

# Homework 9

Due Thursday, 4/22/21 on Canvas.

1. This problem is an extension of Problem 4 from HW8.

- (a) Use REML + GLS to fit the model you proposed in part (c)(i) of Problem 4. Test the null hypothesis that there is no patient  $\times$  treatment interaction. Make sure to report your null and alternative models.
- (b) Use the model you fitted above to estimate the parameters in the model you proposed part (c)(ii) of Problem 4. Which drug would you recommend an insomniac take if they slept an average of 2 hours a night without treatment? How about one that slept an average of 6 hours a night?

2. (REML is maximum quasi-likelihood on the residuals of  $\mathbf{Y}$ ) Suppose  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ , where  $\mathbf{X} \in \mathbb{R}^{n \times p}$  is full rank,  $\boldsymbol{\beta} \in \mathbb{R}^p$ ,  $\mathbb{E}(\boldsymbol{\epsilon}) = \mathbf{0}_n$  and  $\text{Var}(\mathbf{Y}) = \mathbf{V}(\boldsymbol{\theta})$  for  $\boldsymbol{\theta} \in \mathbb{R}^b$  and some function  $\mathbf{V} : \mathbb{R}^b \rightarrow \mathbb{R}^{n \times n}$ . Throughout this problem, assume that  $\mathbf{V}(\boldsymbol{\theta})$  is positive definite and let  $\mathbf{A}_X \in \mathbb{R}^{n \times (n-p)}$  be any matrix whose columns form a basis for  $\ker(\mathbf{X}^T)$ . Recall that  $\mathbf{A}_X$  is not unique.

$$\begin{aligned} \mathbb{E}(\mathbf{Y}\sim) &= \mathbb{E}(\mathbf{Y}) - \\ &= \mathbb{E}(\dots\mathbf{e}) - \mathbb{E}(\mathbf{X}\mathbf{B}) \\ &= \mathbb{E}(\mathbf{X}\mathbf{B}) - \mathbb{E}(\mathbf{X}\mathbf{B}) = \\ &= \mathbf{0} \end{aligned}$$

$$\begin{aligned} \text{Var}(\mathbf{Y}\sim) &= \\ \text{Var}(\mathbf{Q}\mathbf{Y}) &= \\ \mathbf{Q}\text{Var}(\mathbf{Y})\mathbf{Q}^T &= \\ = \mathbf{Q}\mathbf{V}(\boldsymbol{\theta})\mathbf{Q} & \text{ (bc} \\ \mathbf{Q} \text{ is symmetric)} & \end{aligned}$$

- (a) Let  $\mathbf{Q} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ . Show that  $\mathbf{A}_X(\mathbf{A}_X^T\mathbf{A}_X)^{-1}\mathbf{A}_X^T = \mathbf{Q}$ . multiply both sides with  $\mathbf{A}$
- (b) Let  $\tilde{\mathbf{Y}} = \mathbf{Q}\mathbf{Y} \in \mathbb{R}^n$  be the residuals after you regress out  $\mathbf{X}$  from  $\mathbf{Y}$ . Find expressions for  $\mathbb{E}(\tilde{\mathbf{Y}})$  and  $\text{Var}(\tilde{\mathbf{Y}})$ . Show that  $\tilde{\mathbf{Y}}$  is a degenerate distribution (i.e. the support of  $\tilde{\mathbf{Y}}$  is a strict subspace of  $\mathbb{R}^n$ ).
- (c) Given that  $\tilde{\mathbf{Y}}$  is degenerate, the usual normal likelihood is not an appropriate quasi-likelihood for  $\tilde{\mathbf{Y}}$ . To circumvent this, define

$$\tilde{\ell}(\boldsymbol{\theta}) = -\frac{1}{2} \log [\det_+ \{\mathbf{Q}\mathbf{V}(\boldsymbol{\theta})\mathbf{Q}\}] - \frac{1}{2} \tilde{\mathbf{Y}}^T \{\mathbf{Q}\mathbf{V}(\boldsymbol{\theta})\mathbf{Q}\}^\dagger \tilde{\mathbf{Y}}$$

to be the likelihood of the degenerate normal distribution, where  $\mathbf{B}^\dagger$  is the usual Moore-Penrose pseudoinverse of  $\mathbf{B} \in \mathbb{R}^{n \times n}$  and  $\det_+(\mathbf{B})$ , called the **pseudo determinant** of  $\mathbf{B}$ , is defined to be

$$\det_+(\mathbf{B}) = \lim_{\epsilon \rightarrow 0} \frac{\det(\mathbf{B} + \epsilon \mathbf{I}_n)}{\epsilon^{n - \text{rank}(\mathbf{B})}}.$$

- (i) Show that  $\{\mathbf{Q}\mathbf{V}(\boldsymbol{\theta})\mathbf{Q}\}^\dagger = \mathbf{A}_X \{\mathbf{A}_X^T \mathbf{V}(\boldsymbol{\theta}) \mathbf{A}_X\}^{-1} \mathbf{A}_X^T$ .
- (ii) If  $\mathbf{B} \in \mathbb{R}^n$  is symmetric positive semi-definite, show that  $\det_+(\mathbf{B})$  is the product of the non-zero eigenvalues of  $\mathbf{B}$ . Use this to show that  $\det_+ \{\mathbf{Q}\mathbf{V}(\boldsymbol{\theta})\mathbf{Q}\} = \det(\mathbf{A}_X^T \mathbf{A}_X)^{-1} \det \{\mathbf{A}_X^T \mathbf{V}(\boldsymbol{\theta}) \mathbf{A}_X\}$ .

**Remark.** The function  $\tilde{\ell}$  is the log-likelihood of  $\tilde{\mathbf{Y}}$  if we assume  $\tilde{\mathbf{Y}} \sim N(\mathbf{0}, \mathbf{Q}\mathbf{V}(\boldsymbol{\theta})\mathbf{Q})$ . Note that implicit in  $\tilde{\ell}$  is the requirement that  $\tilde{\mathbf{Y}} \in \ker(\mathbf{X}^T)$ . The likelihood, by definition, is 0 for  $\tilde{\mathbf{Y}} \notin \ker(\mathbf{X}^T)$ . For the probability purists,  $\tilde{\ell}$  is the log of the density of  $\tilde{\mathbf{Y}}$  (ignoring constants that do not depend on  $\boldsymbol{\theta}$ ) with respect to the Lebesgue measure on  $\ker(\mathbf{X}^T)$ .

(d) Define the usual REML objective function to be

$$\ell_{REML}(\theta) = -\frac{1}{2} \log \left[ \det \left\{ \mathbf{A}_X^T \mathbf{V}(\theta) \mathbf{A}_X \right\} \right] - \frac{1}{2} \left( \mathbf{A}_X^T \mathbf{Y} \right)^T \left\{ \mathbf{A}_X^T \mathbf{V}(\theta) \mathbf{A}_X \right\}^{-1} \left( \mathbf{A}_X^T \mathbf{Y} \right).$$

Show that  $\ell_{REML}(\theta) = \tilde{\ell}(\theta) + C$ , where  $C$  is a constant that does not depend on  $\theta$ . Use this to show that  $\hat{\theta}_{REML} = \arg \max_{\theta \in \Theta} \ell_{REML}(\theta) = \arg \max_{\theta \in \Theta} \tilde{\ell}(\theta)$ . Why does this imply that  $\hat{\theta}_{REML}$  is invariant to the choice of  $\mathbf{A}_X$ ?

(e) Lastly, suppose  $\theta = \sigma^2 > 0$  and that  $\mathbf{V}(\theta) = \sigma^2 \mathbf{I}_n$  (i.e. the usual assumptions on the variance in standard linear regression). Using part (d) (i.e. that the REML estimator for  $\sigma^2$  can be computed by maximizing  $\tilde{\ell}$ ), show that  $\hat{\sigma}_{REML}^2 = MSE$ , where  $MSE = (n - p)^{-1} \tilde{\mathbf{Y}}^T \tilde{\mathbf{Y}}$  is the usual mean squared error.

3. The data Age1.txt and Age3.txt contain measurements of the concentration of bilirubin, a ubiquitous small molecule metabolite, in the blood plasma of children at age 1 and 3, respectively. The goal is to understand relationship between the concentration of bilirubin and recurrent wheeze, the latter of which is a diagnosis of  $\geq 4$  wheezing episodes in that year. The complete list of covariates is given below.

- **IndividualID:** A unique ID given to each child. There should be 33 children whose bilirubin concentration was measured at both age 1 and age 3. You may assume that data collected on different individuals are independent.
- **Bilirubin:** The log-concentration of bilirubin. Treat this as the dependent variable.
- **Wheeze:** A factor variable with 3 levels. level 0: no wheezing episodes in that year; level 1: 1-3 wheezing episodes in that year; level 2:  $\geq 4$  wheezing episodes in that year.
- **Diet:** A factor variable with 2 levels. level 0: exclusively breastfed for the first six months of life; level 1: not exclusively breastfed for the first six months of life.
- **Daycare:** A factor variable with 2 levels. level 0: did not attend daycare in the first year of life; level 1: attended daycare in the first year of life.
- **Sex:** A factor variable with 2 levels. level 0: male; level 1: female.

(a) By performing two separate linear regressions (one at age 1 and one at age 3) and treating all of the above-mentioned covariates in your model as additive fixed effects, estimate the expected difference in the log-concentration of bilirubin **between recurrent wheezers (those with  $\geq 4$  wheezing episodes in that year) and healthy controls (those who did not wheeze in that year)** at ages 1 and 3. Report 95% confidence interval for both expected differences.

(b) The inference you performed in part (a) presumably relied on approximating test statistics with a normal or t-distribution. Do you trust this approximation? Give an argument as to why you do or do not. Include plots if necessary.

Z test to compare?

- (c) Let  $\beta_j$  be the expected difference in the log-concentration of bilirubin between recurrent wheezers and healthy controls at age  $j$ , and let  $\hat{\beta}_j$  be its estimate you determined in part (a). Your collaborator has reason to believe that  $\beta_1 = \beta_3 = \beta$ . To estimate  $\beta$ , suppose you decide to meta-analyze the results at ages 1 and 3 by modelling  $\hat{\beta}_1, \hat{\beta}_3$  as

$$\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_3 \end{pmatrix} \sim N(\mathbf{1}_2\beta, \text{diag}(\hat{v}_1, \hat{v}_3)),$$

where  $\hat{v}_j$  is the estimate for the variance of  $\hat{\beta}_j$  determined in part (a). Assuming this model is correct, report a point estimate and 95% confidence interval for  $\beta$ .

- (d) Do you think the model in part (c) is appropriate? Why or why not? Given your answer, do you suspect the confidence interval determined in part (c) is too narrow, too wide, or accurate? Justify your answer.
- (e) Let  $y_{ji}$  be the log-bilirubin concentration at age  $j$  in individual  $i = 1, \dots, n_j$ , and let  $\mathbf{x}_{ji}$  be the covariates for individual  $i$  at age  $j$  used in part (a). Assume that  $\mathbb{E}(y_{ji}) = \mathbf{x}_{ji}^T \gamma_j$  and

$$\text{Cov}(y_{ji}, y_{j'i'}) = \begin{cases} \sigma_j^2 & \text{if } (j, i) = (j', i') \quad \text{diff person diff age} \\ \phi & \text{if } j \neq j' \text{ and } i = i' \quad \text{same person diff age} \Rightarrow \text{sign should be } > 0 \\ 0 & \text{otherwise} \end{cases}$$

- (i) What do you expect the sign of  $\phi$  to be?
- (ii) Based on the results from part (a), do you think it would be appropriate to let  $\sigma_1^2 = \sigma_3^2$ ? Justify your answer.
- (f) Using whatever method you deem most appropriate, estimate  $\sigma_j^2$  and  $\phi$ .
- (g) Let  $\mathbf{Y}_j = (y_{j1}, \dots, y_{jn_j})^T$  and define  $z_{ji}$  to be 1 if individual  $i$  at age  $j$  was diagnosed with recurrent wheeze, and 0 otherwise. Let  $\mathbf{Z}_j = (z_{j1}, \dots, z_{jn_j})^T$ . Suppose you parametrize the mean model for  $\mathbf{Y}_j$  as

$$\mathbb{E}(\mathbf{Y}_j) = \mathbf{Z}_j \beta_j + \tilde{\mathbf{X}}_j \tilde{\gamma}_j,$$

where  $\beta_j$  is defined in part (c) and  $\tilde{\mathbf{X}}_j$  contains the other covariates at age  $j$  (that are not of interest). Show that for  $\hat{\beta}_j$  the ordinary least squares estimate for  $\beta_j$  defined in part (c) and  $\tilde{\mathbf{Q}}_j = \mathbf{I}_{n_j} - \tilde{\mathbf{X}}_j (\tilde{\mathbf{X}}_j^T \tilde{\mathbf{X}}_j)^{-1} \tilde{\mathbf{X}}_j^T$ ,

$$\hat{\beta}_j = (\mathbf{Z}_j^T \tilde{\mathbf{Q}}_j \mathbf{Z}_j)^{-1} \mathbf{Z}_j^T \tilde{\mathbf{Q}}_j \mathbf{Y}_j.$$

- (h) Using parts (f) and (g), estimate  $\mathbf{V} = \text{Var}\left\{\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_3 \end{pmatrix}\right\}$ . Comment on the off-diagonal elements of this matrix.
- (i) Assume  $\hat{\beta}_1, \hat{\beta}_3$  are jointly normal and that  $\beta_1 = \beta_3 = \beta$ . Use your estimate for  $\mathbf{V}$  from part (h) to provide a point estimate and 95% confidence interval for  $\beta$ .