

SF1901: Sannolikhetslära och statistik  
Föreläsning 6.  
Normalfördelning (del 2), Centrala gränsvärdessatsen  
(CGS), Approximationer

Jan Grandell & Timo Koski

17.09.2008



$X \in N(0,1)$  har tätheten

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

$X \in N(\mu, \sigma)$  har tätheten

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

där  $\mu$  godtycklig konstant och  $\sigma > 0$ .

## Sats

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

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

$X \in N(\mu, \sigma)$  där  $\mu$  reell och  $\sigma > 0$ ,

$$Z = \frac{X - \mu}{\sigma} \in N(0, 1)$$

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

## 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.

## Extra: Lognormalfördelning

Antag att  $Z$  är en positiv stokastisk variabel, dvs.  $P(Z > 0) = 1$ . Vi sätter  $Y = \ln Z$  och ANTAR att  $Y \in N(\mu, \sigma)$ . Vad är fördelningen för  $Z$ ?  
Tag  $z > 0$ .

$$\begin{aligned} F_Z(z) &= P(Z \leq z) = P(\ln Z \leq \ln z) \\ &= P(Y \leq \ln z) = P\left(\frac{Y - \mu}{\sigma} \leq \frac{\ln z - \mu}{\sigma}\right) \\ &= \Phi\left(\frac{\ln z - \mu}{\sigma}\right) \end{aligned}$$

Derivering ger

$$\begin{aligned} f_Z(z) &= \frac{d}{dz} F_Z(z) = \varphi\left(\frac{\ln z - \mu}{\sigma}\right) \cdot \frac{1}{\sigma z} = \\ &= \frac{1}{\sigma z \sqrt{2\pi}} e^{-(\ln z - \mu)^2 / 2\sigma^2} \end{aligned}$$



# Nytt: Lognormalfördelning

$Z$  är en positiv stokastisk variabel, d.v.s.  $P(Z > 0) = 1$  och  
 $Y = \ln Z \in N(\mu, \sigma) \Rightarrow$

$$f_Z(z) = \begin{cases} \frac{1}{\sigma z \sqrt{2\pi}} e^{-(\ln z - \mu)^2 / 2\sigma^2} & \text{för } z \geq 0, \\ 0 & \text{för } z < 0. \end{cases}$$

Vi säger att  $Z$  är lognormalfördelad,  $Z \in LN(\mu, \sigma)$ . Hur kan man motivera denna fördelning? (Svaret ges senare i denna föreläsning).



$Z \in LN(\mu, \sigma)$ . Observera att  $\mu$  och  $\sigma$  INTE är väntevärde och standardavvikelse för  $Z$ . Man kan i själva verket härleda att

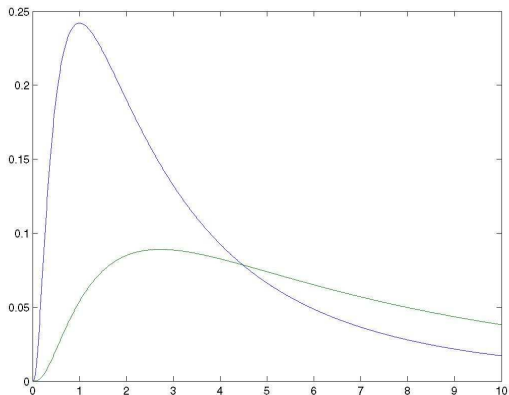
$$E(Z) = e^{\mu + \sigma^2/2}, V(Z) = e^{2\mu + \sigma^2} \cdot (e^{\sigma^2} - 1)$$



# Lognormalfördelning

$$Z \in LN(1,1) \Rightarrow E(Z) = 4.48, V(Z) = 34.5$$

$$Z \in LN(2,1) \Rightarrow E(Z) = 12.18, V(Z) = 255.0$$



## Sats

Om  $X$  är  $N(\mu_X, \sigma_X)$ -fördelad,  $Y$  är  $N(\mu_Y, \sigma_Y)$ -fördelad och  $X$  och  $Y$  är oberoende så gäller att

$$X + Y \text{ är } N\left(\mu_X + \mu_Y, \sqrt{\sigma_X^2 + \sigma_Y^2}\right) \text{-fördelad}$$

och

$$X - Y \text{ är } N\left(\mu_X - \mu_Y, \sqrt{\sigma_X^2 + \sigma_Y^2}\right) \text{-fördelad.}$$

Återkalla i minnet att Låt  $X$  och  $Y$  vara två oberoende (okorrelerade ränder) s.v. Då gäller

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

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

Det nya är att vi kan ge fördelningen för summan av oberoende normalfördelade variabler

## Sats

Låt  $X_1, \dots, X_n$  vara oberoende och  $N(\mu_1, \sigma_1), \dots, N(\mu_n, \sigma_n)$ . Då gäller att

$$\sum_{k=1}^n c_k X_k \text{ är } N\left(\sum_{k=1}^n c_k \mu_k, \sqrt{\sum_{k=1}^n c_k^2 \sigma_k^2}\right)\text{-fördelad.}$$

*Allmän regel:* Linjärkombinationer av oberoende normalfördelade stokastiska variabler är normalfördelade med rätt väntevärde och rätt standardavvikelse.

## Sats

Låt  $X_1, \dots, X_n$  vara oberoende (okorrelerade räcker) 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

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

Det nya är att vi för linjärkombinationer av oberoende normalfördelade stokastiska variabler kan ge hela fördelningen.

## Följdsats

Låt  $X_1, X_2, \dots, X_n$  vara oberoende och  $N(\mu, \sigma)$ -fördelade s.v. Då gäller att

$$\bar{X} \text{ är } N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)\text{-fördelad.}$$

Återkalla i minnet:

## Sats

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

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

# Centrala gränsvärdessatsen

Vi har sett några exempel på att normalfördelningen har trevliga statistiska egenskaper. Detta skulle vi inte ha så stor glädje av, om normalfördelningen inte dessutom var vanligt förekommande. Centrala gränsvärdessatsen CGS, som är den huvudsakliga motiveringen för normalfördelningen, kan utan vidare sägas vara ett av sannolikhetssteorins och statistikens allra viktigaste resultat.

## Sats

**(CGS)** Låt  $X_1, X_2, \dots$  vara oberoende och lika fördelade s.v. med väntevärde  $\mu$  och standardavvikelse  $\sigma$ . Då gäller att

$$P\left(\frac{\sum_{i=1}^n X_i - n\mu}{\sigma\sqrt{n}} \leq x\right) \rightarrow \Phi(x) \quad \text{då } n \rightarrow \infty.$$

Ofta uttrycker man slutsatsen i CGS som att

$$\frac{\sum_{i=1}^n X_i - n\mu}{\sigma\sqrt{n}}$$
 är approximativt  $N(0, 1)$ -fördelad

eller att

$$\sum_{i=1}^n X_i$$
 är approximativt  $N(n\mu, \sigma\sqrt{n})$ -fördelad.



En, för statistiken mycket vanlig användning av CGS är följande:

### Följdsats

Låt  $X_1, X_2, \dots$  vara oberoende och lika fördelade s.v. med väntevärde  $\mu$  och standardavvikelse  $\sigma$ . Då gäller att

$$P(a < \bar{X} \leq b) \approx \Phi\left(\frac{b - \mu}{\sigma/\sqrt{n}}\right) - \Phi\left(\frac{a - \mu}{\sigma/\sqrt{n}}\right)$$

om  $n$  är tillräckligt stort.

## Sats

Låt  $X_1, X_2, \dots, X_n$  vara oberoende och likafördelade s.v. och positiva dvs. för all  $i$  har vi  $P(X_i > 0) = 1$ . Om  $n$  är tillräckligt stort, så har

$$Y = X_1 \cdot X_2 \cdots X_n$$

*approximativt en lognormalfördelning.*

## Bevis.

$$\ln Y = \ln X_1 + \ln X_2 + \dots + \ln X_n$$

Således är  $\ln Y$  en summa av oberoende likafördelade s.v. er (dvs.  $\ln X_i$ na) och därmed enligt CGS approximativt normalfördelad om  $n$  är tillräckligt stort, och därför är  $Y$  approximativt lognormalfördelad, om  $n$  är tillräckligt stort.



Det är tyvärr inte möjligt att ge några generella och enkla tumregler om hur stort  $n$  måste vara för att normalapproximationen ska vara användbar. Detta beror på hur "normalliknande" de enskilda variablerna  $X_k$  är. Om  $X_k$ na är normalfördelade så "gäller" ju CGS för alla  $n$ . En tumregel är att om  $X_k$ na är någorlunda symmetriskt fördelade så räcker ganska små  $n$ , säg något tiotal. Om  $X_k$ na är påtagligt skevt fördelade så behöver  $n$  var något eller i värsta fall några hundratal.

Det är svårt att formulera strikt, men det räcker i CGS att  $X_k$ na är någorlunda oberoende och någorlunda lika fördelade. Med "någorlunda lika fördelade" menas framförallt att det inte finns vissa  $X_k$  som är mycket dominerande. Detta innebär att mätfel i välgjorda försök kan anses vara approximativt normalfördelade. I mindre välgjorda försök kan det däremot mycket väl finnas någon dominerande felkälla som inte alls behöver vara approximativt normalfördelad.

# Binomialfördelningen

Antag att vi gör ett försök där en händelse  $A$ , med sannolikheten  $p = P(A)$ , kan inträffa. Vi upprepar försöken  $n$  gånger, där försöken är oberoende. Sätt

$X$  = antalet gånger som  $A$  inträffar i de  $n$  försöken.

Vi säger då att  $X$  är binomialfördelad med parametrarna  $n$  och  $p$ , eller kortare att  $X$  är  $\text{Bin}(n, p)$ -fördelad.

Vi har

$$p_X(k) = \binom{n}{k} p^k q^{n-k}, \text{ för } k = 0, \dots, n,$$

där  $q = 1 - p$ .



$$X \in \text{Bin}(n, p), X = U_1 + \dots + U_n$$

Låt  $U_1, \dots, U_n$  vara s.v. definierade av

$$U_i = \begin{cases} 0 & \text{om } A^* \text{ inträffar i försök nummer } i, \\ 1 & \text{om } A \text{ inträffar i försök nummer } i. \end{cases}$$

Lite eftertanke ger att  $U_1, \dots, U_n$  är oberoende och att

$$X = U_1 + \dots + U_n.$$



$$X \in \text{Bin}(n, p), X = U_1 + \dots + U_n$$

Då

$$E(U_i) = 0 \cdot (1 - p) + 1 \cdot p = p$$

och

$$V(U_1) = E(U_i^2) - E(U_i)^2 = E(U_i) - E(U_i)^2 = p - p^2 = p(1 - p)$$

så följer

$$E(X) = nE(U_i) = np \quad \text{och} \quad V(X) = nV(U_i) = npq.$$



# Bin( $n, p$ ) approximativt $N(np, \sqrt{npq})$ -fördelad

Av  $X$ s representation som en summa följer att CGS kan tillämpas.

## Sats

Om  $X$  är Bin( $n, p$ )-fördelad med  $npq \geq 10$  så är  $X$  approximativt  $N(np, \sqrt{npq})$ -fördelad.

Detta innebär att

$$\left. \begin{array}{l} P(X \leq k) \\ P(X < k) \end{array} \right\} \approx \Phi\left(\frac{k - np}{\sqrt{npq}}\right).$$

Med *halvkorrektion* menas att vi använder följande approximation:

$$P(X \leq k) \approx \Phi\left(\frac{k + \frac{1}{2} - np}{\sqrt{npq}}\right),$$

$$P(X < k) \approx \Phi\left(\frac{k - \frac{1}{2} - np}{\sqrt{npq}}\right).$$





Vi påminner om urnmodellerna. Vi hade en urna med kulor av två slag:  $v$  vita och  $s$  svarta. Vi drog  $n$  kulor ur urnan slumpmässigt.

Sätt  $A =$  "Man får  $k$  vita kulor i urvalet".

*Dragning utan återläggning:*

$$P(A) = \frac{\binom{v}{k} \binom{s}{n-k}}{\binom{v+s}{n}}.$$

*Dragning med återläggning:*

$$P(A) = \binom{n}{k} \left(\frac{v}{v+s}\right)^k \left(\frac{s}{v+s}\right)^{n-k}.$$

# Hypergeometrisk fördelning

Antag att vi har  $N$  enheter, där proportionen  $p$ , dvs  $Np$  stycken, har egenskapen  $A$ . Drag ett slumpmässigt urval om  $n$  stycken enheter. Sätt  $X =$  antalet enheter i urvalet med egenskapen  $A$ .

I termer av urnmodellen för dragning utan återläggning gäller  $Np = v$  och  $N(1 - p) = s$  om  $A =$  "vit kula". Således fås

$$p_X(k) = P(X = k) = \frac{\binom{Np}{k} \binom{N(1-p)}{n-k}}{\binom{N}{n}},$$

för  $0 \leq k \leq Np$  och  $0 \leq n - k \leq N(1 - p)$ .

Man säger att  $X$  är  $\text{Hyp}(N, n, p)$ -fördelad.

# Hypergeometrisk fördelning

$X$  är  $\text{Hyp}(N, n, p)$ -fördelad, man kan visa att

$$E(X) = np \quad \text{och} \quad V(X) = \frac{N-n}{N-1} np(1-p).$$



# Approximationer, Hyp( $N, n, p$ )

Om  $n/N$  är någolunda liten, så verkar det troligt att det inte spelar så stor roll om vi drar med återläggning eller ej.

Vi har

$$\begin{aligned}\frac{\binom{Np}{k} \binom{N(1-p)}{n-k}}{\binom{N}{n}} &= \frac{Np!}{k!(Np-k)!} \frac{N(1-p)!}{(n-k)![N(1-p)-(n-k)]!} \frac{n!(N-n)!}{N!} \\ &= \frac{n!}{k!(n-k)!} \frac{Np!(N(1-p)!(N-n)!}{(Np-k)![N(1-p)-(n-k)]!N!} \\ &\approx \frac{n!}{k!(n-k)!} \frac{(Np)^k (N(1-p))^{n-k}}{N^n} = \binom{n}{k} p^k q^{n-k}.\end{aligned}$$

# Approximationer, Hyp( $N, n, p$ )

## Sats

Om  $X$  är Hyp( $N, n, p$ )-fördelad med  $n/N \leq 0.1$  så är  $X$  approximativt Bin( $n, p$ )-fördelad.

## Sats

Om  $X$  är Bin( $n, p$ )-fördelad med  $npq \geq 10$  så är  $X$  approximativt  $N(np, \sqrt{npq})$ -fördelad.

Detta innebär att

$$\left. \begin{array}{l} P(X \leq k) \\ P(X < k) \end{array} \right\} \approx \Phi\left(\frac{k - np}{\sqrt{npq}}\right).$$



# Approximation av $\text{Hyp}(N, n, p)$

Av detta följer att  $\text{Hyp}(N, n, p) \approx N(np, \sqrt{npq})$  om  $n/N \leq 0.1$  och  $npq \geq 10$ . Det räcker dock att kräva  $\frac{N-n}{N-1}np(1-p) \geq 10$ .



## Definition

En diskret s.v.  $X$  säges vara Poissonfördelad med parameter  $\mu$ ,  $Po(\mu)$ -fördelad, om

$$p_X(k) = \frac{\mu^k}{k!} e^{-\mu}, \text{ för } k = 0, 1, \dots$$

Vi påminner om att om  $X$  är  $Po(\mu)$ -fördelad, så gäller

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

# Summan av två oberoende Poissonfördelade

Poissonfördelningen är den viktigaste diskreta fördelningen, och har t.ex. följande trevliga egenskaper.

## Sats

Om  $X$  och  $Y$  vara oberoende  $Po(\mu_X)$ - resp.  $Po(\mu_Y)$ -fördelade s.v. Då gäller att  $X + Y$  är  $Po(\mu_X + \mu_Y)$ -fördelad.

## Bevis.

$$\begin{aligned}P(X + Y = k) &= \sum_{i=0}^k P(X = i)P(Y = k - i) = \sum_{i=0}^k \frac{\mu_X^i}{i!} e^{-\mu_X} \frac{\mu_Y^{(k-i)}}{(k-i)!} e^{-\mu_Y} \\&= e^{-(\mu_X + \mu_Y)} \sum_{i=0}^k \frac{\mu_X^i \mu_Y^{(k-i)}}{i!(k-i)!} \\&= e^{-(\mu_X + \mu_Y)} \frac{(\mu_X + \mu_Y)^k}{k!} \underbrace{\sum_{i=0}^k \binom{k}{i} \left(\frac{\mu_X}{\mu_X + \mu_Y}\right)^i \left(\frac{\mu_Y}{\mu_X + \mu_Y}\right)^{(k-i)}}_{= 1, \text{ enl. Binomialsatsen.}}\end{aligned}$$



# Approximation med $Po(np)$

Vi införde ju Poissonfördelningen som en approximation av binomialfördelningen. Detta kan vi formalisera till följande sats.

## Sats

*Om  $X$  är  $\text{Bin}(n, p)$ -fördelad med  $p \leq 0.1$  så är  $X$  approximativt  $Po(np)$ -fördelad.*

I vår approximation antog vi även att  $n$  var stor. Detta är inte nödvändigt, men vårt enkla resonemang fungerar inte utan denna extra förutsättning. Man kan visa att om  $X$  är  $\text{Bin}(n, p)$  och  $Y$  är  $Po(np)$  så gäller att

$$|P(X = k) - P(Y = k)| \leq np^2.$$



# Approximation med $Po(np)$ , ett exempel

I en datalänk växlar de binära siffrorna eller *bitarna* (0 och 1) polaritet dvs. 0 övergår i 1 eller omvänt med sannolikheten  $10^{-7}$ . En viss bit i en bitgrupp växlar polaritet oberoende av alla de andra bitarna. Protokollet i datalänken kontrollerar en bitgrupp som omfattar 8000000 ( $\approx 1$  Mbyte) bitar.

Felkontrollen upptäcker om fem eller flera bitar fått sina polariteter omkastade i bitgruppen. Om fem eller flera polaritetsväxlingar har upptäckts, begärs av protokollet en upprepad överföring av denna bitgrupp. Vi säger att *en överföring lyckas* om en upprepad överföring av bitgruppen *inte* kommer att begäras. De olika överföringarna av en och samma bitgrupp antas vara oberoende av varandra.



# Approximation med $Po(np)$ , ett exempel: forts.

Betrakta den stokastiska variabeln

$X$  = antalet överföringar av en bitgrupp om 8000000 bitar som behövs för att överföringen skall lyckas för första gången, den lyckade överföringen medräknad.

Sökt: väntevärdet för  $X$ .



## Exempel: forts.

Låt  $p_s$  beteckna sannolikheten för att en ny överföring av en bitgrupp om 8000000 bitar INTE begärs av protokollet dvs.  $p_s =$  sannolikheten för högst fyra växlingar av polaritet.

Den stokastiska variabeln  $X$  har *ffg* ( $p_s$ ) - fördelningen eller "För-första-gången"-fördelningen p.g.a. att de olika överföringarna är oberoende och p.g.a. att den lyckade överföringen medräknas. Dvs.

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

En enkel kalkyl ger väntevärdet  $E(X) = \frac{1}{p_s}$  (se kursens formelsamling). Därmed behöver vi  $p_s$ .



## Exempel: forts.

Låt  $Z$  vara antalet bitar som får sin polaritet omkastad, när en bitgrupp om 8000000 bitar överförs. Då fås

$$p_s = P(Z \leq 4).$$

P.g.a. att bitarna i en bitgrupp växlar polaritet oberoende av varandra är vi ledda till den statistiska modellen  $Z \in \text{Bin}(8000000, 10^{-7})$ , dvs. att  $Z$  är binomialfördelad med parametrarna  $n = 8000000$  och  $p = 10^{-7}$ .



## Exempel: forts.

Vi beräknar  $p_s$  med en approximation av binomialfördelningen  $\text{Bin}(8000000, 10^{-7})$  med Poissonfördelningen med parametern  $n \cdot p = 0.8$ . Denna approximation är rimlig, ty  $p = 10^{-7} \leq 0.1$ . Vi har alltså att  $Z$  är approximativt  $\text{Po}(0.8)$ -fördelad. Tabellsamlingen ger med Poissonfördelningen  $\text{Po}(0.8)$  approximationen  $p_s = P(Z \leq 4) \approx 0.99859$ .  
SVAR:  $E(X) \approx \frac{1}{0.99859} \approx 1.0014$ .

## Exempel: uppskattningen ovan

$X \in \text{Bin}(8000000, 10^{-7})$  och  $Y$  är  $\text{Po}(0.8)$  så gäller att

$$|P(X = k) - P(Y = k)| \leq np^2.$$

$$|P(X = k) - P(Y = k)| \leq 8000000 \cdot 10^{-14} = 8.0000e - 08$$



# Exempel: check mot MATLAB Statistics Toolbox

```
>> help binocdf
```

BINOCDF Binomial cumulative distribution function.

$Y = \text{BINOCDF}(X, N, P)$  returns the binomial cumulative distribution function with parameters  $N$  and  $P$  at the values in  $X$ .

The size of  $Y$  is the common size of the input arguments. A scalar input functions as a constant matrix of the same size as the other inputs.

The algorithm uses the cumulative sums of the binomial masses.

See also `binofit`, `binoinv`, `binopdf`, `binornd`, `binostat`, `cdf`.

```
>> binocdf(4, 8000000, 10-7)
```

```
ans = 0.9986
```





# Approximation av $Po(\mu)$

Om bägge villkoren  $p \leq 0.1$  och  $npq \geq 10$  är uppfyllda kan vi välja om vi vill Poissonapproximera eller normalapproximera. Detta är ingen motsägelse, som följande sats visar.

## Sats

*Om  $X$  är  $Po(\mu)$ -fördelad med  $\mu \geq 15$  så är  $X$  approximativt  $N(\mu, \sqrt{\mu})$ -fördelad.*



$$\underbrace{npq \geq 10}_{\approx} N(np, \sqrt{npq})$$

$$\text{Hyp}(N, n, p) \underbrace{n/N \leq 0.1}_{\approx} \text{Bin}(n, p)$$

$$\underbrace{\frac{N-n}{N-1} np(1-p) \geq 10}_{\approx}$$

$$\underbrace{p \leq 0.1}_{\approx} \text{Po}(\underbrace{np}_{=\mu}) \underbrace{\mu \geq 15}_{\approx} N(\mu, \sqrt{\mu})$$

$$N(np, \sqrt{npq})$$

# Berry-Esseens olikhet

Följande hastighet för CGS har tagits fram av Carl-Gustav Esseen (1942) (och av A.C. Berry)

## Sats

Låt  $X_1, X_2, \dots$  vara oberoende och lika fördelade s.v. med väntevärde 0 och standardavvikelse  $\sigma$ . Låt  $Y_n = \frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i$  och låt

$$F_n(x) = P(Y_n \leq x)$$

Då gäller att om  $E(|X_i|^3) = c < \infty$ , så finns en konstant  $K > 0$  så att

$$|F_n(x) - \Phi(x)| \leq \frac{Kc}{\sigma^3\sqrt{n}}$$

Alltså är konvergensthastigheten  $n^{-1/2}$ . Carl-Gustav Esseen visade att  $K \leq 7.59$ . Det är numera känt att  $K \leq 0.7655$  (Shiganov 1986).

