

---

# Class Notes

ESTIMATION IN HIGH-DIMENSIONAL SPACES - ECON231C

**Mauricio Vargas-Estrada**  
Master in Quantitative Economics  
University of California - Los Angeles

---

## 1 Markov's Inequality

Being  $X$  a random variable such that  $X \geq 0$ , then:

$$\mathbb{P}(X \geq t) \leq \frac{\mathbb{E}[X]}{t}, \quad \forall t > 0$$

*Proof.* We can rewrite the left-hand side of the inequality using the indicator function:

$$X \geq X\mathbf{1}_{\{X \geq t\}}$$

The left-hand side would be greater when  $X < t$  and equal when  $X \geq t$ . Given that:

$$X \geq X\mathbf{1}_{\{X \geq t\}} \geq t\mathbf{1}_{\{X \geq t\}}$$

In this case,  $X\mathbf{1}_{\{X \geq t\}} > t\mathbf{1}_{\{X \geq t\}}$  when  $X > t$ , and  $X\mathbf{1}_{\{X \geq t\}} = t\mathbf{1}_{\{X \geq t\}}$  when  $X \leq t$  because  $X = t$  or the indicator function is zero.

Taking the expectation of the inequality:

$$\begin{aligned}\mathbb{E}[X] &\geq \mathbb{E}[X\mathbf{1}_{\{X \geq t\}}] \geq \mathbb{E}[t\mathbf{1}_{\{X \geq t\}}] \\ \mathbb{E}[X] &\geq \mathbb{E}[X\mathbf{1}_{\{X \geq t\}}] \geq t\mathbb{E}[\mathbf{1}_{\{X \geq t\}}] \\ \frac{\mathbb{E}[X]}{t} &\geq \frac{\mathbb{E}[X\mathbf{1}_{\{X \geq t\}}]}{t} \geq \mathbb{E}[\mathbf{1}_{\{X \geq t\}}]\end{aligned}$$

But  $\mathbb{E}[\mathbf{1}_{\{X \geq t\}}] = \mathbb{P}(X \geq t)$ , so:

$$\frac{\mathbb{E}[X]}{t} \geq \mathbb{P}(X \geq t)$$

□

## 2 Chebyshev's Inequality

Given a random variable  $X$  with mean  $\mu$  and variance  $\sigma^2$ , then:

$$\mathbb{P}(|X - \mu| \geq t) \leq \frac{\sigma^2}{t^2}, \quad \forall t > 0$$

*Proof.* We are going to use the fact that a strictly increasing function of a random variable does not change the probability of an event. Let  $Y = (X - \mu)^2$ . Then,

$$\mathbb{P}(|X - \mu|^2 \geq t^2) = \mathbb{P}(Y \geq t^2)$$

Using Markov's inequality, we have:

$$\mathbb{P}(Y \geq t^2) \leq \frac{\mathbb{E}[Y]}{t^2}$$

Given a random variable  $Z$ , the variance of  $Z$  is  $\text{Var}(Z) = \mathbb{E}[(Z - \mathbb{E}[Z])^2]$ . Therefore,

$$\begin{aligned} \mathbb{P}((X - \mu)^2 \geq t^2) &\leq \frac{\mathbb{E}[(X - \mu)^2]}{t^2} \\ \mathbb{P}((X - \mu)^2 \geq t^2) &\leq \frac{\text{Var}(X)}{t^2} \\ \mathbb{P}((X - \mu)^2 \geq t^2) &\leq \frac{\sigma^2}{t^2} \end{aligned}$$

□

## 3 Weak Law of Large Numbers

Given a collection of i.i.d. random variables  $\{X_i\}_{i=1}^n$ , with mean  $\mu$  and variance  $\sigma^2$ .

Defining  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ , the weak law of large numbers states that for any  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\bar{X}_n - \mu| \geq \epsilon) = 0$$

or equivalently,

$$\bar{X}_n \xrightarrow{\mathbb{P}} \mu$$

*Proof.* Calculating the variance of  $\bar{X}_n$ ,

$$\begin{aligned}\text{Var}(\bar{X}_n) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) \\ &= \frac{1}{n^2} \cdot n\sigma^2 \\ &= \frac{\sigma^2}{n}\end{aligned}$$

Then, by Chebyshev's inequality,

$$\begin{aligned}\mathbb{P}(|\bar{X}_n - \mu| \geq \epsilon) &\leq \frac{\text{Var}(\bar{X}_n)}{\epsilon^2} \\ &= \frac{\sigma^2}{n\epsilon^2}\end{aligned}$$

Taking the limit as  $n \rightarrow \infty$ ,

$$\begin{aligned}\lim_{n \rightarrow \infty} \mathbb{P}(|\bar{X}_n - \mu| \geq \epsilon) &\leq \lim_{n \rightarrow \infty} \frac{\sigma^2}{n\epsilon^2} \\ \lim_{n \rightarrow \infty} \mathbb{P}(|\bar{X}_n - \mu| \geq \epsilon) &\leq 0\end{aligned}$$

□

## 4 Hoeffding's Inequality

If  $\{X_i\}_{i=1}^n$  is a random sample from a distribution with mean  $\mu$  such that, for a number <sup>1</sup>  $a > 0$ , we have:

$$|X_i - \mu| \leq a, \quad \forall i = 1, 2, \dots, n$$

Then, for any  $t > 0$ , the following inequality holds:

$$\mathbb{P}(|\bar{X}_n - \mu| \geq t) \leq 2 \exp\left(-\frac{nt^2}{2a^2}\right)$$

where  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  is the sample mean.

*Proof.* Lets define  $Z_i = X_i - \mu, \forall i = 1, 2, \dots, n$ . Then, we have:

---

<sup>1</sup>In the field of statistics, it is common to consider bounded random variables, which naturally leads to the assumption that all moments exist. However, in econometrics, it's often more pragmatic to soften this assumption, focusing instead on the existence of only a select subset of moments. This approach allows for greater flexibility in dealing with real-world data, where the behavior of economic variables can't always be neatly bounded, and full moment conditions may not hold.

$$\bar{Z}_n = \frac{1}{n} \sum_{i=1}^n Z_i = \bar{X}_n - \mu$$

and,

$$|Z_i| \leq a, \forall i = 1, 2, \dots, n$$

Consider the events:

$$\begin{aligned} A &= \{\bar{Z}_n \geq t\} \\ B &= \{\bar{Z}_n \leq -t\} \end{aligned}$$

then, the probability of a event  $C = \{|\bar{Z}_n| \geq t\}$ , can be written as:

$$\mathbb{P}(C) = \mathbb{P}(A \cup B) \leq \mathbb{P}(A) + \mathbb{P}(B)$$

by the union bound. Now, we can write:

$$\begin{aligned} \mathbb{P}(|\bar{Z}_n| \geq t) &\leq \mathbb{P}(\bar{Z}_n \geq t) + \mathbb{P}(\bar{Z}_n \leq -t) \\ \mathbb{P}(|\bar{Z}_n| \geq t) &\leq \mathbb{P}\left(\sum_{i=1}^n Z_i \geq nt\right) + \mathbb{P}\left(\sum_{i=1}^n Z_i \leq -nt\right) \end{aligned}$$

for any  $\lambda > 0$ , we have <sup>2</sup>:

$$\mathbb{P}(|\bar{Z}_n| \geq t) \leq \mathbb{P}\left(\lambda \sum_{i=1}^n Z_i \geq \lambda nt\right) + \mathbb{P}\left(\lambda \sum_{i=1}^n Z_i \leq -\lambda nt\right)$$

and, by Markov's inequality:

$$\mathbb{P}(|\bar{Z}_n| \geq t) \leq \frac{\mathbb{E}[\exp(\lambda \sum_{i=1}^n Z_i)]}{\exp(\lambda nt)} + \frac{\mathbb{E}[\exp(-\lambda \sum_{i=1}^n Z_i)]}{\exp(-\lambda nt)}$$

since  $Z_i$  are independent and identically distributed, we can write:

---

<sup>2</sup>Given that  $f(x) = \lambda x$  is a monotonically increasing function when  $\lambda > 0$

$$\begin{aligned}
\mathbb{P}(|\bar{Z}_n| \geq t) &\leq \frac{\prod_{i=1}^n \mathbb{E}[\exp(\lambda Z_i)]}{\exp(\lambda nt)} + \frac{\prod_{i=1}^n \mathbb{E}[\exp(-\lambda Z_i)]}{\exp(-\lambda nt)} \\
\mathbb{P}(|\bar{Z}_n| \geq t) &\leq \frac{\prod_{i=1}^n \mathbb{E}[\exp(\lambda Z_i)]}{\exp(\lambda nt)} + \frac{\prod_{i=1}^n \exp(-1) \mathbb{E}[\exp(\lambda Z_i)]}{\exp(-\lambda nt)} \\
\mathbb{P}(|\bar{Z}_n| \geq t) &\leq \frac{2 \prod_{i=1}^n \mathbb{E}[\exp(\lambda Z_i)]}{\exp(\lambda nt)}
\end{aligned}$$

Applying the Hoeffding's lemma <sup>3</sup> to the above expression, we have:

$$\mathbb{P}(|\bar{Z}_n| \geq t) \leq \frac{2 \prod_{i=1}^n \exp\left(\frac{\lambda^2 a^2}{2}\right)}{\exp(\lambda nt)}$$

simplifying the above expression:

$$\begin{aligned}
\mathbb{P}(|\bar{Z}_n| \geq t) &\leq \frac{2 \exp\left(\frac{n \lambda^2 a^2}{2}\right)}{\exp(\lambda nt)} \\
\mathbb{P}(|\bar{Z}_n| \geq t) &\leq 2 \exp\left(\frac{n \lambda^2 a^2}{2} - \lambda nt\right)
\end{aligned}$$

Because the above inequality holds for any  $\lambda > 0$ , we can optimize the right-hand side with respect to  $\lambda$ .

$$\lambda^* = \arg \min_{\lambda > 0} \left\{ \frac{n \lambda^2 a^2}{2} - \lambda nt \right\}$$

Calculating the F.O.C. with respect to  $\lambda$ , we get:

$$na^2 \lambda^* - nt = 0 \Rightarrow \lambda^* = \frac{t}{a^2}$$

Substituting  $\lambda^*$  back into the inequality:

$$\begin{aligned}
\mathbb{P}(|\bar{Z}_n| \geq t) &\leq 2 \exp\left(\frac{n \left(\frac{t}{a^2}\right)^2 a^2}{2} - \frac{t}{a^2} nt\right) \\
\mathbb{P}(|\bar{Z}_n| \geq t) &\leq 2 \exp\left(\frac{nt^2}{2a^2} - \frac{nt^2}{a^2}\right) \\
\mathbb{P}(|\bar{Z}_n| \geq t) &\leq 2 \exp\left(-\frac{nt^2}{2a^2}\right)
\end{aligned}$$

---

<sup>3</sup>If  $X$  is a random variable such that  $X \leq a$ , then for any  $\lambda > 0$ , we have:

$$\mathbb{E}[\exp(\lambda X)] \leq \exp\left(\frac{\lambda^2 a^2}{2}\right)$$

replacing  $\bar{Z}_n$  by  $\bar{X}_n - \mu$ :

$$\mathbb{P}(|\bar{X}_n - \mu| \geq t) \leq 2 \exp\left(-\frac{nt^2}{2a^2}\right)$$

□

## 5 Maximal Inequality

Being  $\{X_i\}_{i=1}^n$  a random sample, where  $\dim X_i = p$ , from a distribution with mean  $\mu = [\mu_1, \dots, \mu_p]'$ , such that:

$$|\mu_{i,j} - \mu_j| \leq a, \quad \forall i = 1, \dots, n, \quad \forall j = 1, \dots, p, \quad \forall a > 0$$

then:

$$\mathbb{P}\left(\max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n X_{i,j} - \mu_j \right| \geq a\sqrt{2} \sqrt{\frac{\ln p + \ln \frac{2}{\varepsilon}}{n}}\right) \leq \varepsilon, \quad \forall \varepsilon \in (0, 1)$$

$X_i = [X_{i,1}, \dots, X_{i,p}]'$ , where  $X_{i,j}$  is the  $j$ -th component of the  $i$ -th random vector.

The term  $\left| \frac{1}{n} \sum_{i=1}^n X_{i,j} - \mu_j \right|$  represent the deviation of the sample mean of the  $j$ -th component from the population mean of the  $j$ -th component.

The maximal inequality is a generalization of the weak law of large numbers, and it is used to bound the probability of the maximum deviation of the sample mean from the population mean.

*Proof.* Applying the union bound, we have:

$$\mathbb{P}\left(\max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n X_{i,j} - \mu_j \right| \geq a\sqrt{2} \sqrt{\frac{\ln p + \ln \frac{2}{\varepsilon}}{n}}\right) \leq \sum_{j=1}^p \mathbb{P}\left(\left| \frac{1}{n} \sum_{i=1}^n X_{i,j} - \mu_j \right| \geq a\sqrt{2} \sqrt{\frac{\ln p + \ln \frac{2}{\varepsilon}}{n}}\right)$$

By the Hoeffding's inequality, we have:

$$\sum_{j=1}^p \mathbb{P}\left(\left| \frac{1}{n} \sum_{i=1}^n X_{i,j} - \mu_j \right| \geq a\sqrt{2} \sqrt{\frac{\ln p + \ln \frac{2}{\varepsilon}}{n}}\right) \leq \sum_{j=1}^p 2 \exp\left(-\ln p + \ln \frac{2}{\varepsilon}\right)$$

Simplifying the expression, we have:

$$\begin{aligned} \mathbb{P}\left(\max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n X_{i,j} - \mu_j \right| \geq a\sqrt{2} \sqrt{\frac{\ln p + \ln \frac{2}{\varepsilon}}{n}}\right) &\leq \sum_{j=1}^p 2 \exp(-\ln p) \exp\left(-\ln \frac{2}{\varepsilon}\right) \\ &\leq \sum_{j=1}^p 2 \exp\left(\ln \frac{1}{p}\right) \exp\left(\ln \frac{\varepsilon}{2}\right) \\ &\leq \frac{2p\varepsilon}{2p} \end{aligned}$$

$$\mathbb{P} \left( \max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n X_{i,j} - \mu_j \right| \geq a\sqrt{2} \sqrt{\frac{\ln p + \ln \frac{2}{\varepsilon}}{n}} \right) \leq \varepsilon$$

□

We have proved that the rate of convergence of the maximal inequality is  $\sqrt{\ln p/n}$ , or in big-O notation:

$$\max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n X_{i,j} - \mu_j \right| = O_p \left( \sqrt{\frac{\ln p}{n}} \right)$$

## 6 High-Dimensional Linear Regression

Consider the regression:

$$Y = \mathbf{X}'\beta + \epsilon$$

Where  $Y$  is a scalar,  $\mathbf{X}$  is a  $p$ -dimensional vector of regressors,  $\beta$  is a  $p$ -dimensional vector of coefficients, and  $\epsilon$  is a scalar error term. And suppose we have a random sample:

$$(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n) \sim \text{i.i.d.}(\mathbf{X}, Y)$$

We are interested in case where  $p$  is large, meaning that:  $p \sim n$ ,  $p > n$  or  $p \gg n$ .

**Claim.** OLS linear estimator does not exist when  $p > n$ .

*Proof.* The OLS linear regression estimator is given by:

$$\hat{\beta} = (\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\mathbf{Y}$$

Where  $\mathbb{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)'$  is the  $n \times p$  design matrix, and  $\mathbf{Y} = (Y_1, \dots, Y_n)'$  is the  $n \times 1$  vector of dependent variables.

**Lemma.** If  $\mathbb{A}$  is degenerated, then  $\mathbb{A}^{-1}$  does not exist.

**Definition.** A matrix  $\mathbb{A}$  is degenerated if there exists a non-zero vector  $\mathbf{v}$  such that  $\mathbb{A}\mathbf{v} = \mathbf{0}$ .

*Proof.* Assume that  $\mathbb{A}$  is invertible, then there exist  $\mathbb{A}\mathbb{A}^{-1} = \mathbb{I}_p$ . If  $\mathbb{A}$  is degenerated, then there exists a non-zero vector  $\mathbf{v}$  such that  $\mathbb{A}\mathbf{v} = \mathbf{0}$ . Multiplying both sides by  $\mathbb{A}^{-1}$  we get:

$$\mathbb{A}^{-1}\mathbb{A}\mathbf{v} = \mathbb{A}^{-1}\mathbf{0} \implies \mathbf{v} = \mathbf{0}$$

Which is a contradiction. Therefore,  $\mathbb{A}$  is not invertible. □

**Lemma.** The matrix  $\mathbb{A} = \mathbb{X}'\mathbb{X}$  is degenerated when  $p > n$ .

*Proof.* Consider the linear system of equations:

$$\mathbb{X}\mathbf{b} = \mathbf{0}, \quad \mathbf{b} \in \mathbb{R}^p$$

Where  $\mathbb{X}$  is a  $n \times p$  matrix. This system has a non unique solution when  $p > n$ . therefore, there exists a non-zero vector  $\mathbf{b}$  such that  $\mathbb{X}\mathbf{b} = \mathbf{0}$ , meaning that  $\mathbb{X}'\mathbb{X}$  is degenerated.  $\square$

From the previous lemma, we know that  $\mathbb{X}'\mathbb{X}$  is degenerated when  $p > n$ . Therefore,  $\mathbb{X}'\mathbb{X}$  is not invertible, and the OLS linear estimator does not exist.  $\square$

## 6.1 Expected Value of the OLS Linear Regression Estimator when $p > n$

Assuming that the linear model is homoscedastic, that is:

$$\mathbb{E} [\epsilon^2 | \mathbf{X}] = \sigma^2$$

then:

$$\mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \left( \mathbf{X}'_i (\hat{\beta}^{OLS} - \beta) \right)^2 \right] = \frac{p\sigma^2}{n}$$

*Proof.* Given the linear regression in matrix form:

$$\mathbf{Y} = \mathbb{X}\beta + \varepsilon$$

Where  $\mathbf{Y} = (Y_1, \dots, Y_n)'$  and  $\varepsilon = (\epsilon_1, \dots, \epsilon_n)'$  are  $n \times 1$  column vectors;  $\beta = (\beta_1, \dots, \beta_p)'$  is a  $p \times 1$  column vector; and  $\mathbb{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)'$  is a  $n \times p$  matrix.

Recall that the OLS linear estimator is given by:

$$\hat{\beta}^{OLS} = (\mathbb{X}'\mathbb{X})^{-1} \mathbb{X}'\mathbf{Y}$$

Substituting the linear regression in matrix form into the OLS linear estimator we get:

$$\begin{aligned} \hat{\beta}^{OLS} &= (\mathbb{X}'\mathbb{X})^{-1} \mathbb{X}'(\mathbb{X}\beta + \varepsilon) \\ \hat{\beta}^{OLS} - \beta &= (\mathbb{X}'\mathbb{X})^{-1} \mathbb{X}'\varepsilon \\ \mathbb{X}(\hat{\beta}^{OLS} - \beta) &= \mathbb{X}(\mathbb{X}'\mathbb{X})^{-1} \mathbb{X}'\varepsilon \end{aligned}$$

Where:

$$\mathbf{e} = \hat{\mathbf{Y}} - \mathbb{E}[\mathbf{Y} | \mathbf{X}] = \mathbb{X}(\hat{\beta}^{OLS} - \beta)$$

Therefore, the expected value of the mean squared error for the OLS estimator is:



$$\begin{aligned}
\mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \left( \mathbf{X}'_i (\hat{\beta}^{OLS} - \beta) \right)^2 \right] &= \frac{1}{n} \mathbb{E} [\mathbf{e}' \mathbf{e}] \\
&= \frac{1}{n} \mathbb{E} [(\mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\varepsilon)'(\mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\varepsilon)] \\
&= \frac{1}{n} \mathbb{E} [\varepsilon'\mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\varepsilon]
\end{aligned}$$

$$\mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \left( \mathbf{X}'_i (\hat{\beta}^{OLS} - \beta) \right)^2 \right] = \frac{1}{n} \mathbb{E} [\varepsilon'\mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\varepsilon] \quad (1)$$

**Lemma.** *If  $\mathbb{A}$  is a  $n \times p$  matrix, and  $\mathbb{B}$  is a  $p \times n$  matrix, then:*

$$tr(\mathbb{A}\mathbb{B}) = tr(\mathbb{B}\mathbb{A})$$

The expected value in equation 1 is an scalar, therefore:

$$\frac{1}{n} \mathbb{E} [\varepsilon'\mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\varepsilon] = \frac{1}{n} tr \left( \mathbb{E} [\varepsilon'\mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\varepsilon] \right)$$

Because the trace and the expected value are linear operators, we interchange the order of the trace and the expected value:

$$\begin{aligned}
\frac{1}{n} \mathbb{E} [\varepsilon'\mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\varepsilon] &= \frac{1}{n} \mathbb{E} [tr (\varepsilon'\mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\varepsilon)] \\
&= \frac{1}{n} \mathbb{E} [tr (\varepsilon\varepsilon'\mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}')] \\
&= \frac{1}{n} tr \left( \mathbb{E} [\varepsilon\varepsilon'\mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'] \right)
\end{aligned}$$

By the law of iterated expectations:

$$\begin{aligned}
\frac{1}{n} \mathbb{E} [\varepsilon' \mathbb{X} (\mathbb{X}' \mathbb{X})^{-1} \mathbb{X}' \varepsilon] &= \frac{1}{n} \text{tr} \left( \mathbb{E} \left[ \mathbb{E} [\varepsilon \varepsilon' \mathbb{X} (\mathbb{X}' \mathbb{X})^{-1} \mathbb{X}' | \mathbb{X}] \right] \right) \\
&= \frac{1}{n} \text{tr} \left( \mathbb{E} \left[ \mathbb{E} [\varepsilon \varepsilon' | \mathbb{X}] \mathbb{X} (\mathbb{X}' \mathbb{X})^{-1} \mathbb{X}' \right] \right) \\
&= \frac{1}{n} \text{tr} \left( \mathbb{E} [\sigma^2 \mathbb{I}_n \mathbb{X} (\mathbb{X}' \mathbb{X})^{-1} \mathbb{X}'] \right) \\
&= \frac{\sigma^2}{n} \text{tr} \left( \mathbb{E} [\mathbb{X} (\mathbb{X}' \mathbb{X})^{-1} \mathbb{X}'] \right) \\
&= \frac{\sigma^2}{n} \mathbb{E} [\text{tr} (\mathbb{X} (\mathbb{X}' \mathbb{X})^{-1} \mathbb{X}')] \\
&= \frac{\sigma^2}{n} \mathbb{E} [\text{tr} (\mathbb{X}' \mathbb{X} (\mathbb{X}' \mathbb{X})^{-1})] \\
&= \frac{\sigma^2}{n} \mathbb{E} [\text{tr} (\mathbb{I}_p)] \\
&= \frac{\sigma^2}{n} \mathbb{E} [p]
\end{aligned}$$

$$\frac{1}{n} \mathbb{E} [\varepsilon' \mathbb{X} (\mathbb{X}' \mathbb{X})^{-1} \mathbb{X}' \varepsilon] = \frac{p \sigma^2}{n}$$

□

Meaning that the expected value of the mean squared error of the OLS linear estimator does not converge to zero when  $p \gg n$ .

## 7 Sparsity

Given the linear regression model:

$$Y = \mathbf{X}'\beta + \epsilon, \quad \mathbb{E}[\epsilon | \mathbf{X}] = 0$$

where  $Y$  is the response variable,  $\mathbf{X}$  is the  $p$ -dimensional feature vector,  $\beta$  is the  $p$ -dimensional coefficient vector, and  $\epsilon$  is the error term. We assume that  $\epsilon$  is independent of  $\mathbf{X}$ .

**Definition** (Sparsity Index).

$$\begin{aligned}
S &= \sum_{j=1}^p \mathbf{1}_{\{\beta_j \neq 0\}} \\
S &= \|\beta\|_0
\end{aligned}$$

*It is often called the  $l_0$  norm of  $\beta$ , even though does not satisfy all the properties of a norm<sup>4</sup>.*

---

<sup>4</sup>See Norms in the appendix for more details.

We are going to assume that  $S$  is small, i.e.,  $S \ll p$ . This assumption is often referred to as the *sparsity assumption*.

A more relaxed version of the sparsity assumption is the *approximate sparsity assumption*, where most  $\beta_j$  are close to zero, but not exactly zero.

## 7.1 Best Subset Selection

$$\hat{\beta} = \arg \min_{b \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{X}_i' b)^2 + \lambda \|b\|_0$$

where  $\lambda$  is the penalty parameter, and  $\lambda \|b\|_0$  is the penalty term that encourages sparsity in the solution.

In order to solve this optimization problem, we need to consider all possible subset of features  $\binom{n}{k}$  for  $k = 0, 1, \dots, p$ . This is computationally infeasible for large  $p$ , and the  $\ell_0$  penalty is non-convex, which makes the optimization problem even harder.

## 7.2 Least Absolute Shrinkage and Selection Operator (LASSO)

Introduced by Tibshirani (1996), the LASSO estimator is defined as:

$$\hat{\beta} = \arg \min_{b \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{X}_i' b)^2 + \lambda \|b\|_1$$

### 7.2.1 Main LASSO Result

If  $\lambda$  is chosen appropriately, and some regularization conditions are satisfied (including sparsity), then the LASSO estimator  $\hat{\beta}^{LASSO}$  satisfies:

$$\|\hat{\beta}^{LASSO} - \beta\|_2 \leq C \sqrt{\frac{S \ln P}{n}}$$

and,

$$\|\hat{\beta}^{LASSO} - \beta\|_1 \leq CS \sqrt{\frac{\ln P}{n}}$$

with probability approaching 1, where  $C$  is some constant. Intuitively, the LASSO estimator is consistent in terms of the  $\ell_2$  norm, if  $\frac{S \ln p}{n} \rightarrow 0$ .

## A Big O Notation

The *Big O* notation in statistics deals with the convergence of sets of random variables, where convergence is in the sense of convergence in probability. The notation is used to describe the rate of convergence of a sequence of random variables to a limit.

For a set of random variables  $X_n$ , and a corresponding set of constants  $a_n$  (both indexed by  $n$ ), the notation:

$$X_n = O_p(a_n), \quad \text{as } n \rightarrow \infty$$

means that the set of values  $\frac{X_n}{a_n}$  is stochastically bounded, That means:

$$\forall \varepsilon > 0, \exists M \in (0, \infty), \exists N_\varepsilon \in (0, \infty) : \forall n > N_\varepsilon \left( n > N_\varepsilon \implies \mathbb{P} \left( \left| \frac{X_n}{a_n} \right| > M \right) < \varepsilon \right)$$

equivalently, we can rewrite the above expression as:

$$\forall \varepsilon > 0, \exists \delta_\varepsilon \in (0, \infty), \exists N_\varepsilon \in (0, \infty) : \forall n > N_\varepsilon (n > N_\varepsilon \implies \mathbb{P}(|X_n| > \delta_\varepsilon) < \varepsilon)$$

## B Norms

**Definition.** A norm on a metric vector space  $V$  over a field  $F$  (usually  $\mathbb{R}$  or  $\mathbb{C}$ ) is a function  $\|\cdot\|: V \rightarrow [0, \infty)$  that satisfies the following properties for all vectors  $u, v \in V$  and all scalars  $\alpha \in F$ :

1.  $\|v\| = 0$  if and only if  $v = 0$ .
2.  $\|\alpha v\| = |\alpha| \|v\|$  (Absolute scalability).
3.  $\|u + v\| \leq \|u\| + \|v\|$  (Triangle inequality).

Some common norms include:

- **Zero norm** (or  $l_0$  norm):

$$\|v\|_0 = \sum_{i=1}^n \mathbf{1}_{\{v_i \neq 0\}}$$

- **Manhattan norm** (or  $l_1$  norm):

$$\|v\|_1 = \sum_{i=1}^n |v_i|$$

- **Euclidean norm** (or  $l_2$  norm):

$$\|v\|_2 = \sqrt{\sum_{i=1}^n |v_i|^2}$$

- **General  $l_p$  norm:**

$$\|v\|_p = \left( \sum_{i=1}^n |v_i|^p \right)^{1/p}$$

- **Maximum norm** (or  $l_\infty$  norm):

$$\|v\|_\infty = \max_{i=1}^n |v_i|$$