

4.4 Testi võimsus

Mudelite võrdlemisel F-testi abil eeldasime, et keerukam mudel on kindlasti õige. Lihtsam kahest võrreldavast mudelist võis aga ei pruukinud olla õige. Nullhüpotees väitis, et ka lihtsam mudel on õige. Mis juhtub aga siis, kui lihtsam mudel pole õige? Kas ja millise tõenäosusega suudame F-testi abil tõestada, et lihtsam mudel pole õige? Alljärgnevalt arvutamegi, kui hea on F-testi võime tõestada alternatiivset hüpoteesi ehk milline on F-testi võimsus.

Dramatis personae:

- $\mathbf{M}_0 : \mathbf{y} = \mathbf{X}_0\boldsymbol{\beta}_0 + \boldsymbol{\varepsilon}$.
Lihtsaim mudelitest. Nullhüpoteesi väide. Siin peatükis osutub valeks valikuks.
- $\mathbf{M}_1 : \mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\varepsilon}$.
Rikkam meie poolt kaalutav mudel. Õige mudel (kuid võib sisaldada mittevajalikke parameetreid, osad parameetervektori $\boldsymbol{\beta}_1$ elemendid võivad olla ka nullid).

Meenutame varem tõesatud teoreemi 3.1:

Olgu $\mathbf{y} \sim N(\boldsymbol{\mu}; \mathbf{I})$, ja olgu maatriks \mathbf{A} ortogonaalne projektor (sümmeetriline ja idempotentne). Siis $\mathbf{y}^T \mathbf{A} \mathbf{y} \sim \chi^2(r, ncp)$, kus $r := \text{rank}(\mathbf{A})$, $ncp := \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}$.

Ülalmainitud teoreemist järeldus (vaata näiteks valemit 3.1), et

$$\mathbf{y}^T (\mathbf{I} - \mathbf{P}_{\mathbf{X}_1}) \mathbf{y} / \sigma^2 \sim \chi_{df=n-p_1}^2$$

ja juhul kui kehtib keerukam mudel, $E(\mathbf{y}) = \mathbf{X}_1\boldsymbol{\beta}_1$:

$$\begin{aligned} (\mathbf{y}/\sigma)^T (\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}) (\mathbf{y}/\sigma) &\sim \chi_{df=p_1-p_0; ncp=E\mathbf{y}^T (\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}) \mathbf{E}\mathbf{y} / \sigma^2} \\ \mathbf{y}^T (\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}) \mathbf{y} / \sigma^2 &\sim \chi_{df=p_1-p_0; ncp=\boldsymbol{\beta}_1^T \mathbf{X}_1^T (\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}) \mathbf{X}_1 \boldsymbol{\beta}_1 / \sigma^2} \end{aligned}$$

Juhul, kui õige on vaid keerukam mudel, siis antud avaldises mittetsentraalsuse parameeter ei muutu nulliks. Selgitame siinkohal veidi lähemalt mittetsentraalsuse parameetri arvutuseeskirja. Esmalt asendame kõik vaatlused oma andmestikus nende keskväärtustega. Seejärel leiame sellise vigadeta andmestiku pealt prognoosid samadele vaatlustele mõlema mudeli abil — nii lihtsama kui ka keerukama (õige) mudeli abil. Saadud prognooside erinevuste ruutude summa jagame jääkide dispersiooniga ja saamegi mittetsentraalsuse parameetri väärtuse.

Definitsioon 4.1 *Olgu*

$$X \sim \chi_{df_1; ncp}^2; \quad Y \sim \chi_{df_2}^2$$

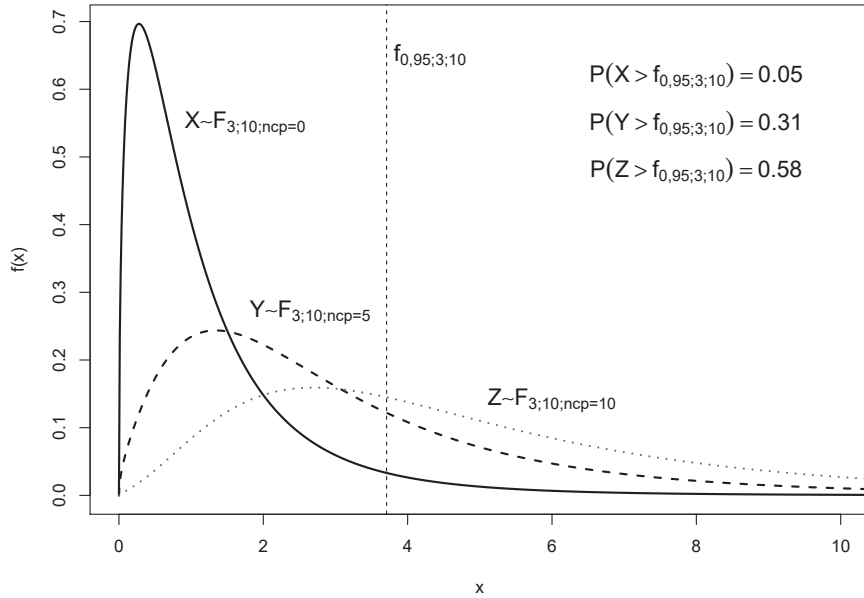
siis juhusliku suuruse

$$Z = \frac{X/df_1}{Y/df_2}$$

jaotuseks on mittetsentraalne F -jaotus vabadusastmete arvudega df_1 , df_2 ja mittetsentraalsuse parameetriga ncp .

Tsentraalse ja mittetsentraalse F -jaotuse tihedusfunktsioonid on toodud ka joonisel 4.2.

Joonis 4.2: Tsentraalne ja mittetsentraalne F -jaotus



Seega tegelik F -statistiku jaotus juhul, kui õige on alternatiivne hüpotees (keerukam mudel) on

$$F \sim F_{df_1=p_1-p_0; df_2=n-p_1; ncp=\beta_1^T \mathbf{X}_1^T (\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}) \mathbf{X}_1 \beta_1 / \sigma^2}.$$

Kuna testi võimsus on tõenäosus kummutada nullhüpotees (ehk saada kriitilisest väärtusest suurem teststatistiku väärtus), siis võime testi võimsust leida valemiga

$$1 - F_{df_1=p_1-p_0; df_2=n-p_1; ncp=\beta_1^T \mathbf{X}_1^T (\mathbf{P}_{\mathbf{X}_1} - \mathbf{P}_{\mathbf{X}_0}) \mathbf{X}_1 \beta_1 / \sigma^2} (f_{0,95;p_1-p_0;n-p_1}),$$

kus $f_{0,95;p_1-p_0;n-p_1}$ tähistab tsentraalse F-jaotuse $df_1 = p_1 - p_0 (= \text{rank}(\mathbf{X}_1) - \text{rank}(\mathbf{X}_0))$; $df_2 = n - p_1 (= n - \text{rank}(\mathbf{X}_1))$ 0,95-kvantiili.

Näide 4.4 Raha jätkub süvauuringute tegemiseks kümnele saarlasele, kümnele virumaalasele ja veerandsajale tartlasele (tartlaste puhul ei pea me maksma sõidu- ega ööbimiskulusi). Varasemate uuringute põhjal oletame, et uuritava tunnuse standardhälve on ligikaudu $\sigma \approx 25$. Uuritava tunnuse keskväärtused on ühe vastava ala tippspetsialisti arvates järgmised: $\mu_{saarlased} = 100$; $\mu_{virulased} = 125$; $\mu_{tartlased} = 108$. Milline on tõenäosus, et suudame F-testi abil tuvastada piirkondlike erinevuste olemasolu?

Lahendus:

Lihtsama mudeli $Y_i = \mu + \varepsilon_i$ prognoos kõigile inimestele (kui inimeste keskväärtused oleksid täpselt teada) oleks üldkeskmine,

$$\hat{\mu} = (100 \cdot 10 + 125 \cdot 10 + 108 \cdot 25) / (10 + 10 + 25) = 110.$$

Keerukama mudeli prognoosiks oleks saarlastele saarlaste keskväärtus jne. Prognooside erinevuste ruutude summa tuleks

$$(100 - 110)^2 \cdot 10 + (125 - 110)^2 \cdot 10 + (108 - 110)^2 \cdot 25 = 3350$$

ja mittetsentraalsuse parameetriks saaksime $ncp = 3350/25^2 = 5,36$. Testi võimsuse leidmiseks peaksime veel tabelist või arvutist välja uurima F-jaotuse ($df_1 = 3 - 1$, $df_2 = 45 - 3$) kriitilise väärtuse ($f_{0,95;2;42} = 3,22$) ja võimegi leida, millise tõenäosusega kavandatava uuringu abil on võimalik tõestada alternatiivse hüpoteesi paikapidavust (testi võimsus):

$$1 - F_{df_1=2; df_2=42; ncp=5,36}(3,22) = 0,5027 \dots$$

Kui kasutaksime kavandatud valimimahte, siis oleks regionaalsete erinevuste olemasolu tõestamise šansid ligikaudu 1:1.

Ülesanne

Kavatseme uurida Lõuna-Eestist Kesk-Eestist ja Põhja-Eestist pärit inimesi. Oletame, et uuritava tunnuse keskvärtus on Lõuna-Eestis 100, Kesk-Eestis 110, Põhja-Eestis 120. Uuritava tunnuse standardhälve Lõuna-Eestis on 20, oletame et ka muudes piirkondades võiks uuritava tunnuse standardhälve olla sarnane. Kogutud andmeid kasutades kavatseme hinnata dispersioonanalüüsi mudeli $Y_{ij} = \mu + \alpha_i + \varepsilon_{ij}$ ja kontrollida kas piirkonna mõju on oluline, st võrrelda hinnatud mudelit mudeliga $Y_{ij} = \mu + \varepsilon_{ij}$.

Kas saaksime võimsama testi kui kasutaksime 90 uuritavat või 70 uuritavat? Täpsemalt: kumb alltoodud katseplaanidest võimaldab keskvärtuste erinevust suurema tõenäosusega tuvastada? Leia mõlema katseplaani jaoks testi võimsus!

Plaan A		Plaan B	
n=90		n=70	
Piirkond	inimesi	Piirkond	inimesi
Lõuna-Eesti	30	Lõuna-Eesti	30
Kesk-Eesti	30	Kesk-Eesti	9
Põhja-Eesti	30	Põhja-Eesti	31

Vihjeks: F -jaotuse kvantiile ja jaotusfunktsiooni väärtust saab R-is vastavalt leida käskudega `qf` ja `pf`. Näiteks vabadusastmete arvudega $df_1 = 10$ ja $df_2 = 100$ ning mittetsentraalsuse parameetri väärtusega $n\text{cp} = 3$ F -jaotuse jaotusfunktsiooni kohal 2 saab leida käsuga `pf(2, df1=10, df2=100, ncp=3)`.