

week9: Der zentrale Grenzwertsatz

In den letzten beiden Wochen hatten wir in dem week7.pdf und dem week8.pdf die Monte Carlo Summen

$$\frac{1}{n} \sum_{i=1}^n F(x_i)$$

betrachtet. Dabei waren die x_1, \dots, x_n Zufallszahlen, die alle dieselbe Wahrscheinlichkeitsverteilung $p(x)$ hatten, so dass der Erwartungswert für alle $F(x_i)$ derselbe war, nämlich

$$\mathbb{E}[F(x_i)] = \int_{\mathbb{R}} F(x_i) p(x_i) dx_i = \int_{\mathbb{R}} F(x) p(x) dx = \mathbb{E}[F]$$

Die Varianz war ebenfalls für alle $F(x_i)$ dieselbe, das ist dann der Ausdruck

$$\mathbb{V}[F(x_i)] = \mathbb{E}[[F(x_i)]^2] - (\mathbb{E}[F(x_i)])^2 = \mathbb{E}[F^2] - (\mathbb{E}[F])^2 = \mathbb{V}[F]$$

mit

$$\mathbb{E}[[F(x_i)]^2] = \int_{\mathbb{R}} [F(x_i)]^2 p(x_i) dx_i = \int_{\mathbb{R}} [F(x)]^2 p(x) dx = \mathbb{E}[F^2]$$

Zur Formulierung des zentralen Grenzwertsatzes wollen wir jetzt etwas kompakter schreiben

$$F(x_i) =: X_i$$

mit

$$\mathbb{E}[X_i] = \mathbb{E}[F] =: \mu$$

$$\mathbb{V}[X_i] = \mathbb{V}[F] =: \sigma^2$$

Wir betrachten die Summen und die Mittelwerte (in week7 und 8 hatten wir immer von Monte Carlo Summe gesprochen, in der Terminologie, die wir jetzt benutzen wollen, wären das also genauer Monte Carlo Mittelwerte, die S_n aus week7 und 8 sind hier jetzt M_n 's)

$$S_n := \sum_{i=1}^n X_i$$
$$M_n := \frac{1}{n} \sum_{i=1}^n X_i = \frac{S_n}{n}$$

Wenn man stochastische Größen hat, macht es Sinn, diese zu standardisieren: man zieht den Erwartungswert ab und teilt durch die Standardabweichung, das ist die Wurzel aus der Varianz. Wir wollen also die standardisierten Größen

$$Z_n := \frac{S_n - \mathbb{E}[S_n]}{\sqrt{\mathbb{V}[S_n]}} \quad (1)$$

$$\tilde{Z}_n := \frac{M_n - \mathbb{E}[M_n]}{\sqrt{\mathbb{V}[M_n]}} \quad (2)$$

betrachten. Dann gilt zunächst mal das folgende

Lemma 9.1: Es seien X_1, \dots, X_n unabhängige und identisch verteilte Zufallsgrößen mit

$$\begin{aligned} \mathbb{E}[X_i] &= \mu \\ \mathbb{V}[X_i] &= \sigma^2 \end{aligned}$$

für alle $i = 1, \dots, n$. Dann gilt

a)

$$\begin{aligned} \mathbb{E}[S_n] &= n\mu \\ \mathbb{E}[M_n] &= \mu \end{aligned}$$

und

$$\begin{aligned} \mathbb{V}[S_n] &= n\sigma^2 \\ \mathbb{V}[M_n] &= \frac{\sigma^2}{n} \end{aligned}$$

b) Die standardisierten Größen Z_n und \tilde{Z}_n aus (1) und (2) sind identisch, es gilt

$$Z_n = \tilde{Z}_n = \frac{1}{\sqrt{n}} \left\{ \frac{X_1 - \mu}{\sigma} + \dots + \frac{X_n - \mu}{\sigma} \right\}$$

Beweis: a) Das haben wir im wesentlichen im week8.pdf in den Gleichungen (8) und (9) schon gemacht. Schreiben wir die Sachen eben nochmal hin, wir brauchen nur, dass der Erwartungswert und die Kovarianz linear sind:

$$\begin{aligned} \mathbb{E}[S_n] &= \mathbb{E}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \mathbb{E}[X_i] = \sum_{i=1}^n \mu = n\mu \\ \mathbb{E}[M_n] &= \mathbb{E}\left[\frac{S_n}{n}\right] = \frac{1}{n} \mathbb{E}[S_n] = \frac{1}{n} n\mu = \mu \end{aligned}$$

Und für die Varianzen:

$$\begin{aligned} \mathbb{V}[S_n] &= \mathbb{V}\left[\sum_{k=1}^n X_k\right] = \text{Cov}\left[\sum_{k=1}^n X_k, \sum_{\ell=1}^n X_\ell\right] = \sum_{k,\ell=1}^n \text{Cov}[X_k, X_\ell] \\ &\stackrel{\substack{X_k, X_\ell \\ \text{unabhängig}}}{=} \sum_{k=1}^n \text{Cov}[X_k, X_k] = \sum_{k=1}^n \mathbb{V}[X_k] = \sum_{k=1}^n \sigma^2 = n\sigma^2 \end{aligned}$$

und damit

$$\mathbb{V}[M_n] = \mathbb{V}\left[\frac{S_n}{n}\right] = \frac{1}{n^2} \mathbb{V}[S_n] = \frac{1}{n^2} n\sigma^2 = \frac{\sigma^2}{n}$$

b) Wir haben

$$\begin{aligned} Z_n &= \frac{S_n - \mathbb{E}[S_n]}{\sqrt{\mathbb{V}[S_n]}} = \frac{S_n - n\mu}{\sqrt{n\sigma^2}} = \frac{\frac{1}{n}S_n - \mu}{\frac{1}{n}\sqrt{n\sigma^2}} = \frac{\frac{1}{n}S_n - \mu}{\sqrt{\frac{\sigma^2}{n}}} \\ &= \frac{M_n - \mathbb{E}[M_n]}{\sqrt{\mathbb{V}[M_n]}} = \tilde{Z}_n, \end{aligned}$$

die beiden Grössen sind also identisch, wir können beides mit Z_n bezeichnen und brauchen kein \tilde{Z}_n . Wir können das Z_n dann auch folgendermassen schreiben:

$$Z_n = \frac{S_n - n\mu}{\sqrt{n\sigma^2}} = \frac{1}{\sqrt{n}} \frac{X_1 + \dots + X_n - n\mu}{\sigma} = \frac{1}{\sqrt{n}} \left\{ \frac{X_1 - \mu}{\sigma} + \dots + \frac{X_n - \mu}{\sigma} \right\}$$

und das Lemma ist bewiesen. ■

Definieren wir die Zufallsgrössen

$$Y_i := \frac{X_i - \mu}{\sigma} = \frac{X_i - \mathbb{E}[X_i]}{\sqrt{\mathbb{V}[X_i]}}$$

dann gilt

$$\mathbb{E}[Y_i] = \mathbb{E}\left[\frac{X_i - \mu}{\sigma}\right] = \frac{\mathbb{E}[X_i] - \mu}{\sigma} = \frac{\mu - \mu}{\sigma} = 0$$

und

$$\mathbb{V}[Y_i] = \mathbb{V}\left[\frac{X_i - \mu}{\sigma}\right] = \frac{1}{\sigma^2} \mathbb{V}[X_i - \mu] = \frac{1}{\sigma^2} \mathbb{V}[X_i] = \frac{1}{\sigma^2} \sigma^2 = 1$$

Die Y_i sind also standardisiert, sie haben Mittelwert 0 und Standardabweichung 1. Mit Hilfe der Y_i können wir das Z_n dann auch einfach schreiben als

$$Z_n = \frac{Y_1 + \dots + Y_n}{\sqrt{n}} \tag{3}$$

Wären die Y_i jetzt alle normalverteilt, dann können wir exakt vorhersagen, welche Verteilung (also wie die Histogramme aussehen wenn man viele von den Z_n erzeugt) die Z_n hätten, sie wären dann nämlich ebenfalls wieder normalverteilt. Das ergibt sich aus dem folgenden allgemeinen Theorem, welches wir auf dem neuen Übungsblatt 9 durch eine geeignete R-Simulation überprüfen:

Theorem 9.1 (Summen von normalverteilten Zufallszahlen): Es seien ϕ_1, \dots, ϕ_n normalverteilte Zufallsvariablen mit Mittelwerten μ_1, \dots, μ_n und Standardabweichungen $\sigma_1, \dots, \sigma_n$. Dann gilt: Die Summe

$$\phi := \phi_1 + \phi_2 + \dots + \phi_n \tag{4}$$

ist normalverteilt mit Mittelwert

$$\mu = \mu_1 + \mu_2 + \dots + \mu_n \tag{5}$$

und Varianz

$$\sigma^2 = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_n^2 \quad (6)$$

Das heisst genauer, für eine beliebige Funktion $f : \mathbb{R} \rightarrow \mathbb{R}$ gilt

$$\int_{\mathbb{R}^n} f(\phi_1 + \dots + \phi_n) \prod_{k=1}^n e^{-\frac{(\phi_k - \mu_k)^2}{2\sigma_k^2}} \frac{d\phi_k}{\sqrt{2\pi\sigma_k^2}} = \int_{\mathbb{R}} f(\phi) e^{-\frac{(\phi - \mu)^2}{2\sigma^2}} \frac{d\phi}{\sqrt{2\pi\sigma^2}} \quad (7)$$

mit μ und σ^2 gegeben durch (5) und (6).

Wählen wir für das f dann etwa die Indikator-Funktion

$$f(\phi) := \chi(x \leq \phi \leq x + dx) := \begin{cases} 1 & \text{falls } x \leq \phi \leq x + dx \\ 0 & \text{sonst} \end{cases}$$

dann ergibt sich aus (7), wenn wir das dx sehr klein wählen,

$$\begin{aligned} & \text{Prob} \left[\phi_1 + \dots + \phi_n \in [x, x + dx] \right] \\ &= \int_{\mathbb{R}^n} \chi(x \leq \phi_1 + \dots + \phi_n \leq x + dx) \prod_{k=1}^n e^{-\frac{(\phi_k - \mu_k)^2}{2\sigma_k^2}} \frac{d\phi_k}{\sqrt{2\pi\sigma_k^2}} \\ &\stackrel{\text{Thm. 9.1}}{=} \int_{\mathbb{R}} \chi(x \leq \phi \leq x + dx) e^{-\frac{(\phi - \mu)^2}{2\sigma^2}} \frac{d\phi}{\sqrt{2\pi\sigma^2}} \\ &\stackrel{dx \text{ klein}}{\approx} e^{-\frac{(x - \mu)^2}{2\sigma^2}} \frac{dx}{\sqrt{2\pi\sigma^2}} \end{aligned}$$

die Summe $\phi_1 + \dots + \phi_n$ ist also wieder normalverteilt mit Mittelwert μ und Standardabweichung σ gegeben durch die Formeln (5) und (6).

Kehren wir zu den Z_n aus (1), (2) und (3) zurück,

$$\begin{aligned} Z_n &= \frac{S_n - n\mu}{\sqrt{n\sigma^2}} = \frac{M_n - \mu}{\sqrt{\sigma^2/n}} \\ &= \frac{1}{\sqrt{n}} \left\{ \frac{X_1 - \mu}{\sigma} + \dots + \frac{X_n - \mu}{\sigma} \right\} = \frac{Y_1 + \dots + Y_n}{\sqrt{n}} \end{aligned} \quad (8)$$

Wären die X_i nun alle normalverteilt, dann wären die Y_i alle standard-normalverteilt und nach dem Theorem 9.1 wäre dann auch das Z_n normalverteilt mit Varianz

$$\mathbb{V}[Z_n] = \left(\frac{1}{\sqrt{n}} \right)^2 \mathbb{V}[Y_1 + \dots + Y_n] = \frac{1}{n} \left(\underbrace{\sigma_1^2}_{=1} + \dots + \underbrace{\sigma_n^2}_{=1} \right) = \frac{1}{n} n = 1$$

die Z_n wären also standard-normalverteilte Zufallszahlen. Wenn die X_i nun eine beliebige Verteilung haben, etwa alle X_i sind gleichverteilt auf dem Intervall $[0, 1]$, dann ist klar, dass

die Z_n nicht normalverteilt sein können, weil die S_n 's dann zum Beispiel immer im endlichen Intervall $[0, n]$ liegen und die Z_n 's dann auch immer in einem endlichen Intervall liegen müssen. Bei normalverteilten Zahlen kommt da, zumindest theoretisch, aber immer ganz \mathbb{R} in Frage.

Also die exakte Verteilung der Z_n kann irgendetwas komplizierteres sein, was sich möglicherweise auch nur sehr schwer oder überhaupt nicht explizit berechnen lassen tut. Nun ist es aber so, dass für grössere Werte von n (n gleich 20 oder 50 reicht da gelegentlich schon aus, werden wir uns dann konkret in der nächsten Veranstaltung in R anschauen) sich die Verteilung der Z_n durch die Standard-Normalverteilung approximieren lässt, das ist jetzt die Aussage des zentralen Grenzwertsatzes:

Theorem 9.2 (Zentraler Grenzwertsatz): Es seien X_1, \dots, X_n eine Folge von unabhängigen, identisch verteilten Zufallsvariablen mit Erwartungswert μ und Standardabweichung σ und die normierten oder standardisierten Zufallsgrössen Z_n seien gegeben durch die Gleichung (8). Dann gilt: Im Limes $n \rightarrow \infty$ sind die Z_n standard-normalverteilt. Das heisst genauer, für ein beliebiges $f : \mathbb{R} \rightarrow \mathbb{R}$ gilt:

$$\lim_{n \rightarrow \infty} \mathbb{E}[f(Z_n)] = \int_{\mathbb{R}} f(z) e^{-\frac{z^2}{2}} \frac{dz}{\sqrt{2\pi}} \quad (9)$$

Beweis: Werden wir in der nächsten Veranstaltung mit einer R-Simulation verifizieren. ■

Wählen wir für das f wieder die Indikator-Funktion

$$f(z) := \chi(x \leq z \leq x + dx)$$

dann ergibt sich aus (9)

$$\lim_{n \rightarrow \infty} \text{Prob} \left[Z_n \in [x, x + dx] \right] = e^{-\frac{x^2}{2}} \frac{dx}{\sqrt{2\pi}}$$

Im Limes $n \rightarrow \infty$ sind die Z_n also standard-normalverteilt. Wir notieren uns noch die folgende

Folgerung 9.1: Mit $X_i = F(x_i)$ und

$$M_n = \frac{1}{n} \sum_{i=1}^n F(x_i)$$

und

$$Z_n = \frac{M_n - \mu}{\sqrt{\sigma^2/n}} = \frac{M_n - \mathbb{E}[F]}{\sqrt{\mathbb{V}[F]/n}}$$

können wir auch schreiben:

$$M_n = \mathbb{E}[F] + \sqrt{\frac{\mathbb{V}[F]}{n}} Z_n \quad (10)$$

Da standard-normalverteilte Zufallszahlen typischerweise so zwischen -5 und +5 liegen, oder, machen wir es etwas genauer (wir schauen uns das auch nochmal in Aufgabe 3 auf dem Übungsblatt 9 an),

$$\lim_{n \rightarrow \infty} \mathbf{Prob} \left[Z_n \in [-5, +5] \right] = \int_{-5}^{+5} e^{-\frac{x^2}{2}} \frac{dx}{\sqrt{2\pi}} = \Phi(+5) - \Phi(-5) \approx 0.9999994$$

lässt sich aus der Gleichung (10) also ablesen, dass der Monte Carlo Fehler typischerweise von der Grössenordnung

$$\text{Monte Carlo Fehler} = O\left(\sqrt{\frac{\mathbf{V}[F]}{n}}\right)$$

ist.