

Föreläsning 9-10. Punktskattningar och konfidensintervaller

Väntevärdesriktig punktskattning av μ

Vi antar att vektor (X_1, X_2, \dots, X_n) består av n oberoende slumpvariabler med samma fördelning som har okänt populations parameter μ och samma populations standardavvikelse σ . Ett stickprov (x_1, x_2, \dots, x_n) är en given realisering (bland flera möjliga) av slumpvektorn (X_1, X_2, \dots, X_n) .

En naturlig punktskattning av μ är stickprovs medelvärde

$$\bar{x} = \frac{x_1 + \dots + x_n}{n},$$

som är en realisering (bland flera möjliga) av slumpvariabeln

$$\bar{X} = \frac{X_1 + \dots + X_n}{n}.$$

Punktskattning \bar{x} är väntevärdesriktig punktskattning av μ eftersom

$$E(\bar{X}) = E\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{E(X_1) + \dots + E(X_n)}{n} = \frac{\mu + \dots + \mu}{n} = \mu.$$

Väntevärdesriktig betyder att det finns ingen systematiskt fel (unbiased estimate).

Dock som regel finns det ett slumpfel $\bar{x} \neq \mu$ som orsakas av att n stickprovs värden väljs på måfå. Olika stickprovs utfall ger olika värde för \bar{x} . Medelfel

$$\sigma_{\bar{X}} = \sqrt{\text{Var}(\bar{X})}$$

mäter slumpfels storlek: ju större $\sigma_{\bar{X}}$ desto större avvikelser $\bar{x} - \mu$ förväntas. Enligt oberoende antagandet

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{\text{Var}(X_1) + \dots + \text{Var}(X_n)}{n^2} = \frac{\sigma^2 + \dots + \sigma^2}{n^2} = \frac{\sigma^2}{n}.$$

Det innebär att

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

ju större är stickprovs storlek n desto mindre medelfel.

Populations andel p och stickprovs andel \hat{p}

Population består av N individer. Det finns två slags individer

$$N = N_1 + N_0, \quad p = \frac{N_1}{N}, \quad 1 - p = \frac{N_0}{N}$$

och vi är intresserade i populations andel p . Parameter p är okänd och i praktiken omöjligt att räkna exakt.

Vi tar ett stickprov av storlek n och för $i = 1, \dots, n$ skriver $x_i = 1$ om person nummer i tillhör grupp 1 och $x_i = 0$ om personen tillhör grupp 0. Vi tillämpar två olika modeller för att beskriva slumpexperiment för stickprovs samling:

Modell 1: dragning med återläggning (oberoende slumpvariabler X_1, X_2, \dots, X_n),

Modell 2: dragning utan återläggning (beroende slumpvariabler X_1, X_2, \dots, X_n).

Modell 1 tolkas som oändlig population och Modell 2 som ändlig population.

Frågor

1. vilken fördelning har slumpvariabel X_i ?
2. vilken fördelning har slumpvariabel $Y = X_1 + \dots + X_n$?
3. förklara varför modell 1 kan tolkas som oändlig population?
4. vad är väntevärde och varians av Y ?

Vi skattar p med hjälp av stickprovs andel

$$\hat{p} = \bar{x}.$$

Stickprovs andel \hat{p} är en väntevärdesriktig punktskattning av populations andel p eftersom $X_i \sim \text{Ber}(p)$ och därför

$$\mu = p, \quad \sigma^2 = p(1-p).$$

Medelfelet är

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$$

om populationen är oändlig, och

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n} \frac{N-n}{N-1}}$$

om populationen är ändlig.

Skattad medelfel

Vi antar igen att vektor (X_1, X_2, \dots, X_n) består av n oberoende slumpvariabler med samma fördelning som har okänt populations parameter μ och samma populations standardavvikelse σ . Om σ är okänd då medelfelet

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

kan inte räknas fram. För ett givet stickprov (x_1, \dots, x_n) beräknar man (skattad) medelfel enligt formel

$$s_{\bar{x}} = \frac{s}{\sqrt{n}}, \quad s = \sqrt{\frac{(x_1 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n-1}},$$

där s kallas stickprovs standardavvikelse. Man delar med $n-1$ istället för n eftersom

$$\begin{aligned} E(S^2) &= E\left(\frac{(X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2}{n-1}\right) = \frac{nE(X_1 - \bar{X})^2}{n-1} = \frac{n\text{Var}(X_1 - \bar{X})}{n-1} \\ &= \frac{\text{Var}((n-1)X_1 - X_2 - X_3 - \dots - X_n)}{n(n-1)} = \frac{(n-1)^2\sigma^2 + \sigma^2 + \sigma^2 + \dots + \sigma^2}{n(n-1)} = \sigma^2. \end{aligned}$$

Skattad medelfel för stickprovs andel är

$$s_{\hat{p}} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n-1}}$$

om populationen är oändlig, och

$$s_{\hat{p}} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n-1} \frac{N-n}{N}}$$

om populationen är ändlig.

Approximativa konfidensintervaller för p

För dragning med återläggning, om stickprovsstorlek n är stor, då

$$\sigma_{\hat{p}} \approx s_{\hat{p}} \approx \sqrt{\frac{\hat{p}(1-\hat{p})}{n}},$$
$$\hat{P} \approx N(p, \sigma_{\hat{p}}) \approx N\left(p, \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right).$$

Det ger formeln

$$I_p = \hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

för en approximativt $100(1 - \alpha)\%$ konfidensintervall för p i oändlig population.

För dragning utan återläggning då stickprovsstorlek n är stor, vi får

$$I_p = \hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n} \frac{N-n}{N}}$$

som en approximativt $100(1 - \alpha)\%$ konfidensintervall för p i ändlig population.

Exakta konfidensintervaller för μ

Vi antar igen att vektor (X_1, X_2, \dots, X_n) består av n oberoende slumpvariabler med samma normalfördelning $N(\mu, \sigma)$. Om σ är känd då

$$\bar{X} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right).$$

Det ger formel

$$I_\mu = \bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

för ett exakt konfidensintervall för μ med känd σ .

Om σ är okänd då

$$T = \frac{\bar{X} - \mu}{s_{\bar{x}}} \sim t_{n-1}$$

har så kallade t-fördelning med $n - 1$ frihetsgrader. Det innebär

$$P\left(-t_{\alpha/2}(n-1) < \frac{\bar{X} - \mu}{s_{\bar{x}}} < t_{\alpha/2}(n-1)\right) = 1 - \alpha,$$

som ger formeln

$$I_\mu = \bar{x} \pm t_{\alpha/2}(n-1) \frac{s}{\sqrt{n}}$$

för ett exakt konfidensintervall för μ med okänd σ .