

Métodos de reducción de varianza

1. Eficiencia de un estimador

La eficiencia de un estimador $\hat{\theta}(\mathbf{X})$ de un parámetro θ de un vector aleatorio \mathbf{X} es

$$\text{Eff}(\hat{\theta}(\mathbf{X})) \stackrel{\text{def}}{=} \frac{1}{\text{MSE}(\hat{\theta}(\mathbf{X})) \text{Coste}(\hat{\theta}(\mathbf{X}))}$$

donde

$$\begin{aligned} \text{MSE}(\hat{\theta}(\mathbf{X})) &\stackrel{\text{def}}{=} \mathbb{E}[(\hat{\theta}(\mathbf{X}) - \theta)^2] = \text{Var}(\hat{\theta}(\mathbf{X})) + \text{Sesgo}(\hat{\theta}(\mathbf{X}), \theta)^2 \\ \text{Coste}(\hat{\theta}(\mathbf{X})) &\stackrel{\text{def}}{=} \text{tiempo esperado de cálculo de } \hat{\theta}(\mathbf{X}) \end{aligned}$$

Dados dos estimadores insesgados, $\hat{\theta}_1(\mathbf{X})$ y $\hat{\theta}_2(\mathbf{X})$, que requieren el mismo tiempo de cálculo, se tiene que

$$\text{Eff}(\hat{\theta}_1(\mathbf{X})) > \text{Eff}(\hat{\theta}_2(\mathbf{X})) \iff \text{Var}(\hat{\theta}_1(\mathbf{X})) < \text{Var}(\hat{\theta}_2(\mathbf{X}))$$

Para $\hat{\mu}_n(\mathbf{X})$ el estimador de Montecarlo de $\mu = \mathbb{E}[g(\mathbf{X})]$ se tiene que

$$\text{Var}(\hat{\mu}_n(\mathbf{X})) = \frac{\sigma^2}{n} \quad \text{Sesgo}(\hat{\mu}_n(\mathbf{X}), \mu) = 0 \quad \text{Coste}(\hat{\mu}_n(\mathbf{X})) = cn$$

Por tanto, $\text{Eff}(\hat{\mu}_n(\mathbf{X})) = \frac{1}{c\sigma^2}$. Es decir, al aumentar el tamaño de la muestra, la reducción en varianza se compensa exactamente con el incremento del tiempo de computación, por lo que no hay ganancia en eficiencia.

Para encontrar un estimador más eficiente se necesita, por tanto, reducir la varianza en un factor superior al incremento en tiempo de computación.

No hay una regla fija sobre cuán grande debe ser la mejora de la eficiencia proporcionada por un método de reducción de varianza para que valga la pena usarlo. En algunos escenarios, como la representación de gráficos por ordenador para películas animadas, donde miles de CPUs se mantienen ocupadas durante meses, una mejora del 10 % aporta un ahorro significativo. En otros escenarios, como un cálculo único, una ganancia de factor 60 que convierte un minuto en un segundo de espera puede que no justifique el coste de la programación de un método más complicado.

La ganancia de eficiencia necesaria para justificar el uso de un método es menor si el esfuerzo de programación puede amortizarse en muchas aplicaciones. El umbral es alto para un programa de un solo uso, más bajo por algo que añadimos a nuestra biblioteca personal, más bajo aún para código que se comparte y aún más bajo para código que se incluye en una biblioteca o herramienta de simulación de uso general.

2. Método del muestreo estratificado

Sean \mathbf{X} un vector aleatorio (con d componentes) que sigue una distribución de probabilidad dada por la función de densidad f y $g: \mathbb{R}^d \rightarrow \mathbb{R}$ una función medible. Entonces,

$$\mu = \mathbb{E}[g(\mathbf{X})] = \int_D g(\mathbf{x})f(\mathbf{x}) d\mathbf{x}$$

donde $D = \{\mathbf{x} \in \mathbb{R}^d \mid f(\mathbf{x}) > 0\}$.

Consideremos una partición de D en los estratos D_1, \dots, D_J tal que para cada $j = 1, \dots, J$ se conoce $p_j = \mathbb{P}(\mathbf{X} \in D_j)$ y se sabe cómo generar valores aleatorios de la variable condicionada $\mathbf{X} \mid \mathbf{X} \in D_j$ (que sigue una distribución de probabilidad dada por la función de densidad $f_j(\mathbf{x}) = 1/p_j f(\mathbf{x}) \mathbb{1}_{D_j}(\mathbf{x})$).

Para estimar el valor de μ mediante el método del muestreo estratificado, basta entonces calcular

$$\hat{\mu}_{\text{str},n} = \sum_{j=1}^J p_j \left(\frac{1}{n_j} \sum_{i=1}^{n_j} g(\mathbf{x}_{ij}) \right)$$

donde $\mathbf{x}_{1j}, \dots, \mathbf{x}_{n_jj} \stackrel{\text{i. i. d.}}{\sim} \mathbf{X} \mid \mathbf{X} \in D_j$ y $n_1 + \dots + n_J = n$.

Se tiene que, en efecto, $\hat{\mu}_{\text{str},n}$ es un estimador insesgado de μ :

$$\begin{aligned} \mathbb{E}[\hat{\mu}_{\text{str},n}] &= \sum_{j=1}^J p_j \mathbb{E} \left[\frac{1}{n_j} \sum_{i=1}^{n_j} g(\mathbf{x}_{ij}) \right] \\ &= \sum_{j=1}^J p_j \mathbb{E}[g(\mathbf{X}) \mid \mathbf{X} \in D_j] \\ &= \sum_{j=1}^J p_j \int_{\mathbb{R}^d} g(\mathbf{x}) f_j(\mathbf{x}) d\mathbf{x} \\ &= \sum_{j=1}^J \int_{D_j} g(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} \\ &= \int_D g(\mathbf{x}) f(\mathbf{x}) d\mathbf{x} \\ &= \mathbb{E}[g(\mathbf{x})] \end{aligned}$$

Denotemos $\mu_j = \mathbb{E}[g(\mathbf{X}) \mid \mathbf{X} \in D_j]$ y $\sigma_j^2 = \text{Var}(g(\mathbf{X}) \mid \mathbf{X} \in D_j)$. Entonces la varianza del estimador viene dada por

$$\text{Var}(\hat{\mu}_{\text{str},n}) = \sum_{j=1}^J p_j^2 \frac{\sigma_j^2}{n_j}$$

Esto proporciona una guía para elegir los estratos, ya que cuanto menor sean las σ_j , es decir, cuanto más constante sea la función g en cada estrato, menor será el error estándar $\sqrt{\text{Var}(\hat{\mu}_{\text{str},n})}$ de la estimación de μ por el método del muestreo estratificado.

Para construir un intervalo de confianza, podemos estimar el error estándar como sigue:

$$\begin{aligned}\hat{\mu}_j &= \frac{1}{n_j} \sum_{i=1}^{n_j} g(\mathbf{x}_{ij}) \\ \hat{\sigma}_j^2 &= \frac{1}{n_j - 1} \sum_{i=1}^{n_j} (g(\mathbf{x}_{ij}) - \hat{\mu}_j)^2 \\ \widehat{\text{Var}}(\hat{\mu}_{\text{str},n}) &= \sum_{j=1}^J p_j^2 \frac{\hat{\sigma}_j^2}{n_j}\end{aligned}$$

Obsérvese que para que $\hat{\mu}_{\text{str},n}$ esté bien definido es necesario que $n_j \geq 1$, para cada $j = 1, \dots, J$, y para que $\widehat{\text{Var}}(\hat{\mu}_{\text{str},n})$ esté bien definido es necesario, además, que $n_j \geq 2$, para cada $j = 1, \dots, J$.

El estimador $\hat{\mu}_{\text{str},n}$ depende del tamaño de cada uno de los estratos. Una manera simple de elegir esos tamaños es mediante una *asignación proporcional*, $n_j = p_j n$ (si, además, $p_j = 1/J$ para todo $j = 1, \dots, J$, la asignación se dice que es *sistemática*). Entonces,

$$\hat{\mu}_{\text{str,prop},n} = \sum_{j=1}^J p_j \left(\frac{1}{n_j} \sum_{i=1}^{n_j} g(\mathbf{x}_{ij}) \right) = \frac{1}{n} \sum_{j=1}^J \sum_{i=1}^{n_j} g(\mathbf{x}_{ij})$$

y

$$\text{Var}(\hat{\mu}_{\text{str,prop},n}) = \sum_{j=1}^J p_j^2 \frac{\sigma_j^2}{n_j} = \frac{1}{n} \sum_{j=1}^J p_j \sigma_j^2$$

Por otra parte, por la ley de la varianza total se tiene que

$$\begin{aligned}\text{Var}(g(\mathbf{X})) &= \mathbb{E} \left[\text{Var}(g(\mathbf{X}) \mid \mathbf{X} \in D_j) \right] + \text{Var} \left(\mathbb{E}[g(\mathbf{X}) \mid \mathbf{X} \in D_j] \right) \\ &= \sum_{j=1}^J p_j \sigma_j^2 + \sum_{j=1}^J p_j (\mu_j - \mu)^2\end{aligned}$$

Es decir, si denotamos $\sigma_W^2 = \sum_{j=1}^J p_j \sigma_j^2$ a la varianza en los estratos y $\sigma_B^2 = \sum_{j=1}^J p_j (\mu_j - \mu)^2$ a la varianza entre los estratos, entonces

$$\text{Var}(\hat{\mu}_n) = \frac{1}{n} (\sigma_W^2 + \sigma_B^2) \quad \text{y} \quad \text{Var}(\hat{\mu}_{\text{str,prop},n}) = \frac{1}{n} \sigma_W^2$$

Por lo tanto, siempre se tiene que $\text{Var}(\hat{\mu}_{\text{str,prop},n}) \leq \text{Var}(\hat{\mu}_n)$. Además, cuanto mayor sea σ_B^2 con respecto a σ_W^2 , mayor será la reducción de varianza conseguida mediante el método del muestreo estratificado con asignación proporcional.

La asignación proporcional no es necesariamente la mejor. Por ejemplo, dados dos estratos del mismo tamaño, pero con diferente varianza condicional, es más beneficioso generar menos valores en el estrato de menor varianza.

Asumiendo que el coste de generar cada valor es el mismo en cada estrato, la asignación óptima viene dada por

$$n_j = \frac{n p_j \sigma_j}{\sum_{k=1}^J p_k \sigma_k}, \quad j = 1, \dots, J$$

siendo $\text{Var}(\hat{\sigma}_{\text{str}, \text{opt}, n}) = \frac{1}{n} \left(\sum_{j=1}^J p_j \sigma_j \right)^2$.

La aplicación práctica de esta asignación óptima se ve dificultada por el hecho de que los σ_j no se conocen usualmente. Se puede usar un procedimiento en dos etapas:

1. En una primera etapa de tanteo se usa una asignación proporcional para generar valores $\mathbf{x}_{1j}, \dots, \mathbf{x}_{n_j j} \stackrel{\text{i. i. d.}}{\sim} \mathbf{X} \mid \mathbf{X} \in D_j$ que se usan para estimar σ_j^2 mediante la cuasivarianza muestral:

$$\hat{\mu}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} g(\mathbf{x}_{ij})$$

$$\hat{\sigma}_j^2 = \frac{1}{n_j - 1} \sum_{i=1}^{n_j} (g(\mathbf{x}_{ij}) - \hat{\mu}_j)^2$$

2. En la etapa de producción se usan las estimaciones $\hat{\sigma}_j^2$ obtenidas para calcular el tamaño de muestra óptimo de cada estrato y realizar la estimación $\hat{\mu}_{\text{str}, \text{opt}, n}$.

La cantidad total de valores generados en la fase de tanteo debería ser pequeña, para que el coste no sea excesivo, pero lo suficientemente grande para obtener buenas estimaciones de σ_j . En caso contrario, podría ocurrir que la varianza de la estimación obtenida superara al de la estimación con asignación proporcional, e incluso al de la estimación por el método de Montecarlo directo.

3. Método del muestreo por importancia

Sean \mathbf{X} un vector aleatorio con función de densidad *nominal* f_1 y $g: \mathbb{R}^d \rightarrow \mathbb{R}$ una función medible. Entonces

$$\mu = \mathbb{E}_{f_1}[g(\mathbf{X})] = \int_{\mathbb{R}^d} g(\mathbf{x}) f_1(\mathbf{x}) d\mathbf{x}$$

El método de muestreo por importancia consiste en considerar una función de densidad *instrumental* f_2 para \mathbf{X} , de tal manera que f_1 sea absolutamente continua con respecto a f_2 (es decir, $f_2(\mathbf{x}) \neq 0$ para todo \mathbf{x} tal que $f_1(\mathbf{x}) \neq 0$). Si denotamos el soporte de una función f como $\text{sop}_f = \{\mathbf{x} \in \mathbb{R}^d \mid f(\mathbf{x}) \neq 0\}$ y definimos $L(\mathbf{x}) = \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})}$ como la *razón de verosimilitud*, se tiene que

$$\begin{aligned} \mathbb{E}_{f_2}[g(\mathbf{x})L(\mathbf{x})] &= \int_{\mathbb{R}^d} g(\mathbf{x})L(\mathbf{x})f_2(\mathbf{x}) d\mathbf{x} \\ &= \int_{\text{sop}_{f_2}} g(\mathbf{x})L(\mathbf{x})f_2(\mathbf{x}) d\mathbf{x} \\ &= \int_{\text{sop}_{f_2}} g(\mathbf{x})f_1(\mathbf{x}) d\mathbf{x} \end{aligned}$$

$$\begin{aligned}
&= \int_{\text{sop}_{f_1}} g(\mathbf{x})f_1(\mathbf{x}) d\mathbf{x} + \int_{\text{sop}_{f_2} \setminus \text{sop}_{f_1}} g(\mathbf{x})f_1(\mathbf{x}) d\mathbf{x} \\
&= \int_{\text{sop}_{f_1}} g(\mathbf{x})f_1(\mathbf{x}) d\mathbf{x} \\
&= \int_{\mathbb{R}^d} g(\mathbf{x})f_1(\mathbf{x}) d\mathbf{x} \\
&= \mathbb{E}_{f_1}[g(\mathbf{x})] \\
&= \mu
\end{aligned}$$

Se tiene entonces que el estimador

$$\hat{\mu}_{\text{is},f_2,n}(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n g(\mathbf{x}_i)L(\mathbf{x}_i)$$

donde $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{\text{i. i. d.}}{\sim} f_2$, es un estimador insesgado de μ .

Además,

$$\begin{aligned}
\text{Var}_{f_2}(\hat{\mu}_{\text{is},f_2,n}(\mathbf{X})) &= \frac{1}{n} \text{Var}_{f_2}(g(\mathbf{X})L(\mathbf{X})) \\
&= \frac{1}{n} (\mathbb{E}_{f_2}[g(\mathbf{X})^2 L(\mathbf{X})^2] - \mathbb{E}_{f_2}[g(\mathbf{X})L(\mathbf{X})]^2) \\
&= \frac{1}{n} \left(\int_{\mathbb{R}^d} g(\mathbf{x})^2 L(\mathbf{x})^2 f_2(\mathbf{x}) d\mathbf{x} - \mu^2 \right) \\
&= \frac{1}{n} \left(\int_{\mathbb{R}^d} g(\mathbf{x})^2 L(\mathbf{x}) f_1(\mathbf{x}) d\mathbf{x} - \mu^2 \right) \\
&= \frac{1}{n} (\mathbb{E}_{f_1}[g(\mathbf{X})^2 L(\mathbf{X})] - \mu^2)
\end{aligned}$$

Como $\text{Var}_{f_1}(\hat{\mu}_n(\mathbf{X})) = \frac{1}{n} (\mathbb{E}_{f_1}[g(\mathbf{X})^2] - \mu^2)$, se tiene que $\text{Var}_{f_2}(\hat{\mu}_{\text{is},f_2,n}(\mathbf{X})) < \text{Var}_{f_1}(\hat{\mu}_n(\mathbf{X}))$ si y solo si $\mathbb{E}_{f_1}[g(\mathbf{X})^2 L(\mathbf{X})] < \mathbb{E}_{f_1}[g(\mathbf{X})^2]$.

Es posible determinar la densidad instrumental óptima: si $0 < \int_{\mathbb{R}^d} |g(\mathbf{x})| f_1(\mathbf{x}) d\mathbf{x} < +\infty$, entonces la elección de f_2 que minimiza el valor de $\text{Var}_{f_2}(\hat{\mu}_{\text{is},f_2,n}(\mathbf{X}))$ es

$$f_2^*(\mathbf{x}) = \frac{|g(\mathbf{x})| f_1(\mathbf{x})}{\int_{\mathbb{R}^d} |g(\mathbf{x})| f_1(\mathbf{x}) d\mathbf{x}}$$

y la mínima varianza es

$$\left(\int_{\mathbb{R}^d} |g(\mathbf{x})| f_1(\mathbf{x}) d\mathbf{x} \right)^2 - \left(\int_{\mathbb{R}^d} g(\mathbf{x}) f_1(\mathbf{x}) d\mathbf{x} \right)^2$$

Este resultado de optimalidad es, en realidad, un resultado formal: por ejemplo, cuando g es siempre positiva o negativa se alcanzaría varianza 0, pero esta elección óptima requeriría conocer $\int_{\mathbb{R}^d} g(\mathbf{x}) f_1(\mathbf{x}) d\mathbf{x}$, que es la integral que se quiere calcular.

No obstante, desde un punto de vista práctico el resultado sugiere probar con densidades instrumentales que sean lo más similares a $|g(\mathbf{x})| f_1(\mathbf{x})$ que sea posible. Es decir, las

tuplas \mathbf{x} que deberían generarse con mayor probabilidad (las más «importantes») deberían ser aquellas con valores altos de $|g(\mathbf{x})|f_1(\mathbf{x})$. También hay que tener en cuenta que la varianza del método es menor cuanto menor sea $\mathbb{E}_{f_1}[g(\mathbf{X})^2 L(\mathbf{X})]$. Esto quiere decir que, puesto que el denominador de L es f_2 , una densidad instrumental mal elegida puede dar lugar a una estimación de gran varianza, incluso infinita. Para evitarlo, los valores cercanos a cero de esta densidad instrumental deberían cancelarse con valores cercanos a cero de la función g o, en su defecto, la densidad instrumental f_2 debería acercarse a cero más lentamente que la densidad nominal f_1 .