

Denne uken: kap. 6.1-6.2-6.3:

Introduksjon til statistisk inferens

- Konfidensintervall
- Hypotesetesting
- P-verdier
- Statistisk signifikans

Statistisk inferens

- Mål: Trekke konklusjoner fra data
- Tidligere har vi undersøkt data og trukket uformelle konklusjoner
- Formell statistisk inferens:
 - Basere konklusjoner på sannsynlighetsberegninger
 - Tar hensyn til usikkerhet/variasjon

Tilfeldig plassering av trær?

- Statistisk analyse:
 - Så mye klustering/sammenklumping (eller enda mer) vil kun skje i 4% av tilfeller med tilfeldig plassering
 - Observert mønster er altså svært usannsynlig med tilfeldig plassering
 - Konkluderer at det er klustering, og ikke tilfeldig plassering

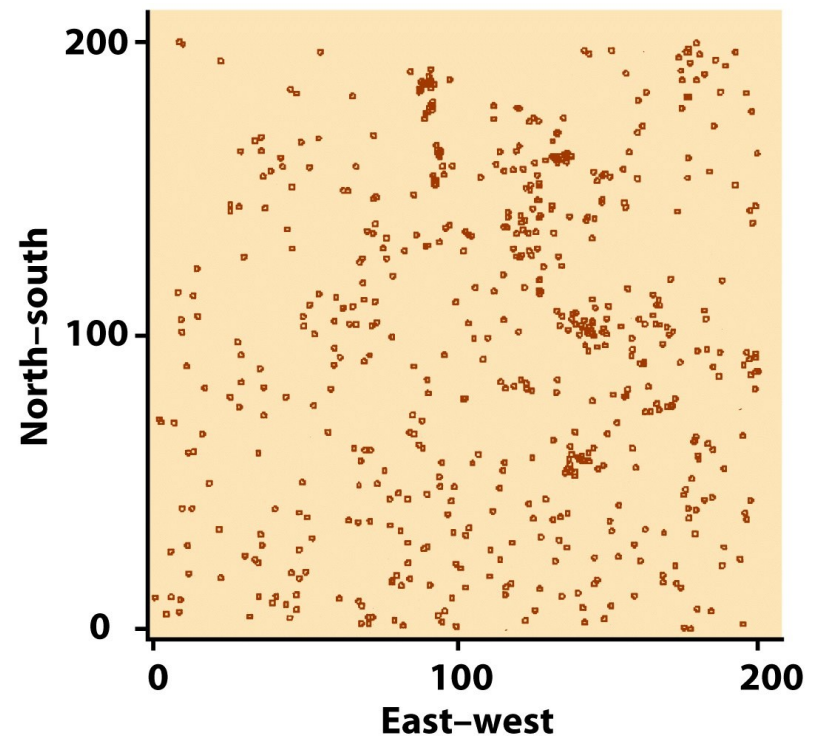


Figure 6-1
Introduction to the Practice of Statistics, Fifth Edition
© 2005 W. H. Freeman and Company

Effektiv ny medisin?

- Gir ny medisin til 20 pasienter, 12 viser bedring av tilstand
- Gir placebo til 20 andre pasienter, 8 viser bedring av tilstand
- Er den nye medisinen mer effektiv enn placebo?
- Kanskje, men en forskjell som den observerte, eller større, mellom resultatene for de to gruppene vil skje i 20% av slike forsøk bare pga tilfeldig variasjon
- En effekt som så lett kan skje bare ved tilfeldighet er ikke overbevisende nok til å konkludere at den nye medisinen er bedre enn placebo

Formell statistisk inferens

- To viktige metoder
 - *Konfidensintervall*
 - *Signifikanstester*
- Basert på *fordelingen* til en statistikk (observator)
- Krever *sannsynlighetsmodell* for dataene
- Statistisk inferens baserer seg på at dataene kommer fra et tilfeldig utvalg eller et randomisert eksperiment
 - Viktig å huske på!

Kap. 6: Kjent σ

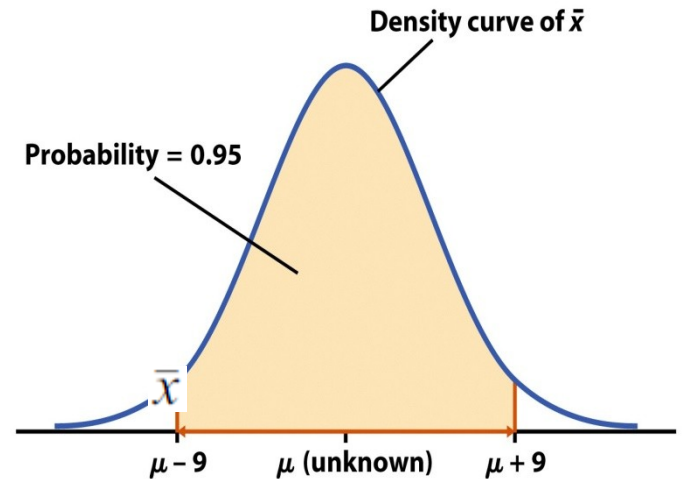
- Hensikten denne uken er å beskrive tankegangen bak statistisk inferens
- Vi skal se på noen enkle metoder for statistisk inferens, som krever en urealistisk antagelse: at vi kjenner det teoretiske standardavviket σ
- Fra og med kap. 7 og utover vil vi se på mer realistiske metoder som kan brukes for de fleste typer av data vi har sett på tidligere
- I dag kap 6.1: *Konfidensintervall*

Konfidensintervall intuitivt

- SATM-poeng:
 - SATM er den matematiske delen av SAT-testen som tester hvor godt man egner seg for college i USA
 - Ønsker å estimere forventet SATM-score for alle sisteårs high school-elever i California
 - Bare 38% av sisteårs high school-elever i California velger å ta SAT. Disse elevene planlegger å gå på college, og er nok ikke representative for alle sisteårs high school-elever
 - Derfor gjør man sitt beste for å gi SAT-testen til et enkelt tilfeldig utvalg på 500 sisteårs high school-elever fra California
 - Gjennomsnittlig SATM-score for disse er $\bar{x} = 485$

Konfidensintervall intuitivt

- SATM-poeng:
 - Dersom poeng for individene i populasjonen er $N(\mu, \sigma)$ -fordelt, vet vi at gjennomsnittet \bar{x} er $N(\mu, \sigma/\sqrt{n})$ -fordelt
 - Antar at vi vet at $\sigma=100$ for populasjonen av sisteårs high school-elever i California. For $n=500$ er da $\sigma/\sqrt{n}=4.5$
 - 68-95-99.7-regelen:



Konfidensintervall intuitivt

- ATM-poeng:
 - Dersom peng for individene i populasjonen er $N(\mu, \sigma)$ -fordelt, vet vi at gjennomsnittet \bar{x} er $N(\mu, \sigma/\sqrt{n})$ -fordelt
 - Antar at vi vet at $\sigma=100$. For $n=500$ er da $\sigma/\sqrt{n}=4.5$
 - 68-95-99.7-regelen: \bar{x} er i $[\mu-2\sigma/\sqrt{n}, \mu+2\sigma/\sqrt{n}] = [\mu-9, \mu+9]$ med ca 95% sannsynlighet
 - Å si at \bar{x} er 9 poeng mindre eller større enn μ er det samme som å si at μ er 9 poeng fra \bar{x}
 - Altså vil den sanne verdien av μ i 95% av utvalg vil ligge i intervallet:

Konfidensintervall intuitivt

- SAT-poeng:
 - $[\bar{x} - 9, \bar{x} + 9]$ er et 95% *konfidensintervall* for μ
 - Har observeret et utvalg med $n=500$ der $\bar{x} = 485$
 - Vi er 95% sikre (confident) på at den ukjente forventningen μ ligger i intervallet:

Konfidensintervall: Tolkning

- Vi *vet ikke* om vårt utvalg er et av de 95% av utvalgene som har konfidensintervall som fanger den ukjente μ , *eller* om det er et av de 5% av utvalgene som har konfidensintervall som ikke fanger μ
- Utsagnet «Vi er 95% sikre (confident) på at den ukjente forventningen μ ligger i intervallet [476, 494]» betyr at vi kom frem til dette intervallet ved å bruke en metode som gir korrekt resultat i 95% av tilfellene

25 utvalg gir 25 forskjellige konfidensintervall, ett dekker ikke μ (dvs 4%). Dersom vi tok veldig mange slike utvalg vil 95% av dem gi intervall som dekker den ukjente parameteren

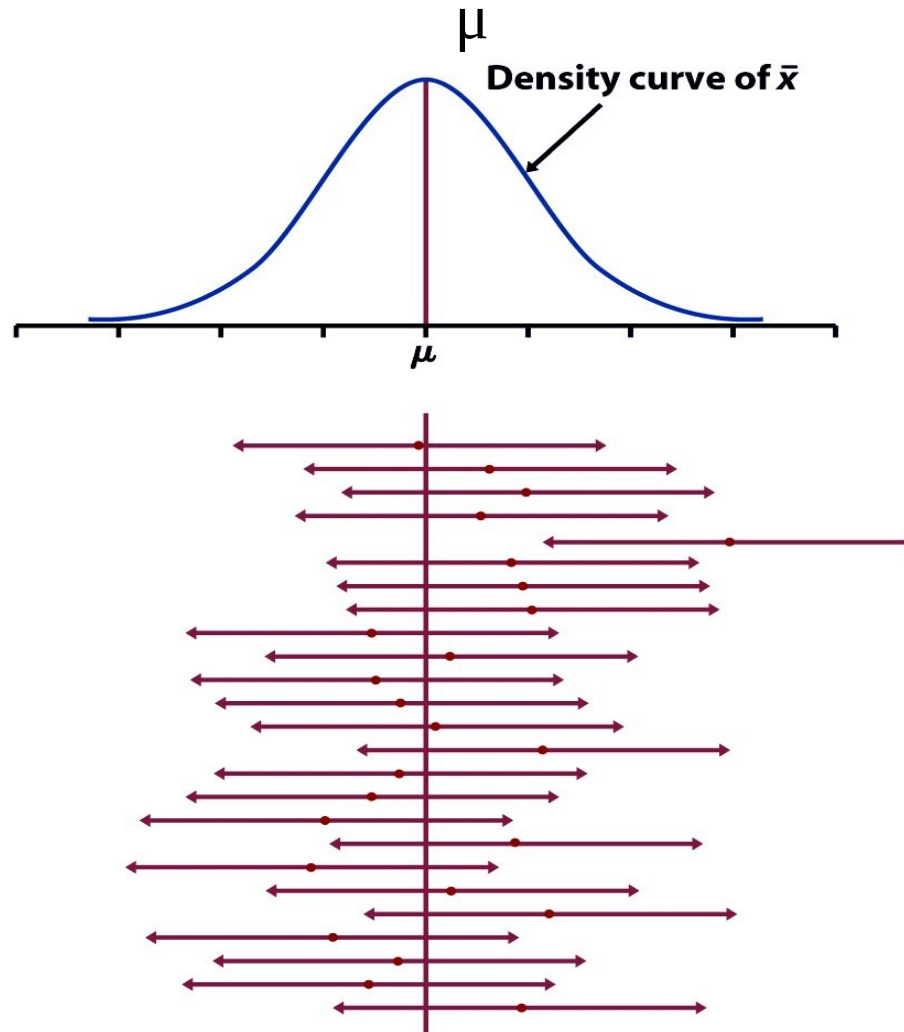


Figure 6-3
Introduction to the Practice of Statistics, Fifth Edition
© 2005 W. H. Freeman and Company

Konfidensintervall og konfidensnivå

- Man kan velge konfidensnivået, 99, 95 og 90 % er det vanligste
- I SAT-eksempelet var 95%-intervallet $\bar{x} \pm 9$
- Generelt konfidensintervall: Estimat \pm feilmargin
- Estimat: Senter i intervall
- Feilmargin: Presisjon av estimat

Konfidensintervall og konfidensnivå

- Generelt konfidensintervall: Estimat \pm feilmargin
- Feilmarginen avhenger av konfidensnivå
 - 68-95-99.7 regel: \bar{x} er i $[\mu - 3\sigma/\sqrt{n}, \mu + 3\sigma/\sqrt{n}]$ med 99.7% sannsynlighet
 - \bar{x} er i $[\mu - 3\sigma/\sqrt{n}, \mu + 3\sigma/\sqrt{n}]$ er ekvivalent med at μ er i $[\bar{x} - 3\sigma/\sqrt{n}, \bar{x} + 3\sigma/\sqrt{n}]$
 - μ er i $[\bar{x} - 3\sigma/\sqrt{n}, \bar{x} + 3\sigma/\sqrt{n}]$ i 99.7% av utvalg
 - $[\bar{x} - 3\sigma/\sqrt{n}, \bar{x} + 3\sigma/\sqrt{n}]$ er et 99.7% *konfidensintervall* for μ
 - Høyere konfidensnivå gir større intervall (større feilmargin)
- *Konfidensnivå C (95/99.7)*: Hvor sikre vi er på at konfidensintervallet inneholder sann parameter

Konfidensnivå for forventning

- Normalfordelte data: \bar{x} er eksakt $N(\mu, \sigma/\sqrt{n})$ -fordelt
- Sentralgrenseteorem for store utvalg: \bar{x} er tilnærmet $N(\mu, \sigma/\sqrt{n})$ -fordelt
- Vi så at vi kunne finne et omtrentlig konfidensintervall for μ ved å bruke 68-95-99.7-regelen
- Skal nå se hvordan vi lager mer presise konfidensintervall for μ
- Må starte med å finne feilmarginene for et bestemt konfidensnivå C
- Går veien om standard normalfordeling: Da kan vi finne generelle feilmarginer som alltid kan brukes for konfidensnivå C når gjennomsnittet \bar{x} er normalfordelt

Konfidensnivå for forventning

- Når Z er $N(0,1)$ -fordelt:
- Velg z^* slik at arealet under kurven mellom $-z^*$ og z^* er C , dvs $P(-z^* < Z < z^*) = C$
- For tre verdier av C (fra nederste rad i Table D):

| | | | |
|-------|-------|-------|-------|
| z^* | 1.645 | 1.960 | 2.576 |
| C | 90% | 95% | 99% |

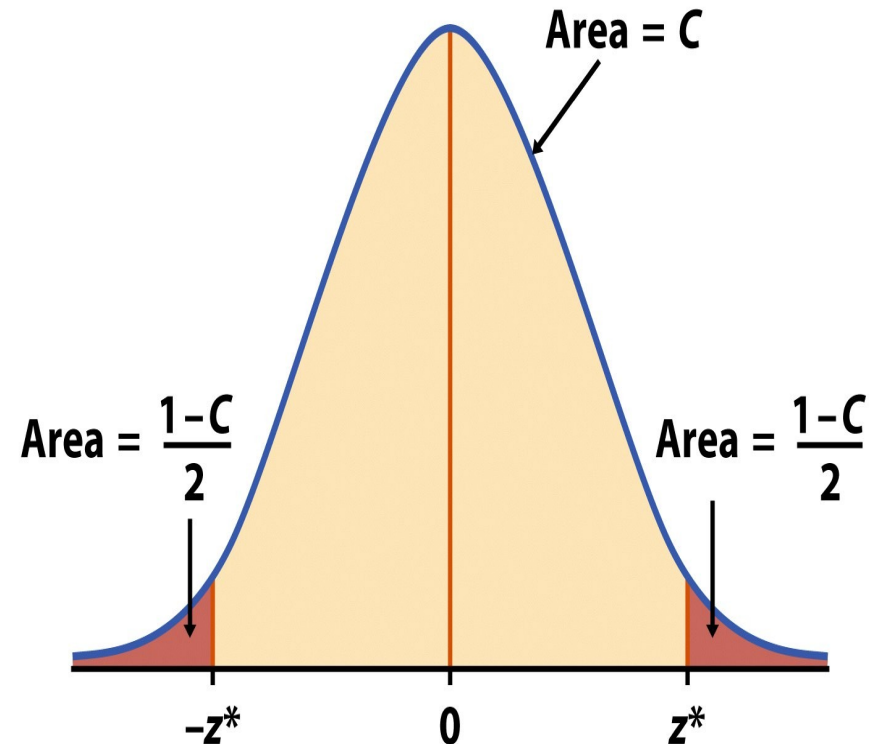


Figure 6-4
Introduction to the Practice of Statistics, Fifth Edition
© 2005 W. H. Freeman and Company

Konfidensintervall

- \bar{x} er $N(\mu, \sigma/\sqrt{n})$ -fordelt, dvs $(\bar{x} - \mu)/(\sigma/\sqrt{n})$ er $N(0,1)$ -fordelt
- Fant z^* slik at $P(-z^* < Z < z^*) = C$ når Z er $N(0,1)$ -fordelt
- $P(-z^* < (\bar{x} - \mu)/(\sigma/\sqrt{n}) < z^*) = C$
- $P(-z^* \sigma/\sqrt{n} < \bar{x} - \mu < z^* \sigma/\sqrt{n}) = C$
- $P(\bar{x} - z^* \sigma/\sqrt{n} < \mu < \bar{x} + z^* \sigma/\sqrt{n}) = C$
- **Konfidensintervall** for μ med nivå C :
- **Feilmargin:**
- Konfidensintervall for μ med nivå 95%:

CONFIDENCE INTERVAL FOR A POPULATION MEAN

Choose an SRS of size n from a population having unknown mean μ and known standard deviation σ . The **margin of error** for a level C confidence interval for μ is

$$m = z^* \frac{\sigma}{\sqrt{n}}$$

Here z^* is the value on the standard normal curve with area C between the critical points $-z^*$ and z^* . The level C **confidence interval** for μ is

$$\bar{x} \pm m$$

This interval is exact when the population distribution is normal and is approximately correct for large n in other cases.

Definition, pg 388

Introduction to the Practice of Statistics, Fifth Edition

© 2005 W.H. Freeman and Company

Konfidensintervall: Eksempel

- I en undersøkelse ble et utvalg på 1280 tidligere studenter spurt om hvor stort studielån de hadde
- Gjennomsnittet \bar{x} var \$18900 og standardavviket s var \$49000
- Klart skjev fordeling, men pga stort utvalg er \bar{x} tilnærmet normalfordelt
- Ukjent teoretisk populasjons-standardavvik σ , men later som om det er kjent lik det observerte standardavviket s , dvs \bar{x} er tilnærmet $N(\mu, \sigma/\sqrt{1280})$ -fordelt, med $\sigma = \$49.000$
- Ønsker å lage et 95%-konfidensintervall for den ukjente forventningen μ
- Vet at $z^* = 1.960$ for $C = 95\%$

Konfidensintervall: Eksempel

- I en undersøkelse ble et utvalg på 1280 tidligere studenter spurt om hvor stort studielån de hadde
- x er tilnærmet $N(\mu, \sigma)$ -fordelt, med $\sigma = \$49000$
- Vet at $z^* = 1.960$ for $C = 95\%$
- Feilmarginen er:
- 95%-konfidensintervall for den ukjente forventningen μ er

Generell form

- Konfidensintervall $[\bar{x} - z^* \sigma / \sqrt{n}, \bar{x} + z^* \sigma / \sqrt{n}]$

$$\sigma_{\bar{x}} = \sigma / \sqrt{n}$$

$$\text{Konfidensintervall} = \text{estimat} \pm z^* \sigma_{\text{estimat}}$$

- *Generell formel* for utvalgsestimater som er normalfordelte (eller tilnærmet normalfordelte)

Hvordan oppfører konfidensintervallene seg?

- Feilmargin $m = z^* \sigma / \sqrt{n}$
- Avhenger av
 - z^* (som avhenger av C): Kan redusere m med mindre C
 - Typisk $C=95\%$ (90% og 99% brukes også)
 - σ : Kan redusere m med data med mindre variasjon
 - n : Kan redusere m med flere målinger

Betydning av C (z^*)

99% og 95% konfidensintervaller for samme utvalgsstørrelsen ($n=1280$)

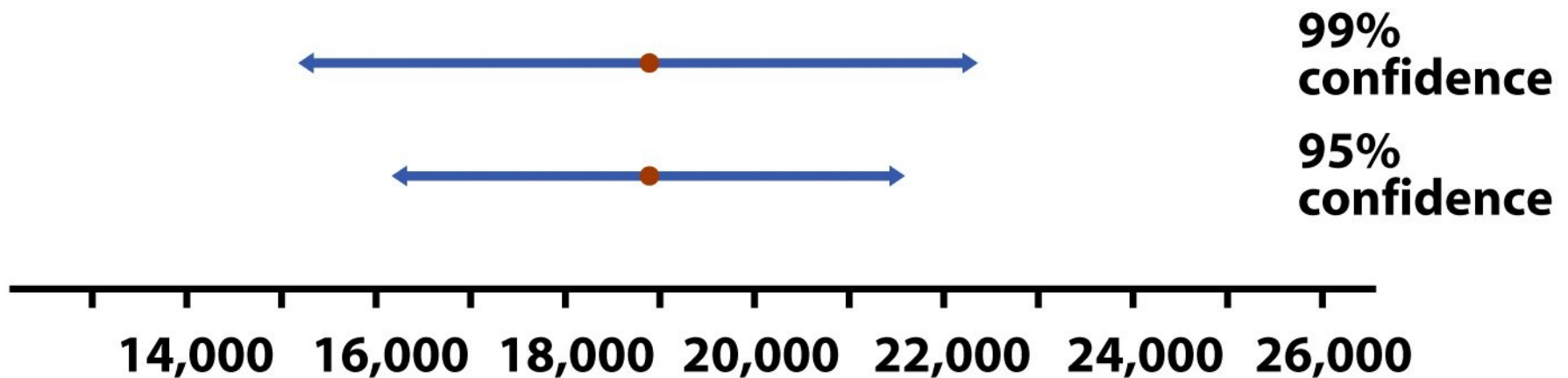


Figure 6-6

Introduction to the Practice of Statistics, Fifth Edition

© 2005 W. H. Freeman and Company

Betydning av n

95% konfidensintervall for to ulike verdier av utvalgsstørrelsen n
(later som om studielånsundersøkelsen var gjort for et utvalg på $n=320$ istedet for de virkelige $n=1280$)

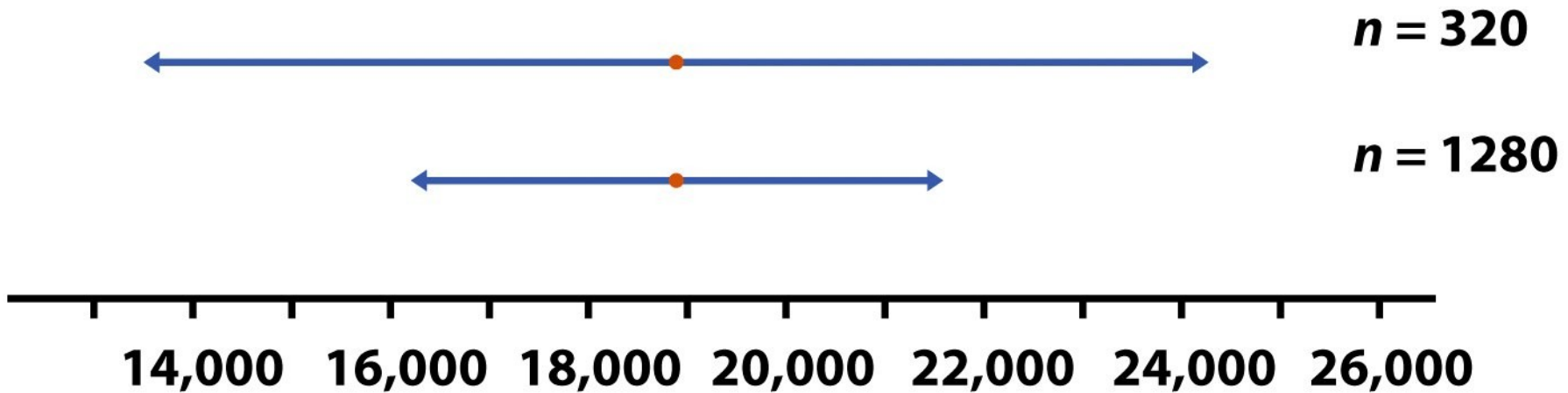


Figure 6-5
Introduction to the Practice of Statistics, Fifth Edition
© 2005 W. H. Freeman and Company

Valg av utvalgsstørrelse

- Design av eksperiment med bestemt feilmargin m :
 - Ønsker feilmargin m med konfidens C
 - C gir deg verdien av z^*
 - Velg utvalgsstørrelse n slik at $m = z^* \sigma / \sqrt{n}$
 - Løser ut for n :

Noen forsiktighetsregler

- Data bør være fra et *enkelt randomisert utvalg* (SRS) av populasjonen
 - Viktig med *uavhengige* observasjoner fra populasjonen
- Andre, korrigerede formler for mer kompliserte design
- Følsomt for *uteliggere*
- Lite robust for små n (bør ha $n > 15$) når data ikke er normalfordelte
- Må kjenne σ . Senere skal vi se på hva vi gjør når σ er ukjent
 - Hvis n stor, kan vi bruke $[\bar{x} - z*s/\sqrt{n}, \bar{x} + z*s/\sqrt{n}]$ (som da er et *tilnærmet* konfidensintervall for μ)