

SF1901: Sannolikhetslära och statistik

Föreläsning 5.

Kovarians, korrelation, väntevärde och varians för
summor av s.v.:er, normalfördelning (del 1)

Jan Grandell & Timo Koski

15.09.2008



Väntevärdet av $g(X, Y)$

Sats

Låt (X, Y) vara en tvådimensionell s.v. Då gäller

$$E[g(X, Y)] = \begin{cases} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f_{X, Y}(x, y) dx dy & \text{för } (X, Y) \text{ kontinuerlig,} \\ \sum_{k=0}^{\infty} \sum_{j=0}^{\infty} g(k, j) f_{X, Y}(k, j) & \text{för } (X, Y) \text{ diskret.} \end{cases}$$



Väntevärdet av $g(X, Y) = X + Y$

Sats

Låt (X, Y) vara en tvådimensionell s.v. Då gäller

$$E[X + Y] = E[X] + E[Y]$$

Bevis. Låt (X, Y) vara en kontinuerlig tvådimensionell s.v.. Den föregående satsen visar med $g(x, y) = x + y$ att

$$\begin{aligned} & \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x + y) f_{X,Y}(x, y) dx dy = \\ & \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f_{X,Y}(x, y) dx dy + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f_{X,Y}(x, y) dx dy \\ & = \int_{-\infty}^{\infty} x \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy dx + \int_{-\infty}^{\infty} y \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx dy \\ & = \int_{-\infty}^{\infty} x f_X(x) dx + \int_{-\infty}^{\infty} y f_Y(y) dy = E[X] + E[Y] \end{aligned}$$

Låt (X, Y) vara en tvådimensionell s.v. där vi är intresserade av sambandet mellan X s och Y s variation. Det kan vara naturligt att betrakta variablerna

$$X - \mu_X \quad \text{och} \quad Y - \mu_Y.$$

Vi skiljer på fallen då X och Y "samvarierar" resp. "motverkar varandra", dvs. då

ett stort/litet värde på X gör ett stort/litet värde på Y troligt resp.

ett stort/litet värde på X gör ett litet/stort värde på Y troligt.



Kovariansen mellan X och Y

Betraktar vi nu variabeln

$$(X - \mu_X)(Y - \mu_Y),$$

så innebär detta att den i första fallet, eftersom $+\cdot+=+$ och $-\cdot-=+$, att den har en tendens att vara positiv. På motsvarande sätt, eftersom $-\cdot+=-$ och $+\cdot=-$, har den i andra fallet en tendens att vara negativ. Det som vi, lite slarvigt, har kallat tendens, kan vi ersätta med väntevärde. Vi leds då till följande definition.

Definition

Kovariansen mellan X och Y är

$$C(X, Y) = E[(X - \mu_X)(Y - \mu_Y)],$$

där $\mu_X = E(X)$ och $\mu_Y = E(Y)$.

Sats

Kovariansen mellan X och Y är

$$C(X, Y) = E[XY] - \mu_X \cdot \mu_Y,$$

där $\mu_X = E(X)$ och $\mu_Y = E(Y)$.

Bevis. :

$$\begin{aligned} C(X, Y) &= E[(X - \mu_X)(Y - \mu_Y)] = E[XY - X\mu_Y - \mu_X Y + \mu_X \mu_Y] \\ &= E[XY] - E[X\mu_Y] - E[\mu_X Y] - \mu_X \mu_Y \\ &= E[XY] - \mu_Y E[X] - \mu_X E[Y] - \mu_X \mu_Y = E[XY] - \mu_Y \mu_X - \mu_X \mu_Y + \mu_X \mu_Y \\ &= E[XY] - \mu_Y \mu_X \end{aligned}$$



Korrelationskoefficienten mellan X och Y

Kovariansen kan sägas ha fel sort. Det verkar rimligt att ett mått på ett så abstrakt begrepp som samvariation skall vara "sortfritt". Det vanligaste måttet är korrelationskoefficienten.

Definition

Korrelationskoefficienten mellan X och Y är

$$\rho = \rho(X, Y) = \frac{C(X, Y)}{D(X)D(Y)}.$$

Man kan visa att $|\rho| \leq 1$, där $\rho = \pm 1$ betyder att det finns ett perfekt linjärt samband, dvs. $Y = aX + b$.

Sats

Om X och Y är oberoende så är de okorrelerade, dvs. $\rho(X, Y) = 0$.

Omvändningen gäller ej, dvs. okorrelerade variabler kan vara beroende.

Exempel

Den simultana sannolikhetsfunktionen för stokastiska variablerna X och Y med $p_{X,Y}(j, k)$

X/Y	0	1	2	3
0	0.2	0	0	0
1	0	0.1	0.1	0
2	0	0.1	0.1	0
3	0	0	0	0.4

- Marginalfördelning för X :

$$p_X(0) = 0.2 + 0 + 0 + 0 = 0.2,$$

$$p_X(1) = 0 + 0.1 + 0.1 + 0 = 0.2$$

$$p_X(2) = 0 + 0.1 + 0.1 + 0 = 0.2,$$

$$p_X(3) = 0 + 0 + 0 + 0.4 = 0.4$$



- På samma sätt fås marginalfördelning för Y :

$$p_Y(0) = 0.2, p_Y(1) = 0.2$$

$$p_Y(2) = 0.2, p_Y(3) = 0.4.$$

- X och Y är INTE oberoende, ty, t.ex.,

$$p_X(0) \cdot p_Y(0) = 0.2 \cdot 0.2 = 0.04 \neq 0.2 = p_{X,Y}(0,0)$$

Exempel (forts.): BETINGADE FÖRDELNINGAR

$$p_{X|Y=k}(j) \stackrel{\text{def}}{=} \frac{p_{X,Y}(j,k)}{p_Y(k)}, j = 0, 1, 2, 3$$

$$p_{Y|X=j}(k) \stackrel{\text{def}}{=} \frac{p_{X,Y}(j,k)}{p_X(j)}, k = 0, 1, 2, 3$$

I detta exempel

$$p_{Y|X=2}(0) = \frac{p_{X,Y}(2,0)}{p_X(2)} = \frac{0}{0.2} = 0$$

$$p_{Y|X=2}(1) = \frac{p_{X,Y}(2,1)}{p_X(2)} = \frac{0.1}{0.2} = \frac{1}{2}$$

$$p_{Y|X=2}(2) = \frac{p_{X,Y}(2,2)}{p_X(2)} = \frac{0.1}{0.2} = \frac{1}{2}$$

$$p_{Y|X=2}(3) = \frac{p_{X,Y}(2,3)}{p_X(2)} = \frac{0}{0.2} = 0$$



Exempel (forts.): BETINGADE FÖRDELNINGAR

och

$$p_{X|Y=0}(0) = \frac{p_{X,Y}(0,0)}{p_Y(0)} = \frac{0.2}{0.2} = 1$$

$$p_{X|Y=0}(1) = p_{X|Y=0}(2) = p_{X|Y=0}(3) = 0$$



Utifrån definitionerna följer vidare att vi har

$$p_X(k) = \sum_{j=0}^3 p_Y(j) p_{X|Y=j}(k) = \sum_{j=0}^3 p_{X,Y}(k,j).$$

$$p_Y(j) = \sum_{k=0}^3 p_X(k) p_{Y|X=k}(j) = \sum_{k=0}^3 p_{X,Y}(k,j).$$

Exempel (forts.):

$$E(X) = 0 \cdot 0.2 + 1 \cdot 0.2 + 2 \cdot 0.2 + 3 \cdot 0.4 = 1.8.$$

samt

$$E(Y) = 1.8.$$

$$E(X^2) = 0^2 \cdot 0.2 + 1^2 \cdot 0.2 + 2^2 \cdot 0.2 + 3^2 \cdot 0.4 = 4.6.$$

$$V(X) = E(X^2) - E(X)^2 = 4.6 - 1.8^2 = 1.36 = V(Y).$$



Exempel (forts.): Kovarians och korrelation



$$E(X \cdot Y) = 0 \cdot 0 \cdot 0.2 + 1 \cdot 1 \cdot 0.1 \\ + 1 \cdot 2 \cdot 0.1 + 2 \cdot 1 \cdot 0.1 + 2 \cdot 2 \cdot 0.1 + 3 \cdot 3 \cdot 0.4 = 4.5$$

- Kovariansen

$$C(X, Y) = E(X \cdot Y) - E(X) E(Y) = 4.5 - 1.8 \cdot 1.8 = 1.26$$

- Korrelationskoefficienten

$$\rho(X, Y) = \frac{C(X, Y)}{\sqrt{V(X)}\sqrt{V(Y)}} = \frac{1.26}{\sqrt{1.36}\sqrt{1.36}} \\ = \frac{1.26}{1.36} = 0.926.$$

Vad säger detta ?



Sats

Låt (X, Y) vara en tvådimensionell s.v. Då gäller

$$(1) \quad E(aX + bY) = aE(X) + bE(Y);$$

$$(2) \quad V(aX + bY) = a^2V(X) + b^2V(Y) + 2abC(X, Y).$$

Bevis. (1) har visats ovan.

(2) fås av följande

$$\begin{aligned}V(aX + bY) &= E[(aX + bY - a\mu_X - b\mu_Y)^2] = E[(aX - a\mu_X + bY - b\mu_Y)^2] \\&= E[a^2(X - \mu_X)^2 + b^2(Y - \mu_Y)^2 + 2ab(X - \mu_X)(Y - \mu_Y)] \\&= a^2V(X) + b^2V(Y) + 2abC(X, Y).\end{aligned}$$



Följdsats

Låt X och Y vara två oberoende (okorrelerade ränder) s.v. Då gäller

$$E(X + Y) = E(X) + E(Y) \quad V(X + Y) = V(X) + V(Y)$$

$$E(X - Y) = E(X) - E(Y) \quad V(X - Y) = V(X) + V(Y).$$

Detta går att utvidga till godtyckligt många variabler:

Sats

Låt X_1, \dots, X_n vara oberoende (okorrelerade räkcker) s.v. och sätt

$$Y = c_1 X_1 + \dots + c_n X_n.$$

Då gäller

$$E(Y) = c_1 E(X_1) + \dots + c_n E(X_n)$$

och

$$V(Y) = c_1^2 V(X_1) + \dots + c_n^2 V(X_n)$$

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

Sats

Låt X_1, X_2, \dots, X_n vara oberoende och likafördelade s.v. med väntevärde μ och standardavvikelse σ . Då gäller att

$$E(\bar{X}) = \mu, \quad V(\bar{X}) = \frac{\sigma^2}{n} \quad \text{och} \quad D(\bar{X}) = \frac{\sigma}{\sqrt{n}}.$$

Uttrycket " X_1, X_2, \dots, X_n är likafördelade" betyder att de stokastiska variablernas fördelningar, dvs. att de stokastiska variablernas *statistiska egenskaper*, är identiska. *Utfallen* av variablerna varierar dock.



Sats

(Tjebysjovs olikhet)

För varje $\varepsilon > 0$ gäller

$$P(|X - \mu| > \varepsilon) \leq \frac{V(X)}{\varepsilon^2}.$$

(Ersätter vi ε med $k\sigma$ fås formuleringen i Blom m.fl.)

Bevis. Detta är den enda riktigt djupa satsen i kursen som vi kan bevisa.

Vi nöjer oss med det kontinuerliga fallet. Vi har

$$\begin{aligned} V(X) &= \int_{-\infty}^{\infty} (x - \mu)^2 f_X(x) dx \geq \int_{|x - \mu| > \varepsilon} (x - \mu)^2 f_X(x) dx \\ &\geq \varepsilon^2 \int_{|x - \mu| > \varepsilon} f_X(x) dx = \varepsilon^2 P(|X - \mu| > \varepsilon). \end{aligned}$$



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

Sats

Stora talens lag *För varje $\varepsilon > 0$ gäller*

$$P(|\bar{X} - \mu| > \varepsilon) \rightarrow 0 \quad \text{då } n \rightarrow \infty.$$

Bevis. Enl. Tjebysjovs olikhet gäller

$$P(|\bar{X} - \mu| > \varepsilon) \leq \frac{V(\bar{X})}{\varepsilon^2} = \frac{\sigma^2}{n\varepsilon^2} \rightarrow 0$$

då $n \rightarrow \infty$.

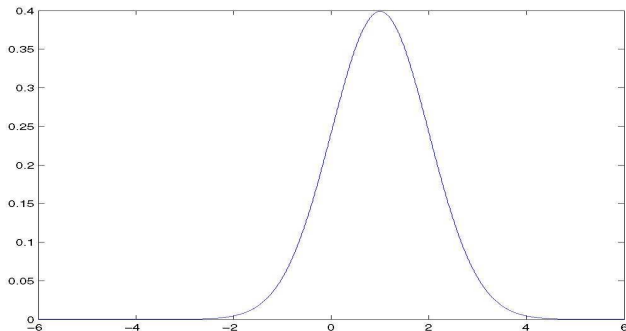


Normalfördelning

Vi gör en mer detaljerad studie av $X \in N(\mu, \sigma)$ med

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

där μ godtycklig konstant och $\sigma > 0$. I figuren för $f_X(x)$ har vi $\mu = 1, \sigma = 1$



Normalfördelning (även känd som Gaussfördelning efter C.F. Gauss, 1777-1855)



Definition

En s.v. Z säges vara standardiserad normalfördelad om den är $N(0, 1)$ -fördelad, dvs. om den har täthetsfunktionen

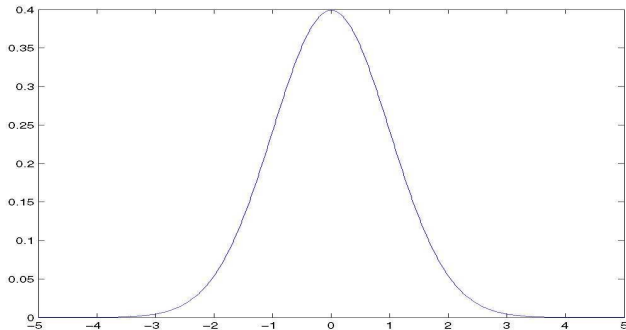
$$\varphi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}.$$

Dess fördelningsfunktion betecknas med $\Phi(z)$, dvs.

$$\Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx.$$

Standardiserad normalfördelning, $N(0, 1)$

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



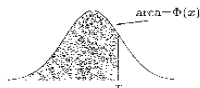
Ett problem är att fördelningsfunktionen inte kan ges på en analytisk form. Det är dock lätt att numeriskt beräkna fördelningsfunktionen och vi använder programvara för beräkning av $\Phi(x)$. (T.ex. normcdf.m i MATLAB Statistics Toolbox.) I våra tentor använder vi dock en tabell över $\Phi(x)$, som återfinns i kursens formelsamling.

Tabellen för $\Phi(x)$ ur kursens formelsamling

Tabell 1. Standard normalfördelning.

$\Phi(x) = P\{X \leq x\}$, där $X \in N(0, 1)$.

För negativa x , utnyttja att $\Phi(-x) = 1 - \Phi(x)$



x	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
0.1	.5398	.5438	.5478	.5517	.5557	.5596	.5636	.5675	.5714	.5753
0.2	.5793	.5832	.5871	.5910	.5948	.5987	.6026	.6064	.6103	.6141
0.3	.6179	.6217	.6255	.6293	.6331	.6368	.6406	.6443	.6480	.6517
0.4	.6554	.6591	.6628	.6664	.6700	.6736	.6772	.6808	.6844	.6879
0.5	.6915	.6950	.6985	.7019	.7054	.7088	.7123	.7157	.7190	.7224
0.6	.7257	.7291	.7324	.7357	.7389	.7422	.7454	.7486	.7517	.7549
0.7	.7580	.7611	.7642	.7673	.7704	.7734	.7764	.7794	.7823	.7852
0.8	.7881	.7910	.7939	.7967	.7995	.8023	.8051	.8078	.8106	.8133
0.9	.8159	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.8461	.8485	.8508	.8531	.8554	.8577	.8599	.8621

$$\Phi(-z) = 1 - \Phi(z).$$

Vi observerar att $\varphi(-z) = \varphi(z)$. $\Phi(z)$ är tabulerad i kursens formelsamling endast för $z \geq 0$. Vi har dock

$$\begin{aligned}\Phi(-z) &= \int_{-\infty}^{-z} \varphi(x) dx = [y = -x] = - \int_{\infty}^z \varphi(-y) dy \\ &= \int_z^{\infty} \varphi(y) dy = 1 - \Phi(z).\end{aligned}$$

Sats

$$\Phi(-z) = 1 - \Phi(z).$$



$$Z \in N(0,1), E(Z), V(Z)$$

Om Z är $N(0,1)$ -fördelad, så kan man visa att

$$E(Z) = 0 \quad (\text{ty } \varphi(-z) = \varphi(z))$$

$$V(Z) = 1.$$



När vi kommer till statistikdelen av kursen behöver vi ofta lösa ekvationer av följande slag:

Bestäm z så att vi för givet α har

- $P(Z \leq z) = 1 - \alpha;$
- $P(Z > z) = 1 - \alpha;$
- $P(-z < Z \leq z) = 1 - \alpha.$

För att lösa sådana ekvationer inför vi α -kvantilen λ_α definierad av $P(Z > \lambda_\alpha) = \alpha$ eller

$$\alpha = 1 - \Phi(\lambda_\alpha).$$



$$\alpha = 1 - \Phi(\lambda_\alpha).$$

Det är då bra att observera att

$$1 - \alpha = 1 - \Phi(\lambda_{1-\alpha})$$

$$\Leftrightarrow$$

$$\alpha = \Phi(\lambda_{1-\alpha})$$

$$\Leftrightarrow$$

$$\alpha = 1 - \Phi(-\lambda_{1-\alpha}),$$

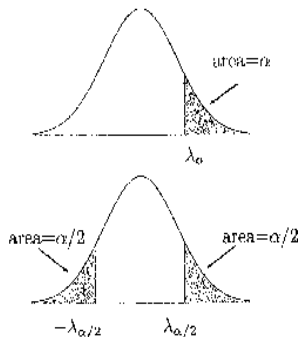
vilket ger

$$\lambda_{1-\alpha} = -\lambda_\alpha.$$

Tab 2. Normalfördelningens kvantiler

$P(X > \lambda_\alpha) = \alpha$ där $X \in N(0, 1)$

α	λ_α	α	λ_α
0.10	1.2816	0.001	3.0902
0.05	1.6449	0.0005	3.2905
0.025	1.9600	0.0001	3.7190
0.010	2.3263	0.00005	3.8906
0.005	2.5758	0.00001	4.2649



Definition

En s.v. X säges vara $N(\mu, \sigma)$ -fördelad, där μ reell och $\sigma > 0$, om

$$Z = \frac{X - \mu}{\sigma} \text{ är } N(0, 1)\text{-fördelad.}$$

Sats

Låt X vara $N(\mu, \sigma)$ -fördelad. Då gäller

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

och

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

Bevis. Vi har

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

Derivation ger $f_X(x) = \frac{1}{\sigma} \varphi\left(\frac{x - \mu}{\sigma}\right)$. □

Beviset ovan innehåller en viktig räkneregel. Om X är $N(\mu, \sigma)$ -fördelad, så gäller det att

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

Man kan m.a.o. använda tabellen för $\Phi(x)$ även för att beräkna $F_X(x)$ för $X \in N(\mu, \sigma)$.

Sats

Om X är $N(\mu, \sigma)$ -fördelad så gäller

$$E(X) = \mu \quad \text{och} \quad V(X) = \sigma^2.$$

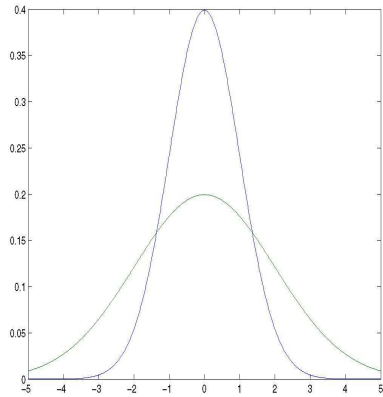
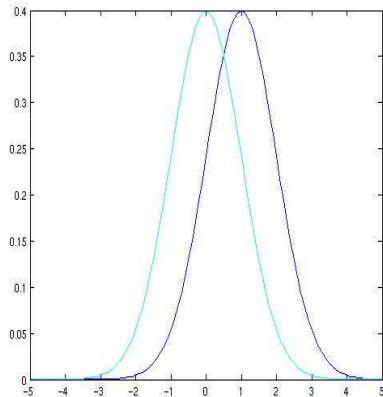
Bevis. Vi ska nu se hur listig definitionen är!

$$X = \sigma Z + \mu$$

$$E(X) = \sigma E(Z) + \mu = 0 + \mu = \mu$$

$$V(X) = \sigma^2 V(Z) + 0 = \sigma^2.$$

Täthetsfunktionerna för $N(0,1)$ och $N(1,1)$ och $N(0,1)$ och $N(0,2)$ (från vänster till höger)



Sats

Låt X vara $N(\mu, \sigma)$ -fördelad och sätt $Y = aX + b$. Då gäller det att

Y är $N(a\mu + b, |a|\sigma)$ -fördelad.

Bevis. Från definitionen följer att $X = \mu + \sigma Z$ där Z är $N(0, 1)$ -fördelad. Detta ger

$$Y = aX + b = a(\mu + \sigma Z) + b = a\mu + b + a\sigma Z$$

$$\frac{Y - (a\mu + b)}{a\sigma} = Z.$$

Om $a > 0$ följer satsen. Om $a < 0$ utnyttjar vi att Z och $-Z$ har samma fördelning.



Hur vet vi att $\varphi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$ är en sannolikhetstäthet ? Dvs. varför gäller det att

$$\int_{-\infty}^{\infty} \varphi(z) dz = 1.$$

Svaret ges t.ex. i Eike Petermann: *Analytiska metoder II*, Studentlitteratur 2002, sid. 235, Ex. 9.14, Anmärkning 9.6 eller bilagan nedan

Bilaga (ur Eike Petermann: *Analytiska metoder II*) :

$$\int_{-\infty}^{\infty} \varphi(z) dz = 1.$$

Exempel 9.14:

$$\begin{aligned} \iint_{\mathbb{R}^2} e^{-(x^2+y^2)} dx dy &= \left[\begin{array}{l} x = r \cos v, \quad 0 \leq r \\ y = r \sin v, \quad 0 \leq v < 2\pi, \quad \det \left(\frac{d(x,y)}{d(r,v)} \right) = r \end{array} \right] = \\ &= \iint_{\substack{0 \leq r \\ 0 \leq v < 2\pi}} e^{-r^2} r dr dv = \int_0^{2\pi} \left[\int_0^{\infty} e^{-r^2} r dr \right] dv = \int_0^{\infty} e^{-r^2} r dr \cdot \int_0^{2\pi} dv = \\ &= \left[-\frac{1}{2} e^{-r^2} \right]_0^{\infty} \cdot \left[v \right]_0^{2\pi} = \frac{1}{2} \cdot 2\pi = \pi. \quad \blacksquare \end{aligned}$$



$$\int_{-\infty}^{\infty} \varphi(z) dz = 1.$$

9.6

Anmärkning 9.6:

Använder man sats 9.11 utan föregående substitution på integralen i exemplet ovan, så får man:

$$\begin{aligned} \iint_{\mathbb{R}^2} e^{-(x^2+y^2)} dx dy &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} e^{-x^2-y^2} dx \right] dy = \int_{-\infty}^{\infty} e^{-x^2} dx \cdot \int_{-\infty}^{\infty} e^{-y^2} dy = \\ &= \left(\int_{-\infty}^{\infty} e^{-x^2} dx \right)^2 \end{aligned}$$

Eftersom dubbelintegralens värde beräknades till π , får man det anmärkningsvärda resultatet:

$$\int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi}$$

Trots att e^{-x^2} saknar elementär primitiv funktion²⁰ så kan alltså den speciella integralen $\int_{-\infty}^{\infty} e^{-x^2} dx$ beräknas exakt!

$$\int_{-\infty}^{\infty} \varphi(z) dz = 1.$$

Sambandet är av intresse bland annat för normalfördelningen i sannolikhetsläran (jämför exempel 9.12): Konstanten k i frekvensfunktionen

$$f(x) = k \cdot e^{-x^2/2\sigma^2}$$

måste väljas så att

$$\int_{-\infty}^{\infty} k \cdot e^{-x^2/2\sigma^2} dx = 1$$

Efter variabelsubstitutionen $x = \sigma\sqrt{2}t$ får man

$$\frac{1}{k} = \sigma\sqrt{2} \int_{-\infty}^{\infty} e^{-t^2} dt = \sigma\sqrt{2\pi}$$

dvs.

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