

**3.2.1. Beispiel: Taxiproblem.** In einer Stadt gebe es  $N$  Taxis mit den vom Straßenrand aus lesbaren Nummern  $1, \dots, N$ . Ein Passant stehe eine gewisse Zeit lang an einer viel befahrenen Straße und notiere sich die Nummern  $x_1, \dots, x_k$  der vorbeifahrenden Taxis. Es sei angenommen, daß <sup>1</sup>  $x_1 < \dots < x_k$  und daß der Passant ein mehrmals vorbeifahrendes Taxi nur einmal zählt. Unter der Annahme, daß im Beobachtungszeitraum alle Taxis in Betrieb sind, ist die Anzahl  $N$  aller Taxis zu schätzen.

Als *statistisches Modell* <sup>2</sup> kann <sup>3</sup>  $(\mathcal{X}_k, \mathcal{G}_k; \mathbf{P}_{k;N}, N \in \mathbb{N}, N \geq k)$  mit

$$\mathcal{X}_k = \text{Menge der } k\text{-elementigen Teilmengen von } \mathbb{N}^4,$$

$$\mathcal{G}_k = \text{Pot}(\mathcal{X}_k),$$

$$\mathbf{P}_{k;N} = \text{Gleichverteilung auf der Menge der } k\text{-elementigen Teilmengen von } \{1, \dots, N\}, \quad N \in \mathbb{N}, N \geq k^5,$$

gewählt werden. Dieser Ansatz führt zur Likelihood-Funktion

$$L_{(k; x_1, \dots, x_k)}(N) = \begin{cases} 6 \binom{N}{k}^{-1}, & \text{falls } x_k \leq N, \\ 7 0, & \text{falls } x_k > N, \end{cases}$$

zur Beobachtung von  $k$  Taxis mit den Nummern  $x_1 < x_2 < \dots < x_k$ . Da für jedes  $k$  die Funktion  $\{k, k+1, \dots\} \ni N \rightarrow \binom{N}{k}^{-1}$  monoton fällt, ist <sup>8</sup>

$$S_1 = x_k$$

der *Maximum-Likelihood-Schätzer* für die Gesamtzahl  $N$  der Taxis.

Der Maximum-Likelihood-Schätzer  $S_1$  ist in der vorliegenden Situation unbefriedigend, da offensichtlich immer  $S_1 \leq N$  gilt, d.h., die „wahre“ Anzahl aller Taxis wird systematisch unterschätzt.

**Alternative Schätzer.** Mit heuristischen Argumenten können zwei weitere, evtl. <sup>9</sup> plausible Schätzer vorgeschlagen werden.

- Aus „Symmetriegründen“ sollte <sup>10</sup>  $x_1 - 1 \approx N - x_k$  gelten <sup>11</sup>. Als Schätzer für  $N$  ergibt sich dann:

$$S_2 = x_k + x_1 - 1.$$

- Es wäre auch sinnvoll, den Ansatz <sup>10 12</sup>

$$N - x_k \approx \frac{1}{k} \sum_{r=1}^k (x_r - x_{r-1} - 1) = \frac{1}{k} (x_k - k),$$

<sup>1</sup>Die Nummern der vorbeifahrenden Taxis werden in aufsteigender Reihenfolge notiert.

<sup>2</sup>Vgl. Abschnitt 3.1.

<sup>3</sup>Die Anzahl  $k$  der beobachteten Taxis wird nicht als eine Beobachtungsgröße, die zu den statistischen Schlussfolgerungen herangezogen wird, betrachtet. Nach dem Ende der Beobachtungen steht  $k$  fest und wird dann vor dem eigentlichen Beginn der statistischen Überlegungen als eine deterministische, d.h. nicht zufällige Zahl festgehalten.

<sup>4</sup>Beachte, daß  $\mathcal{X}_k$  abzählbar ist.

<sup>5</sup>Hier geht die Annahme ein, daß alle Taxis gleichmäßig im Stadtgebiet im Einsatz sind.

<sup>6</sup>In der Menge  $\{1, \dots, N\}$  existieren  $\binom{N}{k}$  Teilmengen mit  $k$  Elementen.

<sup>7</sup>Offensichtlich muß die Anzahl  $N$  aller Taxis  $\geq$  der höchsten beobachteten Nummer  $x_k$  sein.

<sup>8</sup>Der Maximum-Likelihood-Schätzer für die Gesamtzahl aller Taxis ist somit die größte der beobachteten Nummern.

<sup>9</sup>Dies ist natürlich Ansichtssache.

<sup>10</sup>Diese Vermutung sollte zumindest „im Mittel bei vielen Beobachtungsreihen“ gelten.

<sup>11</sup>Die Lücke bis zur kleinsten beobachteten Nummer  $x_1$ , bzw. die Lücke nach der größten beobachteten Nummer  $x_k$  sollten in etwa gleich sein.

<sup>12</sup>Es wird hier  $x_0 = 0$  gesetzt.

zu wählen <sup>13</sup>. Diese Überlegung führt nun zu <sup>14</sup>

$$S_3 = x_k + \lceil (x_k - k)/k \rceil$$

als Schätzer für  $N$  <sup>15</sup>.

#### LITERATUR

- [1] U. Krengel. Einführung in die Wahrscheinlichkeitstheorie und Statistik, 7. Auflage. Vieweg, 2003.

---

<sup>13</sup>Hierbei wird die Größe der Lücke nach der größten beobachteten Nummer  $x_k$  durch die „mittlere Größe aller Lücken“ geschätzt.

<sup>14</sup>Es sollte sinnvollerweise  $N$  durch eine ganze Zahl geschätzt werden.

<sup>15</sup>Die drei Schätzer  $S_1$ ,  $S_2$  und  $S_3$  für die Gesamtzahl  $N$  der Taxis besitzen unterschiedliche Eigenschaften, vgl. [1], Abschnitte 4.2 - 4.4. Zunächst kann nachgewiesen werden, daß  $S_2$  und  $S_3$  *erwartungstreue Schätzer* sind, d.h., für  $i = 2, 3$  gilt:

$$\mathbf{E}_{k;N}[S_i] := \sum_{l=k}^{\infty} l \cdot \mathbf{P}_{k;N}[S_i = l] = N, \quad N \in \mathbb{N}, N \geq k. \quad (*)$$

Andererseits ist  $S_1$  nicht erwartungstreu, d.h.,  $S_1$  erfüllt (\*) nicht. „Im Mittel“ wird daher durch die Schätzer  $S_2$  und  $S_3$  der wahre Wert von  $N$  gefunden. Hingegen wird durch  $S_1$  „im Mittel“ ein falscher Wert geschätzt.

Beim Vergleich von  $S_2$  und  $S_3$  zeigt sich, daß der *mittlere quadratische Fehler* für  $S_3$  kleiner als für  $S_2$  ist, d.h.,

$$\begin{aligned} \mathbf{E}_{k;N}[(S_3 - \mathbf{E}_{k;N}[S_3])^2] &= \sum_{l=k}^{\infty} (l - \mathbf{E}_{k;N}[S_3])^2 \cdot \mathbf{P}_{k;N}[S_3 = l] \\ &< \mathbf{E}_{k;N}[(S_2 - \mathbf{E}_{k;N}[S_2])^2], \quad N \in \mathbb{N}, N \geq k. \end{aligned}$$

Der Schätzer  $S_3$  schwankt daher „im quadratischen Mittel“ weniger als  $S_2$  um den wahren Wert von  $N$ .

Zusammenfassend ist also der Schätzer  $S_3$  gegenüber den beiden anderen Schätzern zu bevorzugen.