

Något om kovarians och korrelation

Definition: Kovariansen mellan två stokastiska variabler X och Y definieras

$$C(X, Y) = E[(X - \mu_x)(Y - \mu_y)].$$

Variabeln

$$(X - \mu_x)(Y - \mu_y)$$

mäter om X och Y tenderar att variera åt samma eller motsatt håll. Om X är stor/liten (relativt μ_x) och Y samtidigt är stor/liten (relativt μ_y) så är kovariansen positiv eftersom $+\cdot+=+$ och $-\cdot-=+$. Analogt, om X och Y varierar åt motsatt håll är kovariansen negativ, $+\cdot-=-$ och $-\cdot+=-$. Kovariansen mäter *linjär* samvariation.

Korrelationen (korrelationskoefficienten) är det dimensionslösa linjära samvariationsmättet. Den definieras och betecknas

$$\rho_{X,Y} = \frac{C(X, Y)}{\sqrt{V(X)V(Y)}} = \frac{C(X, Y)}{D(X)D(Y)}.$$

Korrelationskoefficienten uppfyller alltid

$$-1 \leq \rho_{X,Y} \leq 1$$

och $|\rho_{X,Y}| = 1$ om $Y = aX + b$.

Kovariansen kan även räknas ut som

$$C(X, Y) = E[XY] - E[X]E[Y].$$

Notera att $V(X) = C(X, X)$ så en beräkningsformel för varianser är $V(X) = E[X^2] - (E[X])^2$.

Normalfördelningen

En stokastisk variabel X sägs vara normalfördelad med parametrar μ och $\sigma > 0$ om

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

för alla x . Läs gärna beviset i boken att detta är en giltig täthet, dvs. integrerar sig till 1. Kodbeteckning X är $N(\mu, \sigma)$.

Fördelningen är symmetrisk runt μ så alltså är parametern μ väntevärdet, $E[X] = \mu$. Man kan visa (se boken!) att $\sigma = D(X)$.

Fördelningsfunktionen för en normalfördelad stokastisk variabel X

$$F_X(t) = P(X \leq t) = \int_{-\infty}^t f_X(x) dx = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2} dx$$

finns inte på någon sluten form. Vi skall se hur vi kan klara oss med fallet $\mu = 0$, $\sigma = 1$.

Sats. Låt X_1, \dots, X_n vara en sekvens av oberoende, normalfördelade stokastiska variabler, och a_1, \dots, a_n och b konstanter. Då är

$$Y = a_1X_1 + \dots + a_nX_n + b$$

normalfördelad, Y är $N(\mu, \sigma)$, med väntevärde

$$\mu = E[Y] = E[a_1X_1 + \dots + a_nX_n + b] = a_1E[X_1] + \dots + a_nE[X_n] + b$$

och varians

$$\sigma^2 = V(Y) = a_1^2 V(X_1) + \dots + a_n^2 V(X_n).$$

Speciellt medför detta att om X är $N(\mu, \sigma)$ så har

$$Z = \frac{X - \mu}{\sigma}$$

väntevärde $E[Z] = E\left[\frac{X - \mu}{\sigma}\right] = \frac{1}{\sigma}(E[X] - \mu) = 0$ och varians $V(Z) = \frac{1}{\sigma^2}V(X) = 1$. Satsen säger att Z är normalfördelad och

$$f_Z(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2} = \varphi(x).$$

En normalfördelad s.v. med väntevärde 0 och varians (standardavvikelse) 1 sägs vara *standardnormalfördelad*. Dess täthetsfunktion betecknas av Blom med $\varphi(x)$ och dess fördelningsfunktion med $\Phi(x) = P(Z \leq x)$.

Notera att

$$F_X(x) = P(X \leq x) = P\left(\frac{X - \mu}{\sigma} \leq \frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - \mu}{\sigma}\right).$$

Så för att räkna ut fördelningsfunktionens värde i en punkt för en godtycklig normalfördelning översätter vi den till motsvarande punkt för en standardnormalfördelning. Funktionen Φ finns tabulerad i formelsamlingen tillsammans med några kvantiler.

Sats (Centrala gränsvärdessatsen (CGS)). Låt X_1, X_2, \dots vara en sekvens av oberoende, likafördelade stokastiska variabler med väntevärde μ och standardavvikelse σ . Då gäller att

$$\lim_{n \rightarrow \infty} P\left(\frac{1}{\sigma\sqrt{n}}\left(\sum_{k=1}^n X_k - n\mu\right) \leq x\right) = \Phi(x)$$

Man använder konvergensen för att säga att för stora n är

$$\frac{\sum_{k=1}^n X_k - n\mu}{\sigma\sqrt{n}}$$

approximativt $N(0, 1)$, eller

$$\sum_{k=1}^n X_k \text{ approx är } N(n\mu, \sqrt{n}\sigma).$$

Hur stort n skall vara för att approximationen skall vara bra beror på vad som avses med "bra" och fördelningen för de stokastiska variablerna. Symmetriska fördelningar konvergerar snabbare än assymetriska.

Binomialfördelning och dess släktingar

Genomgående modellsituation: Betrakta oberoende försök där vid varje försök det finns en sannolikhet p att en händelse A inträffar.

Låt X beteckna antalet gånger man får göra försöket tills man ser att A inträffar för första gången. Då är X för första gången-fördelad, skrivet X är ffg(p), om

$$P(X = k) = (1 - p)^{k-1}p, \quad k = 1, 2, 3, \dots$$

för $k = 1, 2, 3, \dots$

Sedan tidigare är $E[X] = 1/p$ och man kan visa att $V(X) = (1-p)/p^2$.

Den omvända frågeställningen är, av n gjorda försök, låt X beskriva antalet gånger A inträffar. Då är X binomialfördelad, det vill säga

$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}$$

för $k = 0, 1, \dots, n$.

Det är ofta vettigt att modellera situationen på följande sätt. Låt U_1, \dots, U_n vara oberoende, likafördelade stokastiska variabler, där

$$U_k = \begin{cases} 1 & \text{med sannolikhet } p \\ 0 & \text{med sannolikhet } 1-p \end{cases}$$

för $k = 1, \dots, n$. Att $U_k = 1$ betyder att A inträffade i försök k och $U_k = 0$ att A inte gjorde det. Notera att

$$E[U_k] = \sum_{i=0}^1 i P(U_k = i) = 0 \cdot (1-p) + 1 \cdot p = p$$

och

$$E[U_k^2] = \sum_{i=0}^1 i^2 P(U_k = i) = 0^2 \cdot (1-p) + 1^2 \cdot p = p,$$

och alltså $V(U_k) = p - p^2 = p(1-p)$.

Med $X = \sum_{k=1}^n U_k$ får vi att X beskriver antalet lyckade försök. Då är X binomialfördelad och

$$E[X] = E\left[\sum_{k=1}^n U_k\right] = \sum_{k=1}^n E[U_k] = np$$

och

$$V(X) = V\left(\sum_{k=1}^n U_k\right) = \{\text{oberoende}\} = \sum_{k=1}^n V(U_k) = np(1-p)$$

Med modellen $X = \sum_{k=1}^n U_k$ säger Centrala gränsvärdesatsen att för stora n är X approximativt normalfördelad med parametrar $\mu = np$ och $\sigma = \sqrt{np(1-p)}$. Vi kräver dock för att approximationen skall fungera väl att p inte är för liten eller för stor. Vi formulerar det gemensamma kravet på n och p som

$$np(1-p) = V(X) \geq 10.$$

Hypergeometrisk fördelning

Exempel: Opinionsundersökning (ändlig population): I en population om N individer finns det s sympatisörer, dvs. $p = s/N$ är andelen sympatisörer. Välj ut n på måfå utan återläggning och låt X beskriva antalet sympatisörer. Då är

$$P(X = k) = \frac{\binom{s}{k} \binom{N-s}{n-k}}{\binom{N}{n}} = \frac{\binom{Np}{k} \binom{N(1-p)}{n-k}}{\binom{N}{n}}, \quad \max(n - (N-s), 0) \leq k \leq \min(n, s)$$

(Notera: om urvalet sker med återläggning är X binomialfördelad, X är $\text{Bin}(n, p)$.)

$$E[X] = np \quad V(X) = np(1-p) \frac{N-n}{N-1}$$

Om populationen är stor i förhållande till urvalet kan den hypergeometrisk fördelningen approximeras av binomialfördelningen. Tumregel: $n/N < 0.10$.