

## Bias und Biaskorrektur

Bei einem erwartungstreuen Parameterschätzer  $\hat{\theta}$  ist der Erwartungswert der Schätzfunktion  $E(\hat{\theta})$  gleich dem wirklichen Parameter  $\theta$ . Somit ist die Differenz  $E(\hat{\theta}) - \theta = 0$ .

### BIAS

Ist ein Schätzer nicht erwartungstreu, so ist entweder  $E(\hat{\theta}) > \theta$  oder aber  $E(\hat{\theta}) < \theta$ .

Ist  $E(\hat{\theta}) > \theta$ , also  $E(\hat{\theta}) - \theta > 0$ , dann wird der Parameter tendenziell *überschätzt*, man spricht von einem *positiven Bias*. (Vorsicht: Einzelne Schätzwerte können auch in diesem Fall zu klein geraten, im Durchschnitt kommt es aber zur Überschätzung.)

Analog bedeutet ein *negativer Bias*, dass der Parameter systematisch *unterschätzt* wird ( $E(\hat{\theta}) < \theta \Leftrightarrow E(\hat{\theta}) - \theta < 0$ ).

### Ein Beispiel zur Veranschaulichung

$\eta = 100$  sei ein bestimmter Parameter der Gesamtpopulation,  $\hat{\lambda}$ ,  $\hat{\epsilon}$  und  $\hat{\phi}$  seien drei verschiedene Schätzfunktionen für  $\eta$ .

Ich ziehe nun 10 unabhängige Stichproben und ermittle für jede von ihnen mit allen drei Methoden einen Schätzwert für  $\eta$ :

$i$	$\hat{\lambda}_i$	$\hat{\epsilon}_i$	$\hat{\phi}_i$
1	104	81	134
2	97	102	142
3	93	74	120
4	99	88	151
5	108	78	138
6	101	112	144
7	102	93	98
8	98	72	123
9	100	79	139
10	104	94	128
$\sum_i$	1006	873	1317

Es ergeben sich folgende durchschnittliche Schätzwerte:

$$\begin{aligned}\hat{\lambda} &= 100.6 \\ \hat{\varepsilon} &= 87.3 \\ \hat{\phi} &= 131.7\end{aligned}$$

Die Schätzfunktion  $\hat{\lambda}$  schätzt den Parameter  $\eta$  offenbar am besten. Die Abweichungen sind vermutlich nicht überzufällig, also nicht groß genug, um statistisch aussagekräftig („signifikant“) zu sein.

Im Gegensatz dazu wirft  $\hat{\varepsilon}$  tendenziell Werte aus, die  $\eta$  eindeutig (signifikant) unterschätzen ( $\rightarrow$  negativer Bias), während  $\hat{\phi}$  vermutlich einen positiven Bias aufweist.

### ***Bias der Stichprobenvarianz:***

Ein Beispiel für einen nicht erwartungstreuen Schätzer ist die Stichprobenvarianz  $s_X^2$ . Dies kann man zeigen, indem man ihren Erwartungswert umformt und mit der Populationsvarianz  $\sigma_X^2$  vergleicht.

Wir ziehen immer wieder unabhängige Stichproben und berechnen jeweils  $\bar{x}$  und  $s_X^2$ , woraus sich die Zufallsvariablen  $\bar{X}$  und  $S_X^2$  ergeben.

Wir wissen:

$$s_X^2 = \frac{1}{n} \cdot \sum_i x_i^2 - \bar{x}^2$$

Für die entsprechenden Zufallsvariablen gilt demnach:

$$S_X^2 = \frac{1}{n} \cdot \sum_i X_i^2 - \bar{X}^2$$

Nun berechnen wir auf beiden Seiten den Erwartungswert:

$$E(S_X^2) = E\left(\frac{1}{n} \cdot \sum_i X_i^2 - \bar{X}^2\right)$$

$$E(S_X^2) = \frac{1}{n} \cdot \sum_i E(X_i^2) - E(\bar{X}^2)$$

*Einschub 1:*

Da alle  $X_i$  dieselbe Verteilung wie die Variable  $X$  haben, gilt:

$$\begin{aligned}\frac{1}{n} \cdot \sum_i E(X_i^2) &= \frac{1}{n} \cdot \sum_i E(X^2) = \\ &= \frac{n}{n} \cdot E(X^2) = \\ &= E(X^2)\end{aligned}$$

*Ende Einschub 1*

*Einschub 2:*

$$\begin{aligned}E(\bar{X}^2) &= E\left[\left(\frac{1}{n} \cdot \sum_i X_i\right)^2\right] = \\ &= E\left[\left(\frac{1}{n}\right)^2 \cdot \left(\sum_i X_i\right)^2\right] = \\ &= \frac{1}{n^2} \cdot E\left[\left(\sum_i X_i\right)^2\right] = \\ &= \frac{1}{n^2} \cdot E\left[\left(\sum_j X_j\right) \cdot \left(\sum_l X_l\right)\right] = \\ &= \frac{1}{n^2} \cdot E\left[\sum_j \sum_l (X_j \cdot X_l)\right] =\end{aligned}$$

Wir ziehen nun aus der Doppelsumme jene Produkte heraus, bei denen der Index gleich ist (*der Einfachheit halber schreiben wir für die Doppelsumme nur noch ein einzelnes Summenzeichen*):

$$\begin{aligned}&= \frac{1}{n^2} \cdot E\left[\sum_i X_i^2 + \sum_{j \neq l} (X_j \cdot X_l)\right] = \\ &= \frac{1}{n^2} \cdot E\left(\sum_i X_i^2\right) + \frac{1}{n^2} \cdot E\left[\sum_{j \neq l} (X_j \cdot X_l)\right] = \\ &= \frac{1}{n^2} \cdot \sum_i E(X_i^2) + \frac{1}{n^2} \cdot \left[\sum_{j \neq l} E(X_j \cdot X_l)\right] =\end{aligned}$$

Da je zwei Variablen  $X_j$  und  $X_l$  unabhängig sind, gilt für die Erwartungswerte  $E(X_j \cdot X_l) = E(X_j) \cdot E(X_l)$ :

$$= \frac{1}{n^2} \cdot [n \cdot E(X^2)] + \frac{1}{n^2} \cdot \sum_{j \neq l} [E(X_j) \cdot E(X_l)] =$$

Alle  $X_j$  und  $X_l$  sind wie  $X$  verteilt, somit gilt  $E(X_j) = E(X_l) = E(X)$ :

$$= \frac{1}{n} \cdot E(X^2) + \frac{1}{n^2} \cdot \sum_{j \neq l} [E(X)]^2 =$$

Insgesamt gäbe es  $n \cdot n$  Kombinationen von  $j$  und  $l$ , wir haben aber alle  $n$  gleichen Index-Paare herausgezogen, somit bleiben  $n \cdot n - n = n \cdot (n - 1)$ .

$$\begin{aligned} &= \frac{1}{n} \cdot E(X^2) + \frac{1}{n^2} \cdot n \cdot (n - 1) \cdot [E(X)]^2 \\ E(\bar{X}^2) &= \frac{1}{n} \cdot E(X^2) + \frac{n - 1}{n} \cdot [E(X)]^2 \end{aligned}$$

*Ende Einschub 2*

Zurück zur Gleichung:

$$E(S_X^2) = \frac{1}{n} \cdot \sum_i E(X_i^2) - E(\bar{X}^2)$$

Wir setzen unsere Ergebnisse aus *Einschub 1* und *Einschub 2* ein:

$$\begin{aligned} E(S_X^2) &= E(X^2) - \frac{1}{n} \cdot E(X^2) - \frac{n - 1}{n} \cdot [E(X)]^2 \\ E(S_X^2) &= \frac{n}{n} \cdot E(X^2) - \frac{1}{n} \cdot E(X^2) - \frac{n - 1}{n} \cdot [E(X)]^2 \\ E(S_X^2) &= \frac{n - 1}{n} \cdot E(X^2) - \frac{n - 1}{n} \cdot [E(X)]^2 \\ E(S_X^2) &= \frac{n - 1}{n} \cdot [E(X^2) - [E(X)]^2] \end{aligned}$$

Zum Vergleich:

$$\sigma_X^2 = E(X^2) - [E(X)]^2$$

also:

$$E(S_X^2) = \frac{n-1}{n} \cdot \sigma_X^2$$

Wir sehen, dass die Schätzer durchschnittlich kleiner sind als der zu schätzende Parameter. Die Stichprobenvarianz als Schätzfunktion der Populationsvarianz hat also einen negativen Bias.

## BIASKORREKTUR

Natürlich ist man daran interessiert, derartige systematische Fehler auszumerzen. Ein Schätzer, der stets zu kleine oder zu große Werte liefert, ist wenig bis gar nicht brauchbar. Man versucht, eine derartige Verzerrung formal in den Griff zu bekommen, um die Schätzfunktion korrigieren zu können.

Dies lässt sich sehr gut am (oben errechneten) Bias der Stichprobenvarianz demonstrieren.

Wir haben gezeigt:

$$E(S_X^2) = \frac{n-1}{n} \cdot \sigma_X^2$$

Um auf die gesuchte Populationsvarianz zu kommen, genügt eine kleine Umformung:

$$\sigma_X^2 = \frac{n}{n-1} \cdot E(S_X^2)$$

Wir können also die verzerrte Schätzung korrigieren, indem wir  $s_X^2 = \frac{1}{n} \cdot \sum_i (x_i - \bar{x})^2$  mit  $\frac{n}{n-1}$  multiplizieren. So erhalten wir folgende erwartungstreue Schätzfunktion für  $\sigma_X^2$ :

$$s_{korr}^2 = \frac{n}{n-1} \cdot s_X^2$$

$$s_{korr}^2 = \frac{1}{n-1} \cdot \sum_i (x_i - \bar{x})^2$$