

week8a: Kapitel 5: Die Maximum Likelihood Methode 5.1: Grundidee und Beispiele, Teil1

Die Maximum Likelihood Methode ist eine Standard-Methode in der Statistik, um die Modell-Parameter eines statistischen Modells zu bestimmen. Maximum-Likelihood-Schätzer haben in der Regel gute statistische Eigenschaften. Die Funktionsweise der Methode versteht man am besten, wenn man ein paar Beispiele durchrechnet. Wir betrachten die folgenden 3 Beispiele:

- In Beispiel 1 sind n Zufallszahlen gegeben, die von einer Normalverteilung generiert worden sind, es sind aber nicht der Mittelwert und die Standardabweichung der Normalverteilung bekannt. Diese Parameter sollen aus den gegebenen Daten zurückgewonnen werden.
- In Beispiel 2 sind n Poisson-verteilte Zufallszahlen gegeben, der Parameter λ der Poisson-Verteilung ist aber nicht bekannt und soll geschätzt werden. Konzeptionell ist das nicht wesentlich anders als das Beispiel 1.
- In Beispiel 3 sind n Zeitreihendaten gegeben, die durch einen zeitdiskreten Ornstein-Uhlenbeck oder kurz OU-Prozess generiert worden sind. Der OU-Prozess, das schauen wir uns dann noch an, hat 3 Parameter, einen Mean-Reversion Level μ , eine Mean-Reversion Speed κ und einen Volatilitätsparameter σ . Diese 3 Parameter sollen aus den Zeitreihendaten zurückgewonnen werden.

Zufallszahlen, die von stochastischen Zeitreihenmodellen erzeugt werden, sind typischerweise nicht mehr unabhängig, deshalb ist das Beispiel 3 konzeptionell ein bisschen anders als die ersten beiden Beispiele 1 und 2. Alle drei Beispiele haben aber die schöne Eigenschaft, dass die Maximum-Likelihood-Schätzer analytisch in geschlossener Form berechnet werden können und nicht, wie es etwa bei ARCH- oder GARCH-Zeitreihenmodellen der Fall ist, numerisch durch Maximieren der log-Likelihood-Funktion bestimmt werden müssen. Ok, schauen wir uns das jetzt im Detail an:

Beispiel 1: Gegeben seien n Zufallszahlen x_1, x_2, \dots, x_n . Wir wissen, dass diese Zahlen mit Hilfe einer Normalverteilung generiert worden sind, kennen aber nicht den Mittelwert μ und die Standardabweichung σ . Welche Werte von μ und σ passen ‘am besten’ zu den gegebenen Zahlen x_1, \dots, x_n ? Und welches Kriterium nehmen wir für ‘am besten’?

Lösung 1: Die Aussage: “Die x_i sind normalverteilt mit Mittelwert μ und Standardabweichung σ ” ist äquivalent zu

$$\text{Prob} \left[x_i \in [\tilde{x}_i, \tilde{x}_i + d\tilde{x}_i] \right] = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\tilde{x}_i - \mu)^2}{2\sigma^2}} d\tilde{x}_i \quad (1)$$

Wir nehmen an, dass wir unabhängige Zufallszahlen haben, und bekommen dann aus (1)

$$\begin{aligned}
 \text{Prob} \left[x_1 \in [\tilde{x}_1, \tilde{x}_1 + d\tilde{x}_1), x_2 \in [\tilde{x}_2, \tilde{x}_2 + d\tilde{x}_2), \dots, x_n \in [\tilde{x}_n, \tilde{x}_n + d\tilde{x}_n) \right] &= \\
 \stackrel{\text{unabh.}}{=} \text{Prob} \left[x_1 \in [\tilde{x}_1, \tilde{x}_1 + d\tilde{x}_1) \right] \text{Prob} \left[x_2 \in [\tilde{x}_2, \tilde{x}_2 + d\tilde{x}_2) \right] \cdots \text{Prob} \left[x_n \in [\tilde{x}_n, \tilde{x}_n + d\tilde{x}_n) \right] & \\
 \stackrel{(1)}{=} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\tilde{x}_i - \mu)^2}{2\sigma^2}} d\tilde{x}_i & \\
 = (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (\tilde{x}_i - \mu)^2 \right\} d\tilde{x}_1 \cdots d\tilde{x}_n & \quad (2)
 \end{aligned}$$

Jetzt haben wir ja konkrete Werte für die Zufallszahlen, nämlich x_1, \dots, x_n . Diese Werte können wir auf der rechten Seite von (2) für die \tilde{x}_i einsetzen. Die $d\tilde{x}_i$ sind Intervallbreiten, etwa $d\tilde{x}_i = 0.1$ oder $d\tilde{x}_i = 0.01$. Wir werden gleich sehen, dass diese Intervallbreiten für die eigentliche Rechnung keine Rolle spielen, aber um eine konkrete Vorstellung zu haben, wählen wir etwa $d\tilde{x}_i = 0.01$ und wir lassen die Tilde weg, schreiben einfach dx_i . Dann bekommen wir eine Funktion, die nur noch von μ und σ abhängt,

$$L(\mu, \sigma) := (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\} dx_1 \cdots dx_n \quad (3)$$

Diese Funktion heisst **Likelihood-Funktion**, sie gibt also die W'keit an, dass, unter der Annahme, die Zufallszahlen sind von einer Normalverteilung mit Mittelwert μ und Standardabweichung σ generiert worden, dass sich dann konkret Werte in den Intervallen $[x_i, x_i + dx_i)$ realisieren. Also ist es plausibel, das μ und das σ dann so zu wählen, dass diese W'keit maximal wird,

$$L(\mu, \sigma) \stackrel{!}{\rightarrow} \max \quad \Rightarrow \quad \mu = \mu_{\text{ML}}, \sigma = \sigma_{\text{ML}} \quad (4)$$

Die so gewonnenen Werte von μ und σ heissen dann die **Maximum-Likelihood-Schätzer** von μ und σ .

In unserem konkreten Beispiel können wir das Maximum der Likelihood-Funktion (3) analytisch in geschlossener Form berechnen, das machen wir gleich in dem Theorem 5.1.1 weiter unten. In vielen Anwendungssituationen ist jedoch eine analytische Berechnung in geschlossener Form nicht möglich und man ist auf numerische Prozeduren angewiesen. Anstatt die Likelihood-Funktion L selber zu maximieren, betrachtet man dann typischerweise den Logarithmus der Likelihood-Funktion $\log L$. Da der Logarithmus eine streng monoton steigende Funktion ist, gilt

$$(\mu, \sigma) \text{ maximiert } L(\mu, \sigma) \quad \Leftrightarrow \quad (\mu, \sigma) \text{ maximiert } \log L(\mu, \sigma)$$

Das macht man jetzt aus folgendem Grund: Das L ist typischerweise durch ein Produkt gegeben und die Anzahl der Faktoren in diesem Produkt ist gleich der Anzahl der verfügbaren Daten. Wenn man etwa eine Finanz-Zeitreihe hat mit täglichen Schlusskursen der letzten 10 Jahre, dann sind das etwa $N = 2500$ Daten und das L ist dann im wesentlichen

$$L \sim \text{prob}^N = \text{prob}^{2500}$$

Auf dem Computer ist die Zahl prob^{2500} entweder 0 oder unendlich, je nachdem, ob das prob kleiner oder grösser ist als 1 (probieren Sie das ruhig mal aus), also das ist numerisch nicht handhabbar. In Gegensatz dazu ist

$$\log L \sim \log[\text{prob}^N] = N \log[\text{prob}] = 2500 \log[\text{prob}]$$

und das ist numerisch unproblematisch.

Kehren wir jetzt zu unserem konkreten Beispiel 1 zurück. Auch für die analytische Rechnung ist es günstig, den Logarithmus zu betrachten,

$$\begin{aligned} \log L(\mu, \sigma) &= \log \left[(2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\} dx_1 \cdots dx_n \right] \\ &= -\frac{n}{2} \log[2\pi\sigma^2] - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 + \log[dx_1 \cdots dx_n] \\ &= -n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 + \text{const} \end{aligned} \quad (5)$$

wobei die Grösse

$$\text{const} := -\frac{n}{2} \log[2\pi] + \log[dx_1 \cdots dx_n]$$

eine von μ und σ unabhängige Zahl ist. Es ist also hinreichend, die Funktion

$$\log \tilde{L}(\mu, \sigma) := -n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \quad (6)$$

zu betrachten. Wie oben bereits angekündigt, stellen wir also fest, dass die Intervallbreiten dx_1, \dots, dx_n für die eigentliche Rechnung irrelevant sind. Es gilt das folgende

Theorem 5.1.1: Das $\log \tilde{L}(\mu, \sigma)$ aus Gleichung (6) wird maximiert durch

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i =: \mu_{\text{ML}}(x_1, \dots, x_n) \quad (7)$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 =: \sigma_{\text{ML}}^2(x_1, \dots, x_n) \quad (8)$$

mit $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i = \mu_{\text{ML}}(x_1, \dots, x_n)$.

Beweis: Notwendige Bedingung für ein Maximum ist

$$\frac{\partial}{\partial \mu} \log \tilde{L}(\mu, \sigma) = 0 \quad (9)$$

$$\frac{\partial}{\partial \sigma} \log \tilde{L}(\mu, \sigma) = 0 \quad (10)$$

Nun ist

$$\frac{\partial}{\partial \mu} \log \tilde{L}(\mu, \sigma) = + \frac{2}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)$$

also folgt aus Gleichung (9)

$$\sum_{i=1}^n (x_i - \mu) = \sum_{i=1}^n x_i - n\mu = 0$$

oder

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

Für die Ableitung nach σ erhalten wir

$$\frac{\partial}{\partial \sigma} \log \tilde{L}(\mu, \sigma) = -\frac{n}{\sigma} + \frac{2}{2\sigma^3} \sum_{i=1}^n (x_i - \mu)^2$$

und Gleichung (10) liefert

$$\frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 = \frac{n}{\sigma}$$

oder

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

Damit ist das Theorem bewiesen. ■

Machen wir gleich das Beispiel 2, wie gesagt, das geht fast genauso:

Beispiel 2: Gegeben seien n Zufallszahlen x_1, x_2, \dots, x_n . Wir wissen, dass diese Zahlen mit Hilfe einer Poisson-Verteilung generiert worden sind, kennen aber nicht den Parameter λ der Poisson-Verteilung. Welcher Wert von λ passt am besten, im Sinne der Maximum Likelihood Methode, zu den gegebenen Zahlen x_1, \dots, x_n ?

Lösung 2: Die Aussage: “Die x_i sind Poisson-verteilt mit Parameter $\lambda \in \mathbb{R}$ ” ist äquivalent zu

$$\text{Prob}[x_i = k_i] = \frac{\lambda^{k_i}}{k_i!} e^{-\lambda} \quad (11)$$

mit natürlichen Zahlen $k_i \in \mathbb{N}$. Wir nehmen an, dass wir unabhängige Zufallszahlen haben, und bekommen dann aus (11)

$$\begin{aligned} & \text{Prob}[x_1 = k_1, x_2 = k_2, \dots, x_n = k_n] \\ & \stackrel{\text{unabh.}}{=} \text{Prob}[x_1 = k_1] \text{Prob}[x_2 = k_2] \cdots \text{Prob}[x_n = k_n] \\ & \stackrel{(11)}{=} \prod_{i=1}^n \left\{ \frac{\lambda^{k_i}}{k_i!} e^{-\lambda} \right\} = \frac{\lambda^{k_1 + \dots + k_n}}{k_1! \cdots k_n!} e^{-n\lambda} \end{aligned} \quad (12)$$

Jetzt haben wir ja konkrete Werte für die Zufallszahlen, nämlich x_1, \dots, x_n . Diese Werte können wir auf der rechten Seite von (12) für die k_i einsetzen. Dann bekommen wir wieder eine Funktion, die nur noch von dem gesuchten λ abhängt,

$$L(\lambda) := \frac{\lambda^{x_1 + \dots + x_n}}{x_1! \dots x_n!} e^{-n\lambda} \quad (13)$$

und das ist dann die Likelihood-Funktion für dieses konkrete Beispiel. Wir wollen das Maximum finden und betrachten dazu wieder den Logarithmus der Likelihood-Funktion,

$$\begin{aligned} \log L(\lambda) &= \log \left[\frac{\lambda^{x_1 + \dots + x_n}}{x_1! \dots x_n!} e^{-n\lambda} \right] \\ &= \log[\lambda^{x_1 + \dots + x_n}] - \log[x_1! \dots x_n!] + \log[e^{-n\lambda}] \\ &= (x_1 + \dots + x_n) \log \lambda - n\lambda + \text{const} \end{aligned} \quad (14)$$

wobei die Zahl

$$\text{const} := -\log[x_1! \dots x_n!]$$

wieder nicht von dem gesuchten Parameter, in diesem Fall das λ , abhängt, so dass sie für das Maximieren vernachlässigt werden kann, da sie etwa beim Ableiten nach λ wegfällt. Also:

$$\frac{d}{d\lambda} \log L(\lambda) = \frac{x_1 + \dots + x_n}{\lambda} - n \stackrel{!}{=} 0$$

und wir bekommen

$$\lambda = \frac{x_1 + \dots + x_n}{n} =: \hat{\lambda}_{\text{ML}}(x_1, \dots, x_n) \quad (15)$$

Die Funktion $\hat{\lambda}_{\text{ML}} : \mathbb{R}^n \rightarrow \mathbb{R}$ gegeben durch (15) heisst dann der Maximum Likelihood Schätzer für das λ . Dass das in diesem Fall gerade der Mittelwert der x_i ist, ist natürlich nicht weiter erstaunlich, da für eine mit Parameter λ Poisson-verteilte Zufallszahl x ja die Gleichung $E[x] = \lambda$ gilt.

Das Beispiel 3 diskutieren wir dann in der nächsten Veranstaltung.