

Markovkedjor

- tillståndsrum $S = \{s_1, \dots, s_r\}$ och övergångsmatris $P = (p_{ij})_{i,j=1,\dots,r}$ där
 $p_{ij} = P(X_{n+1} = s_j | X_n = s_i)$ för godtyckligt $n \in \mathbb{N}_0$, $\sum_{j=1}^r p_{ij} = 1$ (P radstokastisk)
- flerstegs övergångsslh. $P(X_n = s_j | X_0 = s_i) = (P^n)_{ij}$ (element (i,j) i matrispotensen P^n)
- multiplikation från höger! $u^{(n)} := (P(X_n = s_1), P(X_n = s_2), \dots, P(X_n = s_r))^T \in [0,1]^r$
 $\rightarrow u^{(n)} = u^{(0)} \cdot P^n$
- MK är irreducibel om $(P^n)_{ij} > 0$ för godtyckliga i,j och något n
MK är regulär om $(P^n)_{ij} > 0$ för något fast $n \in \mathbb{N}_0$ och alla i,j
- För en regulär MK gäller $P^n \rightarrow \Pi = \begin{pmatrix} \pi_1 & & \\ & \ddots & \\ & & \pi_r \end{pmatrix}$ där $\Pi = (\pi_1, \dots, \pi_r)$ kallas gränsfördelning.
- tillstånd s_j absorberande om $p_{jj} = 1$, transient om det har pos. slh. att aldrig återvända
- För en absorberande MK väljer vi kanonisk form $P = \begin{pmatrix} Q & R \\ \text{trans.} & \text{abs.} \\ 0 & I_a \end{pmatrix}$
 $\Rightarrow P^n \rightarrow \begin{pmatrix} 0 & NR \\ 0 & I_a \end{pmatrix}$ där $N = (I - Q)^{-1}$ är den fundamentala matrisen
 - $n_{ij} =$ förväntade antalet besök i s_j med start i s_i (där s_i, s_j är transienta)
 - $d_j = N \cdot \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$ anger: $d_j =$ tiden till absorption efter start i tillstånd s_j (transient)

Genererande funktioner

$$g(x) = \sum_{n=0}^{\infty} a_n x^n \quad \text{till en given talfoljd } (a_n)_{n \in \mathbb{N}_0}$$

viktiga g.f. $g(x) = \sum_{n=0}^{\infty} (cx)^n = \frac{1}{1-cx} \quad [|cx| < 1] \quad \text{geometrisk}$

$$g(x) = \sum_{k=0}^{\infty} \binom{n}{k} x^k = (1+x)^n \quad \text{binomial}$$

$$g(x) = \sum_{n=0}^{\infty} \binom{n+k}{k} x^n = \frac{1}{(1-x)^{k+1}} \quad [|x| < 1]$$

- kan användas till att räkna antalet heltalslösningar till $y_1 + y_2 + \dots + y_k = n$ med givna restriktioner $\left[\text{koeff. framför } x^n \text{ i } g_1(x) \cdots g_k(x) \text{ där } g_i(x) = \sum_{m \text{ tillåtet värde}} x^m \text{ för } y_i \right]$
- kan användas till att lösa rekurrensor (ibland tillsammans med partialbråksuppdelning)
 $\left[\text{gänga rekurrensen med } x^n, \text{ summa över tillåtna } n, \text{ lösa ut } g(x), \text{ skriva om till potensserie} \right]$

Exponentiella gen. fkt.

$$g(x) = \sum_{n=0}^{\infty} a_n \frac{x^n}{n!} \quad \text{till en given talfoljd } (a_n)_{n \in \mathbb{N}_0}$$

- kan användas till att räkna antalet strängar med "bokstäver" ur ett givet alfabet med restriktioner [strategi: summa som ovan fast nu med exp. gen. fkt.]

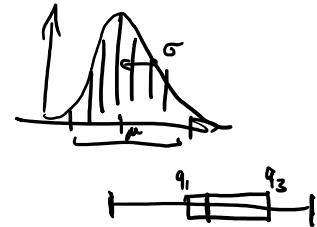
- vanliga: $e^{cx} = \sum_{n=0}^{\infty} \frac{(cx)^n}{n!}, \quad \cosh(x) = \frac{e^x + e^{-x}}{2} = \sum_{n=0}^{\infty} \frac{x^{2n}}{(2n)!}, \quad \sinh(x) = \frac{e^x - e^{-x}}{2} = \sum_{n=0}^{\infty} \frac{x^{2n+1}}{(2n+1)!}$

Statistik

Beskrivande S.

histogram (variationsbredd, rel. frekvenser)

boxplot (median, 1:a och 3:c quartil, iqr)



Inferentiel S.

- Generellt antagande: Stickprov består av n oberoende s.v., alla har samma fördelning som populationen.
- Estimator $\hat{\theta}$ för en parameter θ av populationens fördelning är en s.v. och kan vara väntevärdesriktig ($E\hat{\theta} = \theta$) och konsistent ($\text{var}(\hat{\theta}) \rightarrow 0$ med $n \rightarrow \infty$).
- Stickprovsmedelvärdet $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ är väntevärdesriktig och har varians $\frac{\sigma^2}{n}$ där $\sigma^2 = \text{var}(X_i)$, alltså konsistent.
- Stickprovsvariansen $S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1} = \frac{n \sum_i X_i^2 - (\sum_i X_i)^2}{n(n-1)}$ är väntevärdesriktig ($E S^2 = \sigma^2$) estimator för σ^2 .
- För normalfördelad population (eller stort stickprov, $n \geq 25$) är \bar{X} (approximativt) normalfördelad med väntevärde μ och varians σ^2 .
 $\Rightarrow \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$. (Centrala gränsvärdesl.)
- Ett symmetriskt $(1-\alpha) \cdot 100\%$ konfidensintervall för μ ges då av
 $\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$ där $z_{\alpha/2}$ är sådan att $P(Z > z_{\alpha/2}) = \alpha/2$
ensidigt: $(-\infty, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}})$ resp. $(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \infty)$.
högra vänstra svansen
- För normalfördelad population (eller stor n) är $\frac{(n-1)S^2}{\sigma^2}$ (approximativt) χ^2 -fördelad med $n-1$ frihetsgrader.
- Ett symmetriskt $(1-\alpha) \cdot 100\%$ konfidensintervall för σ^2 ges då av $\left(\frac{(n-1)S^2}{q_{1-\alpha/2}}, \frac{(n-1)S^2}{q_{\alpha/2}} \right)$ där q_α är sådan att $P(X < q_\alpha) = \alpha$.
ensidigt: $[0, \frac{(n-1)S^2}{q_\alpha}]$
- För en normalfördelad population är $\frac{\bar{X} - \mu}{S/\sqrt{n}}$ t-fördelad med $n-1$ frihetsgrader
symmetriskt $(1-\alpha) \cdot 100\%$ k.I. för μ :
 $\bar{X} \pm t_{\alpha/2} \cdot \frac{S}{\sqrt{n}}$ där $t_{\alpha/2}$ är sådan att $P(T > t_{\alpha/2}) = \alpha/2$.

Hypotesprövning

$$\begin{array}{lll} \text{nollhypotes: } H_0: \theta = \theta_0 & , & \theta \geq \theta_0 \text{ eller } \theta \leq \theta_0 \\ \text{alternativ: } H_1: \theta \neq \theta_0 & , & \theta < \theta_0 \quad \theta > \theta_0 \\ & \text{tvåsidigt} & \text{ensidigt} \end{array}$$

- Väljer signifikansnivå och bedömer om nollhypotesen är sannig eller bör förkastas (typ I fel: H_0 förkastas men stämmer, typ II fel: H_0 behålls dock fel)
- Berörer en teststatistik och dess fördelning (givet antagandet att $\theta = \theta_0$) vanligast: $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$ och $T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$ t-fördelad med $n-1$ frihetsgrader (är det gäller $\theta = \mu$)
- om utfallet av Z inte landar i $(-z_{\alpha/2}, z_{\alpha/2})$ (tvåsidigt)
resp. $(-\infty, z_\alpha)$ eller $(-z_\alpha, \infty)$ (ensidigt) förkastar vi H_0 .
Motsvarande regioner för T : $(-t_{\alpha/2}, t_{\alpha/2})$
resp. $(-\infty, t_\alpha)$ eller $(-t_\alpha, \infty)$
- P -värde: minsta signifikansnivå α på vilken H_0 förkastas.
t.ex. för teststatistik Z med utfall z_0 är det

$$P = \begin{cases} 2P(Z \leq z_0) & \text{om } z_0 \leq 0 \\ 2P(Z \geq z_0) & \text{om } z_0 > 0 \end{cases}$$
 (tvåsidigt) resp. $P = \begin{cases} P(Z \leq z_0) & \text{(vänster)} \\ P(Z \geq z_0) & \text{ensidigt} \end{cases}$ (höger)
(samma för andra symmetriskt fördelade)
(teststatistiker, som t.ex. T)

Proportionsskattning.

- X_j : indikator om element j i stickprovet ingår
- $\hat{p} = \frac{\sum X_j}{n}$ är väntevärdesriktig och konsistent estimator av proportionen p .
 - Ett symmetriskt $(1-\alpha)100\%$ konfidensintervall för p är $\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$
 - Vill man ha max. avvikelse d från punktskattningen \hat{p} i detta behövs stickprovsstorlek $n \geq \frac{z_{\alpha/2}^2 p_0(1-p_0)}{d^2}$ med apriori värde p_0 , $n \geq \frac{z_{\alpha/2}^2}{4d^2}$ utan.
 - För hypotesprövning används man teststatistiken $Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \sim N(0,1)$
 - Jämförelse av två populationer:
 $\hat{p}_1 - \hat{p}_2$ approx. $N(\hat{p}_1 - \hat{p}_2, \frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2})$
 \rightarrow symmetriskt $(1-\alpha)100\%-KI$ för $\hat{p}_1 - \hat{p}_2$: $\hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$
 så länge min $\{\hat{p}_1; n(1-\hat{p}_1)\} > 5$

Jämförelse av medelvärden i två populationer

- $\sigma_1 \neq \sigma_2$: $\bar{X}_1 - \bar{X}_2$ (approximativt) $N(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$
 $\Rightarrow \bar{X}_1 - \bar{X}_2 \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$ symmetriskt $(1-\alpha)100\%$ -KI för $\bar{X}_1 - \bar{X}_2$.
- $\sigma_1 = \sigma_2 =: \sigma$ (nöde okänd)
 kombinerad stickprovsvarians $S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1+n_2-2}$
 väntevärdesriktig estimator för σ^2 ,
 $\frac{n_1+n_2-2}{\sigma^2} S_p^2$ χ^2 -fördelad med n_1+n_2-2 frihetsgrader.

$$\frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{S_p^2 (\frac{1}{n_1} + \frac{1}{n_2})}} \quad t\text{-fördelad} \quad \text{---} \quad \text{---} \quad \text{---} \quad \text{---}.$$

Kan användas som teststatistik eller till beräkning av CTT (symmetriskt)
 $(1-\alpha) \cdot 100\%$ -KI för $\mu_1 - \mu_2$: $(\bar{X}_1 - \bar{X}_2) \pm t_{\alpha/2} \sqrt{S_p^2 (\frac{1}{n_1} + \frac{1}{n_2})}$

Linjär regression

X : förklarande variabel, Y : responsvariabel

vill hitta ett linjärt samband $\mu_{Y|X} = \beta_0 + \beta_1 \cdot X$

datamängd $\{(x_i, y_i); i=1 \dots n\}$

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_i x_i^2 - \frac{1}{n} (\sum_i x_i)^2$$

$$S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_i y_i^2 - \frac{1}{n} (\sum_i y_i)^2$$

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_i x_i y_i - \frac{1}{n} (\sum_i x_i)(\sum_i y_i)$$

minsta kvadrat estimatorer

$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$ är en väntevärdesriktig est. för β_1 , $\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{S_{xx}})$

$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ ————— β_0 , $\hat{\beta}_0 \sim N(\beta_0, \frac{\sigma^2}{n S_{xx}} \sum_i x_i^2)$

$$SSE = \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2 = S_{yy} - \hat{\beta}_1 S_{xy} \quad (\text{sum of squared errors})$$

$S^2 = \frac{SSE}{n-2}$ är en väntevärdesriktig est. för σ^2 .

- $Z = \frac{\beta_1 - \beta_1}{\sigma / \sqrt{s_{xx}}} \sim N(0,1)$ eller $T = \frac{\beta_1 - \beta_1}{s / \sqrt{s_{xx}}}$ t-fördelad med $n-2$ frihetsgrader
Kan användas för att bestämma konfidensintervall
 $\beta_1 \pm z_{\alpha/2} \frac{S}{\sqrt{s_{xx}}}$ resp. $\beta_1 \pm t_{\alpha/2} \frac{S}{\sqrt{s_{xx}}}$ eller hypotesprövning.
 - Samma gäller för $Z = \frac{\beta_0 - \beta_0}{\sigma / \sqrt{\sum_i x_i^2 / n \cdot s_{xx}}} \sim N(0,1)$, $T = \frac{\beta_0 - \beta_0}{S \cdot \sqrt{\frac{\sum_i x_i^2}{n \cdot s_{xx}}}}$
 $\Rightarrow \beta_0 \pm z_{\alpha/2} S \sqrt{\frac{\sum_i x_i^2}{n \cdot s_{xx}}}$ resp. $\beta_0 \pm t_{\alpha/2} S \cdot \sqrt{\frac{\sum_i x_i^2}{n \cdot s_{xx}}}$
- t-fördelad med $n-2$ frihetsgrader

Centrala gråhs värdessatsen

X_1, \dots, X_n oberoende, samma fördelning med väntevärde μ och varians σ^2
Då är \bar{X} approximativt normalfördelad med väntevärde μ och varians σ^2/n
alltså $\frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \sim N(0,1)$ ($n \geq 25$).

Chebyshov

För X med $E(X) = \mu$ och $\text{var}(X) = \sigma^2$ och godtyckligt $\varepsilon > 0$ gäller
 $P(|X - \mu| \geq \varepsilon) \leq \frac{\sigma^2}{\varepsilon^2}$

Stora talens lag

X_1, \dots, X_n oberoende, samma väntevärde μ och varians σ^2 , godtyckligt $\varepsilon > 0$

$$\Rightarrow P(|\bar{X} - \mu| \geq \varepsilon) \rightarrow 0 \text{ med } n \rightarrow \infty.$$