y | \mathbf{x}] + \epsilon \\ y &= \mathbf{x}'\beta^* + \epsilon\end{aligned}$$

This is indeed a projection, and we can verify that

$$0 = \mathbb{L}[\epsilon | \mathbf{x}]$$

---

<sup>2</sup>Recall that a projection is a linear operator,  $P$ , such that  $P^2 = P$ .

$$\begin{aligned}\mathbb{L}[\mathbb{L}[y | \mathbf{x}] | \mathbf{x}] &= \mathbf{x}' (\mathbb{E}[\mathbf{x}\mathbf{x}'])^{-1} \mathbb{E}[\mathbf{x} \mathbb{L}[y | \mathbf{x}]] \\ &= \mathbf{x}' (\mathbb{E}[\mathbf{x}\mathbf{x}'])^{-1} \mathbb{E}\left[\underbrace{\mathbf{x} \mathbf{x}' (\mathbb{E}[\mathbf{x}\mathbf{x}'])^{-1} \mathbb{E}[\mathbf{x}y]}_{\mathbb{L}(y | \mathbf{x})}\right] \\ &= \mathbf{x}' (\mathbb{E}[\mathbf{x}\mathbf{x}'])^{-1} \mathbb{E}[\mathbf{x}\mathbf{x}'] (\mathbb{E}[\mathbf{x}\mathbf{x}'])^{-1} \mathbb{E}[\mathbf{x}y] \\ &= \mathbf{x}' (\mathbb{E}[\mathbf{x}\mathbf{x}'])^{-1} \mathbb{E}[\mathbf{x}y] \\ &= \mathbb{L}[y | \mathbf{x}]\end{aligned}$$

alternately stated as

$$0 = \mathbb{E}[\mathbf{x}\epsilon]$$

The last equality says  $\epsilon$  is orthogonal to  $\mathbf{x}$ . Contrast this with the approximation error of the conditional expectation,  $\varepsilon$ , which is not just orthogonal to  $x$  but has zero conditional expectation.

**Theorem 1.** If  $\mathbb{E}[\epsilon | \mathbf{x}] = 0$ , then the linear projection is the conditional expectation.

$$\mathbb{E}[y | \mathbf{x}] = \mathbb{L}[y | \mathbf{x}]$$

Equivalently, if we assume that the conditional expectation is of the form,

$$\mathbb{E}[y | \mathbf{x}] = \mathbf{x}\boldsymbol{\theta} + \varepsilon$$

then the linear projection is the conditional expectation, with

$$\boldsymbol{\beta}^* = \boldsymbol{\theta}$$

## 1.8 Nonlinearity

**Example 1** (Call Option). Consider the value of a call option, denoted  $c$ , and the value of the underlying stock, denoted  $s$ . Suppose call options have the following relationship to underlying stock price:<sup>3</sup>

$$c = f(s) + \epsilon$$

where

- $f$  is a nonlinear equation
- $\epsilon$  is the impact of transaction costs, has zero mean, and is independent of  $s$ .

Then by assumption,  $\mathbb{E}[c | s] = f(s)$ , with  $\mathbb{E}[\epsilon | s] = 0$ . However, the linear projection of  $c$  onto  $\mathbf{x} = \begin{bmatrix} 1 \\ s \end{bmatrix}$  is

$$\begin{aligned} \mathbb{L}[c | \mathbf{x}] &= (\mathbb{E}[\mathbf{x}\mathbf{x}'])^{-1} \mathbb{E}[\mathbf{x}c] \\ &= \mathbb{E}[c] + \frac{\text{cov}(s, c)}{\text{var}(s)} (s - \mathbb{E}[s]) \end{aligned}$$

Thus,

$$c = \mathbb{E}[c] + \frac{\text{cov}(s, c)}{\text{var}(s)} (s - \mathbb{E}[s]) + \epsilon$$

where the only restriction on  $\epsilon$  is that

$$0 = \mathbb{E}[\mathbf{x}\epsilon] \iff 0 = \text{corr}(s, c), 0 = \mathbb{E}[\epsilon]$$

So  $\epsilon$  may predictably be small or large for certain values of  $s$ , since we do not have  $\mathbb{E}[\epsilon | s] = 0$ . However, the ways in which  $s$  predicts  $\epsilon$  are not linear, given the zero correlation.

---

<sup>3</sup>This model is motivated by the Black-Scholes modeling you will see in FINM 33000, but it simplifies a number of things to focus on our point regarding approximation.



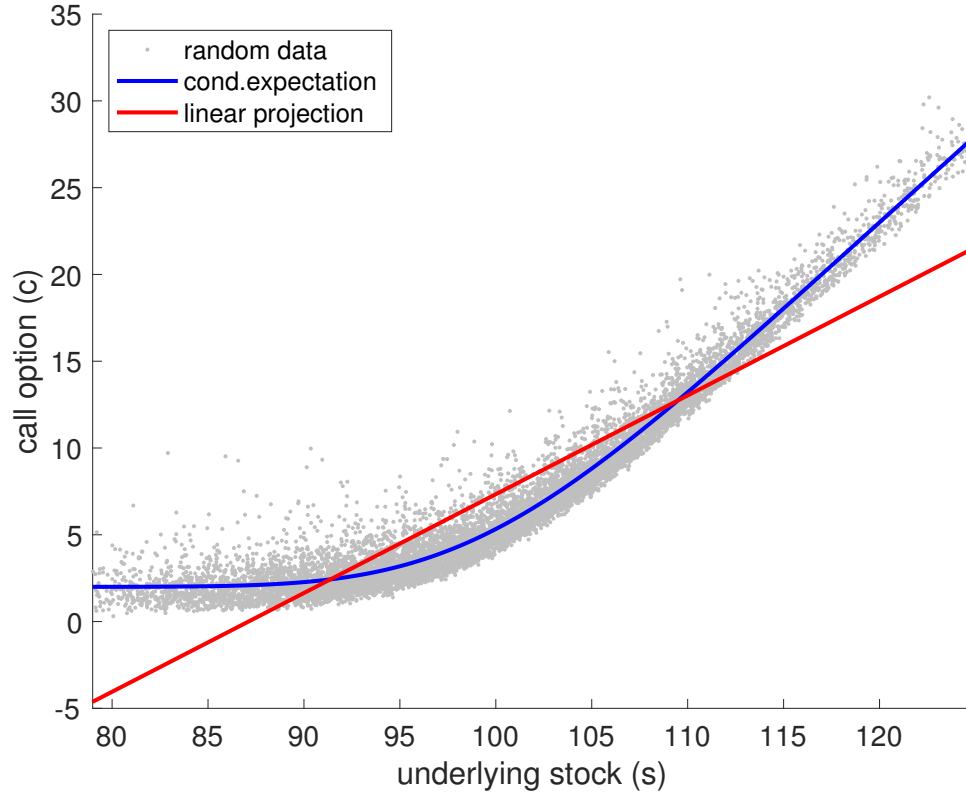


Figure 1: Illustration of Example 1

## 1.9 Partial derivative

The projection coefficient vector,  $\beta^*$  tells us that if we randomly observe an increase in  $\mathbf{x}$ , our optimal linear approximation is impacted by  $\beta^*$ .

However, we are not inferring a causal impact from  $\mathbf{x}$  to  $y$ .

- Suppose some other variable,  $z$ , causes change in  $y$ , but  $z$  is correlated with  $\mathbf{x}$ .
- The projection coefficients on  $x$  may be due to the implied change to  $z$ .
- For the conditional expectation, this indirect channel is wider: even if  $z$  is uncorrelated with  $\mathbf{x}$ , there may be nonlinear relationships between  $z$  from  $\mathbf{x}$  which are useful in approximating  $y$ .

If we do not condition on  $z$ , yet it is correlated with both  $\mathbf{x}$  and  $y$ , then it is known as a **confounding variable**. It causes the projection vector to identify not just causal impact, but inferred impact from  $z$ .

This question of **identification** is simple in theory: we should just condition on  $z$  in addition to  $\mathbf{x}$ . But in realized data, we'll see this may not be possible.

In many areas, researchers are focused on identifying causal impacts:

- $y$  = salary,  $x$  = education,  $z$  = intelligence, responsibility, etc.
- $y$  = crime rate,  $x$  = policing,  $z$  = lawfulness
- $y$  = profitability,  $x$  = leverage,  $z$  = operating margin

However, there are many cases where this is not an issue, including many areas where machine-learning is increasingly used:

- Tax authorities predicting whether someone has unreported income.
- Netflix predicting your rating of a movie.
- Assessing portfolio risk to certain factors.

## 1.10 Confounding variables

**Example 2** (Equity Premia). Suppose we want to understand how a stock's characteristics,  $\mathbf{x}$ , impacts its premium,  $y$ .<sup>4</sup>

$$\mathbb{E}[y | x] = \mathbb{L}[y | x] = x\beta^*$$

Let our population be the stocks of the S&P 500 from Jan 2005 to July 2018. Let  $\mathbf{x}$  be a scalar variable, the stock's dividend-yield. Table 1 gives the projection coefficients.

	cond. on $\mathbf{x}$	cond. on $(\mathbf{x}, z)$
dividend-yield	0.0393	(-0.0009)
volatility		0.0796

Table 1: Projection of all S&P 500 mean excess returns on the stock's own dividend yield and return volatility.

Only using  $\mathbf{x}$ , the vector  $\beta^*$  is our best (linear) estimate of the impact of a stock return. However, it is not causal. It is based on the population's covariance between dividend-yield and other relevant factors. Should a firm take this as a prescription that if it increases the dividend yield by one point, the stock's mean excess return will change as indicated by  $\beta^*$ ?

Let  $z$  denote the stock's return volatility, and include it in the projection:

$$\mathbb{L}[y | \mathbf{x}, z] = \mathbf{x}'\beta^* + z\gamma$$

See Table 1 for the projection coefficients.

<sup>4</sup>In FINM 36700, we discuss the equity premium as the mean excess return and will see that the Fundamental Theorem of Asset Pricing implies that the conditional expectation for  $y$  is linear in a certain set of variables,  $\mathbf{x}$ .

- Now, the estimated impact of dividend-yield on mean stock return almost disappears, while the impact through volatility is large.
- This happens because there is a positive covariance between dividend-yield and volatility. Thus, in the projection without volatility, dividend-yield is measuring its own impact as well as the associated impact via volatility.
- Specifically, for every point of dividend yield, we expect 0.5 points increase in standardized volatility.

From this, it seems that if a firm changes its dividend policy without making any more fundamental changes, it would lead to almost no impact on the mean excess return.

## 2 Stochastic processes

### 2.1 Stationarity

A stochastic process,  $\{z_i\}$  is **stationary** if the joint distribution between  $z_i$  and  $z_{i+h}$  depends only on  $h$ , not on  $i$ . We often only require a weaker version, **covariance-stationary**, which says that the  $\mathbb{E}[z_i]$  does not depend on  $i$  and that the covariances only depend on  $h$ .<sup>5</sup>

$$\gamma_h \equiv \text{cov}[z_i, z_{i-h}] = \text{cov}[z_j, z_{j-h}]$$

Let  $\mathbf{z}$  denote a set of  $n$  observations of the process  $\{z_i\}$ ,

$$\mathbf{z} \equiv [z_i \quad z_{i+1} \quad \dots \quad z_{i+n-1}]'$$

The  $n \times n$  covariance matrix is then restricted to being a Toeplitz matrix,

$$\text{cov}[\mathbf{z}] = \begin{bmatrix} \gamma_0 & \gamma_1 & \gamma_2 & \dots & \gamma_{n-3} & \gamma_{n-2} & \gamma_{n-1} \\ \gamma_1 & \gamma_0 & \gamma_1 & \dots & \gamma_{n-4} & \gamma_{n-3} & \gamma_{n-2} \\ \gamma_2 & \gamma_1 & \gamma_0 & \dots & \gamma_{n-5} & \gamma_{n-4} & \gamma_{n-3} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \gamma_{n-3} & \gamma_{n-4} & \gamma_{n-5} & \dots & \gamma_0 & \gamma_1 & \gamma_2 \\ \gamma_{n-2} & \gamma_{n-3} & \gamma_{n-4} & \dots & \gamma_1 & \gamma_0 & \gamma_1 \\ \gamma_{n-1} & \gamma_{n-2} & \gamma_{n-3} & \dots & \gamma_2 & \gamma_1 & \gamma_0 \end{bmatrix}$$

Thus for a set of  $n$  values from a (covariance) stationary process, the covariance matrix is restricted to  $n$  unique parameters rather than the usual  $n(n+1)/2$ .

### 2.2 Ergodicity

We give only a brief description of ergodicity and stationarity; the technical details require a separate treatment.

A stochastic process,  $\{z_i\}$  is **ergodic** if it is asymptotically independent. Specifically, for any bounded, real-valued functions  $f$  and  $g$ ,

$$\lim_{h \rightarrow \infty} |\mathbb{E}[f(z_i)g(z_{i+h})]| = |\mathbb{E}[f(z_i)]| |\mathbb{E}[g(z_{i+h})]|$$

A few comments,

- Ergodicity allows for dependence among observations and for different observations to come from different distributions.
- Intuitively, ergodicity says that, given enough time, the process is drawn across the entire population—it does not get stuck in a cycle or become forever altered based on earlier realizations.

---

<sup>5</sup>Of course, this requires that the covariances are finite.

- We typically are interested in a slightly stronger version of ergodicity that guarantees that conditional expectations converge to unconditional expectations as the forecast horizon goes to infinity. We skip the technical details.<sup>6</sup>
- This easily generalizes to vector-valued stochastic processes,  $\{\mathbf{z}_i\}$ .

These two give us the following results,

- It implies that the population moments are well approximated by sample averages, for sufficiently large samples.
- Stationarity implies that the population distribution from which  $(y_i, \mathbf{x}_i)$  is drawn is stable.

## 2.3 Consistency

Recall that a sequence of random variables,  $\{z_n\}$  **converges in probability** to  $\theta$ , if  $\forall \epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \Pr(|z_n - \theta| > \epsilon) = 0$$

We denote this,

$$z_n \xrightarrow{P} \theta$$

Suppose we have an estimator,  $A(\{x\}) = A_n$  which is a function of a random sample  $\{x\}$  of size  $n$  and is intended to estimate some parameter of interest,  $\alpha$ .

The estimator generates a sequence of random estimates, indexed by the sample size on which each estimate is based,  $\{A_n\}$ . We say this estimator is **consistent** if its probability limit equals the parameter being estimated,  $\{A_n\} \xrightarrow{P} \alpha$ .

## 2.4 Law of Large Numbers

**Assumption 2** (Ergodic Stationarity). We have a random sample that comes from a stochastic process which is **jointly stationary and ergodic**.

With this assumption, we can invoke perhaps the most important theorem in statistics: **The Law of Large Numbers (LLN)**.

**Theorem 2** (LLN). Suppose the stochastic process  $\{z\}_i$  satisfies Assumption 2, with  $\mathbb{E}[z] = \mu$ . Denote the sample average,

$$\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i$$

---

<sup>6</sup>The extra conditions are often referred to as Gordin's condition, and the main restriction is that the conditional expectation of  $z$  becomes the unconditional expectation in the limit.

$\bar{z}$  is a consistent estimator of the population mean,  $\mu$ .

$$\bar{z} \xrightarrow{P} \mu$$

- This is more general than the more familiar LLN which requires that the stochastic process is i.i.d.
- The LLN is important as it tells us that under Assumption 2, sample averages are consistent estimators for population means, including higher-order moments.

**Example 3** (Estimating a variance). Suppose we have observed  $n$  values of a scalar i.i.d. process  $\{x_i\}$ , with  $\mathbb{E}[x_i] = \mu$  and

$$\text{var}[x_i] = \mathbb{E}[x_i^2] - (\mathbb{E}[x_i])^2 = \sigma^2$$

The LLN tells us that we can consistently estimate this second moment as

$$\begin{aligned} s^2 &\equiv \frac{1}{n} \sum_{i=1}^n x_i^2 - \left( \frac{1}{n} \sum_{i=1}^n x_i \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \\ s^2 &\xrightarrow{P} \sigma^2 \end{aligned}$$

Simply taking the sample average of the targeted moment is a consistent estimator.<sup>7</sup>

**Example 4** (Lazy estimator). Suppose we estimated the mean by simply taking the average of the largest and smallest values of the sample,

$$\hat{x} = x_1$$

Not surprisingly, this is not consistent. As the sample size increases,  $\hat{x}$  ignores the additional information and simply uses the first value seen.

$$\hat{x} \not\xrightarrow{P} \mu$$

---

<sup>7</sup>One might note  $s^2$  usually adjusts for the degrees of freedom by dividing by  $n - 1$  instead of by  $n$ . However, the LLN tells us this is not required for a consistency which is concerned with the limiting behavior, and the usual degrees-of-freedom adjustment does not matter in the limit.

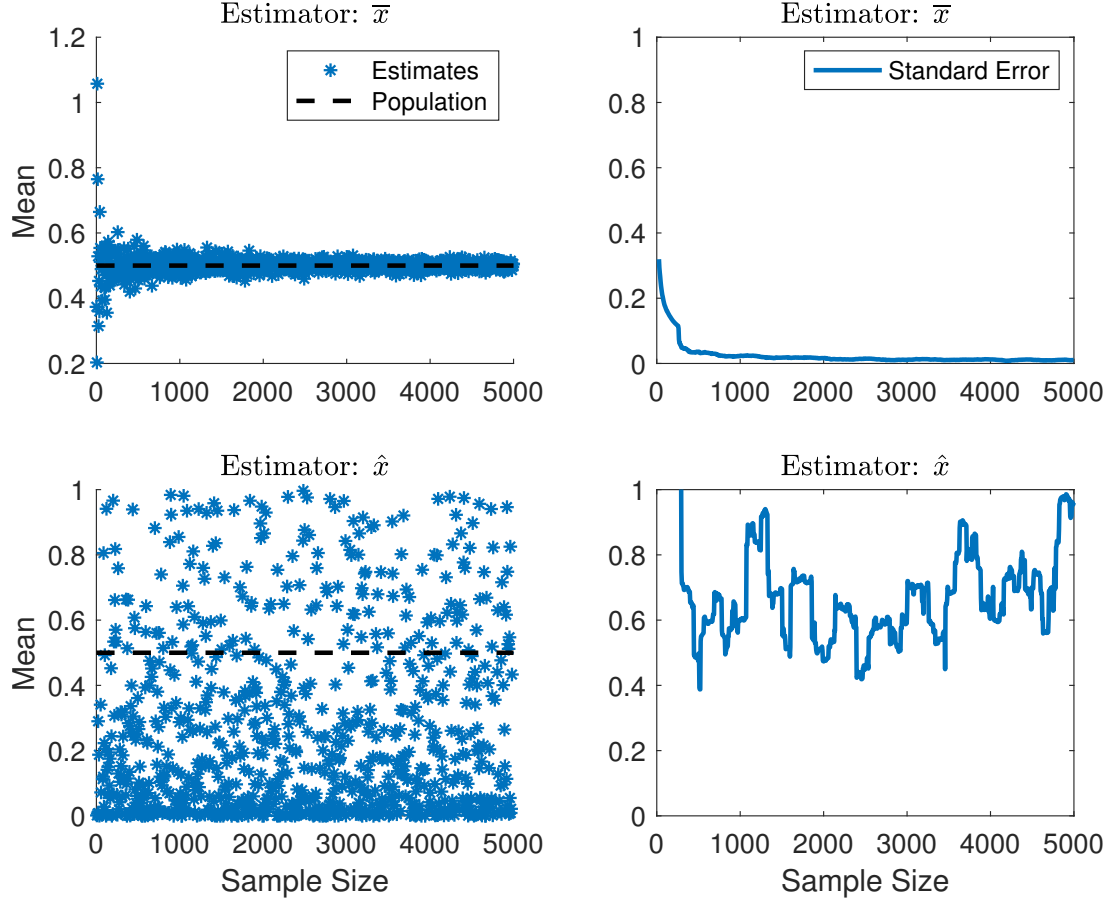


Figure 2: Data is simulated from a gamma distribution, with shape 0.5 and scale 1.0. The convergence is based on samples ranging from  $n = 1$  to  $n = 10,000$ .

## 2.5 Variance of a sum

Suppose we have a set of scalar random variables,  $z_i$ , which do not satisfy Assumption 3 but do satisfy Assumption 2, meaning they are covariance stationary. By the arithmetic of covariances, it is easy to show that

$$\begin{aligned}
 \text{var}(\bar{z}) &= \text{var}\left(\frac{1}{n} \sum_{i=1}^n z_i\right) \\
 &= \frac{1}{n^2} \text{var}\left(\sum_{i=1}^n z_i\right) \\
 &= \frac{1}{n^2} \sum_{h=-(n-1)}^{n-1} (n - |h|) \gamma_h \\
 \gamma_h &\equiv \text{cov}[z_i, z_{i-h}]
 \end{aligned}$$

Extending this for a process of random vectors,  $\{\mathbf{z}_i\}$  leads to,

$$\text{cov} [\bar{\mathbf{z}}] = \frac{1}{n^2} \sum_{h=-(n-1)}^{n-1} (n - |h|) \mathbf{\Gamma}_h \quad (2)$$

$$\mathbf{\Gamma}_h \equiv \text{cov} [\mathbf{z}_i, \mathbf{z}_{i-h}]$$

Note that we are taking advantage of the stationarity to rewrite second moments only as a function of the lag,  $h$ , not the indexing,  $(i, j)$ ,

$$\mathbf{\Gamma}_h = \mathbb{E} [\mathbf{z}_i \mathbf{z}_{i-h}' ] = \mathbb{E} [\mathbf{z}_j \mathbf{z}_{j-h}' ]$$

In fact, we can use the fact that  $\mathbf{\Gamma}_{-h} = \mathbf{\Gamma}_h'$  to write,

$$\text{cov} [\bar{\mathbf{z}}] = \frac{1}{n} \mathbf{\Gamma}_0 + \frac{1}{n^2} \sum_{h=1}^{n-1} (n - h) (\mathbf{\Gamma}_h + \mathbf{\Gamma}_h')$$

## 2.6 Asymptotic Distribution

Recall what is meant by a sequence of variables converging in distribution. For a sequence of random variables,  $\{z_n\}$ , we have a sequence of associated c.d.f.'s, denoted  $\{F_n\}$ . If

$$\{z_n\} \xrightarrow{D} z$$

implies that some variable,  $z$ , with c.d.f  $F$ , and  $\lim_{n \rightarrow \infty} F_n = F$

The **asymptotic distribution** of an estimator refers to the limiting distribution of the sequence of estimators indexed by sample size. The **Central Limit Theorem (CLT)** gives the asymptotic distribution of a sample average under broad conditions.

**Theorem 3 (CLT).** Let  $\{\mathbf{z}_i\}$  be a stochastic process satisfying Assumption 2. Define the following notation, (and assume all the moments exist and are finite.)

$$\mathbb{E} [\mathbf{z}_i] = \boldsymbol{\mu}$$

$$\text{cov} [\mathbf{z}_i \mathbf{z}_{i-h}] = \mathbf{\Gamma}_h, \forall h \geq 0$$

$$\bar{\mathbf{z}} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i$$

Then

$$\sqrt{n} (\bar{\mathbf{z}} - \boldsymbol{\mu}) \xrightarrow{D} \mathcal{N} (\mathbf{0}, \boldsymbol{\Sigma}_{\text{lim}})$$

where the asymptotic covariance must take account of all potential covariances, given that the process is not i.i.d.

$$\boldsymbol{\Sigma}_{\text{lim}} \equiv \sum_{h=-\infty}^{\infty} \mathbf{\Gamma}_h$$



## 2.7 Special case: i.i.d.

ften, the CLT is applied with the additional assumption that the process is **Independently and Identically Distributed (i.i.d.)**

**Assumption 3 (i.i.d.).** A process,  $\{y_i, \mathbf{x}_i\}$  is i.i.d.conditional on  $\{\mathbf{x}_i\}$  if the joint distribution of  $(y_i, \mathbf{x}_i)$  is identical for all  $i$  and independent of the joint distribution for  $(y_j, \mathbf{x}_j)$ , for all  $j \neq i$ .

With Assumption 3, Theorem 3 simplifies as there is no longer any covariation between  $\epsilon_i$  and  $\epsilon_j$ . Thus,  $\mathbf{\Gamma}_i = \mathbf{0}$  for all  $i \neq 0$ .

$$\mathbf{\Sigma}_{\text{lim}} = \mathbf{\Gamma}_0$$

### Example 5. Gamma Distribution

Consider variables,  $z_i$  with an i.i.d.Gamma distribution. Figure 4 shows that the sample average as an estimator of  $\mathbb{E}[z_i]$  converges to a normal distribution with smaller and smaller asymptotic variance.

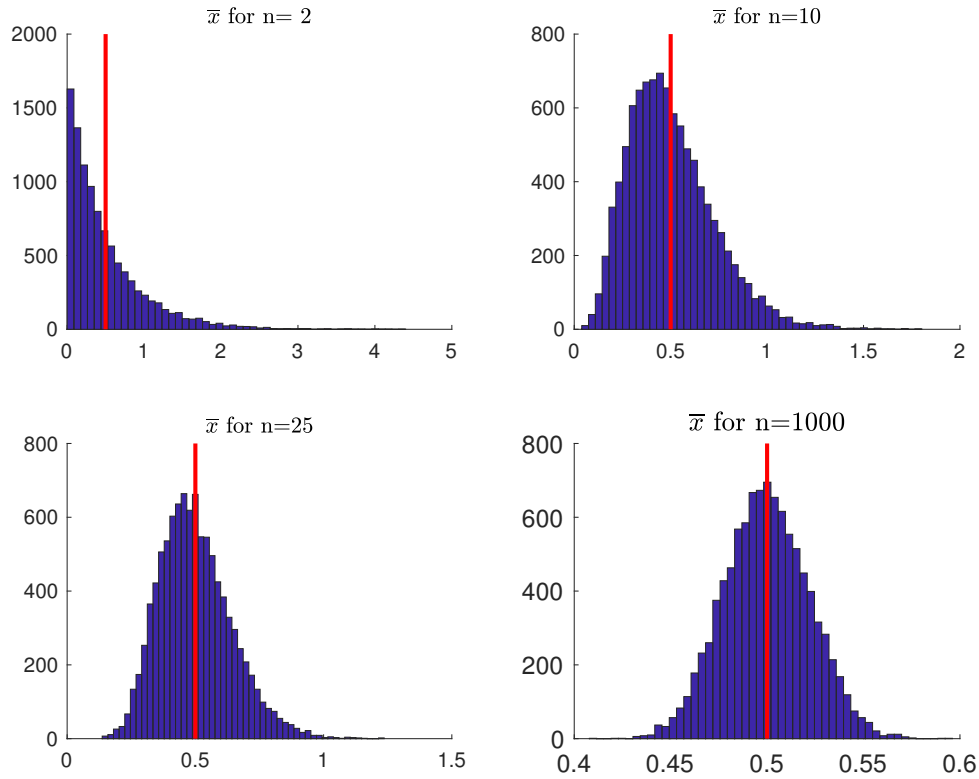


Figure 3: Illustration of the Central Limit Theorem. Sample average converges in distribution to a normal, even if the underlying data comes from a non-normal distribution. This simulation uses a Gamma distribution with shape 0.5 and scale 1.0. Estimated 10,000 times, with each sample mean based on  $n = (2, 10, 25, 1000)$ .

## 3 The Projection Estimate

### 3.1 Estimators

Estimators are functions of random data. Thus, the output of this function, the estimate, is itself a random variable. Its randomness comes from the randomness of the function's input.

For example,  $\bar{x}$  is a function of sample data  $\{x_i\}$ , and the random variation in the different samples will cause random variation in the sample average,  $\bar{x}$ . See Figure 4 to see a histogram of  $\bar{x}$  estimates, each based on samples of  $n$  data points. Note that in some samples, the resulting estimate is far from the true parameter.

Thus, we have several questions to consider when trying to estimate a parameter of interest.

- What estimator, (function,) is being used?
- Do the random estimates center around the true parameter?
- Do the random estimates tightly cluster around the true parameter?
- What is the distribution of the estimates? i.e. What is the probability that the estimate is a certain distance from the true parameter value?

We will be further analyzing each of these issues below, in the context of estimating the projection coefficient vector,  $\beta^*$ .

**Example 6** (Lazy estimator). Consider again the lazy estimator of the mean of  $x_i$ :

$$\hat{x} = x_1$$

Figure 4 shows that this estimator does not center nor bunch tightly around  $\mu$ .

### 3.2 Sample notation

- The linear projection vector in (1) depends on population moments.
- In applications, we rarely know the population; rather we must estimate based on a sample of data.

Suppose we observe  $n$  realizations of the stochastic process,  $\{y_i, \mathbf{x}_i\}$ .

- The observations are denoted  $(y_i, \mathbf{x}_i)$  for  $i = 1 \dots n$ .

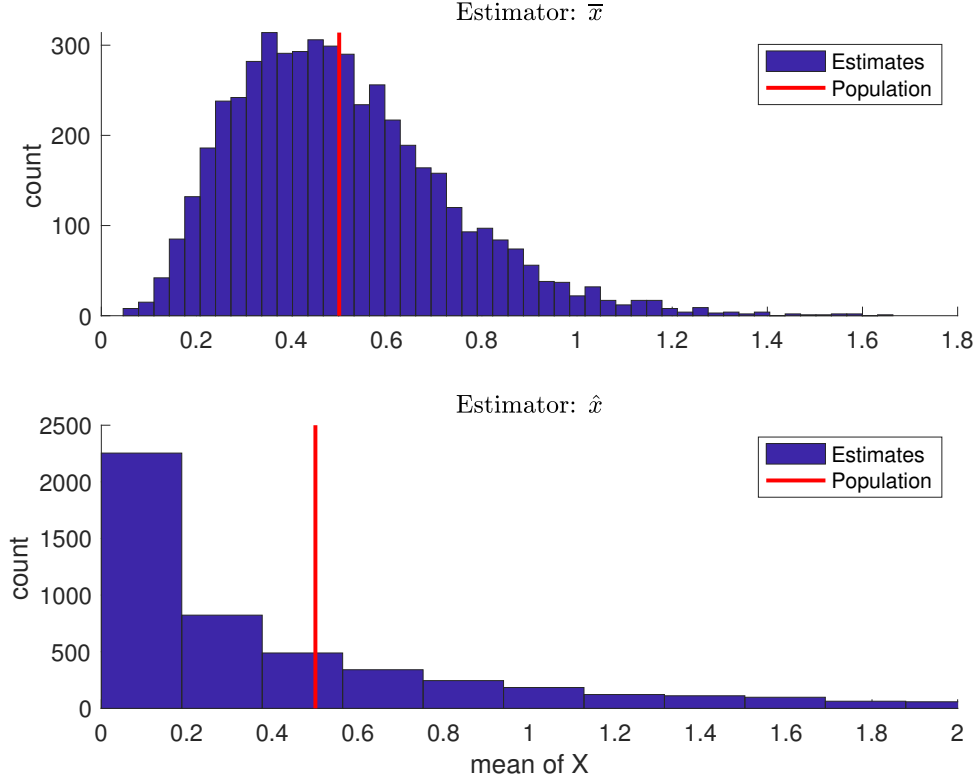


Figure 4: Data is simulated from a gamma distribution, with shape 0.5 and scale 1.0. The histogram is based on 5000 estimates of  $\bar{x}$ , which each are estimated from a sample size of 10.

- The total sample is denoted with the  $n \times 1$  vector  $\mathbf{y}$  and the  $n \times k$  vector  $\mathbf{X}$ , where each row of  $\mathbf{X}$  is an observation,  $\mathbf{x}_i'$ .

$$\underbrace{\mathbf{Y}}_{n \times 1} = \underbrace{\mathbf{X}}_{n \times k} \underbrace{\mathbf{b}}_{k \times 1} + \underbrace{\mathbf{e}}_{n \times 1}$$

Consider estimating the projection from this sample:

- An estimate of  $\beta^*$  with some coefficient vector,  $\mathbf{b}$ , will decompose the sample,  $\{y, \mathbf{x}\}_n$  into a sample of  $\mathbf{x}'\mathbf{b}$  and  $\mathbf{e}$ .
- The projection errors,  $\epsilon_i$ , associated with the sample values will be unobserved and unknown.

Consider the sample notation applied to Example 2:

$$\underbrace{\begin{bmatrix} \text{mean return}_1 \\ \text{mean return}_2 \\ \vdots \\ \text{mean return}_n \end{bmatrix}}_{\mathbf{Y}} = \underbrace{\begin{bmatrix} \text{div}_1 & \text{vol}_1 \\ \text{div}_2 & \text{vol}_2 \\ \vdots & \vdots \\ \text{div}_n & \text{vol}_n \end{bmatrix}}_{\mathbf{X}} \underbrace{\begin{bmatrix} -.0009 \\ .0796 \end{bmatrix}}_{\beta^*} + \underbrace{\begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}}_{\text{unobserved error}}$$

Contrast this with the notation for a single observation:

$$\underbrace{\text{mean return}_i}_{y_i} = \underbrace{\begin{bmatrix} \text{div}_i & \text{vol}_i \end{bmatrix}}_{(\mathbf{x}_i)'} \underbrace{\begin{bmatrix} -.0009 \\ .0796 \end{bmatrix}}_{\boldsymbol{\beta}^*} + \epsilon_i$$

### 3.3 Projection estimate

Consider estimating the projection coefficient vector,  $\mathbf{b}^*$ , using the same formula but simply replacing the population moments with their sample averages:<sup>8</sup>

$$\mathbf{b}^* = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$$

Assumption 1 is required to imply that  $\mathbf{X}'\mathbf{X}$  is nonsingular—that it has a non-zero determinant. This assumption will rarely be a problem mathematically; if it is violated, then we likely just need to drop the linearly dependent element of  $\mathbf{X}$ .<sup>9</sup> While this assumption will not be a problem mathematically, we will find that “nearly” violating it causes our biggest statistical issues.

This estimator ensures the resulting in-sample decomposition is indeed a projection,

$$\begin{aligned} \mathbf{Y} &= \mathbf{X}\mathbf{b}^* + \mathbf{e} \\ \mathbf{Y} &= \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} + \mathbf{e} \\ \mathbf{Y} &= \mathbf{P}\mathbf{Y} + \mathbf{e} \end{aligned}$$

where the projection matrix,  $\mathbf{P}$  is

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

and it is easy to check that it indeed is a projection, maintaining its own map. Furthermore, one can easily show that the sample residual is uncorrelated to  $\mathbf{X}$  and to  $\mathbf{P}\mathbf{Y}$ .

$$\mathbf{P}^2 = \mathbf{P} \quad \mathbf{P}\mathbf{e} = 0 \quad \mathbf{X}'\mathbf{e} = 0$$

Notice that the estimated errors,  $\mathbf{e}$ , will be different than the population errors,  $\boldsymbol{\epsilon}$ .

### 3.4 $\mathbf{b}^*$ as a random variable

The estimated linear projection vector is itself a random variable, which can be written as the sum of the population projection vector plus an error term: Rewrite the random variable,  $\mathbf{b}^*$ ,

$$\begin{aligned} \mathbf{b}^* &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} \\ &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'(\mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\epsilon}) \\ &= \underbrace{\boldsymbol{\beta}^*}_{\text{constant}} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\boldsymbol{\epsilon} \end{aligned} \tag{3}$$

<sup>8</sup>Below we motivate this approach by showing that under certain conditions it leads to consistent estimators. Generally, this approach is known as the method of moments.

<sup>9</sup>Technically, Assumption 1 only guarantees the condition in the population, not in any particular finite sample. Still, this means it will hold almost surely as the sample size goes to infinity.

Thus, the estimate is a constant plus a random variable estimation error:

$$\underbrace{\mathbf{b}^*}_{\text{random vector}} = \underbrace{\beta^*}_{\text{unknown constant vector}} + \underbrace{(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\epsilon}_{\text{estimation error}}$$

- $\epsilon$  is the  $n \times 1$  vector of projection errors associated with the sample data,  $\{y_i, \mathbf{x}_i\}$ . They are unobserved.
- $\mathbf{e}$  is the  $n \times 1$  vector of estimated projection errors associated with the sample data and the sample estimate,  $\mathbf{b}^*$ .
- By construction,  $\epsilon$  is orthogonal to  $\mathbf{x}$  in population, while  $\mathbf{e}$  is orthogonal to the sample  $\mathbf{X}$ .
- Even though  $\epsilon$  is orthogonal to  $\mathbf{x}$  in population, it almost surely will not be orthogonal to the sample  $\mathbf{X}$ , due to random variation.
- The estimation error,  $(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\epsilon$  is an unobserved random variable, given that it is the function of the unobserved  $n \times 1$  vector,  $\epsilon$ .

### 3.5 Consistent projection

**Theorem 4.** Given Assumptions 1 and 2, the estimator  $\mathbf{b}^*$  is a consistent estimator of  $\beta^*$ .

$$\mathbf{b}^* \xrightarrow{P} \beta^*$$

The proof is direct.<sup>10</sup> In fact, we can view  $\mathbf{b}^*$  as being based on consistent estimates for the two second

---

10

*Proof.* Recall from Equation (3),

$$\mathbf{b}^* = \beta^* + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\epsilon$$

Thus, the consistency of  $\mathbf{b}^*$  depends on two things. First, that  $\lim_{n \rightarrow \infty} (\mathbf{X}'\mathbf{X})^{-1}$  is well defined (finite), which is guaranteed by Assumption 1. It also depends on

$$\mathbf{X}'\epsilon \xrightarrow{P} \mathbf{0}$$

Given Assumption 2, the sample average converges to the population moment,

$$\mathbf{X}'\epsilon \xrightarrow{P} \mathbb{E}[\mathbf{x}\epsilon]$$

By construction, linear projection ensures a decomposition with an  $\epsilon$  such that,

$$\mathbb{E}[\mathbf{x}\epsilon] = \mathbf{0}$$

□

moments in a projection:

$$\begin{aligned}\frac{1}{n} (\mathbf{X}'\mathbf{X}) &\xrightarrow{P} \mathbb{E} [\mathbf{x}\mathbf{x}'] \\ \frac{1}{n} (\mathbf{X}'\mathbf{Y}) &\xrightarrow{P} \mathbb{E} [\mathbf{x}y] \\ (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} &\xrightarrow{P} (\mathbb{E} [\mathbf{x}\mathbf{x}'])^{-1} \mathbb{E} [\mathbf{x}y] = \boldsymbol{\beta}^*\end{aligned}$$

### 3.6 Variability of the estimation error

Though  $\mathbf{b}^*$  has arbitrary accuracy in estimating  $\boldsymbol{\beta}^*$  as  $n \rightarrow \infty$ , we want to know how much variation there is in  $\mathbf{b}^*$  for a fixed sample size  $n$ , and how fast this variation shrinks as  $n$  grows.

Consider the second moment of the estimation error.<sup>11</sup> Define the notation,

$$\mathbf{v}_i \equiv \mathbf{x}_i \epsilon_i$$

Define the sample averages estimating second moments,

$$\begin{aligned}\mathbf{S}_x &\equiv \frac{1}{n} \mathbf{X}'\mathbf{X} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' \\ \bar{\mathbf{v}} &\equiv \frac{1}{n} \mathbf{X}'\boldsymbol{\epsilon} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \epsilon_i\end{aligned}$$

We can write  $\mathbf{b}^*$  as,

$$\mathbf{b}^* = \boldsymbol{\beta}^* + \mathbf{S}_x^{-1} \bar{\mathbf{v}}$$

Then the variation in  $\mathbf{b}^*$  is

$$\begin{aligned}\boldsymbol{\Sigma}_b &= \mathbb{E} [(\mathbf{b}^* - \boldsymbol{\beta}^*)(\mathbf{b}^* - \boldsymbol{\beta}^*)' | \mathbf{X}] \\ &= \mathbb{E} [\mathbf{S}_x^{-1} \bar{\mathbf{v}} \bar{\mathbf{v}}' \mathbf{S}_x^{-1} | \mathbf{X}]\end{aligned}$$

Rewrite this as

$$\begin{aligned}\boldsymbol{\Sigma}_b &= \mathbf{S}_x^{-1} \boldsymbol{\Sigma}_v \mathbf{S}_x^{-1} \\ \boldsymbol{\Sigma}_v &\equiv \mathbb{E} [\bar{\mathbf{v}} \bar{\mathbf{v}}' | \mathbf{X}]\end{aligned}\tag{4}$$

### 3.7 Asymptotic distribution of $\mathbf{b}^*$

**Theorem 5** (Limiting distribution of  $\mathbf{b}^*$ ). Under Assumptions 1 and 2,

$$\sqrt{n} (\mathbf{b}^* - \boldsymbol{\beta}^*) \xrightarrow{D} \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_b^{\text{lim}})\tag{5}$$

---

<sup>11</sup>Throughout this section we discuss uncentered second moments rather than the covariances. This is because we want to avoid assuming that  $\epsilon$  has zero mean conditional on  $x$ , which is a much stronger assumption that guarantees our linear projection is the conditional expectation.

where we define notation,

$$\begin{aligned}\boldsymbol{\Sigma}_{\mathbf{x}} &\equiv \mathbb{E} [\mathbf{x}\mathbf{x}'] \\ \boldsymbol{\Sigma}_{\mathbf{b}}^{\text{lim}} &\equiv \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \boldsymbol{\Sigma}_{\mathbf{v}} \boldsymbol{\Sigma}_{\mathbf{x}}^{-1}\end{aligned}$$

Note that the asymptotic covariance is the same as the finite sample covariance in (4), but replacing the sample estimates  $\mathbf{S}_{\mathbf{x}}$  with the population moment  $\mathbb{E} [\mathbf{x}\mathbf{x}']$ .

We refer to  $\boldsymbol{\Sigma}_{\mathbf{b}}^{\text{lim}}$  as the asymptotic covariance matrix of  $\mathbf{b}^*$ . In applications with large  $n$ , use  $\boldsymbol{\Sigma}_{\mathbf{b}}$  as the covariance matrix of the estimation error.<sup>12</sup>

The theorem is important as it tells us how the random estimate,  $\mathbf{b}^*$  is distributed around the value we are estimating,  $\boldsymbol{\beta}^*$ . This is necessary in order to test hypotheses on  $\boldsymbol{\beta}^*$

---

<sup>12</sup>Careful readers will note that in finite samples we have not shown  $\boldsymbol{\Sigma}_{\mathbf{b}}$  is a proper covariance matrix but only a second-moment. The consistency of  $\mathbf{b}^*$  guarantees that in the limit,  $\boldsymbol{\Sigma}_{\mathbf{b}}^{\text{lim}}$  is a covariance.

## 4 Measuring variability of $\mathbf{b}^*$

This section considers how to estimate  $\Sigma_{\mathbf{b}}$  under both an i.i.d.  $\{\epsilon_i\}$  process as well as a non-i.i.d. process.

### 4.1 Errors that are i.i.d.

Let  $\mathcal{I}$  denote the identity matrix. Under Assumption 3, we have

$$\begin{aligned}\mathbb{E}[\epsilon_i \epsilon_j | \mathbf{x}] &= 0, \quad i \neq j \\ \mathbb{E}[\epsilon_i \epsilon_i | \mathbf{x}] &= \gamma_0^2, \quad i = 1, \dots, n\end{aligned}$$

This simplifies Equation (4),

$$\begin{aligned}\Sigma_{\mathbf{v}} &= \mathbb{E}[\mathbf{x}_i \epsilon_i \epsilon_i \mathbf{x}_i'] \\ &= \Sigma_{\mathbf{x}} \Sigma_{\epsilon} \\ \Sigma_{\epsilon} &\equiv \mathbb{E}[\epsilon \epsilon' | \mathbf{x}] = \gamma_0^2 \mathcal{I}\end{aligned}$$

Thus,

$$\Sigma_{\mathbf{b}} = \gamma_0^2 \Sigma_{\mathbf{x}}^{-1}$$

Estimate these moments with the sample averages,

$$\begin{aligned}s_0^2 &\equiv \frac{1}{n} \mathbf{e}' \mathbf{e} \\ \mathbf{S}_{\mathbf{b}} &= s_0^2 (\mathbf{X}' \mathbf{X})^{-1}\end{aligned}\tag{6}$$

### 4.2 Simplifying to conditional homoskedasticity

**Homoskedasticity** refers to having constant second moments.<sup>13</sup> **Unconditional homoskedasticity** says that second moments are constant:

$$\mathbb{E}[\epsilon_i \epsilon_j] = \sigma_{i,j}$$

Under the stationarity of Assumption 2,  $\{\epsilon_i\}$  is unconditionally homoskedastic:

$$\mathbb{E}[\epsilon_i \epsilon_j] = \gamma_{i-j}$$

For estimation of  $\Sigma_{\mathbf{v}}$ , we are interested in whether  $\{\epsilon_i\}$  is **conditionally homoskedastic**, which is not guaranteed by Assumption 2.

---

<sup>13</sup>This is often stated in terms of covariances, but uncentered second moments allows for the case that the conditional mean is nonzero.



**Assumption 4** (Conditional Homoskedasticity). The process  $\{\epsilon_i\}$  is homoskedastic conditional on  $\{\mathbf{x}_i\}$ .

$$\mathbb{E}[\epsilon_i \epsilon_j | \mathbf{x}_i, \mathbf{x}_j] = \sigma_{i,j}$$

Under Assumption 4,  $\Sigma_v$  simplifies to

$$\Sigma_v = \mathbf{X}' \Sigma_\epsilon \mathbf{X}$$

Estimating  $\Sigma_x$  is simple, but estimating  $\Sigma_\epsilon$  is harder given that we are assuming only conditional homoskedasticity, not i.i.d.. Thus,  $\epsilon_i$  and  $\epsilon_j$  can covary, but this covariance cannot depend on  $\mathbf{x}$ .

- To directly estimate  $\Sigma_\epsilon$ , we have only  $n$  realizations with which to estimate the  $n(n+1)/2$  unique elements.
- We could simply use the cross-products of the sample residuals to calculate each element,  $\gamma_{i,j}$ , as  $e_i e_j$ , but this is estimating making an estimate on each second moment based on a single sample cross product,  $e_i e_j$ .
- We need to add more structure in order to get any power to this estimate.

Under Assumption 2, the joint distribution of  $(\epsilon_i, \epsilon_j)$  depends only on  $h = i - j$ , not on  $i, j$ . Then for any pair of points in the sample,

$$\mathbb{E}[\epsilon_i \epsilon_j] = \mathbb{E}[\epsilon_{i+h} \epsilon_{j+h}]$$

This structure restricts the covariance matrix of  $n$  realizations of  $\{\epsilon_i\}$  to

$$\Sigma_\epsilon = \begin{bmatrix} \gamma_0 & \gamma_1 & \gamma_2 & \dots & \gamma_{n-3} & \gamma_{n-2} & \gamma_{n-1} \\ \gamma_1 & \gamma_0 & \gamma_1 & \dots & \gamma_{n-4} & \gamma_{n-3} & \gamma_{n-2} \\ \gamma_2 & \gamma_1 & \gamma_0 & \dots & \gamma_{n-5} & \gamma_{n-4} & \gamma_{n-3} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \gamma_{n-3} & \gamma_{n-4} & \gamma_{n-5} & \dots & \gamma_0 & \gamma_1 & \gamma_2 \\ \gamma_{n-2} & \gamma_{n-3} & \gamma_{n-4} & \dots & \gamma_1 & \gamma_0 & \gamma_1 \\ \gamma_{n-1} & \gamma_{n-2} & \gamma_{n-3} & \dots & \gamma_2 & \gamma_1 & \gamma_0 \end{bmatrix}$$

Given a sample size of  $n$ , we can then estimate each  $\gamma_h$  with  $s_h$ :

$$s_h = \frac{1}{n-h} \sum_{i=h+1}^n e_i e_{i-h} \quad (7)$$

Let  $\mathbf{S}_e$  denote the  $n \times n$  matrix of these estimates,  $s_h$ , for  $0 \leq h \leq n-1$ .

Then use the estimators,

$$\begin{aligned} \mathbf{S}_u &= \mathbf{X}' \mathbf{S}_e \mathbf{X} \\ \mathbf{S}_b &= \mathbf{S}_x^{-1} \mathbf{S}_u \mathbf{S}_x^{-1} \end{aligned} \quad (8)$$

Estimating  $s_h$  for  $h$  near  $n$  will still have low power and possibly give poor estimates. Thus, an additional assumption is often added:

**Assumption 5** (Orthogonality). Error terms,  $\epsilon_i$  and  $\epsilon_{i-h}$  are conditionally orthogonal beyond for large  $h$ .

$$\mathbb{E}[\epsilon_i \epsilon_j] = 0, \quad \forall |i - j| > H$$

Suppose Assumption 5 holds for  $H = 2$ ,

$$\mathbf{\Sigma}_\epsilon = \begin{bmatrix} \gamma_0 & \gamma_1 & \gamma_2 & \dots & & & \\ \gamma_1 & \gamma_0 & \gamma_1 & \dots & & & \\ \gamma_2 & \gamma_1 & \gamma_0 & \dots & & & \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ & & & \dots & \gamma_0 & \gamma_1 & \gamma_2 \\ & \mathbf{0} & & \dots & \gamma_1 & \gamma_0 & \gamma_1 \\ & & & \dots & \gamma_2 & \gamma_1 & \gamma_0 \end{bmatrix}$$

- This does not increase our power for estimating  $\gamma_h$  for  $h = 0, 1, 2$ .
- However, it eliminates the estimation of the cross products for which we had the lowest power, (due to fewer high-lag cross products.)
- So this simplification makes sense when we have theoretical reasons to believe in orthogonality past a certain lag or simply a case where we do not trust the noisily estimated high-lag terms.

### 4.3 Estimation with conditional heteroskedasticity

Without Assumption 3 nor 4, we have the general Equation (4):

$$\begin{aligned} \mathbf{\Sigma}_b &= \mathbf{S}_x^{-1} \mathbf{\Sigma}_v \mathbf{S}_x^{-1} \\ \mathbf{\Sigma}_v &\equiv \mathbb{E}[\bar{\mathbf{v}} \bar{\mathbf{v}}' | \mathbf{X}] \end{aligned}$$

Note that  $\mathbf{\Sigma}_v$  is the second moment matrix of a sample average,  $\bar{\mathbf{v}}$ . Thus, Equation (9) indicates that in general,

$$\begin{aligned} \mathbf{\Sigma}_v &= \frac{1}{n^2} \sum_{h=-(n-1)}^{n-1} (n - |h|) \mathbf{\Gamma}_h^v \\ &= \frac{1}{n^2} \left[ n \mathbf{\Gamma}_0^v + 2 \sum_{h=1}^{n-1} (n - h) (\mathbf{\Gamma}_h^v + (\mathbf{\Gamma}_h^v)') \right] \\ \mathbf{\Gamma}_h^v &\equiv \mathbb{E}[\mathbf{v}_i \mathbf{v}_{i-h}'] \end{aligned}$$

This suggests an estimator for  $\mathbf{S}_u$ , with elements that use sample averages to replace the covariance moments above.

$$\begin{aligned} \mathbf{S}_u &= \frac{1}{n^2} \sum_{h=-(n-1)}^{n-1} (n - |h|) \mathbf{S}_h \\ &= \frac{1}{n^2} \left[ n \mathbf{S}_0 + \sum_{h=1}^{n-1} (n - h) (\mathbf{S}_h + \mathbf{S}_h') \right] \end{aligned}$$

Define the estimator of the mixed second-moments,

$$\mathbf{S}_h \equiv \frac{1}{n-h} \sum_{i=h+1}^n \mathbf{u}_i \mathbf{u}'_{i-h}$$

Substituting, we have the well-known spectral density estimate or  $\Sigma_v$ ,

$$\mathbf{S}_u = \frac{1}{n^2} \left[ \sum_{i=1}^n \mathbf{u}_i \mathbf{u}'_i + \sum_{h=1}^{n-1} \sum_{i=h+1}^n (\mathbf{u}_i \mathbf{u}'_{i-h} + \mathbf{u}_{i-h} \mathbf{u}'_i) \right] \quad (9)$$

#### 4.3.1 Restricting for large $h$

The **Hansen-Hodrick correction** implements Equation (9) with additionally making Assumption 5.

$$\mathbf{S}_h = \begin{cases} \frac{1}{n-h} \sum_{i=h+1}^n \mathbf{u}_i \mathbf{u}'_{i-h} & \text{for } h < H \\ \mathbf{0} & \text{for } h \geq H \end{cases}$$

Thus the restriction limits the estimate to considering  $H$  total sample second-moments,

$$\mathbf{S}_u = \frac{1}{n^2} \left[ \sum_{i=1}^n \mathbf{u}_i \mathbf{u}'_i + \sum_{h=1}^{H-1} \sum_{i=h+1}^n (\mathbf{u}_i \mathbf{u}'_{i-h} + \mathbf{u}_{i-h} \mathbf{u}'_i) \right] \quad (10)$$

Of course, setting  $H = n$  leaves us with the unrestricted Equation (9).

#### 4.3.2 Ensuring the covariance is positive definite

Unfortunately, Equation (10) is not guaranteed to be positive definite, and often is not. The **Newey-West estimator** remedies this by putting less weight on  $\mathbf{S}$  for covariances with larger lags,  $h$ :

$$\mathbf{S}_u = \frac{1}{n^2} \left[ n \mathbf{S}_0 + \sum_{h=1}^{H-1} (n-h) w_h (\mathbf{S}_h + \mathbf{S}'_h) \right]$$

$$w_h \equiv 1 - \frac{h}{H}$$

Thus,

$$\mathbf{S}_u = \frac{1}{n^2} \left[ \sum_{i=1}^n \mathbf{u}_i \mathbf{u}'_i + \sum_{h=1}^{H-1} \left( 1 - \frac{h}{H} \right) \sum_{i=h+1}^n (\mathbf{u}_i \mathbf{u}'_{i-h} + \mathbf{u}_{i-h} \mathbf{u}'_i) \right] \quad (11)$$

The Newey West Estimator in (11) is widely implemented in computational statistics libraries.

### 4.3.3 Assuming no serial correlation

While the estimator in Equation (11) is very general, it requires trusting that the estimates for high-order serial correlation are better than a theoretical restriction that these equal zero.

- This is the usual tradeoff between assumptions and statistical power.
- In particular, financial return data often has very low serial correlation, and supported by basic financial theory.
- Thus, we might want an estimator that allows for heteroskedasticity without allowing for serial correlation.

**White's Estimator** is a popular way to do this, and is simply a restricted version of (10) and (11):

$$\mathbf{S}_u = \frac{1}{n^2} \left[ \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i' \right] \quad (12)$$

### 4.3.4 Putting it together to get $\Sigma_b$

Using the estimators in Equations (9), (10), or (11), we finally have a consistent estimator for  $\Sigma_b$ :

$$\begin{aligned} \Sigma_b &= \mathbf{S}_x^{-1} \mathbf{S}_u \mathbf{S}_x^{-1} \\ &= (\mathbf{X}\mathbf{X})^{-1} \left[ \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i' + \sum_{h=1}^{n-1} \sum_{i=h+1}^n (\mathbf{u}_i \mathbf{u}_{i-h}' + \mathbf{u}_{i-h} \mathbf{u}_i') \right] (\mathbf{X}\mathbf{X})^{-1} \end{aligned}$$

where we note that the  $\frac{1}{n^2}$  in  $\mathbf{S}_u$  cancels with the  $\frac{1}{n}$  inside each  $\mathbf{S}_x$ .

## 4.4 Estimating the asymptotic variance

Theorem 5 contains an expression for the asymptotic covariance of  $\mathbf{b}^*$ , denoted  $\Sigma_b^{\text{lim}}$ . Estimating this asymptotic covariance given a sample size of  $n$  is done by estimating  $\Sigma_b$  using Equation (11), (8), or (6) according to the assumptions used.

As a corollary, the consistency of  $\mathbf{b}^*$  ensures the consistency of using sample averages to estimate  $\Sigma_x$ ,  $\Sigma_v$ , and  $\Sigma_b$ .

$$s_0^2 \xrightarrow{P} \gamma_0^2 \quad \mathbf{S}_e \xrightarrow{P} \Sigma_e \quad \mathbf{S}_u \xrightarrow{P} \Sigma_u \quad \mathbf{S}_b \xrightarrow{P} \Sigma_b$$

## 5 Inference

### 5.1 Notation

Define the  $i$ -coordinate vector,

$$\boldsymbol{\iota}_i \equiv \begin{bmatrix} 0 & 0 & \dots & \underbrace{1}_{\text{element } i} & \dots & 0 \end{bmatrix}'$$

For any vector,  $\mathbf{z}$ , let  $\mathbf{z}_{[i]}$  denote element  $i$  of  $\mathbf{z}$ .

$$\mathbf{z}_{[i]} \equiv (\boldsymbol{\iota}_i)' \mathbf{z}$$

### 5.2 $\mathbf{z}$ -statistic

We wish to calculate probabilities of observing the estimate  $\mathbf{b}_{[i]}^*$  conditional on  $\beta_{[i]}^*$  equaling some stated value,  $\beta_i$ .

$$\Pr \left( \text{estimate for } \beta^* \leq \mathbf{b}_{[i]}^* \mid \beta_{[i]}^* = \beta_i \right) \quad (13)$$

Recall that  $\mathbf{b}^*$  has asymptotic distribution given by Equation (5),

$$\sqrt{n} (\mathbf{b}^* - \beta^*) \xrightarrow{D} \mathcal{N} \left( \mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{b}}^{\text{lim}} \right)$$

then for sample sizes with large  $n$ , we approximate,

$$\mathbf{b}^* \approx \mathcal{N} \left( \beta^*, \frac{1}{n} \boldsymbol{\Sigma}_{\mathbf{b}} \right)$$

Define the  $\mathbf{z}$ -statistic as

$$\mathbf{z}^* = \frac{\mathbf{b}_{[i]}^* - \beta_i}{\sigma_{[i]}^{\mathbf{b}}}$$

Then in terms of  $\mathbf{z}$ , the expression in (13) can be rewritten as,

$$\Pr (\mathbf{z} \leq \mathbf{z}^*) = \Phi_{\mathcal{N}} (\mathbf{z}^*)$$

- $\left( \sigma_{[i]}^{\mathbf{b}} \right)^2$  is element  $i$  of  $\text{diag} (\boldsymbol{\Sigma}_{\mathbf{b}}^{\text{lim}})$ .

$$\sigma_{[i]}^{\mathbf{b}} \equiv \sqrt{\boldsymbol{\iota}_i' \boldsymbol{\Sigma}_{\mathbf{b}} \boldsymbol{\iota}_i}$$

- $\Phi_{\mathcal{N}}$  denotes the cdf for a standard normal, and  $\mathbf{t}^* \rightarrow \mathbf{z}^*$ .

Define the  $\mathbf{p}$ -statistic as the complementary probability, which gives the probability, conditional on the assumed value for  $\beta_{[i]}^*$ , that an estimate would have been larger than the realized  $\mathbf{b}_{[i]}^*$ :

$$\mathbf{p} = 1 - \Phi_{\mathcal{N}} (\mathbf{z}^*)$$

### 5.3 t-statistic

The z-statistic above relies on  $\Sigma_{\mathbf{b}}$ , which depends on population moments which we do not know. Using the estimator,  $\mathbf{S}_{\mathbf{b}}$  changes the Normal distribution to a t-Distribution with  $n - k$  degrees of freedom.

Thus estimating (13) can be done as,

$$\begin{aligned} \mathbf{s}_{[i]}^{\mathbf{b}} &\equiv \sqrt{\boldsymbol{\iota}_i' \mathbf{S}_{\mathbf{b}} \boldsymbol{\iota}_i} \\ \mathbf{t}^* &= \frac{\mathbf{b}_{[i]}^* - \beta_i}{\mathbf{s}_{[i]}^{\mathbf{b}}} \\ \mathbf{p} &= 1 - \Phi_{\mathbf{t}}(\mathbf{t}^*, n - k) \end{aligned}$$

where  $\Phi_{\mathbf{t}}$  denotes the cdf for a t-distribution.

Notice that as  $n \rightarrow \infty$ , the t-Distribution converges to a standard normal distribution.

### 5.4 Hypothesis Testing

Being able to calculate the probability in (13) through the t-statistic gives us a way to test the hypothesis that a particular element of  $\beta^*$  equals a particular number,  $\beta$ .

The implementation depends on whether we are testing against the right-tail alternative, the left-tail alternative, or the two-sided alternative.

#### 5.4.1 Right-tailed test

$$\begin{aligned} \mathcal{H}_0 &: \beta_{[i]}^* = \beta_i \\ \mathcal{H}_1 &: \beta_{[i]}^* > \beta_i \end{aligned}$$

To implement the test, we must choose a **significance level**, which is

$$\Pr(\text{reject } \mathcal{H}_0 \mid \mathcal{H}_0 = \text{true}) = \hat{\mathbf{p}}$$

The test result is simply,

$$\begin{aligned} \mathbf{p} \geq \hat{\mathbf{p}} & \text{ Do not reject } \mathcal{H}_0 \\ \mathbf{p} < \hat{\mathbf{p}} & \text{ Reject } \mathcal{H}_0 \end{aligned}$$

Or, in terms of the t-statistic, it is

$$\begin{aligned} \mathbf{t}^* \leq \mathbf{t}_{\hat{\mathbf{p}}} & \text{ Do not reject } \mathcal{H}_0 \\ \mathbf{t}^* > \mathbf{t}_{\hat{\mathbf{p}}} & \text{ Reject } \mathcal{H}_0 \\ \mathbf{t}_{\hat{\mathbf{p}}} & \equiv \Phi_{\mathbf{t}}^{-1}(1 - \hat{\mathbf{p}}, n - k) \end{aligned}$$

### 5.4.2 Left-tailed test

Simplify modify to

$$\begin{aligned}\mathcal{H}_1 : \beta_{[i]}^* &< \beta_i \\ \mathbf{p} &\equiv \Phi_{\mathbf{t}}(\mathbf{t}^*, n - k) \\ \mathbf{t}_{\hat{\mathbf{p}}} &\equiv \Phi_{\mathbf{t}}^{-1}(\hat{\mathbf{p}}, n - k)\end{aligned}$$

### 5.4.3 Two-sided test

Suppose we test  $\mathcal{H}_0$  against a two-sided alternative,

$$\begin{aligned}\mathcal{H}_0 : \beta_{[i]}^* &= \beta_i \\ \mathcal{H}_1 : \beta_{[i]}^* &\neq \beta_i\end{aligned}$$

Then we simply modify the test as follows,

$$\begin{aligned}|\mathbf{t}^*| \leq \mathbf{t}_{\hat{\mathbf{p}}} &\text{ Do not reject } \mathcal{H}_0 \\ |\mathbf{t}^*| > \mathbf{t}_{\hat{\mathbf{p}}} &\text{ Reject } \mathcal{H}_0 \\ \mathbf{t}_{\hat{\mathbf{p}}} &\equiv \Phi_{\mathbf{t}}^{-1}\left(1 - \frac{\hat{\mathbf{p}}}{2}, n - k\right)\end{aligned}$$

## 5.5 Errors and Power

A **Type 1 error** is the probability that we reject  $\mathcal{H}_0$  even though it is true.

- Thus,  $\hat{\mathbf{p}}$  is the probability that using the test will lead to a Type 1 error.
- We can reduce this simply by choosing a lower significance level for the test.
- The significance level is commonly chosen to be 0.05.

A **Type 2 error** is the probability that we do not reject  $\mathcal{H}_0$  even though it is false.

- The more cautious we are about Type 1 errors, the more we make Type 2 errors.

The **Power** of the test is the probability that we reject  $\mathcal{H}_0$  when it is false.

- Power is the complement of the Type 2 error

$$\text{Power} = 1 - \text{Type 2 error}$$

- A consistent test has

$$\text{Power} \lim_{n \rightarrow \infty} 1$$

## 6 Descriptive statistics

### 6.1 R-squared

The **R-squared**, or coefficient of determination, in a regression is defined as

$$\begin{aligned} R_{y,x}^2 &= \frac{\text{regression sum of squares}}{\text{total sum of squares}} \\ &= 1 - \frac{\text{error sum of squares}}{\text{total sum of squares}} \end{aligned}$$

Algebraically, this is

$$\begin{aligned} R_{y,x}^2 &= \frac{\mathbf{b}^* \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \\ &= 1 - \frac{\mathbf{e}'\mathbf{e}}{\sum_{i=1}^n (y_i - \bar{y})^2} \end{aligned}$$

Intuitively, the R-squared is the square of the correlation between  $y$  and the projection of  $y$  onto  $\mathbf{x}$ .

$$R_{y,\mathbf{x}}^2 = [\text{corr}(\mathbf{Y}, \mathbf{PY})]^2$$

In a univariate regression of  $y$  on  $x$ ,

$$R_{y,x}^2 = [\text{corr}(y, x)]^2$$

#### 6.1.1 Caveat: Regressing on a constant

The interpretation and formula for R-squared does not hold if there is no constant regressor.

- Without a constant, the R-squared will not necessarily be between 0 and 1.
- Without a constant, the R-squared will not necessarily be the square of the correlation between the sample  $\mathbf{Y}$  and the projected  $Y$  values.
- Without a regressor, the fit can be improved simply by shifting the sample  $\mathbf{Y}$  data by a constant.

### 6.2 Problems with multicollinearity

If regressor  $j$  is highly correlated with the other regressors, then the variance of the coefficient estimate can be written as



$$\text{var}[b_j | \mathbf{x}] = \left( \frac{1}{1 - R_{x_j, \mathbf{x}_{[-j]}}^2} \right) \frac{\sigma^2}{\sum_{i=1}^n (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_j)^2}$$

where

$$R_{x_j, \mathbf{x}_{[-j]}}^2$$

denotes the R-squared from regressing  $x_j$  on the remaining regressors, (all columns except column  $j$  of  $\mathbf{X}$ .)

The term

$$\frac{1}{1 - R_{x_j, \mathbf{x}_{[-j]}}^2}$$

is known as the **variance inflation factor**.

- The VIF is closely related to the **condition number** of  $\mathbf{X}'\mathbf{X}$ .
- The condition number in linear algebra measures the sensitivity of inverting a matrix.
- It compares the largest and smallest eigenvalues.
- Most software packages will warn the user if the condition number (VIF) is too big.

### 6.3 Variation in regressors

The sample variance of  $x_j$  reduces the variance of the OLS estimate  $b_j$ .

- Again write

$$\text{var}[b_j | \mathbf{x}] = \left( \frac{1}{1 - R_{x_j, \mathbf{x}_{[-k]}}^2} \right) \frac{\sigma^2}{\sum_{i=1}^n (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_j)^2}$$

- The denominator of the second term is the scaled sample variance.
- If there is not enough variation in the regressor data, then the OLS estimation can not precisely estimate  $\beta_j^*$ .

The standard error of the OLS estimator depends on the variation in the regressors.

- The standard error of  $\mathbf{b}^*$  decreases as the variation in  $\mathbf{X}$  increases, holding other things equal.
- Recall that net interest margin refers to the spread in the lending and borrowing of banks.
- Consider using this as a regressor, for some financial or economic data.

## References

- [1] Fumio Hayashi, *Econometrics*. 2000.
- [2] John Cochrane, *Asset Pricing*. 2005.
- [3] Jeffrey Wooldridge, *Econometric Analysis of Cross Section and Panel Data*. 2010.
- [4] James Hamilton *Time Series Analysis*. 1994
- [5] Ruey Tsay, *Analysis of Financial Time Series*. 2010