

- Oletetaan malli $\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$ ($r(\mathbf{X}) = p$).
- Tyypillisin esimerkki luottamusväleistä liittyy parametrivektorin $\boldsymbol{\beta}$ yksittäisiin komponentteihin β_j ($1 \leq j \leq p$).
- Tämä on erikoistapaus parametrivektorin $\boldsymbol{\beta}$ lineaarikombinaatiosta $\mathbf{a}'\boldsymbol{\beta} = a_1\beta_1 + \dots + a_p\beta_p$, jossa $\mathbf{a} \neq \mathbf{0}$ ($p \times 1$) on tunnettu
 - $\mathbf{a}' = [0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0] \Rightarrow \mathbf{a}'\boldsymbol{\beta} = \beta_j$
- Muita tyypillisiä erikoistapauksia lineaarikombinaatiosta $\mathbf{a}'\boldsymbol{\beta}$:
 - $\mathbf{a}' = [x_1^* \ \dots \ x_p^*] \Rightarrow \mathbf{a}'\boldsymbol{\beta} = \beta_1 x_1^* + \dots + \beta_p x_p^* = Y$:n odotusarvo, kun selittävälle muuttujille annetaan arvot x_1^*, \dots, x_p^*
 - odotusarvojen erotus $\mu_1 - \mu_2$ kahden riippumattoman normaalisen otoksen mallissa tai
 - vastaavat erotukset $\mu_j - \mu_k$ ($j \neq k$) eli ns. *kontrastit* yleisemmässä yksisuuntaisessa varianssianalyysimallissa
- Kuten tilastollisen päättelyn kurssilla todetaan, voidaan luottamusväljä muodostaa testien avulla.

- T -testi nollahypoteesille

$$H: \mathbf{a}'\boldsymbol{\beta} = \mathbf{a}'\boldsymbol{\beta}_0, \quad \boldsymbol{\beta}_0 (p \times 1) \text{ tunnettu}$$

saadaan kuten jaksossa 3.2 tapauksessa $\mathbf{a}' = [0 \cdots 0 \ 1 \ 0 \cdots 0]$:

$$T(\mathbf{Y}) = \frac{\mathbf{a}'\hat{\boldsymbol{\beta}} - \mathbf{a}'\boldsymbol{\beta}_0}{S\sqrt{\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a}}} \stackrel{H}{\sim} t_{n-p}.$$

- Vaihtoehtoa $\mathbf{a}'\boldsymbol{\beta} \neq \mathbf{a}'\boldsymbol{\beta}_0$ vastaava hylkäysalue merk. tasolla α :

$$C_\alpha(\mathbf{a}'\boldsymbol{\beta}_0) = \{\mathbf{y} : |T(\mathbf{y})| \geq t_{n-p}(\alpha/2)\}, \quad P(|T_{n-p}| \geq t_{n-p}(\frac{\alpha}{2})) = \alpha$$

- Vastaava hyväksymisalue muodostuu aineistoista, joilla

$$-t_{n-p}(\alpha/2) < T(\mathbf{y}) < t_{n-p}(\alpha/2)$$

tai yhtäpitävästi (käyttäen $T(\mathbf{y})$:n lauseketta)

$$\mathbf{a}'\hat{\boldsymbol{\beta}} - t_{n-p}(\frac{\alpha}{2})s\sqrt{\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a}} < \mathbf{a}'\boldsymbol{\beta}_0 < \mathbf{a}'\hat{\boldsymbol{\beta}} + t_{n-p}(\frac{\alpha}{2})s\sqrt{\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a}}$$

Luottamusväli lineaarikombinaatiolle

- Todennäköisyys (P_{β_0, σ^2}), että (nyt satunnaisiksi tulkitut) epäyhtälöt

$$\mathbf{a}'\hat{\boldsymbol{\beta}} - t_{n-p}(\frac{\alpha}{2})S\sqrt{\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a}} < \mathbf{a}'\boldsymbol{\beta}_0 < \mathbf{a}'\hat{\boldsymbol{\beta}} + t_{n-p}(\frac{\alpha}{2})S\sqrt{\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a}}$$

pätevät kaikilla $\boldsymbol{\beta}_0$ ja σ^2 on $1 - \alpha$.

- Siis, lineaarikombinaation $\mathbf{a}'\boldsymbol{\beta}$ luottamusväli luottamustasolla $1 - \alpha$ on

$$\mathbf{a}'\hat{\boldsymbol{\beta}} \pm t_{n-p}(\alpha/2) s\sqrt{\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a}}.$$

- Koska $\text{Var}(\mathbf{a}'\hat{\boldsymbol{\beta}}) = \sigma^2\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a}$ (Lause 2.1(i)), on $s\sqrt{\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a}}$ yllä estimaattorin $\mathbf{a}'\hat{\boldsymbol{\beta}}$ keskivirhe.

- Jos $\mathbf{a}'\boldsymbol{\beta} = \beta_j$, saadaan luottamusväli

$$\hat{\beta}_j \pm t_{n-p}(\alpha/2) s\sqrt{m^{jj}},$$

jossa $m^{jj} = \left[(\mathbf{X}'\mathbf{X})^{-1} \right]_{jj}$ on matriisin $(\mathbf{X}'\mathbf{X})^{-1}$ j . diagonaalialkio.

- Johdetaan luottamusjoukko parametrivektorille β kokonaisuudessaan.
- Tarkastellaan nollahypoteesia $H : \beta = \beta_0$, jossa β_0 ($p \times 1$) annettu.
- F -testisuure ($\mathbf{A} = \mathbf{I}_p$ ja $\mathbf{c} = \beta_0$)

$$F = (\hat{\beta} - \beta_0)' \mathbf{X}' \mathbf{X} (\hat{\beta} - \beta_0) / pS^2 \stackrel{H}{\sim} F_{p, n-p}.$$

- Jos $F_{p, n-p}(\alpha)$ toteuttaa $P(F_{p, n-p} \geq F_{p, n-p}(\alpha)) = \alpha$, kaikilla β_0 ja σ^2 pätee

$$P_{\beta_0, \sigma^2} \left((\hat{\beta} - \beta_0)' \mathbf{X}' \mathbf{X} (\hat{\beta} - \beta_0) / pS^2 < F_{p, n-p}(\alpha) \right) = 1 - \alpha.$$

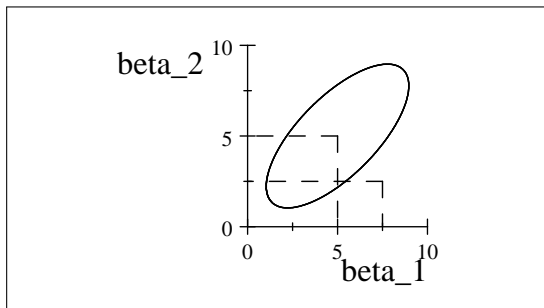
- β :n luottamusjoukko luottamustasolla $1 - \alpha$ on siten

$$\left\{ \beta \in \mathbb{R}^p : (\hat{\beta} - \beta)' \mathbf{X}' \mathbf{X} (\hat{\beta} - \beta) / ps^2 < F_{p, n-p}(\alpha) \right\}.$$

- β :n luottamusjoukko luottamustasolla $1 - \alpha$ on

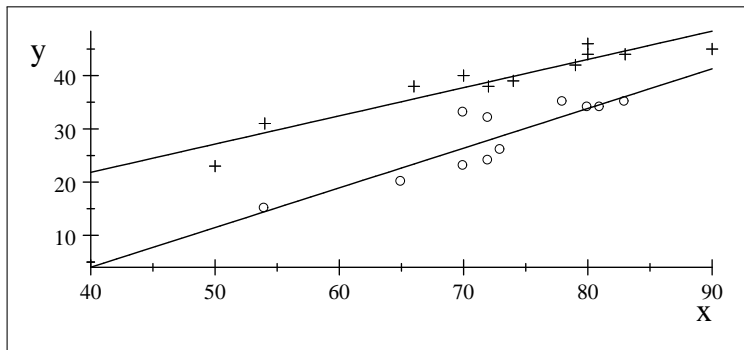
$$\left\{ \beta \in \mathbb{R}^p : (\hat{\beta} - \beta)' \mathbf{X}' \mathbf{X} (\hat{\beta} - \beta) / ps^2 < F_{p, n-p}(\alpha) \right\}$$

- Sen rajoittama pinta on \mathbb{R}^p :n ellipsoidi, jonka keskipiste on $\hat{\beta}$ ja muodon määrää matriisi $\mathbf{X}' \mathbf{X}$.



Kuva 4.1. $\beta = (\beta_1, \beta_2)$:n luottamusellipsi. PNS-estimaatti $\hat{\beta} = (5, 5)$.

- **Aineisto:** 22 satunnaisesti valittua karitsaa
- **Koeryhmä:** $i = 1, \dots, 11$; uusi ruokavalio
- **Kontrolliryhmä:** $i = 11, \dots, 22$; vanha ruokavalio
- **Kysymys:** Onko ruokavalioiden välillä eroa ja riippuuko mahdollinen ero ruokinta-ajan pituudesta, joka vaihteli 50 ja 90 päivän välillä?
- **Selitettävä muuttuja y :** Karitsan painon nousu koeajan aikana (mitattuna nauloissa)
- **Selittävät muuttujat:** Koeajan pituus x sekä ryhmää osoittavat indikaattorimuuttujat d_1 (koeryhmä) ja d_2 (kontrolliryhmä)



Kuva 5.1. Karitsojen painonnouluun (y) ja koeajan pituuteen (x) liittyv \ddot{a} aineisto koeryhm \ddot{a} n (+) ja kontrolliryhm \ddot{a} n (o) mukaan luokiteltuna sek \ddot{a} kummallekin ryhm \ddot{a} lle PNS:ll \ddot{a} estimoidut regressiosuorat

Koeryhm \ddot{a} :
$$y_i = \underset{(4.68)}{0.64} + \underset{(0.064)}{0.53} x_i + \hat{\varepsilon}_i, \quad s_1^2 = 5.98$$

Kontrolliryhm \ddot{a} :
$$y_i = \underset{(9.78)}{-25.80} + \underset{(0.13)}{0.75} x_i + \hat{\varepsilon}_i, \quad s_2^2 = 12.23.$$

Koeryhmä:
$$y_i = \underset{(4.68)}{0.64} + \underset{(0.064)}{0.53} x_i + \hat{\varepsilon}_i, \quad s_1^2 = 5.98$$

Kontrolliryhmä:
$$y_i = \underset{(9.78)}{-25.80} + \underset{(0.13)}{0.75} x_i + \hat{\varepsilon}_i, \quad s_2^2 = 12.23.$$

- **Kiinnostavia kysymyksiä:** Ovatko regressiosuorat yhdensuuntaisia eli voidaanko estimaattien 0.53 ja 0.75 välinen ero tulkita pelkästään satunnaisvaihtelusta johtuvaksi.
- Jos voidaan, on ruokavalion vaikutuksen mahdollinen ero riippumaton ruokinta-ajan pituudesta, jolloin regressiosuorien vakioiden ero mittaa sitä kaikilla selittävän muuttujan arvoilla.
- Jos regressiosuoria voidaan pitää yhdensuuntaisina, on seuraava kiinnostava kysymys siten voidaanko regressiosuorien vakioiden 0.64 ja -25.80 ero tulkita pelkästään satunnaisvaihtelusta johtuvaksi.

- Kiinnostavien kysymysten tutkiminen perustetaan lineaariseen malliin

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_{11} \\ Y_{12} \\ \vdots \\ Y_{22} \end{bmatrix} = \begin{bmatrix} 1 & 0 & x_1 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & x_{11} & 0 \\ 0 & 1 & 0 & x_{12} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 1 & 0 & x_{22} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \gamma_1 \\ \gamma_2 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_{11} \\ \varepsilon_{12} \\ \vdots \\ \varepsilon_{22} \end{bmatrix},$$

jossa $\varepsilon_1, \dots, \varepsilon_{22} \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$.

- PNS $\Rightarrow \hat{\alpha}_1 = 0.64, \hat{\gamma}_1 = 0.53, \hat{\alpha}_2 = -25.80, \hat{\gamma}_2 = 0.75$ ja $s^2 = 9.10$ (vrt. aikaisempi).
- Ensimmäinen testattava hypoteesi on $\gamma_1 = \gamma_2$ ja, jos se jää voimaan, testataan hypoteesia $\alpha_1 = \alpha_2$.
- Ennen testaamista on kuitenkin syytä tutkia ovatko mallista tehdyt oletukset realistiset.

Mallin oletusten tarkistaminen

- Mallin oletusten tarkistamisessa residuaalit ovat keskeisiä. Tarkastellaan residuaalivektoria

$$\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}, \quad \hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y},$$

jolla on esitys

$$\hat{\boldsymbol{\varepsilon}} = (\mathbf{I}_n - \mathbf{P}) \mathbf{y}, \quad \mathbf{P} = [p_{ij}] = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \quad (\text{projektio}).$$

- Mallin ollessa oikein spesifioitu pätee vastaavalle sv:lle

$$E(\hat{\boldsymbol{\varepsilon}}) = \mathbf{0} \quad \text{ja} \quad \text{Cov}(\hat{\boldsymbol{\varepsilon}}) = \sigma^2 (\mathbf{I}_n - \mathbf{P}) \neq \sigma^2 \mathbf{I}_n = \text{Cov}(\boldsymbol{\varepsilon}).$$

- Residuaalit siis oikeinkin mallin kyseessä ollessa korreloituneita ja niiden varianssit vaihtelevat havaintoyksiköstä toiseen. Lisäksi, jos normaalisuusoletus $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ pätee, noudattaa residuaalivektori singulaarista normaalijakaumaa ($r(\text{Cov}(\hat{\boldsymbol{\varepsilon}})) = n - p$).

Mallin oletusten tarkistaminen

- Koska $\text{Var}(\hat{\varepsilon}_i) = \sigma^2(1 - p_{ii})$, on luontevaa tarkastella standardoituja residuaaleja

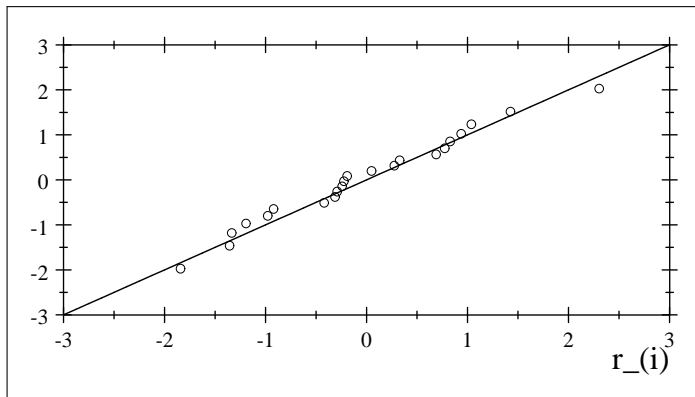
$$r_i = \hat{\varepsilon}_i / s\sqrt{1 - p_{ii}}, \quad s^2 = \frac{1}{n - p} \hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}},$$

joiden varianssi on likimain 1.

- Järjestetään standardoidut residuaalit $r_{(1)} \leq r_{(2)} \leq \dots \leq r_{(n)}$.
- Jos $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$, pätee (Φ^{-1} on $N(0, 1)$ -jakauman kf:n käänteisfunktio)

$$E(R_{(i)}) \approx \Phi^{-1} \left(\frac{i - 0.5}{n} \right) \approx r_{(i)}.$$

- Normaalisuusetusta voidaan siis tutkia tutkimalla ovatko pisteet $\left(r_{(i)}, \Phi^{-1} \left(\frac{i - 0.5}{n} \right) \right)$, $i = 1, \dots, n$, likimain 45° :n suoralla.



Kuva 5.2. Mallin järjestetyt standardoidut residuaalit arvoja $\Phi^{-1}((i - 0.5) / n)$ ($i = 1, \dots, 22$) vastaan eli ns. normaalipaperipiirros (normal probability plot tai normal QQ plot).

Tämän perusteella normalisuusoletus näyttää kohtuulliselta.

Mallin oletusten tarkistaminen

- Muista oletuksista riippumattomuusoletus tuntuu perustellulta, koska kysymyksessä on satunnaisotanta, mutta vaihtelee $\text{Var}(\varepsilon_i)$ ryhmien mukaan?
- Koska virheiden $\varepsilon_1, \dots, \varepsilon_{22}$ riippumattomuus voidaan olettaa, on $S_1^2 \perp\!\!\!\perp S_2^2$ ja (Lause 2.1(ii))

$$\frac{S_1^2}{S_2^2} \stackrel{H_0}{\sim} \frac{1}{9}\chi_{1,9}^2 / \frac{1}{9}\chi_{2,9}^2 \sim F_{9,9} \quad \left(\chi_{i,9}^2 \sim \chi_9^2 \text{ ja } \chi_{1,9}^2 \perp\!\!\!\perp \chi_{2,9}^2 \right).$$

Testisuureen arvoksi saadaan

$$\frac{s_1^2}{s_2^2} = \frac{5.98}{12.23} = 0.49$$

ja P-arvoksi $P(F_{9,9} \leq 0.49) = 0.15$. Vakiovarianssioletus tuntuu siten kohtuulliselta.

- Rajoittamaton malli (malli (1)) on

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_{11} \\ Y_{12} \\ \vdots \\ Y_{22} \end{bmatrix} = \begin{bmatrix} 1 & 0 & x_1 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & x_{11} & 0 \\ 0 & 1 & 0 & x_{12} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 1 & 0 & x_{22} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \gamma_1 \\ \gamma_2 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_{11} \\ \varepsilon_{12} \\ \vdots \\ \varepsilon_{22} \end{bmatrix}, \quad \varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n),$$

josta estimaatit $\hat{\alpha}_1 = 0.64$, $\hat{\gamma}_1 = 0.53$, $\hat{\alpha}_2 = -25.80$, ja $\hat{\gamma}_2 = 0.75$.

- Tarkastellaan hypoteesia $H_0 : \gamma_1 = \gamma_2 = \gamma$.

- Rajoitettu malli on

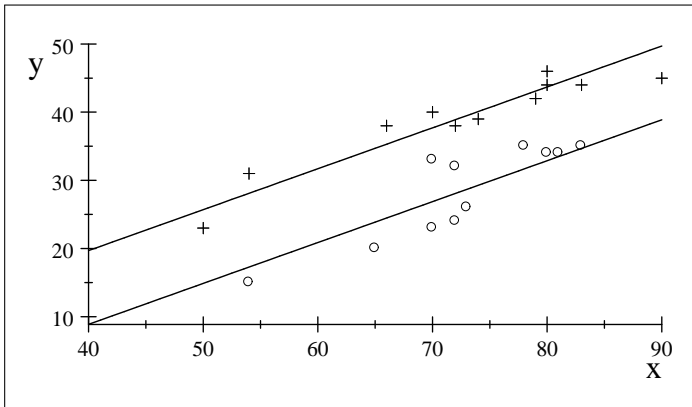
$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_{11} \\ Y_{12} \\ \vdots \\ Y_{22} \end{bmatrix} = \begin{bmatrix} 1 & 0 & x_1 \\ \vdots & \vdots & \vdots \\ 1 & 0 & x_{11} \\ 0 & 1 & x_{12} \\ \vdots & \vdots & \vdots \\ 0 & 1 & x_{22} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \gamma \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_{11} \\ \varepsilon_{12} \\ \vdots \\ \varepsilon_{22} \end{bmatrix}, \quad \varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n).$$

- Virhevarianssin estimaatiksi tulee $s_H^2 = 9.76$ ja $n - p_H = 22 - 3 = 19$.
- Vapaasta mallista saatiin $s^2 = 9.10$ ja $n - p = 22 - 4 = 18$.
- Testisuure saa arvon

$$F = \frac{(19 \times 9.76 - 18 \times 9.10) / 1}{9.10} = 2.37$$

ja vastaavaksi P-arvoksi tulee

$$P = P(F_{1,18} \geq 2.37) = 0.14 \Rightarrow H_0 : \gamma_1 = \gamma_2 \text{ jää voimaan.}$$



Kuva 5.3. Kuvan 5.1 aineistoon piirretyt mallin (2) estimointituloksiin perustuvat samansuuntaiset regressiosuorat.

Estimoitu malli

$$y_i = \underset{(4.98)}{-4.30} d_{i1} - \underset{(4.98)}{15.12} d_{i2} + \underset{(0.067)}{0.60} x_i + \hat{\varepsilon}_i, \quad s^2 = 9.76 \quad (2)$$

- Tarkastellaan nyt hypoteesia $H_0 : \alpha_1 = \alpha_2 = \alpha$ mallissa (2)

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_{11} \\ Y_{12} \\ \vdots \\ Y_{22} \end{bmatrix} = \begin{bmatrix} 1 & 0 & x_1 \\ \vdots & \vdots & \vdots \\ 1 & 0 & x_{11} \\ 0 & 1 & x_{12} \\ \vdots & \vdots & \vdots \\ 0 & 1 & x_{22} \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \gamma \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_{11} \\ \varepsilon_{12} \\ \vdots \\ \varepsilon_{22} \end{bmatrix}, \quad \varepsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n).$$

- Rajoitettu malli on

$$Y_i = \alpha + \gamma x_i + \varepsilon_i, \quad i = 1, \dots, 22.$$

PNS-estimointi \Rightarrow

$$y_i = \underset{(10.16)}{-9.71} + \underset{(0.14)}{0.60} x_i + \hat{\varepsilon}_i, \quad i = 1, \dots, 22, \quad s^2 = 41.46 \quad (3)$$

- Vapaassa mallissa (2)

$$s^2 = 9.76, \quad n - p = 22 - 3 = 19$$

- Rajoitetussa mallissa (3)

$$s^2 = 41.46, \quad n - p_H = 22 - 2 = 20$$

- Testisuure saa arvon

$$F = \frac{(20 \times 41.46 - 19 \times 9.76) / 1}{9.76} = 65.97$$

ja vastaavaksi P-arvoksi tulee

$$P = P(F_{1,19} \geq 65.97) \approx 0 \Rightarrow H_0 : \alpha_1 = \alpha_2 \text{ on syytä hylätä.}$$

- Malli (2)

$$Y_i = \alpha_1 d_{i1} + \alpha_2 d_{i2} + \gamma x_i + \varepsilon_i, \quad i = 1, \dots, 22.$$

- Ruokavalioiden eroa mittaa parametri

$$\alpha_1 - \alpha_2 = [1 \quad -1 \quad 0] \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \gamma \end{bmatrix} = \mathbf{a}'\boldsymbol{\beta}.$$

- $\hat{\alpha}_1 - \hat{\alpha}_2 = -4.30 + 15.12 = 10.82$, $s\sqrt{\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a}} = 1.77$,
 $t_{19}(0.05/2) = 2.09$
- 95%:n luottamusväli erotukselle $\alpha_1 - \alpha_2$ on

$$10.82 - 2.09 \cdot 1.77 < \alpha_1 - \alpha_2 < 10.82 + 2.09 \cdot 1.77$$

eli

$$7.12 < \alpha_1 - \alpha_2 < 14.52$$

- Uusi ruokavalio on (nopean painon nousun kannalta) systemaattisesti vanhaa parempi, sillä sen paremmuus ei riipu tarkasteltujen ruokinta-aikojen pituudesta.
- Kun ruokinta-ajan pituus on n. 50 - 90 päivää, voidaan uudella ruokavaliolla odottaa päästävän 95%:n varmuudella n. 7 - 14.5 naulaa suurempaan karitsan painon nousuun kuin vanhalla ruokavaliolla.

Yksisuuntainen varianssianalyysi

Yksisuuntaisessa varianssianalyysissä tarkastellaan p :tä ryhmää ja niistä satunnaisotannalla poimittuja havaintoja. Mielenkiinnon kohteena on ryhmien odotusarvojen mahdolliset erot. Asetelma on kaaviona

	Havainnot	Keskiarvot
Ryhmä 1	$y_{11}, y_{12}, \dots, y_{1n_1}$	\bar{y}_1
Ryhmä 2	$y_{21}, y_{22}, \dots, y_{2n_2}$	\bar{y}_2
\vdots	\vdots	\vdots
Ryhmä p	$y_{p1}, y_{p2}, \dots, y_{pn_p}$	\bar{y}_p

jossa $\bar{y}_j = (y_{j1} + \dots + y_{jn_j}) / n_j$ ($j = 1, \dots, p$).

Yksisuuntainen varianssianalyysi

- Tilastollista mallia varten oletetaan, että havaintoja vastaavat satunnaismuuttujat $Y_{11}, \dots, Y_{1n_1}, Y_{21}, \dots, Y_{2n_2}, \dots, Y_{p1}, \dots, Y_{pn_p}$ ovat riippumattomia ja $Y_{ji} \sim N(\mu_j, \sigma^2)$
- tai yhtäpitävästi, että

$$Y_{ji} = \mu_j + \varepsilon_{ji}, \quad i = 1, \dots, n_j, \quad j = 1, \dots, p,$$

jossa sm:t $\varepsilon_{ji} \sim N(0, \sigma^2)$ ovat riippumattomia ja $\mu_j \in \mathbb{R}$.

- Käyttäen matriisimerkintöjä saadaan yhtälö

$$\begin{bmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \\ \vdots \\ \mathbf{Y}_p \end{bmatrix} = \begin{bmatrix} \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_p} \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_p \end{bmatrix},$$

josta nähdään, että kysymyksessä on lineaarisen mallin erikoistapaus.

Yksisuuntainen varianssianalyysi

- $Y_{ji} = \mu_j + \varepsilon_{ji}$, $i = 1, \dots, n_j$, $j = 1, \dots, p$,

jossa sm:t $\varepsilon_{ji} \sim N(0, \sigma^2)$ ovat riippumattomia ja $\mu_j \in \mathbb{R}$.

- Mielenkiinnon kohteena on nollahypoteesi

$$H : \mu_1 = \dots = \mu_p \Leftrightarrow \mu_1 - \mu_p = \mu_2 - \mu_p = \dots = \mu_{p-1} - \mu_p = 0,$$

tai matriisein

$$\begin{bmatrix} 1 & 0 & 0 & \dots & \dots & -1 \\ 0 & 1 & 0 & \dots & \dots & -1 \\ \vdots & \vdots & \ddots & & & \vdots \\ \vdots & \vdots & & \ddots & & \vdots \\ 0 & 0 & 0 & & 1 & -1 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \vdots \\ \mu_p \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ \vdots \\ 0 \end{bmatrix}.$$

- Nollahypoteesi on siis lineaarista muotoa $\mathbf{A}\boldsymbol{\mu} = \mathbf{0}$, jossa \mathbf{A} on $(p-1) \times p$ matriisi ja $r(\mathbf{A}) = p-1$.

Yksisuuntainen varianssianalyysi

- Kysymyksessä on lineaarinen malli

$$Y_{ji} = \mu_j + \varepsilon_{ji}, \quad i = 1, \dots, n_j, \quad j = 1, \dots, p, \quad (*)$$

jossa sm:t $\varepsilon_{ji} \sim N(0, \sigma^2)$ ovat riippumattomia ja $\mu_j \in \mathbb{R}$. Testattava (lineaarinen) hypoteesi

$$H: \mu_1 = \dots = \mu_p := \mu_0.$$

- Testisuureen yleinen lauseke on $(r(\mathbf{A}) = p - 1)$

$$F = \frac{(S(\hat{\boldsymbol{\mu}}_H) - S(\hat{\boldsymbol{\mu}})) / (p - 1)}{S(\hat{\boldsymbol{\mu}}) / (n - p)} \stackrel{H}{\sim} F_{p-1, n-p},$$

jossa parametrin $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)$ vapaa PNS-estimaatti $\hat{\boldsymbol{\mu}}$ saadaan mallista (*) ja sidottu PNS-estimaatti $\hat{\boldsymbol{\mu}}_H$ mallista

$$Y_{ji} = \mu_0 + \varepsilon_{ji}, \quad i = 1, \dots, n_j, \quad j = 1, \dots, p.$$

- Sidottu PNS-estimaatti $\hat{\mu}_H$ minimoi jäännöselineliösumman

$$S(\mu_0) = \sum_{j=1}^p \sum_{i=1}^{n_j} (y_{ji} - \mu_0)^2.$$

Siten, $\hat{\mu}_H = (\bar{y}, \dots, \bar{y})$, jossa

$$\bar{y} = \hat{\mu}_0 = \frac{1}{n} \sum_{j=1}^p \sum_{i=1}^{n_j} y_{ji} = \frac{1}{n} \sum_{j=1}^p n_j \bar{y}_j \quad (n = n_1 + \dots + n_p),$$

ja

$$S(\hat{\mu}_H) = \sum_{j=1}^p \sum_{i=1}^{n_j} (y_{ji} - \bar{y})^2.$$

- Vapaa PNS-estimaatti $\hat{\boldsymbol{\mu}}$ minimoi jäännösneliösumman

$$S(\boldsymbol{\mu}) = \sum_{j=1}^p \sum_{i=1}^{n_j} (y_{ji} - \mu_j)^2.$$

- Tämä johtaa estimaatteihin $\hat{\mu}_j = \bar{y}_j = (y_{j1} + \dots + y_{jn_j}) / n_j$, joten $\hat{\boldsymbol{\mu}} = (\bar{y}_1, \dots, \bar{y}_p)$ ja

$$S(\hat{\boldsymbol{\mu}}) = \sum_{j=1}^p \sum_{i=1}^{n_j} (y_{ji} - \bar{y}_j)^2$$

- Testisuureksi saadaan

$$F = \frac{S(\hat{\boldsymbol{\mu}}_H) - S(\hat{\boldsymbol{\mu}}) / (p - 1)}{S(\hat{\boldsymbol{\mu}}) / (n - p)}$$
$$= \frac{\left(\sum_{j=1}^p \sum_{i=1}^{n_j} (Y_{ji} - \bar{Y})^2 - \sum_{j=1}^p \sum_{i=1}^{n_j} (Y_{ji} - \bar{Y}_j)^2 \right) / (p - 1)}{\sum_{j=1}^p \sum_{i=1}^{n_j} (Y_{ji} - \bar{Y}_j)^2 / (n - p)}$$
$$\stackrel{H}{\sim} F_{p-1, n-p},$$

- Tämä riittää, mutta testisuure on perinteisesti esitetty hieman toisessa muodossa.

- Laskemalla nähdään, että F -testisuureen osoittaja voidaan kirjoittaa

$$S(\hat{\boldsymbol{\mu}}_H) - S(\hat{\boldsymbol{\mu}}) = \sum_{j=1}^p \sum_{i=1}^{n_j} (\bar{y}_j - \bar{y})^2 = \sum_{j=1}^p n_j (\bar{y}_j - \bar{y})^2,$$

joten F -testisuureksi saadaan

$$F = \frac{\sum_{j=1}^p n_j (\bar{Y}_j - \bar{Y})^2 / (p-1)}{\sum_{j=1}^p \sum_{i=1}^{n_j} (Y_{ji} - \bar{Y}_j)^2 / (n-p)} \stackrel{H}{\sim} F_{p-1, n-p}.$$

- Usein sanotaan, että testisuure on ryhmien välisen varianssin ja ryhmien sisäisen varianssin suhde, mistä nimitys varianssianalyysi tulee.
- Huomaa, että testisuureen nimittäjässä on $S(\hat{\boldsymbol{\mu}}) / (n-p) = S^2$.