

Robuste Statistik

Musterlösung zu Blatt 2

Aufgabe 2.1:

a)

Um die Äquivarianz- bzw. Invarianz-Eigenschaften des MAD herleiten zu können, ist es essentiell sich zunächst die Äquivarianz- und Invarianz-Eigenschaften des Medians anzuschauen.

Lokationseigenschaften (Median):

Sei $z_n = y_n + l$, $n = 1, \dots, N$, $l \in \mathbb{R}$

\Rightarrow für ungerades N

$$\text{med}(\mathbf{z}) = z_{(\frac{n+1}{2})} = y_{(\frac{n+1}{2})} + l = \text{med}(\mathbf{y}) + l$$

bzw. für gerades N

$$\text{med}(\mathbf{z}) = \left[z_{(\frac{n}{2})}, z_{(\frac{n}{2}+1)} \right] = \left[y_{(\frac{n}{2})} + l, y_{(\frac{n}{2}+1)} + l \right] = \text{med}(\mathbf{y}) + l$$

Diese Gleichheiten gelten, weil sich die Ordnung der Elemente durch die Addition einer Konstanten auf jedes y_n nicht ändert.

\Rightarrow Der Median ist lokations-äquivariant.

Skaleneigenschaften (Median):

Sei $z_n = s \cdot y_n$, $n = 1, \dots, N$, $s > 0$

\Rightarrow für ungerades N

$$\text{med}(\mathbf{z}) = z_{(\frac{n+1}{2})} = s \cdot y_{(\frac{n+1}{2})} = s \cdot \text{med}(\mathbf{y})$$

bzw. für gerades N

$$\text{med}(\mathbf{z}) = \left[z_{(\frac{n}{2})}, z_{(\frac{n}{2}+1)} \right] = \left[s \cdot y_{(\frac{n}{2})}, s \cdot y_{(\frac{n}{2}+1)} \right] = s \cdot \text{med}(\mathbf{y})$$

Diese Gleichheiten gelten, weil sich die Ordnung der Elemente durch die Multiplikation einer positiven Konstanten mit jedem y_n nicht ändert.

\Rightarrow Der Median ist skalen-äquivariant.

Lokationseigenschaften (MAD):

Sei $z_n = y_n + l$, $n = 1, \dots, N$, $l \in \mathbb{R}$

$$\begin{aligned} \text{MAD}(\mathbf{z}) &= \text{med}(|z_1 - \text{med}(\mathbf{z})|, \dots, |z_N - \text{med}(\mathbf{z})|) \\ &= \text{med}(|y_1 + l - \text{med}(\mathbf{y}) - l|, \dots, |y_N + l - \text{med}(\mathbf{y}) - l|) \\ &= \text{med}(|y_1 - \text{med}(\mathbf{y})|, \dots, |y_N - \text{med}(\mathbf{y})|) \\ &= \text{MAD}(\mathbf{y}) \end{aligned}$$

\Rightarrow Der MAD ist lokations-invariant.

Skaleneigenschaften (MAD):

Sei $z_n = s \cdot y_n, n = 1, \dots, N, s > 0$

$$\begin{aligned}MAD(\mathbf{z}) &= \text{med}(|z_1 - \text{med}(\mathbf{z})|, \dots, |z_N - \text{med}(\mathbf{z})|) \\&= \text{med}(|s \cdot y_1 - \text{med}(s \cdot \mathbf{y})|, \dots, |s \cdot y_N - \text{med}(s \cdot \mathbf{y})|) \\&= \text{med}(|s \cdot y_1 - s \cdot \text{med}(\mathbf{y})|, \dots, |s \cdot y_N - s \cdot \text{med}(\mathbf{y})|) \\&= \text{med}(s \cdot |y_1 - \text{med}(\mathbf{y})|, \dots, s \cdot |y_N - \text{med}(\mathbf{y})|) \\&= s \cdot \text{med}(|y_1 - \text{med}(\mathbf{y})|, \dots, |y_N - \text{med}(\mathbf{y})|) \\&= s \cdot MAD(\mathbf{y})\end{aligned}$$

\Rightarrow Der MAD ist skalen-äquivalent.

Somit erfüllt der MAD die Anforderungen an einen Skalenschätzer. Zusätzlich sehen wir, dass der Median die Anforderungen an einen Punktschätzer erfüllt.

b)

Der Modalwert ist derjenige Wert, der in der Stichprobe y_1, \dots, y_N am häufigsten angenommen wird. Werden mehrere Werte gleichhäufig und keiner häufiger angenommen, so ist der Modalwert nicht eindeutig.

Lokationseigenschaften (Modalwert):

Sei $z_n = y_n + l, n = 1, \dots, N, l \in \mathbb{R}$

Die Häufigkeitsverteilung der Werte z_1, \dots, z_N ist identisch zu der Häufigkeitsverteilung der Werte y_1, \dots, y_N , da sich Häufigkeitsverteilungen nicht ändern, wenn auf alle Werte die gleiche Konstante addiert wird. Somit ist $\text{mod}(\mathbf{z}) = \text{mod}(\mathbf{y}) + l$. Der Modalwert ist somit lokations-äquivalent.

Skaleneigenschaften (Modalwert):

Sei $z_n = s \cdot y_n, n = 1, \dots, N, s > 0$

Die Häufigkeitsverteilung der Werte z_1, \dots, z_N ist identisch zu der Häufigkeitsverteilung der Werte y_1, \dots, y_N , da sich Häufigkeitsverteilungen nicht ändern, wenn alle Werte mit der gleichen positiven Konstanten multipliziert werden. Somit ist $\text{mod}(\mathbf{z}) = s \cdot \text{mod}(\mathbf{y})$. Der Modalwert ist somit skalen-äquivalent.

Somit erfüllt der Modalwert die Anforderungen an einen Lokationsschätzer.

Aufgabe 2.2:

a)

Lokationseigenschaften:

Sei $z_n = y_n + l, n = 1, \dots, N, l \in \mathbb{R}$

$$\begin{aligned}\hat{\theta}_1(\mathbf{z}) &= \left(\frac{1}{N} \sum_{n=1}^N |z_n - \bar{z}| \right)^2 \\&= \left(\frac{1}{N} \sum_{n=1}^N |y_n + l - \bar{y} - l| \right)^2 \\&= \left(\frac{1}{N} \sum_{n=1}^N |y_n - \bar{y}| \right)^2 \\&= \hat{\theta}_1(\mathbf{y})\end{aligned}$$

$\Rightarrow \hat{\theta}_1$ ist lokations-invariant.

Skaleneigenschaften:

Sei $z_n = s \cdot y_n, n = 1, \dots, N, s > 0$

$$\begin{aligned}\hat{\theta}_1(\mathbf{z}) &= \left(\frac{1}{N} \sum_{n=1}^N |z_n - \bar{z}| \right)^2 \\ &= \left(\frac{1}{N} \sum_{n=1}^N |s \cdot y_n - s \cdot \bar{y}| \right)^2 \\ &= \left(\frac{1}{N} \sum_{n=1}^N s \cdot |y_n - \bar{y}| \right)^2 \\ &= \left(\frac{s}{N} \sum_{n=1}^N |y_n - \bar{y}| \right)^2 \\ &= s^2 \cdot \left(\frac{1}{N} \sum_{n=1}^N |y_n - \bar{y}| \right)^2 \\ &= s^2 \cdot \hat{\theta}_1(\mathbf{y})\end{aligned}$$

$\Rightarrow \hat{\theta}_1$ ist weder skalen-äquivalent noch skalen-invariant.

Somit erfüllt $\hat{\theta}_1$ weder die Anforderungen an einen Lokationsschätzer noch an einen Skalenschätzer.

b)

Lokationseigenschaften:

Sei $z_n = y_n + l, n = 1, \dots, N, l \in \mathbb{R}$

$$\begin{aligned}\hat{\theta}_2(\mathbf{z}) &= \sqrt{z^2 - \bar{z}^2} \\ &= \sqrt{\frac{1}{N} \sum_{n=1}^N z_n^2 - \left(\frac{1}{N} \sum_{n=1}^N z_n \right)^2} \\ &= \sqrt{\frac{1}{N} \sum_{n=1}^N (y_n + l)^2 - \left(\frac{1}{N} \sum_{n=1}^N (y_n + l) \right)^2} \\ &= \sqrt{\frac{1}{N} \sum_{n=1}^N (y_n^2 + 2ly_n + l^2) - \left(\frac{1}{N} \left(\sum_{n=1}^N y_n + Nl \right) \right)^2} \\ &= \sqrt{\frac{1}{N} \sum_{n=1}^N y_n^2 + \frac{1}{N} \sum_{n=1}^N 2ly_n + \frac{1}{N} \sum_{n=1}^N l^2 - \left(\frac{1}{N} \sum_{n=1}^N y_n + l \right)^2} \\ &= \sqrt{\frac{1}{N} \sum_{n=1}^N y_n^2 + \frac{2l}{N} \sum_{n=1}^N y_n + l^2 - \left(\frac{1}{N} \sum_{n=1}^N y_n \right)^2 - \frac{2l}{N} \sum_{n=1}^N y_n - l^2} \\ &= \sqrt{\frac{1}{N} \sum_{n=1}^N y_n^2 - \left(\frac{1}{N} \sum_{n=1}^N y_n \right)^2} \\ &= \sqrt{y^2 - \bar{y}^2} \\ &= \hat{\theta}_2(\mathbf{y})\end{aligned}$$

$\Rightarrow \hat{\theta}_2$ ist lokations-invariant.

Skaleneigenschaften:

Sei $z_n = s \cdot y_n, n = 1, \dots, N, s > 0$

$$\begin{aligned}
 \hat{\theta}_2(\mathbf{z}) &= \sqrt{z^2 - \bar{z}^2} \\
 &= \sqrt{\frac{1}{N} \sum_{n=1}^N z_n^2 - \left(\frac{1}{N} \sum_{n=1}^N z_n \right)^2} \\
 &= \sqrt{\frac{1}{N} \sum_{n=1}^N (s \cdot y_n)^2 - \left(\frac{1}{N} \sum_{n=1}^N s \cdot y_n \right)^2} \\
 &= \sqrt{\frac{1}{N} \sum_{n=1}^N s^2 \cdot y_n^2 - \left(\frac{s}{N} \sum_{n=1}^N y_n \right)^2} \\
 &= \sqrt{s^2 \cdot \frac{1}{N} \sum_{n=1}^N y_n^2 - s^2 \cdot \left(\frac{1}{N} \sum_{n=1}^N y_n \right)^2} \\
 &= \sqrt{s^2 \cdot \left(\frac{1}{N} \sum_{n=1}^N y_n^2 - \left(\frac{1}{N} \sum_{n=1}^N y_n \right)^2 \right)} \\
 &= s \cdot \sqrt{\frac{1}{N} \sum_{n=1}^N y_n^2 - \left(\frac{1}{N} \sum_{n=1}^N y_n \right)^2} \\
 &= s \cdot \sqrt{y^2 - \bar{y}^2} \\
 &= s \cdot \hat{\theta}_2(\mathbf{y})
 \end{aligned}$$

$\Rightarrow \hat{\theta}_2$ ist skalen-äquivalent.

Somit erfüllt $\hat{\theta}_2$ die Anforderungen an einen Skalenschätzer.

c)

Man kann sich leicht überlegen, dass die in Aufgabe 2.1 gezeigte Lokations- und Skalen-Äquivarianz des Medians auch für allgemeines p -Quantil gilt, siehe dazu auch die Darstellung in 3.1 im Skript zur Vorlesung.

Bonus: Ganz allgemein kann man auch zeigen, dass das p -Quantil äquivalent bezüglich streng monotoner, bijektiver Transformationen ist, d.h. für streng monotonen, bijektives $f : \mathbb{R} \rightarrow \mathbb{R}$ gilt

$$f(\tilde{y})_p = f(\tilde{y}_p). \quad (1)$$

Dabei ist $f(y)$ die koordinatenweise Anwendung der Abbildung auf f auf den Datensatz y . Formel (1) lässt sich wie folgt zeigen:

$$\begin{aligned}
 f(\tilde{y})_p &= \left\{ \mu \in \mathbb{R}; \#\{n; f(y_n) \leq \mu\} \geq pN \text{ und } \#\{n; f(y_n) \geq \mu\} \geq (1-p)N \right\} \\
 &= \left\{ \mu \in \mathbb{R}; \#\{n; y_n \leq f^{-1}(\mu)\} \geq pN \text{ und } \#\{n; y_n \geq f^{-1}(\mu)\} \geq (1-p)N \right\} \\
 &= \left\{ f(\mu) \in \mathbb{R}; \#\{n; y_n \leq \mu\} \geq pN \text{ und } \#\{n; y_n \geq \mu\} \geq (1-p)N \right\} \\
 &= f \left(\left\{ \mu \in \mathbb{R}; \#\{n; y_n \leq \mu\} \geq pN \text{ und } \#\{n; y_n \geq \mu\} \geq (1-p)N \right\} \right) = f(\tilde{y}_p).
 \end{aligned}$$

Wenn wir $f(x) = sx + l$ für $s > 0$ wählen, folgt, dass p -Quantile Lageschätzer sind.

Man beachte, dass der Schätzer $\hat{\theta}_3$ nur sinnvoll ist, falls $y_i > 0$ für beliebige $i = 1, \dots, N$.

Somit gilt insgesamt für $\hat{\theta}_3$:

Lokationseigenschaften:

Sei $z_n = y_n + l$, $n = 1, \dots, N$, $l \in \mathbb{R}$

$$\begin{aligned}\hat{\theta}_3(\mathbf{z}) &= \frac{\tilde{z}_{0.75} - \tilde{z}_{0.25}}{\tilde{z}_{0.5}} \\ &= \frac{(\tilde{y}_{0.75} + l) - (\tilde{y}_{0.25} + l)}{\tilde{y}_{0.5} + l} \\ &= \frac{\tilde{y}_{0.75} - \tilde{y}_{0.25}}{\tilde{y}_{0.5} + l}\end{aligned}$$

$\Rightarrow \hat{\theta}_3$ ist weder lokations-invariant noch lokations-äquivariant.

Skaleneigenschaften:

Sei $z_n = s \cdot y_n$, $n = 1, \dots, N$, $s > 0$

$$\begin{aligned}\hat{\theta}_3(\mathbf{z}) &= \frac{\tilde{z}_{0.75} - \tilde{z}_{0.25}}{\tilde{z}_{0.5}} \\ &= \frac{s \cdot \tilde{y}_{0.75} - s \cdot \tilde{y}_{0.25}}{s \cdot \tilde{y}_{0.5}} \\ &= \frac{s \cdot (\tilde{y}_{0.75} - \tilde{y}_{0.25})}{s \cdot \tilde{y}_{0.5}} \\ &= \frac{\tilde{y}_{0.75} - \tilde{y}_{0.25}}{\tilde{y}_{0.5}} \\ &= \hat{\theta}_3(\mathbf{y})\end{aligned}$$

$\Rightarrow \hat{\theta}_3$ ist skalen-invariant.

Somit erfüllt $\hat{\theta}_3$ weder die Anforderungen an einen Lokationsschätzer noch an einen Skalenschätzer.

Hinweis: Die Lokations- und Äquivarianzeigenschaften können auch gleichzeitig gezeigt werden, z.B.

$$\begin{aligned}\hat{\theta} &\text{ ist ein Lageschätzer, also skalen- und lokationsäquivariant} \\ \Leftrightarrow \forall \mathbf{y} \in \mathbb{R}^N \forall l \in \mathbb{R} \forall s > 0 : \hat{\theta}(s\mathbf{y} + l) &= s\hat{\theta}(\mathbf{y}) + l,\end{aligned}$$

oder z.B.

$$\begin{aligned}\hat{\theta} &\text{ ist ein Skalenschätzer, also skalenäquivariant und lokationsinvariant} \\ \Leftrightarrow \forall \mathbf{y} \in \mathbb{R}^N \forall l \in \mathbb{R} \forall s > 0 : \hat{\theta}(s\mathbf{y} + l) &= s\hat{\theta}(\mathbf{y}).\end{aligned}$$