

## week5: Wahrscheinlichkeitsverteilungen und Zufallszahlen in R

**Allgemeine Logik:** Die Wahrscheinlichkeitsverteilung  $p = p(x)$  einer kontinuierlichen Zufallszahl  $x$  mit  $p(x) \geq 0$  und  $\int_{-\infty}^{+\infty} p(x) dx = 1$  hat die R-Syntax

$$p(x) = \text{dVert}(x, \dots)$$

wobei `Vert` der R-Name der Verteilung ist und die Punkte  $\dots$  stehen für optionale Parameter der Verteilung. Dabei steht das `d` für ‘density’. Eine Übersicht über die in R eingebauten Wahrscheinlichkeitsverteilungen findet man hier:

<https://hsrm-mathematik.de/WS2425/Wirtschaftsmathematik3/W'keitsverteilungen-in-R.pdf>

Neben der Dichte gibt es noch die R-Funktionen

$$\text{pVert}(x, \dots) := \int_{-\infty}^x p(y) dy$$

wobei das `p` etwa für probability steht und es gibt die Umkehr-Funktion von dieser Funktion,

$$\text{pVert}(x, \dots) = y$$

$$\stackrel{\text{Def. qVert}}{\Leftrightarrow} x = \text{qVert}(y, \dots)$$

wobei das `q` dann für Quantil steht. Das `pVert` wird auch als Verteilungsfunktion bezeichnet und das `qVert` als Quantil-Funktion. Schliesslich gibt es noch die Funktion

$$\text{rVert}(n, \dots)$$

die  $n$  Zufallszahlen mit Verteilung  $p(x)$  generiert. Dabei steht das `r` dann für random numbers. Die Funktionen `dVert`, `pVert` und `qVert` sind deterministische Funktionen, nur das `rVert` liefert zufällige Resultate.

**Beispiel 1, uniforme Verteilung:** Die uniforme oder auch Gleichverteilung auf dem Intervall  $[a, b]$  hatten wir im `week2.pdf` definiert durch die Wahrscheinlichkeitsdichte

$$p_{a,b}(x) := \begin{cases} 1/(b-a) & \text{falls } x \in [a,b] \\ 0 & \text{sonst} \end{cases}$$

$$\stackrel{\text{R-Befehl}}{=} \text{dunif}(\mathbf{x}, \text{min} = \mathbf{a}, \text{max} = \mathbf{b})$$

Daraus ergibt sich

$$\text{punif}(\mathbf{x}, \text{min} = \mathbf{a}, \text{max} = \mathbf{b}) := \int_{-\infty}^x p_{a,b}(y) dy$$

$$= \begin{cases} 0 & \text{falls } x \leq a \\ \int_a^x \frac{1}{b-a} dy & \text{falls } a < x < b \\ 1 & \text{falls } x \geq b \end{cases}$$

oder, wir können hier das Integral explizit berechnen,

$$\text{punif}(\mathbf{x}, \text{min} = \mathbf{a}, \text{max} = \mathbf{b}) = \begin{cases} 0 & \text{falls } x \leq a \\ \frac{x-a}{b-a} & \text{falls } a < x < b \\ 1 & \text{falls } x \geq b \end{cases}$$

Die Umkehrfunktion, die Quantil-Funktion  $\text{qunif}(\mathbf{y}, \text{min} = \mathbf{a}, \text{max} = \mathbf{b})$  ergibt sich dann folgendermassen: Gegeben sei ein  $y$  zwischen 0 und 1,  $y \in [0, 1]$ . Das  $y$  hat die Bedeutung einer Wahrscheinlichkeit. Wir müssen die Gleichung

$$y = \int_{-\infty}^x p(\tilde{x}) d\tilde{x} \quad (1)$$

nach  $x$  auflösen, das  $x$  hat dann die folgende Bedeutung: Erzeugen wir  $n$  Zufallszahlen mit Verteilung  $\text{Vert}$ , dann ist der Anteil an Zufallszahlen  $n_{\leq x}$  die kleiner oder gleich  $x$  sind, gegeben durch

$$\frac{n_{\leq x}}{n} = y$$

Hier in unserem konkreten Beispiel reduziert sich die Gleichung (1) auf

$$y = \frac{x-a}{b-a}$$

$$\Leftrightarrow (b-a) \cdot y = x-a$$

$$\Leftrightarrow x = (b-a) \cdot y + a$$

Wenn wir etwa  $y = 0$  setzen, ist  $x = a$  und der Anteil  $n_a/n$  an Zufallszahlen kleiner oder gleich  $a$  ist 0. Setzen wir  $y = 1$ , bekommen wir  $x = b$  und der Anteil  $n_b/n$  an Zufallszahlen kleiner oder gleich  $b$  ist 1. Da alle Zufallszahlen nach Definition von  $p_{a,b}(x)$  zwischen  $a$  und  $b$  liegen, ist klar, dass das dann so sein muss. Also für die uniforme Verteilung ist

$$\text{qunif}(\mathbf{y}, \text{min} = \mathbf{a}, \text{max} = \mathbf{b}) = (b-a) \cdot \mathbf{y} + \mathbf{a}$$

mit  $0 \leq y \leq 1$ .

**Beispiel 2, Normalverteilung:** Die Normalverteilung mit Mittelwert  $\mu$  und Standardabweichung  $\sigma$  ist definiert durch die Wahrscheinlichkeitsdichte

$$\begin{aligned}\varphi_{\mu,\sigma}(x) &:= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \\ &\stackrel{\text{R-Befehl}}{=} \text{dnorm}(x, \text{mean} = \mu, \text{sd} = \sigma)\end{aligned}$$

Die Verteilungsfunktion ist dann gegeben durch

$$\begin{aligned}\Phi_{\mu,\sigma}(x) &:= \int_{-\infty}^x \varphi_{\mu,\sigma}(y) dy = \int_{-\infty}^x e^{-\frac{(y-\mu)^2}{2\sigma^2}} \frac{dy}{\sqrt{2\pi\sigma^2}} \\ &\stackrel{\text{R-Befehl}}{=} \text{pnorm}(x, \text{mean} = \mu, \text{sd} = \sigma)\end{aligned}$$

In diesem Fall lässt sich das Integral nicht weiter vereinfachen. Die Verteilungsfunktion braucht man etwa bei der folgenden typischen Aufgabenstellung:

**Aufgabe:** Eine Zufallszahl  $X$  sei normalverteilt mit Mittelwert  $\mu = 5$  und Standardabweichung  $\sigma = 1$ .

- Mit welcher Wahrscheinlichkeit liegt  $X$  zwischen 3 und 4?
- Überprüfen Sie Ihr Resultat aus (a), indem Sie  $n = 10000$  solcher Zufallszahlen generieren und dann schauen, wie viele davon zwischen 3 und 4 liegen.

**Lösung:** a) Wir bekommen

$$\begin{aligned}\text{Prob}[3 \leq X \leq 4] &= \int_3^4 \varphi_{\mu=5,\sigma=1}(y) dy \\ &= \int_{-\infty}^4 \varphi_{\mu=5,\sigma=1}(y) dy - \int_{-\infty}^3 \varphi_{\mu=5,\sigma=1}(y) dy \\ &= \Phi_{\mu=5,\sigma=1}(4) - \Phi_{\mu=5,\sigma=1}(3) \\ &\stackrel{\text{R-Befehl}}{=} \text{pnorm}(4, \text{mean} = 5, \text{sd} = 1) - \text{pnorm}(3, \text{mean} = 5, \text{sd} = 1) \\ &\approx 0.1587 - 0.0228 = 0.1359\end{aligned}$$

b) → siehe `week5.txt` ■

Die Quantil-Funktion `qnorm` lässt sich in diesem Fall ebenfalls nicht mehr explizit angeben. Sie hat die folgende anschauliche Bedeutung, machen wir gleich konkret im `week5.txt`: Wir generieren etwa  $N = 10000$  (oder  $N = 10001$ , der Grund wird gleich klarer)  $(\mu, \sigma)$ -normalverteilte Zufallszahlen `z` und ordnen sie dann der Grösse nach mit dem Befehl

```
sz = sort(z)
```

so dass also

$$\text{sz}[1] < \text{sz}[2] < \text{sz}[3] < \dots < \text{sz}[9999] < \text{sz}[10000] < \text{sz}[10001]$$

Das Element `sz[5001]` hat dann die folgende Eigenschaft: Genau 5000 Zahlen sind kleiner als `sz[5001]` und 5000 Zahlen sind grösser als `sz[5001]`, deshalb heisst es dann auch der **Median** von der Stichprobe `z`. Entsprechend heissen `sz[2501]` das **erste Quartil**, genau ein viertel aller Zahlen sind kleiner als `sz[2501]` und 3/4 aller Zahlen sind grösser als `sz[2501]`, und `sz[7501]` heisst dann das **dritte Quartil**. Es gilt dann für  $1 \leq i \leq N$ :

$$\text{sz}[i] \stackrel{N \rightarrow \infty}{\approx} \text{qnorm}(i/N, \text{mean} = \mu, \text{sd} = \sigma)$$

Insbesondere gilt also:

$$\text{erstes Quartil} = \text{qnorm}(0.25, \text{mean} = \mu, \text{sd} = \sigma)$$

$$\text{Median} = \text{qnorm}(0.5, \text{mean} = \mu, \text{sd} = \sigma)$$

$$\text{drittes Quartil} = \text{qnorm}(0.75, \text{mean} = \mu, \text{sd} = \sigma)$$

Allgemein, für ein  $0 < \alpha < 1$  heisst die Zahl

$$\text{qnorm}(\alpha, \text{mean} = \mu, \text{sd} = \sigma)$$

ein  $\alpha$ -Quantil. Diese Bezeichnungen gelten nicht nur für die Normalverteilung, sondern für beliebige Verteilungen.