

Prediction

Rasmus Waagepetersen
Department of Mathematics
Aalborg University
Denmark

April 3, 2024

WLS and BLUE (prelude to BLUP)

Suppose that Y has mean $X\beta$ and known covariance matrix V (but Y need not be normal). Then

$$\hat{\beta} = (X^T V^{-1} X)^{-1} X^T V^{-1} Y$$

is a weighted least squares estimate since it minimizes

$$(Y - X\beta)^T V^{-1} (Y - X\beta).$$

It is also the best linear unbiased estimate (BLUE) - that is the unbiased estimate with smallest variance in the sense that

$$\text{Var} \tilde{\beta} - \text{Var} \hat{\beta}$$

is positive semi-definite for any other linear unbiased estimate $\tilde{\beta}$.

BLUE for general parameter and $V = I$

Theorem: Suppose $\mathbb{E}Y = \mu$ is in linear subspace M and $\text{Cov}Y = \sigma^2 I$ and $\psi = A\mu$. Then BLUE of ψ is $\hat{\psi} = A\hat{\mu}$ where $\hat{\mu} = PY$ and P orthogonal projection on M .

Obviously $\hat{\psi}$ is LUE: $\mathbb{E}\hat{\psi} = AP\mu = A\mu$.

Key result:

$$\text{Cov}(\tilde{\psi} - \hat{\psi}, \hat{\psi}) = \mathbb{E}[(\tilde{\psi} - \hat{\psi})\hat{\psi}] = 0$$

for any other LUE $\tilde{\psi} = BY$.

Proof of theorem follows by key result:

$$\text{Var}(\tilde{\psi}) = \text{Var}(\tilde{\psi} - \hat{\psi}) + \text{Var}\hat{\psi} \Rightarrow \text{Var}(\tilde{\psi}) - \text{Var}\hat{\psi} = \text{Var}(\tilde{\psi} - \hat{\psi}) \geq 0.$$

Hence $\hat{\psi}$ is BLUE (here $A \geq B$ means $A - B$ positive semi definite).

Proof of key result:

Assume $\tilde{\psi}$ is LUE. I.e. $\tilde{\psi} = BY$ and $\mathbb{E}\tilde{\psi} = B\mu = A\mu$ for all $\mu \in M$.
Thus for all $w \in \mathbb{R}^p$,

$$(B - AP)Pw = BPw - APw = APw - APw = 0$$

since $Pw \in M$. This implies $(B - AP)P = 0$ which gives

$$\mathbb{Cov}(\tilde{\psi} - \hat{\psi}, \hat{\psi}) = \sigma^2(B - AP)P^T A^T = 0.$$

Recall: for random vectors X and Y and matrices A and B of appropriate dimensions

$$\mathbb{Cov}(AX, BY) = A\mathbb{Cov}(X, Y)B^T$$

BLUE - non-diagonal covariance matrix

Lemma: suppose $\tilde{Y} = KY$ where K is an invertible matrix. If $\hat{\psi} = C\tilde{Y}$ is BLUE of ψ based on data \tilde{Y} then $\hat{\psi} = CKY$ is BLUE based on Y as well.

Corollary: suppose $V = LL^T$ is invertible and $\mu = X\beta$ where X has full rank. Then BLUE of μ is $\hat{\mu}$ where $\hat{\mu} = X(X^T V^{-1} X)^{-1} X^T V^{-1} Y$ is WLS estimate of μ .

Proof: $\tilde{Y} = L^{-1}Y$ has covariance matrix I and mean $\tilde{\mu} = \tilde{X}\beta$ where $\tilde{\mu} = L^{-1}\mu$. Thus by theorem, BLUE of $\mu = L\tilde{\mu}$ is $L\tilde{X}(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \tilde{Y}$. Applying lemma we get BLUE based on Y is $L\tilde{X}(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T L^{-1}Y = \hat{\mu}$.

Remark: $\hat{\mu}$ above is in fact orthogonal projection of Y wrt. inner product $\langle x, y \rangle = x^T V^{-1}y$.

Optimal prediction

X and Y random variables, g real function. General result:

$$\begin{aligned}\text{Cov}(Y - \mathbb{E}[Y|X], g(X)) &= \\ \text{Cov}(\mathbb{E}[Y - \mathbb{E}[Y|X]|X], \mathbb{E}[g(X)|X]) + \\ \mathbb{E}\text{Cov}(Y - \mathbb{E}[Y|X], g(X)|X) &= 0\end{aligned}$$

In particular, for any prediction $\tilde{Y} = f(X)$ of Y :

$$\mathbb{E}[(Y - \mathbb{E}[Y|X])(\mathbb{E}[Y|X] - f(X))] = 0$$

from which it follows that

$$\mathbb{E}(Y - \tilde{Y})^2 = \mathbb{E}(Y - \mathbb{E}[Y|X])^2 + \mathbb{E}(\mathbb{E}[Y|X] - \tilde{Y})^2 \geq \mathbb{E}(Y - \mathbb{E}[Y|X])^2$$

Thus $\mathbb{E}[Y|X]$ minimizes mean square prediction error.

Decomposition of Y :

$$Y = \mathbb{E}[Y|X] + (Y - \mathbb{E}[Y|X])$$

where predictor $\mathbb{E}[Y|X]$ and prediction error $Y - \mathbb{E}[Y|X]$ uncorrelated.

Moreover,

$$\text{Var} Y = \text{Var} \mathbb{E}[Y|X] + \text{Var}(Y - \mathbb{E}[Y|X]) = \text{Var} \mathbb{E}[Y|X] + \mathbb{E} \text{Var}[Y|X]$$

whereby

$$\text{Var}(Y - \mathbb{E}[Y|X]) = \mathbb{E} \text{Var}[Y|X].$$

Prediction variance is equal to the expected conditional variance of Y .

BLUP

Consider random vectors Y and X with mean vectors

$$\mathbb{E}Y = \mu_Y \quad \mathbb{E}X = \mu_X$$

and covariance matrix

$$\Sigma = \begin{bmatrix} \Sigma_Y & \Sigma_{YX} \\ \Sigma_{XY} & \Sigma_X \end{bmatrix}$$

Then the best *linear* unbiased predictor of Y given X is

$$\hat{Y} = \mu_Y + \Sigma_{YX}\Sigma_X^{-1}(X - \mu_X)$$

in the sense that

$$\text{Var}[Y - (a + BX)] - \text{Var}[Y - \hat{Y}]$$

is positive semi-definite for all linear unbiased predictors $a + BX$ and '=' only if $a + BX = \hat{Y}$ (unbiased: $\mathbb{E}[Y - a - BX] = 0$).

Prediction variance/mean square prediction error

Fact:

$$\mathbb{Cov}[Y - \hat{Y}, CX] = 0 \quad \text{for all } C. \quad (1)$$

Thus $\mathbb{Cov}[Y - \hat{Y}, \hat{Y}] = 0$ which implies

$$\mathbb{Var} \hat{Y} = \Sigma_{YX} \Sigma_X^{-1} \Sigma_{XY} = \mathbb{Cov}(Y, \hat{Y})$$

It follows that mean square prediction error is

$$\begin{aligned} \mathbb{Var}[Y - \hat{Y}] &= \mathbb{Var} Y + \mathbb{Var} \hat{Y} - \mathbb{Cov}(Y, \hat{Y}) - \mathbb{Cov}(\hat{Y}, Y) \\ &= \Sigma_Y - \Sigma_{YX} \Sigma_X^{-1} \Sigma_{XY} \end{aligned}$$

Proof of fact:

$$\begin{aligned} \mathbb{Cov}[Y - \hat{Y}, CX] &= \mathbb{Cov}[Y, CX] - \mathbb{Cov}[\hat{Y}, CX] = \\ &\quad \Sigma_{YX} C^T - \Sigma_{YX} \Sigma_X^{-1} \Sigma_X C^T = 0 \end{aligned}$$

Proof of BLUP

By (1), $\text{Cov}[Y - \hat{Y}, CX] = 0$ for all C .

$$\begin{aligned}\mathbb{V}\text{ar}[Y - (a + BX)] &= \mathbb{V}\text{ar}[Y - \hat{Y}] + \mathbb{V}\text{ar}[\hat{Y} - (a + BX)] + \\ &\quad \mathbb{C}\text{ov}[Y - \hat{Y}, \hat{Y} - (a + BX)] + \mathbb{C}\text{ov}[\hat{Y} - (a + BX), Y - \hat{Y}] = \\ &\quad \mathbb{V}\text{ar}[Y - \hat{Y}] + \mathbb{V}\text{ar}[\hat{Y} - (a + BX)]\end{aligned}$$

Hence $\mathbb{V}\text{ar}[Y - (a + BX)] - \mathbb{V}\text{ar}[Y - \hat{Y}] = \mathbb{V}\text{ar}[\hat{Y} - (a + BX)]$
where right hand side is positive semi-definite.

Conditional distribution in multivariate normal distribution

Consider jointly normal random vectors Y and X with mean vector

$$\mu = (\mu_Y, \mu_X)$$

and covariance matrix

$$\Sigma = \begin{bmatrix} \Sigma_Y & \Sigma_{YX} \\ \Sigma_{XY} & \Sigma_X \end{bmatrix}$$

Then (provided Σ_X invertible)

$$Y|X = x \sim N(\mu_Y + \Sigma_{YX}\Sigma_X^{-1}(x - \mu_X), \Sigma_Y - \Sigma_{YX}\Sigma_X^{-1}\Sigma_{XY})$$

Proof: By BLUP

$$Y = \hat{Y} + R$$

where $\hat{Y} = \mu_Y + \Sigma_{YX}\Sigma_X^{-1}(X - \mu_X)$,

$R = Y - \hat{Y} \sim N(0, \Sigma_Y - \Sigma_{YX}\Sigma_X^{-1}\Sigma_{XY})$ and $\text{Cov}(R, X) = 0$. By normality R is independent of X . Given $X = x$, \hat{Y} is constant and distribution of R is not affected. Thus result follows.

Optimal prediction for jointly normal random vectors

By previous result it follows that BLUP of Y given X coincides with $E[Y|X]$ when (X, Y) jointly normal.

Hence for normally distributed (X, Y) , BLUP is optimal prediction.

Prediction in linear mixed model

Let $U \sim N(0, \Psi)$ and $Y|U = u \sim N(X\beta + Zu, \Sigma)$.

Then $\text{Cov}[U, Y] = \Psi Z^T$ and $\text{Var}Y = V = Z\Psi Z^T + \Sigma$.

Thus

$$\hat{U} = \mathbb{E}[U|Y] = \Psi Z^T V^{-1}(Y - X\beta)$$

NB: by Woodbury

$$\Psi Z^T (Z\Psi Z^T + \Sigma)^{-1} = (\Psi^{-1} + Z^T \Sigma^{-1} Z)^{-1} Z^T \Sigma^{-1}$$

- e.g. useful if Ψ^{-1} is sparse (like AR-model).

Similarly

$$\text{Var}[U - \hat{U}] = \mathbb{E}\text{Var}[U|Y] = \Psi - \Psi Z^T V^{-1} Z \Psi^T = (\Psi^{-1} + Z^T \Sigma^{-1} Z)^{-1}$$

One-way anova example at p. 186 in M & T.

IQ example

Y measurement of IQ , U subject specific random effect:

$$Y = \mu + U + \epsilon$$

where standard deviation of U and ϵ are 15 and 5 and $\mu = 100$.

Given $Y = 130$, $\mathbb{E}[\mu + U | Y = 130] = 127$.

Example of shrinkage to the mean.

BLUP as hierarchical likelihood estimates

Maximization of joint density ('hierarchical likelihood')

$$f(y|u; \beta)f(u; \psi)$$

with respect to u gives BLUP (M & T p. 171-172 for one-way anova and p. 183 for general linear mixed model)

Joint maximization wrt. u and β gives Henderson's mixed-model equations (M & T p. 184) leading to BLUE $\hat{\beta}$ and BLUP \hat{u} .

BLUP of mixed effect with unknown β

Assume $\mathbb{E}X = C\beta$ and $\mathbb{E}Y = D\beta$. Given X and β , BLUP of

$$K = A\beta + BY$$

is

$$\hat{K}(\beta) = A\beta + B\hat{Y}(\beta)$$

where BLUP $\hat{Y}(\beta) = D\beta + \Sigma_{YX}\Sigma_X^{-1}(X - C\beta)$.

Typically β is unknown. Then BLUP is

$$\hat{K} = A\hat{\beta} + B\hat{Y}(\hat{\beta})$$

where $\hat{\beta}$ is BLUE (Harville, 1991)

Proof: $\hat{K}(\beta)$ can be rewritten as

$$A\beta + B\hat{Y}(\beta) = [A + BD - B\Sigma_{YX}\Sigma_X^{-1}C]\beta + B\Sigma_{YX}\Sigma_X^{-1}X = T + B\Sigma_{YX}\Sigma_X^{-1}X$$

Note BLUE of $T = [A + BD - B\Sigma_{YX}\Sigma_X^{-1}C]\beta$ is

$$\hat{T} = [A + BD - B\Sigma_{YX}\Sigma_X^{-1}C]\hat{\beta}.$$

Now consider a LUP $\tilde{K} = HX = [H - B\Sigma_{YX}\Sigma_X^{-1}]X + B\Sigma_{YX}\Sigma_X^{-1}X$ of K . By unbiasedness,

$$\tilde{T} = [H - B\Sigma_{YX}\Sigma_X^{-1}]X$$

is LUE of T . Hence $\mathbb{V}\text{ar}[\tilde{T} - T] \geq \mathbb{V}\text{ar}[\hat{T} - T]$. Also note by (1)

$$\mathbb{C}\text{ov}[\tilde{T} - T, \hat{K}(\beta) - K] = 0 \text{ and } \mathbb{C}\text{ov}[\hat{T} - T, \hat{K}(\beta) - K] = 0$$

Using this it follows that

$$\mathbb{V}\text{ar}[\tilde{K} - K] \geq \mathbb{V}\text{ar}[\hat{K} - K]$$

Hint: subtract and add $\hat{K}(\beta)$ both in $\mathbb{V}\text{ar}[\tilde{K} - K]$ and $\mathbb{V}\text{ar}[\hat{K} - K]$.

Application to model assessment

From the mixed model formulation

$$Y = X\beta + ZU + \epsilon$$

we obtain

$$\epsilon = Y - X\beta - ZU$$

It is then easy to see that BLUP of ϵ given Y and β is

$$\hat{\epsilon}(\beta) = Y - X\beta - Z\hat{U}(\beta)$$

where $\hat{U}(\beta)$ is BLUP of U given β . When unknown β is replaced by BLUE $\hat{\beta}$, previous slides give that residual

$$\hat{\epsilon} = Y - X\hat{\beta} - Z\hat{U}(\hat{\beta})$$

is BLUP of ϵ (this is returned by applying `residuals` to `lmer` object).

EBLUP and EBLUE

Typically covariance matrix depends on unknown parameters.

EBLUPS are obtained by replacing unknown variance parameters by their estimates (similar for EBLUE).

Model assessment

Make histograms, qq-plots etc. for EBLUPs of ϵ and U .

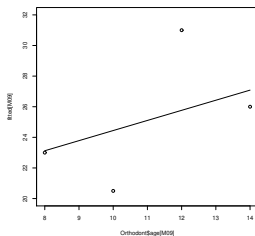
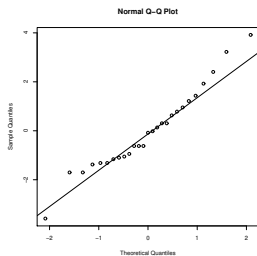
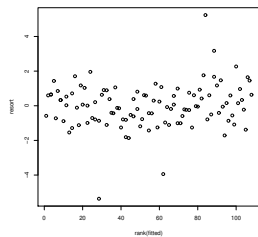
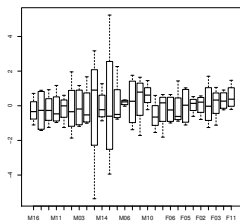
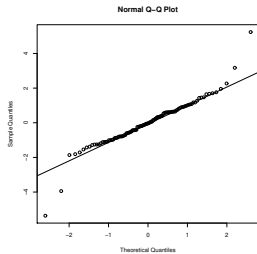
May be advantageous to consider standardized EBLUPS.
Standardized BLUP is

$$[\text{Cov}\hat{U}]^{-1/2}\hat{U}$$

Example: prediction of random intercepts and slopes in orthodont data

```
ort7=lmer(distance~age+factor(Sex)+(1|Subject),data=Orthodont)
#check of model ort7
#residuals
res=residuals(ort7)
qqnorm(res)
qqline(res)
#outliers occur for subjects M09 and M13
#plot residuals against subjects
boxplot(resort~Orthodont$Subject)
#plot residuals against fitted values
fitted=fitted(ort7)
plot(rank(fitted),resort)
```

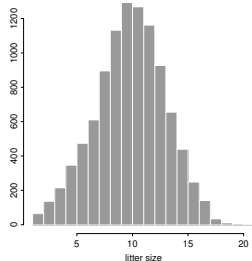
```
#extract predictions of random intercepts
raneffects=ranef(ort7)
#qqplot of random intercepts
qqnorm(ranint[[1]])
qqline(ranint[[1]])
#plot for subject M09
M09=Orthodont$Subject=="M09"
plot(Orthodont$age[M09],fitted[M09],type="l",ylim=c(20,32))
points(Orthodont$age[M09],Orthodont$distance[M09])
```



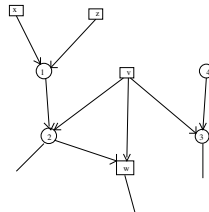
Example: quantitative genetics (Sorensen and Waagepetersen 2003)

X_{ij} size of j th litter of i th pig.

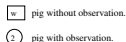
Histogram



Pedigree



Etc.



U_i, \tilde{U}_i random genetic effects influencing size and variability of X_{ij} :

$$X_{ij}|U_i = u_i, \tilde{U}_i = \tilde{u}_i \sim N(\mu_i + u_i, \exp(\tilde{\mu}_i + \tilde{u}_i))$$

$$(U_1, \dots, U_n, \tilde{U}_1, \dots, \tilde{U}_n) \sim N(0, G \otimes A)$$

A: additive genetic relationship (correlation) matrix (depending on pedigree). Correlation structure derived from simple model:

$$U_{\text{offspring}} = \frac{1}{2}(U_{\text{father}} + U_{\text{mother}}) + \epsilon$$

$\Rightarrow Q = A^{-1}$ sparse ! (generalization of AR(1))

$$G = \begin{bmatrix} \sigma_u^2 & \rho\sigma_u\sigma_{\tilde{u}} \\ \rho\sigma_u\sigma_{\tilde{u}} & \sigma_{\tilde{u}}^2 \end{bmatrix}$$

ρ : coefficient of genetic correlation between U_i and \tilde{U}_i .

NB: high dimension $n > 6000$.

Aim: identify pigs with favorable genetic effects

Exercises

1. Fill in the details of the proofs on slides 4-5.
2. Fill in the details of the proof on slide 17.
3. Verify the results on page 186 in M&T regarding BLUPs in case of a one-way anova.

Further results

Rasmus Waagepetersen
Department of Mathematics
Aalborg University
Denmark

April 3, 2024

† Estimable parameters and BLUE

Definition: A linear combination $a^T \beta$ is estimable if it has a LUE $b^T Y$.

Result: $a^T \beta$ is estimable $\Leftrightarrow a^T \beta = c^T \mu$ for some c .

By results on previous slides: If $a^T \beta$ is estimable then BLUE is $c^T \hat{\mu}$.

†Pythagoras and conditional expectation

Space of real random variables with finite variance may be viewed as a vector space with inner product and (L_2) norm

$$\langle X, Y \rangle = \mathbb{E}(XY) \quad \|X\| = \sqrt{\mathbb{E}X^2}$$

Orthogonal decomposition (Pythagoras):

$$\|Y\|^2 = \|\mathbb{E}[Y|X]\|^2 + \|Y - \mathbb{E}[Y|X]\|^2$$

$\mathbb{E}[Y|X]$ may be viewed as projection of Y on X since it minimizes distance

$$\mathbb{E}(Y - \tilde{Y})^2$$

among all predictors $\tilde{Y} = f(X)$.

For zero-mean random variables, orthogonal is the same as uncorrelated.

(Grimmett & Stirzaker, Prob. and Random Processes, Chapter 7.9 good source on this perspective on prediction and conditional expectation)

† BLUP as projection

Y scalar for consistency with slide on L_2 space view.

$X = (X_1, \dots, X_n)^T$. Assume wlog that all variables are centered $\mathbb{E}Y = \mathbb{E}X_i = 0$ (otherwise consider prediction of $Y - \mathbb{E}Y$ based on $X_i - \mathbb{E}X_i$).

BLUP is projection of Y onto *linear* subspace spanned by X_1, \dots, X_n (with orthonormal basis U_1, \dots, U_n where $U = \Sigma_X^{-1/2}X$):

$$\hat{Y} = \sum_{i=1}^n \mathbb{E}[YU_i]U_i = \Sigma_{YX}\Sigma_X^{-1}X$$

(analogue to least squares $\hat{Y} = X(X^T X)^{-1}X^T Y$).

NB: conditional expectation $\mathbb{E}[Y|X]$ projection of Y onto space of *all* variables $Z = f(X_1, \dots, X_n)$ where f real function.

† Conditional simulation using prediction

Suppose Y and X are jointly normal and we wish to simulate $Y|X = x$. By previous result

$$Y|X = x \sim \hat{y} + R$$

where $\hat{y} = \mu_Y + \Sigma_{YX}\Sigma_X^{-1}(x - \mu_X)$. We thus need to simulate R . This can be done by 'simulated prediction': simulate (Y^*, X^*) and compute \hat{Y}^* and $R^* = Y^* - \hat{Y}^*$.

Then our conditional simulation is

$$\hat{y} + R^*$$

Advantageous if it is easier to simulate (Y^*, X^*) and compute \hat{Y}^* than simulate directly from conditional distribution of $Y|X = x$

(e.g. if simulation of (Y, X) easy but $\Sigma_Y - \Sigma_{YX}\Sigma_X^{-1}\Sigma_{XY}$ difficult)