

# Kapitel 7

## Grænsesresultater for stokastiske variable

I dette kapitel behandles to af sandsynlighedsregningens vigtigste resultater, som begge drejer sig om, hvordan gennemsnittet af  $n$  stokastiske variable opfører sig, når  $n$  går mod uendelig.

### 7.1 De store tals lov

Store tals lov siger, at for  $n$  ukorreleerde stokastiske variable med middelværdi og varians vil gennemsnittet med stor sandsynlighed ligge nær middelværdien, når  $n$  er stor. Dette resultat vises ved hjælp af Chebychevs ulighed.

**Sætning 7.1.1 (Chebychevs ulighed).** *Lad  $X$  være en diskret eller kontinuert stokastisk variabel med middelværdi  $\mu$  og varians  $\sigma^2$ . Da gælder, for ethvert  $a > 0$ , at*

$$P(|X - \mu| \geq a) \leq \frac{\sigma^2}{a^2}. \quad (7.1.1)$$

**Bevis:** Definer for  $a > 0$  mængden

$$A = \{x \in \mathbb{R} \mid |x - \mu| \geq a\}.$$

Da er

$$(x - \mu)^2 1_A(x) \geq a^2 1_A(x)$$

for alle  $x \in \mathbb{R}$ . Dermed er

$$\begin{aligned}\sigma^2 &= E((X - \mu)^2) \\ &= E((X - \mu)^2 1_A(X)) + E((X - \mu)^2 1_{\mathbb{R} \setminus A}(X)) \\ &\geq E((X - \mu)^2 1_A(X)) \\ &\geq E(a^2 1_A(X)) = a^2 P(X \in A) = a^2 P(|X - \mu| \geq a).\end{aligned}$$

Det andet ulighedstegn følger af (5.2.6)

□

**Eksempel 7.1.2** Hvis  $X$  har middelværdi 0 og varians 1 er  $P(|X| \geq 10) \leq 0.01$ .

□

Vi kan nu uden større besvær bevise de store tals lov.

**Sætning 7.1.3** (*De store tals lov*). *Lad  $X_1, X_2, \dots$  være en følge af ukorrelerede diskrete eller kontinuerte stokastiske variable, som alle har samme middelværdi og varians. Hvis  $\mu = E(X_i)$ , gælder for ethvert  $\epsilon > 0$  at*

$$P\left(\left|\frac{X_1 + \dots + X_n}{n} - \mu\right| < \epsilon\right) \rightarrow 1 \quad (7.1.2)$$

for  $n \rightarrow \infty$ .

**Bevis:** Gennemsnittet  $\bar{X}_n = \frac{1}{n}(X_1 + \dots + X_n)$  har middelværdi  $\mu$  og varians

$$\text{Var}(\bar{X}_n) = \frac{1}{n^2} (\text{Var}(X_1) + \dots + \text{Var}(X_n)) = \frac{\sigma^2}{n},$$

hvor  $\sigma^2 = \text{Var}(X_i)$ . Ved at anvende Chebychevs ulighed (7.1.1) fås

$$\begin{aligned}1 &\geq P(|\bar{X}_n - \mu| < \epsilon) = 1 - P(|\bar{X}_n - \mu| \geq \epsilon) \\ &\geq 1 - \frac{\text{Var}(\bar{X}_n)}{\epsilon^2} = 1 - \frac{\sigma^2}{n\epsilon^2} \rightarrow 1,\end{aligned}$$

for  $n \rightarrow \infty$ .

□

Den eneste grund til at det i Sætningerne 7.1.1 og 7.1.3 kræves, at de stokastiske variable er diskrete eller kontinuerte, er at vi kun har defineret middelværdi og varians for disse to typer af fordelinger. Chebychevs ulighed og store tals lov gælder for alle typer stokastisk variabel, som har middelværdi og varians.

Sætning 7.1.3 er den enkleste version af de store tals lov. Der findes andre varianter, som vi ikke vil bevise her (se dog Opgave 7.7). F. eks. er det ikke nødvendigt at antage, at de stokastiske variable har varians. Det er nok, at  $E(|X_i|) < \infty$ , hvilket til gengæld klart nok er nødvendigt. I den situation antages uafhængighed i stedet for ukorrelerethed, men de stokastiske variable behøver faktisk ikke være ukorrelerede eller uafhængige; resultatet (7.1.2) gælder, hvis blot afhængigheden er tilstrækkeligt svag.

Vi har nu bevist, at *sandsynlighedsregningens frekvensfortolkning* holder. Lad nemlig  $Y_1, Y_2, \dots$  være en følge af uafhængige stokastiske variable med samme fordeling, og lad  $A$  være en delmængde af  $\mathbb{R}$ . Sæt  $p = P(Y_i \in A)$ . Da er  $E(1_A(Y_i)) = p$  og

$$H_n = \frac{1_A(Y_1) + \dots + 1_A(Y_n)}{n}$$

er den relative hyppighed af udfald i  $A$  blandt de  $n$  første stokastiske variable. Da  $1_A(Y_i)$  er begrænset, har den varians, så ifølge store tals lov gælder

$$P(|H_n - p| < \epsilon) \rightarrow 1$$

for  $n \rightarrow \infty$  for ethvert  $\epsilon > 0$ .

Det skal også lige nævnes, at vi her har stiftet bekendtskab med et af sandsynlighedsregningens vigtigste begreber, *konvergens i sandsynlighed*. Hvis  $Y_1, Y_2, \dots$  er en følge af stokastiske variable, som opfylder, at der findes en stokastisk variabel  $Y$ , så  $P(|Y_n - Y| < \epsilon) \rightarrow 1$  for  $n \rightarrow \infty$  for ethvert  $\epsilon > 0$ , så siger man, at følgen  $\{Y_n\}$  konvergerer i sandsynlighed mod  $Y$ . Store tals lov siger altså, at gennemsnittet af  $n$  uafhængige stokastiske variable med samme middelværdi og varians konvergerer i sandsynlighed mod den fælles middelværdi. Her er  $Y$  lig med en konstant, nemlig middelværdien.

## 7.2 Den centrale grænseværdidisætning

Betrægt en følge af uafhængige diskrete eller kontinuerte stokastiske variable  $X_1, X_2, \dots$ , som alle har samme fordeling med middelværdi  $\mu$  og varians  $\sigma^2$ . De store tals lov fortæller os da, at gennemsnittet af de  $n$  første stokastiske variable

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

er tæt på  $\mu$ , når  $n$  er stor. Gennemsnittet er jo en stokastisk variabel, så det er naturligt at spørge, hvad dets fordeling er, når  $n$  er stor. Variansen er meget lille, så det varierer kun lidt i nærheden af  $\mu$ . Imidlertid kan man studere gennemsnittets stokastiske variation ved at betragte den standardiserede variabel

$$U_n = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}. \quad (7.2.1)$$

Da  $E(\bar{X}_n) = \mu$  og  $\text{Var}(\bar{X}_n) = \sigma^2/n$ , har  $U_n$  middelværdi nul og varians en. Vi har ved at gange med  $\sqrt{n}/\sigma$  forstørret den stokastiske variation af  $\bar{X}_n$ , så vi kan studere den nøjere. Der gælder følgende resultat.

**Sætning 7.2.1** (*Den centrale grænseværdidisætning*). *Lad  $X_1, X_2, \dots$  være en følge af uafhængige, identisk fordelte kontinuerte eller diskrete stokastiske variable med middelværdi  $\mu$  og varians  $\sigma^2$ . Da vil fordelingen af  $U_n$ , defineret ved (7.2.1), konvergere mod en standard normalfordeling, i den forstand at for alle  $u \in \mathbb{R}$  vil*

$$P(U_n \leq u) \rightarrow \Phi(u), \quad (7.2.2)$$

når  $n \rightarrow \infty$ . Som sædvanlig betegner  $\Phi$  fordelingsfunktionen for standard normalfordelingen givet ved (5.3.2).

Vi har på dette kursus ikke de nødvendige tekniske hjælpemidler til at bevise den centrale grænseværdidisætning. Med de rette hjælpemidler er sætningen faktisk ikke svær at vise.

Der er mange gode grunde til, at normalfordelingen dukker op her. Hvis  $X_i$ -erne er normalfordelte, følger det af Sætning 6.3.12, at  $U_n$  er standard normalfordelt for alle  $n$ . Hvis sætningen skal være sand, må

$P(U_n \leq u)$  derfor nødvendigvis konvergere mod  $\Phi(u)$ . Bemærk, at vi kan slutte af Sætning 7.2.1, at standard normalfordelingen er den eneste fordeling med varians, som har den egenskab, at  $U_n$  har samme fordeling som  $X_i$ -erne for alle  $n$ . Det kan faktisk bevises, at den er den eneste fordeling overhovedet med denne egenskab.

Hvis vi definerer  $S_n = X_1 + \dots + X_n$ , kan resultatet (7.2.2) også skrives på formen

$$P\left(\frac{S_n - n\mu}{\sigma\sqrt{n}} \leq u\right) \rightarrow \Phi(u), \quad (7.2.3)$$

når  $n \rightarrow \infty$ . Vi kan altså også opfatte den centrale grænseværdidisætning som et udsagn om fordelingen af summen  $S_n$ , når  $n$  er stor.

Den centrale grænseværdidisætning giver en forklaring på, hvorfor det i praksis så ofte viser sig, at observationer kan antages at være normalfordelte. Vi kan f. eks. tænke os, at man skal måle en størrelse, som har værdien  $\xi$ , men at en række tilfældige forhold indvirker på målingen, så den ikke bliver helt nøjagtig. Det, der faktisk måles, kan derfor antages at være  $Y = \xi + X_1 + \dots + X_n$ , hvor  $X_i$ -erne er virkningen af de forskellige kilder til målefejl. Hvis  $X_i$ -erne opfylder betingelserne i den centrale grænseværdidisætning, kan fordelingen af målingen  $Y$  tilnærmes med en normalfordeling med middelværdi  $\xi + n\mu$  og varians  $n\sigma^2$ . Ved gode målinger er  $\mu$  lig nul eller i hvert fald meget lille.

Der findes andre udgaver af den centrale grænseværdidisætning med væsentligt svagere betingelser, som gør overstående argument for, at målefejl ofte er normalfordelte, mere overbevisende. De stokastiske variable  $X_i$  behøver således ikke at have samme fordeling; ikke engang samme middelværdi og varians. De behøver heller ikke at være uafhængige. Det væsentlige er, at de alle er små i forhold til summen  $S_n$ , og at afhængigheden mellem dem ikke er stor.

Den form for konvergens, som udtrykkes ved (7.2.2) kaldes *konvergens i fordeling*. Et andet eksempel på denne form for konvergens så vi i Sætning 4.1.2. Udforskningen af konvergens i fordeling spiller en central rolle i sandsynlighedsregningen og dens anvendelser i statistik.