

Uporaba statistike na področju kvantitativne genetike

Gregor Gorjanc

Univerza v Ljubljani, Biotehniška Fakulteta, Oddelek za zootehniko, Domžale (Rodica)



IBMI
Ljubljana, Slovenija
11. oktober 2010

UL, BF, Oddelek za zootehniko, Domžale (Rodica)



Izvleček

Kvantitativna genetika skuša odgovoriti na vprašanje kolikšen del fenotipske variabilnosti je povzročen z genetsko variabilnostjo.

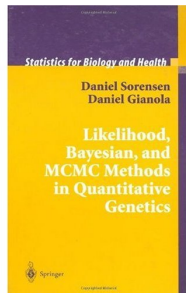
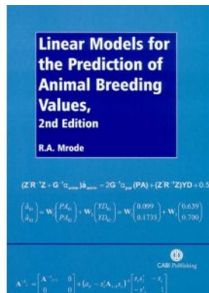
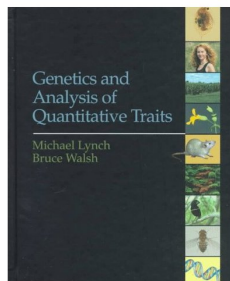
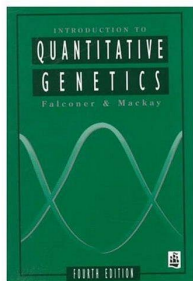
Vloga statističnih metod pri iskanju tega odgovora je ključna in razvoj nekaterih metod (npr. regresija, analiza variance, mešani model, model s pragovi, ...) je potekal hkrati z metodološkim razvojem na področju genetike kvantitativnih lastnosti.

Tekom predavanja si bomo ogledali te metode s stališča statistike in kvantitativne genetike ter njihovo uporabo na področju selekcije domačih živali.

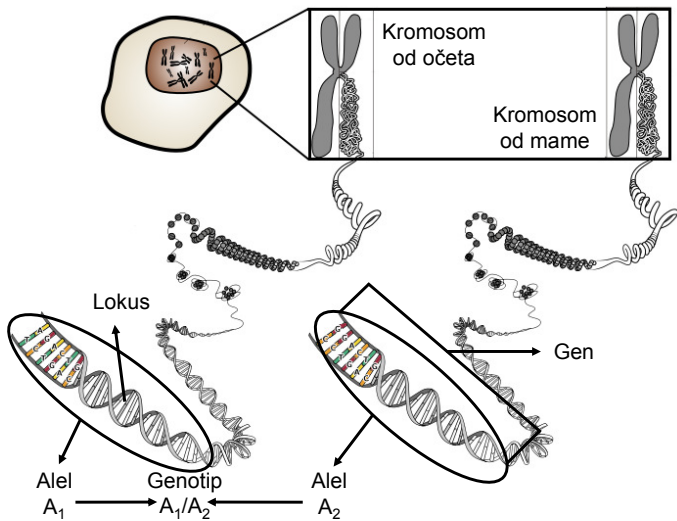
Kazalo

1. Uvod
2. Začetki kvantitativne genetike
3. Uporaba kvantitativne genetike pri selekciji domačih živali

Nekaj literature



Genetika: DNA → beljakovine → fenotip



Področja genetike

- ▶ Populacijska genetika
- ▶ Kvantitativna genetika
- ▶ Evolucijska genetika
- ▶ Molekularna genetika
- ▶ Genetika
 - ▶ človeka
 - ▶ živali
 - ▶ rastlin
 - ▶ mikrobov
 - ▶ ...
- ▶ ...

Podobna “širina” in “prepredanje” kot pri področjih statistike

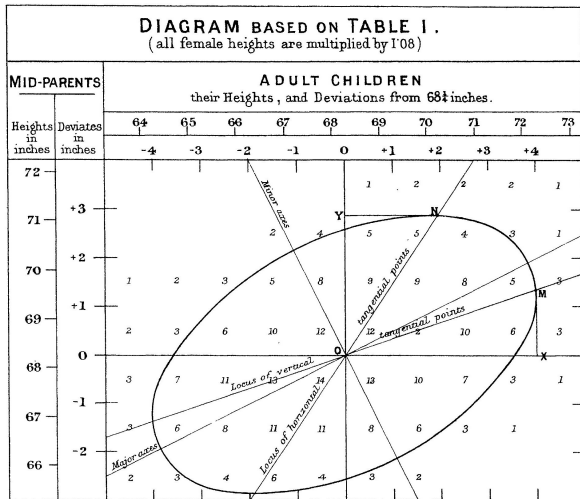
2. Začetki kvantitativne genetike

Začetki kvantitativne genetike

- ▶ Začetki genetike
 - ▶ študije vpliva genov na “enostavne” (diskretne) lastnosti (Mendel, Bateson, ...) → **mendlisti**
 - ▶ študije “kompleksnih” (kvantitativnih) lastnosti (Pearson, Galton, Fisher, Wright, ...) → **biometriki**
- ▶ **Galton** (1869) - regresija
- ▶ Nestrinjanje med mendlisti in biometriki ob koncu 19. stoletja
- ▶ **Fisher** (1918) - infinitezimalni model
- ▶ **Wright, Malecot** - inbriding in sorodstvo
- ▶ ...
- ▶ Razcvet, upad in ponovno razcvet z “genomsko revolucijo”

Galton (1886) - regresija

- ▶ Analiziral povprečno telesno višino staršev in njihovih (odraslih) otrok



Razprava: Zakaj je povezava (korelacija) pozitivna? ▶

Regresija

- ▶ Bivariatna normalna porazdelitev fenotipskih vrednosti
 - ▶ telesna višina očeta - P_f
 - ▶ telesna višina sina - P_s

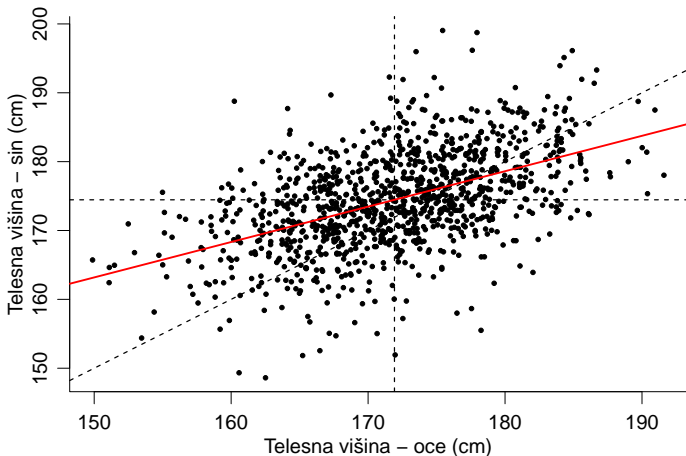
$$p(P_f, P_s) \sim N \left(\begin{array}{c} E(P_f) \\ E(P_s) \end{array}, \begin{array}{cc} \text{Var}(P_f) & \text{Cov}(P_f, P_s) \\ \text{sym.} & \text{Var}(P_s) \end{array} \right)$$

- ▶ Pogojno pričakovanje

$$E(P_s|P_f) = E(P_s) + \text{Cov}(P_s, P_f) (\text{Var}(P_f))^{-1} (P_f - E(P_f))$$

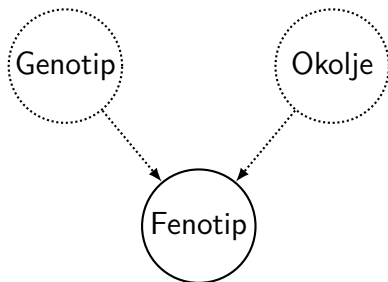
- ▶ Regresijski koeficient: $\text{Cov}(P_s, P_f) (\text{Var}(P_f))^{-1}$

Regresija - Pearsonovi podatki



Razprava: Zakaj “regression/shrinkage” k povprečju?

Fisher (1918) - dekompozicija fenotipske vrednosti



► Koncept vrednosti

$$P = \mu + G + E$$

- P - fenotipska vrednost - **lahko merimo!!!**
- G - genotipska vrednost - **ne moremo meriti!!!**
- E - deviacija zaradi okolja

Fisher (1918) - dekompozicija genotipske vrednosti

▶ Genotipska vrednost

- ▶ En lokus: $G = A + D$
- ▶ Več kot en lokus: $G = A + D + I$
- ▶ **A - additivna genotipska (plemenska) vrednost**
(učinek posameznih alelov)
- ▶ **D - deviacija zaradi dominance**
(interakcije med aleli na enem lokusu)
- ▶ **I - deviacija zaradi epistaze: $A \times A, A \times D, D \times D$**
(interakcije med aleli na različnih lokusih)

$$P = \mu + A + D + I + E$$

▶ Centralni limitni izrek

- ▶ veliko število genov z majhnim učinkom + naključni vplivi iz okolja → normalna porazdelitev

Fisher (1918) - analiza variance

- ▶ Analiza variance za **predpostavljeni** model

$$P \approx \mu + A + D + E$$

- ▶ Vzročne komponente fenotipske variance

$$\sigma_p^2 = \sigma_a^2 + \sigma_d^2 + \sigma_e^2$$

- ▶ Dednostni delež (heritabiliteta)

- ▶ v širšem smislu $h^2 = \sigma_g^2 / \sigma_p^2$
- ▶ v ožjem smislu $h^2 = \sigma_a^2 / \sigma_p^2$

- ▶ Opazovane komponente variance - to lahko ocenimo iz zbranih podatkov

- ▶ varianca med skupinami σ_{med}^2
- ▶ varianca znotraj skupin $\sigma_{znotraj}^2$

Fisher (1918) - analiza variance II

- ▶ Analiza variance skupin sorodnikov → opazovane komponente variance vsebujejo različne vzročne komponente variance

$$\begin{aligned}\sigma_{med}^2 &= r_a \sigma_a^2 + r_d \sigma_d^2 + r_c \sigma_{ec}^2 \\ \sigma_{znotraj}^2 &= \sigma_P^2 - \sigma_{med}^2\end{aligned}$$

- ▶ r_a = verjetnost, da imata sorodnika enake alele
- ▶ r_d = verjetnost, da imata sorodnika enak genotip
- ▶ r_c = "verjetnost", da imata sorodnika "skupno okolje"

Skupine / σ_{med}^2	r_a	r_d	r_c
Enojajčni dvojčki	1	1	1
Dvojajčni dvojčki	1/2	1/4	1
Bratje in sestre (FSIB)	1/2	1/4	1
Pol-bratje in pol-sestre (HSIB)	1/4	0	0

Galton (1886) - revizija

- ▶ Predpostavljeni model za telesno višino sinov (P_s)

$$\begin{aligned}P_s &\approx \mu + A_s + E_s \\ &\approx \mu + 1/2A_f + 1/2A_m + E_s\end{aligned}$$

- ▶ Pogojno pričakovanje na podlagi telesne višine očetov (P_f)

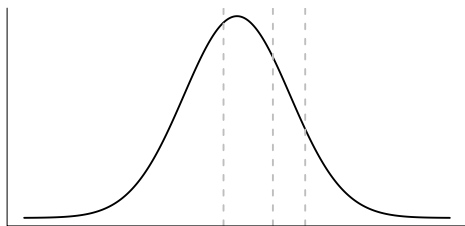
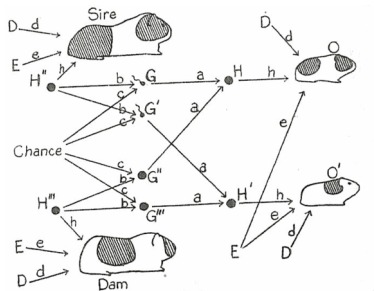
$$E(P_s|P_f) = E(P_s) + \text{Cov}(P_s, P_f) (\text{Var}(P_f))^{-1} (P_f - E(P_f))$$

- ▶ pričakovane vrednosti: $E(P_s) = E(P_f) = \mu$
- ▶ varianca med očeti: $\text{Var}(P_f) = \sigma_p^2$
- ▶ covarianca med sinovi in očeti:
 $\text{Cov}(P_s, P_f) = \text{Cov}(1/2A_f, A_f) = 1/2\sigma_a^2$
- ▶ regresijski koeficient

$$\text{Cov}(P_s, P_f) (\text{Var}(P_f))^{-1} = \frac{1/2\sigma_a^2}{\sigma_p^2} = 1/2h^2 < 1$$

Wright

- ▶ S korelacijami definirali koncept sorodnosti na osnovi informacije iz rodovnikov:
 - ▶ **koef. inbridinga** - korelacija med gametami posameznika
 - ▶ **koef. sorodnosti** - korelacija med gametami posameznikov
- ▶ **Stezna metoda** (ang. path analysis) in **model s pragovi** (1934) (ang. threshold model) = probit model



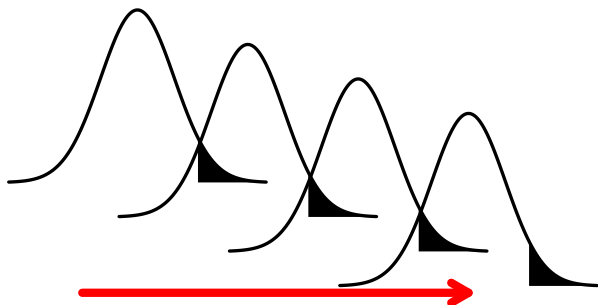
Raziskovalna “področja”

- ▶ **ljudje** → biostatistika, genetika človeka (ang. human genetics)
 - ▶ cilj: splošno znanje in **zdravljenje**
- ▶ **rastline** → genetika rastlin (ang. plant breeding and genetics)
 - ▶ cilj: → splošno znanje in **selekcija** (=žlahtnjenje)
- ▶ **živali** → genetika živali (ang. animal breeding and genetics)
 - ▶ cilj: splošno znanje in **selekcija**
 - ▶ “animal breeding and genetics”
 - = živinoreja + genetika + statistika + ...

3. Uporaba kvantitativne genetike pri selekciji domačih živali

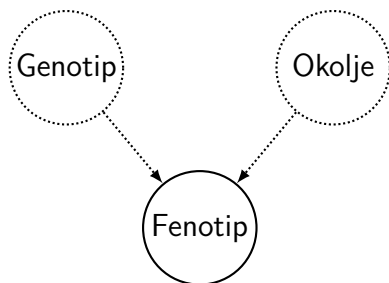
Selekcija

- ▶ Izmerimo **fenotipske vrednosti** kandidatov in izberemo (selekcionirami) tiste z najbolj zaželenimi vrednostmi (ang. **mass/phenotype selection**)
- ▶ Izbrani kandidati bodo starši naslednje generacije



Razprava: Ali se fenotipska vrednost prenaša s staršev na potomcev?

Dekompozicija fenotipske vrednosti



- ▶ **Genetsko vrednotenje** = statistično sklepanje o genotipski (plemnski) vrednosti posameznikov glede na zbrane **podatke** in **predpostavljeni model** (= **BLUP selection**)

Henderson (1949+) - mešani model

- ▶ Mešani model - fiksni/sistematski (**b**) in naključni (**a**) vplivi

$$y_{ijk} = \mu + b_i + a_j + e_{ijk}$$

$$\mathbf{y} = \mathbf{Xb} + \mathbf{Za} + \mathbf{e}$$

- ▶ Predpostavke

$$p(\mathbf{a}) \sim N(\mathbf{0}, \mathbf{G}), \quad \mathbf{G} = \mathbf{A}\sigma_a^2$$

$$p(\mathbf{e}) \sim N(\mathbf{0}, \mathbf{R}), \quad \mathbf{R} = \mathbf{I}\sigma_e^2$$

$$p \begin{pmatrix} \mathbf{y} \\ \mathbf{a} \\ \mathbf{e} \end{pmatrix} \sim N \begin{pmatrix} \mathbf{Xb} & \mathbf{ZGZ}^T + \mathbf{R} & \mathbf{ZG}^T & \mathbf{R} \\ \mathbf{0} & & \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \text{sym.} & & \mathbf{R} \end{pmatrix}$$

A - matrika sorodstva (Wright)

Mešani model - poimenovanje

- ▶ Henderson: model živali, ker je modeliral prirejo živali
 - ▶ model očeta - upoštevamo rodovnik samo med očeti
 - ▶ model očeta in mame - upoštevamo rodovnik samo med očeti
 - ▶ ...
 - ▶ bolj splošno: **mešani model z rodovniki**
- ▶ Mešani model (ang. mixed model)
- ▶ Hierarhični model (ang. hierarchical model)
- ▶ Večnivojski??? model (ang. multilevel model)
- ▶ ...

Henderson (1949+) - sistem enačb

- ▶ Metoda najmanjših kvadratov

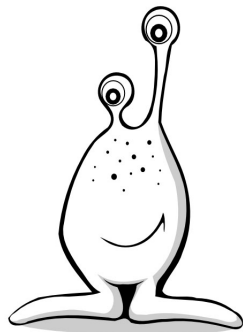
$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$$
$$\left(\mathbf{X}^T \mathbf{R}^{-1} \mathbf{X} \right) \left(\hat{\mathbf{b}} \right) = \left(\mathbf{X}^T \mathbf{R}^{-1} \mathbf{y} \right)$$

- ▶ Mešani model

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{a} + \mathbf{e}$$
$$\hat{\mathbf{a}} = \mathbf{G}\mathbf{Z}^T \mathbf{V}^{-1} \left(\mathbf{y} - \mathbf{X}\hat{\mathbf{b}} \right)$$
$$\left(\begin{array}{cc} \mathbf{X}^T \mathbf{R}^{-1} \mathbf{X} & \mathbf{X}^T \mathbf{R}^{-1} \mathbf{Z} \\ \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{X} & \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{Z} + \mathbf{G}^{-1} \end{array} \right) \left(\begin{array}{c} \hat{\mathbf{b}} \\ \hat{\mathbf{a}} \end{array} \right) = \left(\begin{array}{c} \mathbf{X}^T \mathbf{R}^{-1} \mathbf{y} \\ \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{y} \end{array} \right)$$
$$\mathbf{G}^{-1} = \mathbf{A}^{-1} \sigma_a^{-2} \quad \mathbf{R}^{-1} = \mathbf{I} \sigma_e^{-2}$$

- ▶ $\hat{\mathbf{b}}$ - najboljša linearna nepristranska cenilka (**BLUE**)
- ▶ $\hat{\mathbf{a}}$ - najboljša linearna nepristranska napoved??? (**BLUP**)

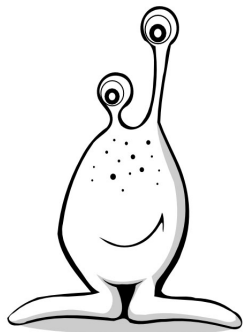
Primer "vesoljcev" - podatki



Posameznik	Oče	Mama	Skupina	Fenotip
1	/	/	/	/
2	/	/	1	103, 106
3	2	1	1	98
4	2	/	2	101
5	4	3	2	106
6	2	3	2	93
7	5	6	/	/
8	5	6	/	/
9	/	/	/	/
10	8	9	1	109

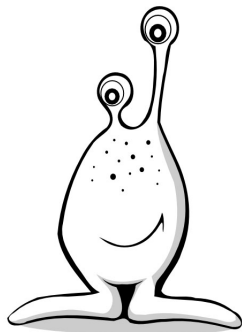
Slika: Jouke

Primer "vesoljcev" - grafični model

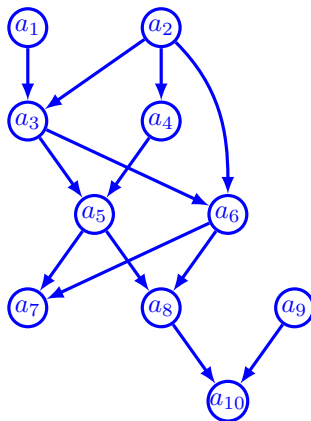


Slika: Jouke

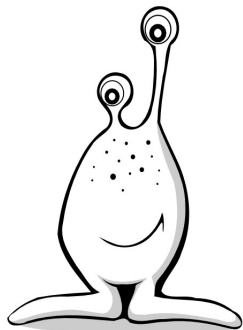
Primer "vesoljcev" - grafični model



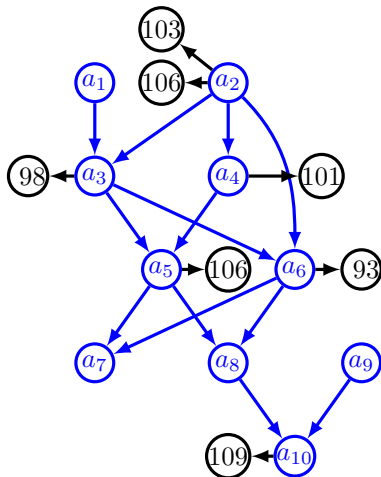
Slika: Jouke



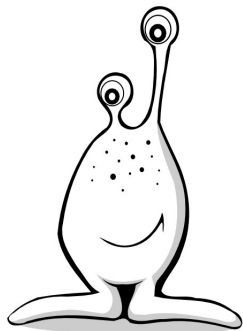
Primer "vesoljcev" - grafični model



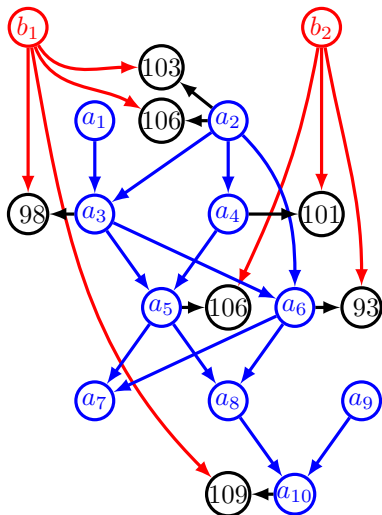
Slika: Jouke



Primer "vesoljcev" - grafični model



Slika: Jouke



Primer “vesoljcev” - R

Demonstracija v R-ju

Programje

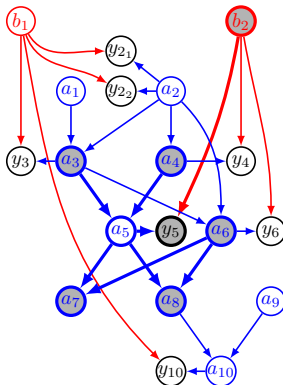
- ▶ Aplikacije lahko zajemajo tudi več 100 tisoč ali milijon živali in toliko ali še več fenotipskih vrednosti
 - ▶ Primer ~150.000 telitev = fenotipskih vrednosti, ~75.000 krav + ~1.000 bikov
- ▶ Specializirani programi (rešitev sistema enačb in/ali ocena komponent variance)
 - ▶ ASREML
 - ▶ BLUPf90
 - ▶ DMU
 - ▶ PEST
 - ▶ SurvivalKit
 - ▶ VCE
 - ▶ WOMBAT
 - ▶ ...

Posebnosti

- ▶ Statistični pomen inverze matrike sorodstva \mathbf{A}^{-1}
- ▶ Bayesovski pogled
- ▶ Komponente variance?

Statistični pomen inverze matrike sorodstva \mathbf{A}^{-1}

- ▶ Pogojna neodvisnost spremenljivk
 - ▶ $\mathbf{A}_{i,j}^{-1} = 0 \rightarrow$ pogojna neodvisnost
 - ▶ $\mathbf{A}_{i,j}^{-1} \neq 0 \rightarrow$ pogojna odvisnost
 - ▶ grafična predstavitev modela = grafični modeli
- ▶ Aditivna genotipska vrednost posameznika
= f(povp. staršev, fenotipsko odstopanje, povp. potomcev)



Bayesovski pogled

- ▶ Model za analizirano spremenljivko

$$p(\mathbf{y}|\mathbf{b}, \mathbf{a}, \mathbf{R}) \sim N(\mathbf{X}\mathbf{b} + \mathbf{Z}\mathbf{a}, \mathbf{R})$$

- ▶ Parametri: \mathbf{b} , \mathbf{a} , σ_a^2 , $\sigma_e^2 \rightarrow$ apriorne porazdelitve

$$p(\mathbf{b}|\boldsymbol{\mu}_b, \mathbf{B}) \sim N(\boldsymbol{\mu}_b, \mathbf{B})$$

$$p(\mathbf{a}|\boldsymbol{\mu}_a, \mathbf{G}) \sim N(\boldsymbol{\mu}_a, \mathbf{G})$$

$$p(\sigma_a^2|\dots), p(\sigma_e^2|\dots) \sim \dots$$

- ▶ Posteriorna porazdelitev za \mathbf{b} in \mathbf{a} če so variance poznane

$$p(\mathbf{b}, \mathbf{a}|\mathbf{R}, \boldsymbol{\mu}_b, \mathbf{B}, \boldsymbol{\mu}_a, \mathbf{G}) \propto p(\mathbf{y}|\mathbf{b}, \mathbf{a}, \mathbf{R}) p(\mathbf{b}|\boldsymbol{\mu}_b, \mathbf{B}) p(\mathbf{a}|\boldsymbol{\mu}_a, \mathbf{G})$$

Bayesovski pogled II

$$p(\mathbf{b}, \mathbf{a} | \mathbf{R}, \boldsymbol{\mu}_b, \mathbf{B}, \boldsymbol{\mu}_a, \mathbf{G}) \sim N(\hat{\boldsymbol{\theta}}, \mathbf{C}^{-1} \sigma_e^2)$$

$$\mathbf{C} \hat{\boldsymbol{\theta}} = \mathbf{r}$$

$$\hat{\boldsymbol{\theta}} = \begin{pmatrix} \hat{\mathbf{b}} \\ \hat{\mathbf{a}} \end{pmatrix}$$

$$\mathbf{C} = \begin{pmatrix} \mathbf{X}^T \mathbf{R}^{-1} \mathbf{X} + \mathbf{B}^{-1} & \mathbf{X}^T \mathbf{R}^{-1} \mathbf{Z} \\ \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{X} & \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{Z} + \mathbf{G}^{-1} \end{pmatrix}$$

$$\mathbf{r} = \begin{pmatrix} \mathbf{X}^T \mathbf{R}^{-1} \mathbf{y} + \mathbf{B}^{-1} \boldsymbol{\mu}_b \\ \mathbf{Z}^T \mathbf{R}^{-1} \mathbf{y} + \mathbf{G}^{-1} \boldsymbol{\mu}_a \end{pmatrix}$$

Ocena komponent variance

- ▶ Metoda največjega verjetja
 - ▶ Maximum Likelihood (ML)

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

- ▶ Restricted/Residual Maximum Likelihood (REML; Patterson & Thompson, 1971)

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - p}$$

- ▶ EM, AI-REML, ...
- ▶ Bayesovski pristopi

Vprašanja?

