

De Wet van Grote Aantallen

In hoofdstuk 5 hebben wij de discrete begrippen uit hoofdstukken 1 en 2 (met aftelbare uitkomstenruimten) gegeneraliseerd naar de uitkomstenruimte \mathbb{R}^d . Hiervoor moesten we hoofdzakelijk bedenken hoe we het \sum -teken konden vervangen door de \int -teken. Bovendien is het een extra eis aan een deelverzameling $A \subseteq \mathbb{R}^d$ dat de integraal

$$\int_A f(x)dx \tag{1}$$

bestaat. In het discrete geval is iedere deelverzameling $A \subseteq \Omega$ een gebeurtenis — we kunnen immers de kans

$$P(A) = \sum_{x \in A} p(x)$$

zonder moeite berekenen — maar om soortgelijk een kans $P(A)$ aan $A \subseteq \mathbb{R}^d$ toe te kennen moet de integraal (1) wel bestaan, en alleen dan noemen we A een gebeurtenis. Meer details hierover in het college ‘maat en integratie’.

Hoofdstuk 3 ging over de stochastische wandeling. Ook hierover kunnen we d.m.v. de uitbreiding naar continue uitkomstenruimten meer uitspraken doen, zie § 6.4.

De inhoud van hoofdstuk 4 was limietstellingen. De centrale limietstelling komen we in hoofdstuk 8 weer tegen, op dit moment beperken we ons tot de wet van grote aantallen. **Opgave 5.11.2** : formuleer en bewijs de stellingen 4.1.1 en 4.1.4 voor continue toevalsvariabelen.

Stelling 5.11.2.A (Zwakke wet van grote aantallen) *Gegeven onafhankelijke en identiek verdeelde (i.i.d.) toevalsvariabelen $X_1, X_2, \dots, X_n : \mathbb{R}^d \rightarrow \mathbb{R}$ met $E(X_1) = \mu \in \mathbb{R}$ en $\text{var}(X_1) = \sigma^2 \in \mathbb{R}$, dus allebei eindig. Dan geldt voor $S_n = X_1 + \dots + X_n$*

$$\forall_{\delta > 0} \quad P\left(\left|\frac{S_n}{n} - \mu\right| > \delta\right) \leq \frac{\sigma^2}{n\delta^2} \xrightarrow{n \rightarrow \infty} 0.$$

Het bewijs is letterlijk hetzelfde als in het discrete geval, vervang alleen de verwijzing naar gevolg 2.3.23 (de ongelijkheid van Chebyshev) door een verwijzing naar opgave 5.11.1.

Lemma 5.11.2.B Zij $p > 0$ en $Z : \mathbb{R}^d \longrightarrow [0, \infty[$ een toevalsvariabele die geen negatieve waarden aanneemt. Dan is

$$E(Z^p) = \int_0^\infty pz^{p-1}P(Z > z) dz .$$

Bewijs : We beperken ons tot continue toevalsvariabelen, schrijven $f = f_Z$ voor de kansdichtheid van Z en berekenen eerst

$$P(Z > z) = \int_z^\infty f(\zeta) d\zeta = \int_{-\infty}^\infty f(\zeta)\mathbf{1}_{z \leq \zeta}(z) d\zeta$$

(want tussen $-\infty$ en z integreren we de nulfunctie). Dan volgt

$$\begin{aligned} \int_0^\infty pz^{p-1}P(Z > z) dz &= \int_0^\infty pz^{p-1} \int_{-\infty}^\infty f(\zeta)\mathbf{1}_{z \leq \zeta}(z) d\zeta dz \\ &= \int_{-\infty}^\infty f(\zeta) \int_0^\infty pz^{p-1}\mathbf{1}_{z \leq \zeta}(z) dz d\zeta \\ &= \int_{-\infty}^\infty f(\zeta) \int_0^\zeta pz^{p-1} dz d\zeta \\ &= \int_{-\infty}^\infty f(\zeta)\zeta^p d\zeta = E(Z^p) . \end{aligned}$$

□

Stelling 5.11.2.C (Algemene zwakke wet van grote aantallen) Gegeven onafhankelijke en identiek verdeelde (i.i.d.) toevalsvariabelen $X_1, X_2, \dots : \mathbb{R}^d \longrightarrow \mathbb{R}$ met $E(X_1) = \mu \in \mathbb{R}$, dus eindig. Dan is

$$\forall_{\delta > 0} \quad P\left(\left|\frac{S_n}{n} - \mu\right| > \delta\right) \xrightarrow{n \rightarrow \infty} 0 .$$

Bewijs : We houden n eerst vast en nemen de limiet pas aan het eind. Dan hebben we te maken met een eindig aantal toevalsvariabelen X_1, \dots, X_n , $S_n = X_1 + \dots + X_n$, en in het geval van eindige variantie kunnen we stelling 5.11.2.A toepassen. Het idee is dan ook om over te stappen naar toevalsvariabelen Y_1, \dots, Y_n die wél eindige variantie hebben. Dit gebeurt d.m.v. truncatie :

$$Y_k(\omega) := X_k(\omega) \cdot \mathbf{1}_{\{|X_k| \leq n\}}(\omega), \quad \text{voor } k = 1, \dots, n \quad (2)$$

en we hebben uiteraard

$$\begin{aligned}
\text{var}(Y_k) &= E(Y_k^2) - (E(Y_k))^2 \leq E(Y_k^2) \\
&\stackrel{5.11.2.B}{=} \int_0^\infty 2xP(|Y_k| > x) \, dx = \int_0^\infty 2xP(|X_k|\mathbf{1}_{\{|X_k| \leq n\}} > x) \, dx \\
&= \int_0^n 2xP(|X_k|\mathbf{1}_{\{|X_k| \leq n\}} > x) \, dx \leq n^2 < \infty,
\end{aligned}$$

maar we hebben een betere afschatting nodig. Immers, de toevalsvariabelen Y_1, \dots, Y_n hangen volgens de definitie (2) direct van n af, we schrijven dan ook liever Y_k^n , en als we hun variantie $\text{var}(Y_k^n) = \sigma_n^2$ niet beter kunnen afschatten dan gaat het quotient $\frac{\sigma_n^2}{n\delta^2}$ uit stelling 5.11.2.A naar oneindig i.p.v. naar nul. Daarom kijken we naar de integrand

$$\begin{aligned}
xP(|X_k| > x) &= xE(\mathbf{1}_{\{|X_k| > x\}}) = E(x\mathbf{1}_{\{|X_k| > x\}}) \leq E(|X_k|\mathbf{1}_{\{|X_k| > x\}}) \\
&= \int_{-\infty}^\infty yf_{|X_k|\mathbf{1}_{\{|X_k| > x\}}}(y) \, dy = \int_x^\infty yf_{|X_k|}(y) \, dy \quad (3) \\
&= E(|X_1|) - \int_{-\infty}^x yf_{|X_k|}(y) \, dy \xrightarrow{x \rightarrow \infty} 0
\end{aligned}$$

(hier gebruiken we aan het eind dat $E(|X_1|)$ eindig is). I.h.b. heeft $xP(|X_k| > x)$ een maximum M . De afchatting $\sigma_n^2 \leq 2Mn$ is al beter, maar nog niet goed genoeg. Het volgende blijkt te zullen voldoen : gegeven $\eta > 0$ bestaat er (vanwege (3)) een $\xi(\eta) \in \mathbb{R}$ met $xP(|X_k| > x) < \eta \forall x \geq \xi(\eta)$ en dus

$$\begin{aligned}
\sigma_n^2 &= \int_0^n 2xP(|X_k|\mathbf{1}_{\{|X_k| \leq n\}} > x) \, dx \leq 2 \int_0^n xP(|X_k| > x) \, dx \\
&\leq 2 \int_0^{\xi(\eta)} xP(|X_k| > x) \, dx + 2 \int_{\xi(\eta)}^n xP(|X_k| > x) \, dx \quad (4) \\
&\leq 2M\xi(\eta) + 2\eta n.
\end{aligned}$$

Om X_k en S_n volledig door Y_k^n en $T_n = Y_1^n + \dots + Y_n^n$ te kunnen vervangen, moeten we ook met $\nu_n = E(Y_k^n)$ werken. Let op dat de truncatie zodanig is gekozen dat $Y_k^n \rightarrow X_k$ als $n \rightarrow \infty$. Hierdoor is

$$\nu_n = \int_{-\infty}^n xf_{X_1}(x) \, dx \xrightarrow{n \rightarrow \infty} \int_{-\infty}^\infty xf_{X_1}(x) \, dx = \mu.$$

Verder geldt

$$\begin{aligned}
 P\left(\left|\frac{S_n}{n} - \nu_n\right| > \delta\right) &= P\left(\left|\frac{S_n}{n} - \nu_n\right| > \delta, S_n = T_n\right) + P\left(\left|\frac{S_n}{n} - \nu_n\right| > \delta, S_n \neq T_n\right) \\
 &\leq P\left(\left|\frac{T_n}{n} - \nu_n\right| > \delta\right) + P\left(\left|\frac{S_n}{n} - \nu_n\right| > \delta, S_n \neq T_n\right) \\
 &\leq \frac{\sigma_n^2}{n\delta^2} + P(S_n \neq T_n)
 \end{aligned}$$

met

$$P(S_n \neq T_n) \leq P(\exists_{k \leq n} X_k \neq Y_k^n) \leq nP(|X_k| > n) \xrightarrow{n \rightarrow \infty} 0$$

vanwege (3) en

$$\frac{\sigma_n^2}{n\delta^2} \leq \frac{2M\xi(\eta)}{n\delta^2} + \frac{2\eta}{\delta^2}$$

vanwege (4).

Om te laten zien dat deze laatste som ook daadwerkelijk naar nul gaat, gebruiken we de definitie van de limiet uit ‘inleiding analyse’. Gegeven $\varepsilon > 0$. Dan kiezen we (nu pas!) $\eta = \frac{1}{6}\delta^2\varepsilon$ en $N_\varepsilon \geq \frac{6M\xi(\eta)}{\delta^2\varepsilon}$. Hiermee volgt voor alle $n \geq N_\varepsilon$

$$\dots \leq \frac{2M\xi(\eta)}{\frac{6M\xi(\eta)}{\delta^2\varepsilon} \cdot \delta^2} + \frac{2 \cdot \frac{\delta^2\varepsilon}{6}}{\delta^2} = \frac{\varepsilon}{3} + \frac{\varepsilon}{3} < \varepsilon.$$

Voor de kleine lettertjes gebruiken we de ‘overgebleven’ $\frac{\varepsilon}{3}$ en vervangen in bovenstaande redenering δ door $\tilde{\delta} = 2\delta$, en kiezen

$$N_\varepsilon \geq \max\left\{\frac{6M\xi(\eta)}{\tilde{\delta}^2\varepsilon}, N_\delta, N_{\frac{\varepsilon}{3}}\right\}.$$

Hier is $N_{\frac{\varepsilon}{3}}$ zodanig gekozen dat

$$\forall_{n \geq N_{\frac{\varepsilon}{3}}} nP(|X_k| > n) < \frac{\varepsilon}{3}$$

en dus

$$P\left(\left|\frac{S_n}{n} - \nu_n\right| > \tilde{\delta}\right) < \varepsilon.$$

Verder is N_δ zodanig gekozen dat

$$\forall_{n \geq N_\delta} |\nu_n - \mu| < \delta$$

en dus

$$\left|\frac{S_n}{n} - \nu_n\right| > \tilde{\delta} \Rightarrow \left|\frac{S_n}{n} - \mu\right| > \delta$$

waarmee uiteindelijk

$$\forall_{n \geq N_\varepsilon} P\left(\left|\frac{S_n}{n} - \mu\right| > \delta\right) \leq \varepsilon. \quad \square$$