
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 μ and variance σ^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 μ and variance σ^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 μ such that, for a number ¹ $a > 0$, 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:

¹In 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 ²:

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

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:

²Given 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 ³ 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 λ .

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

Calculating the F.O.C. with respect to λ , we get:

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

Substituting λ^* back into the inequality:

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

³If 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 Regression

Consider the regression:

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

Where Y is a scalar, \mathbf{X} is a p -dimensional vector of regressors, β is a p -dimensional vector of coefficients, and ϵ 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

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:

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