

Stochastiek 2

Inleiding in de Mathematische Statistiek

Herhaling t-toetsen

Steekproefgemiddelde en -variantie van normale observaties

Stelling.

Laat X_1, \dots, X_n o.o. zijn en $N(\mu, \sigma^2)$ -verdeeld. Dan:

- (i) $\bar{X} \sim N(\mu, \sigma^2/n)$,
- (ii) $(n-1)S_X^2/\sigma^2 \sim \chi_{n-1}^2$,
- (iii) \bar{X} en S_X^2 zijn onafhankelijk (!),
- (iv) $\sqrt{n}(\bar{X} - \mu)/S_X \sim t_{n-1}$.

Eensteekproef-t-toets

Eensteekproef geval: we hebben één steekproef X_1, \dots, X_n , o.o., $N(\mu, \sigma^2)$ -verdeeld, met *zowel μ als σ onbekend*.

Probleem: hypothese over de verwachting μ . Voor gegeven $\mu_0 \in \mathbb{R}$ één van

$$H_0 : \mu \leq \mu_0, \quad H_0 : \mu \geq \mu_0, \quad \text{of} \quad H_0 : \mu = \mu_0.$$

→ vb. 4.30

Tweesteekproevenprobleem

Tweesteekproeven geval: we hebben twee steekproeven, X_1, \dots, X_m o.o. uit verdeling I en Y_1, \dots, Y_n o.o. uit verdeling II.

Probleem:

Algemeen: vergelijken verdelingen I en II. Meestal: is de verwachting van verdeling I gelijk aan die van verdeling II of niet?

Twee deelgevallen:

- **gepaarde waarnemingen**: $m = n$ en data vormen natuurlijke paren $(X_1, Y_1), \dots, (X_n, Y_n)$. Coördinaten binnen een paar mogen afhankelijk zijn, paren onderling onafhankelijk. Typisch voorbeeld: $X_i \sim$ "voor", $Y_i \sim$ "na".
- **ongepaarde waarnemingen**: mogelijk $n \neq m$, meestal X -vector en Y -vector onafhankelijk.

→ vb. 4.33, 4.34, 4.35

Aanpassingstoetsen (goodness-of-fit)

Beoordelen van verdeling van data

Vaak geïntereiseerd in redelijkheid van **aannamen over verdeling van data**.

Grafisch hulpmiddel: **QQ-plot**.

Aanvulling: **aanpassingstoetsen** (goodness-of-fit tests). Formele toetsen voor de nulhypothese dat de data een steekproef uit een gegeven verdeling (of familie van verdelingen) is.

Bijv:

- **Kolmogorov-Smirnov-toets**: toets op basis van empirische verdelingsfunctie (vb. 4.38).
- **Chikwadraat-toets**: toets op basis van verwachte aantal waarnemingen in gegeven intervallen (vb. 4.39).

Empirische verdelingsfunctie

Gegeven: X_1, \dots, X_n o.o., verdelingsfunctie F .

Empirische verdelingsfunctie:

$$\mathbb{F}_n(x) = \frac{1}{n} \#\{i : X_i \leq x\} = \frac{1}{n} \sum_{i=1}^n 1_{\{X_i \leq x\}}, \quad x \in \mathbb{R}.$$

Voor elke $x \in \mathbb{R}$:

$$E\mathbb{F}_n(x) = F(x), \quad \mathbb{F}_n(x) \xrightarrow{\text{b.z.}} F(x) \text{ als } n \rightarrow \infty.$$

Dus voor vaste x is $\mathbb{F}_n(x)$ een **zuivere, consistente** schatter voor $F(x)$.

Kolmogorov-Smirnov-toets (vb 4.38)

Waarnemingen: X_1, \dots, X_n o.o., verdelingsfunctie F .

Hypothesen: $H_0 : F = F_0$, tegen $H_1 : F \neq F_0$, gegeven F_0 .

Toetsingsgrootheid: $T = \sup_{x \in \mathbb{R}} |\mathbb{F}_n(x) - F_0(x)|$.

Toets: “verwerp H_0 als $T \geq c$ ” voor zekere $c > 0$.

Stelling.

Onder H_0 (i.e. als $F = F_0$) heeft T een verdeling die onafhankelijk is van F_0 .

N.B. Er zijn uitbreidingen naar $H_0 : F \in \{F_\theta : \theta \in \Theta\}$.

Belangrijk voorbeeld: $H_0 : F \in \{N(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0\}$.

ZIE OPGAVE 4.46 (4.42)!!!

Chikwadraat-toets

Waarnemingen: X_1, \dots, X_n o.o., verdelingsfunctie F .

Hypothesen: $H_0 : F = F_0$, tegen $H_1 : F \neq F_0$, gegeven F_0 .

Aanpak:

- Verdeel bereik in intervallen (bins) I_1, \dots, I_k .
- Definieer tellingen $N_j = \#\{i : X_i \in I_j\}$.

Stelling.

Onder H_0 (i.e. als $F = F_0$) geldt

$$T = \sum_{j=1}^k \frac{(N_j - np_j)^2}{np_j} \xrightarrow{d} \chi_{k-1}^2$$

Toets: “verwerp H_0 als $T \geq \chi_{k-1, 1-\alpha}^2$ ”.

Likelihood-ratiotoetsen

Likelihood-ratiostatistiek

Setting: waarneming $X \sim p_\theta$, $\theta \in \Theta$, partitie $\Theta = \Theta_0 \cup \Theta_1$,
hypotheses $H_0 : \theta \in \Theta_0$, $H_1 : \theta \in \Theta_1$.

Definitie.

De **likelihood-ratiostatistiek** is

$$\lambda(X) = \frac{\sup_{\theta \in \Theta} p_\theta(X)}{\sup_{\theta_0 \in \Theta_0} p_{\theta_0}(X)}$$

Idee: als $\lambda(X)$ groot is, dan is het sup over Θ_1 groter dan het sup over Θ_0 . Dan ligt de MLS in Θ_1 , wat er op duidt dat H_1 waar is.

→ vb. 4.42

Likelihood-ratiostatistiek - asymptotische verdeling

Nu: o.o. waarnemingen $X_1, \dots, X_n \sim p_\theta$, $\theta \in \Theta \subset \mathbb{R}^k$. Voor alle x , $\theta \mapsto p_\theta(x)$ is differentieerbaar.

Aanname: stel dat voor $\theta_0 \in \Theta_0$:

$$\sqrt{n}(\Theta - \theta_0) \rightarrow \text{"}k\text{-dim. lineaire ruimte"},$$

$$\sqrt{n}(\Theta_0 - \theta_0) \rightarrow \text{"}k_0\text{-dim. lineaire ruimte"}$$

als $n \rightarrow \infty$.

"Stelling"

In deze setting geldt onder regulariteitsvoorwaarden dat onder P_{θ_0} ,

$$2 \log \lambda_n(X_1, \dots, X_n) \xrightarrow{d} \chi_{k-k_0}^2$$

als $n \rightarrow \infty$.

Likelihood-ratiotoets

In de voorgaande setting heeft voor grote n de volgende toets dus bij benadering niveau $\alpha \in (0, 1)$:

“verwerp H_0 als $2 \log \lambda_n(X_1, \dots, X_n) \geq \chi_{k-k_0, 1-\alpha}^2$ ”.

Dit is de **likelihood-ratiotoets**.

N.B.: als niet is voldaan aan de voorwaarden kan de chikwadraat benadering helemaal fout zijn!

→ vb. 4.44, 4.46, 4.47

Toetsingstheorie

Optimale toetsen - Definitie

Definitie.

Een toets met onderscheidend vermogen $\theta \mapsto \pi(\theta; K)$ heet **uniform meest onderscheidend** (UMP) bij niveau α , voor het toetsen van $H_0 : \theta \in \Theta_0$ tegen $H_1 : \theta \in \Theta_1$, als de toets niveau α heeft en voor iedere andere toets van niveau α , met kritiek gebied K' , geldt dat

$$\forall \theta \in \Theta_1 : \pi(\theta; K) \geq \pi(\theta; K').$$

Optimale toetsen - enkele feiten

- UMP toetsen bestaan niet altijd.
- Het kan nodig zijn om **lotingstoetsen** te bekijken. Formeel een afbeelding $\varphi : \mathcal{X} \rightarrow [0, 1]$. Interpretatie: als we $X = x$ observeren, verwerpen we H_0 met kans $\varphi(x)$.
- Als de hypothesen **enkelvoudig** zijn: $H_0 : \theta = \theta_0$ tegen $H_1 : \theta = \theta_1$ voor zekere θ_0, θ_1 , dan bestaat er een UMP (lotings)toets: **Neyman-Pearson** lemma.

Neyman-Pearson

Setting: waarneming $X \sim p_\theta$, $\theta \in \Theta = \{\theta_0, \theta_1\}$.

Statistiek: $L(\theta_1, \theta_0; X) = p_{\theta_1}(X)/p_{\theta_0}(X)$.

Aanname: Voor $\alpha \in (0, 1)$ bestaat er een c_α z.d.d.
 $P_{\theta_0}(L(\theta_1, \theta_0; X) \geq c_\alpha) = \alpha$.

Stelling. (Neyman-Pearson Lemma)

De toets met kritiek gebied $K = \{x : L(\theta_1, \theta_0; x) \geq c_\alpha\}$ is een UMP toets van niveau α voor dit probleem.

→ vb. 6.32

Optimale toetsen - nog meer feiten

- Als de verdelingsfunctie van de statistiek L niet continu is onder H_0 , dan lotingstoets nodig. (Zie Stelling 6.34, vb. 6.35.)
- Voor modellen met een “monotone likelihood-ratio” bestaat er een UMP (lotings)toets voor hypothesen van de vorm $H_0 : \theta \leq \theta_0$ tegen $H_1 : \theta \geq \theta_1$. Speciaal geval: de Gauss-toets is UMP voor het toetsen van deze hypothesen (voor θ de verwachting). (Zie sectie 6.4.2.)

Optimale toetsen - t-toets

- Voor het probleem van de t -toets (normale steekproef, toets voor verwachting bij onbekende variantie) bestaat **geen** UMP toets (voor niveau $\alpha < 1/2$)
- De t -toets is wel uniform meest onderscheidend binnen de klasse van **zuivere** toetsen. (Zie Stelling 6.50 (6.49).)

Definitie.

Een toets heet **zuiver** van niveau α als voor alle $\theta_0 \in \Theta_0$ en $\theta_1 \in \Theta_1$ geldt dat

$$\pi(\theta_0; K) \leq \alpha \leq \pi(\theta_1; K).$$