

# Outline 3

February 20, 2019

## 0.1 Step1:

I was confused about the dimensions of the following equation. I just wanted to add the subscript  $i$  in the equation and specify the equation for each subject.

To calculate the max purity function, first, fit the linear mixed model for the outcome  $\mathbf{y}_i$  and time  $\mathbf{X}_i$ , with baseline covariates  $\mathbf{x}_i$ :

$$\mathbf{y}_i = \mathbf{X}_i(\boldsymbol{\beta} + \mathbf{b}_i + \boldsymbol{\Gamma}(\boldsymbol{\alpha}'\mathbf{x}_i)) + \boldsymbol{\epsilon}_i.$$

where,

- $\mathbf{y}_i$  is the vector of outcome for the  $i$ th subject, i.e., the dimension of  $\mathbf{y}_i$  is  $(n_i, 1)$ .  $n_i$  is the number of observations for subject  $i$ .
- $\mathbf{X}_i$  is the covariance matrix for the  $i$ th subject. The dimension is  $n_i, p$ .  $p$  is the number of covariates. Here in our example  $p = 3$ , which is for  $(1, t, t^2)$ .
- $\boldsymbol{\beta}$  is the coefficient vector for the fixed effects of  $\mathbf{X}_i$ . The dimension is  $(p, 1)$ .
- $\mathbf{x}_i$  is the vector of the baseline covariates for the subject. The dimension is  $(q, 1)$ .
- $\mathbf{b}_i$  is the vector of random effects. The dimension is  $(p, 1)$ .
- $\boldsymbol{\Gamma}$  is the vector of fixed effects of the baseline covariates. Dimension is  $(p, 1)$ .

- $w_i = \alpha' \mathbf{x}_i$  is the combination of the input baseline covariates.  $w_i$  is a scalar.

We can define the covariance matrix of  $\mathbf{X}_i$  as  $\mathbf{z}_i$ . The  $\mathbf{z}_i$  contains both fixed effects and random effects.

$$\mathbf{z}_i = \boldsymbol{\beta} + \mathbf{b}_i + \boldsymbol{\Gamma} w_i$$

We can also write the above equation in the matrix version:

$$Y_i = \begin{bmatrix} 1 & t_1 & t_1^2 \\ 1 & t_2 & t_2^2 \\ \dots & \dots & \dots \\ 1 & t_{n_i} & t_{n_i}^2 \end{bmatrix} \left[ \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} b_0 \\ b_1 \\ b_2 \end{pmatrix} + \begin{pmatrix} \gamma_0 \\ \gamma_1 \\ \gamma_2 \end{pmatrix} w_i \right] + \boldsymbol{\epsilon}_i$$

For different subjects, they have the same  $\boldsymbol{\beta}$  vector and the same  $\boldsymbol{\Gamma}$  vector. Their random effect vector  $\mathbf{b}_i$  can be different.

Sorry I am still a little confused here. How could we combine the  $\mathbf{Y}_i$ ,  $\mathbf{X}_i$  and the matrix of coefficients of each subject together to get an overall equation, like  $\mathbf{Y} = \mathbf{X}$  times some matrix?

## 0.2 Step 2:

Estimate the distribution of  $\mathbf{z}_i$  for drug group and placebo, separately.

$$\begin{aligned} f(\mathbf{z}_i) &= \int_{w_i} f(\mathbf{z}_i, w_i) dw_i \\ &= \int_{w_i} f(\mathbf{z}_i | w_i) g(w_i) dw_i \end{aligned}$$

where,

- the conditional distribution  $f(\mathbf{z}_i | w_i) \sim MVN(\boldsymbol{\beta} + \boldsymbol{\Gamma} w_i, \mathbf{D})$
- $g(w_i)$  is the distribution of the covariates combination  $\alpha' \mathbf{x}_i$ . For example, if the covariates combination only contains "sex", which is binary, then the integral becomes summation.
- $\mathbf{D}_i$  is the covariance matrix of random effects  $\mathbf{b}_i$ .

We could then estimate the  $f(\cdot)$  for drug group and placebo group separately, i.e.  $f_1(\mathbf{z}_i)$  and  $f_2(\mathbf{z}_i)$ .

### 0.3 Step 3:

The purity should be a function of  $w_i$ . The integral of  $z_i$  should be calculated first. Then we can define the purity as,

$$\begin{aligned} P_{w_i} &= \int_{z_i} \frac{[f_1(z_i|w_i) - f_2(z_i|w_i)]^2}{[f_1(z_i|w_i) + f_2(z_i|w_i)]^2} (f_1(z_i|w_i) + f_2(z_i|w_i)) dz_i \\ &= \int_{z_i} \frac{[f_1(z_i|w_i) - f_2(z_i|w_i)]^2}{f_1(z_i|w_i) + f_2(z_i|w_i)} dz_i \end{aligned}$$

where

- the integral of  $z_i$  with high dimension, which is equals to the dimension of  $z_i$ .
- the  $f_1(z_i|w_i)$  and  $f_2(z_i|w_i)$  are pdf of multivariate normal distributions.
- $f_1(z_i|w_i) = (2\pi)^{-\frac{p}{2}} |\mathbf{D}|^{-1/2} \exp(-\frac{1}{2}(\mathbf{z}_i - \hat{\boldsymbol{\beta}}_1 - \hat{\boldsymbol{\Gamma}}_1 w_i)^T \mathbf{D}^{-1} (\mathbf{z}_i - \hat{\boldsymbol{\beta}}_1 - \hat{\boldsymbol{\Gamma}}_1 w_i))$
- $f_2(z_i|w_i) = (2\pi)^{-\frac{p}{2}} |\mathbf{D}|^{-1/2} \exp(-\frac{1}{2}(\mathbf{z}_i - \hat{\boldsymbol{\beta}}_2 - \hat{\boldsymbol{\Gamma}}_2 w_i)^T \mathbf{D}^{-1} (\mathbf{z}_i - \hat{\boldsymbol{\beta}}_2 - \hat{\boldsymbol{\Gamma}}_2 w_i))$

It can be also approximated as the summation form:

$$\sum_{i=1}^n \left[ \frac{f_1(\mathbf{z}_i|w_i) - f_2(\mathbf{z}_i|w_i)}{f_1(\mathbf{z}_i|w_i) + f_2(\mathbf{z}_i|w_i)} \right]^2$$

The summation is also high dimensional, with the same dimension as  $z_i$ .

This is one purity value based on one combination of baseline covariates for subject  $i$ .

Then we can get the subject purity, i.e., the purity for subject  $w_i$ . We need to calculate the integral or summation of  $w_i$  next to get an overall purity for the data.

We would like to find the  $w$  or  $\alpha$  that max the  $P_{\alpha'x}$ , i.e.,

$$\begin{aligned} \hat{\alpha} &= \arg \max_{\alpha} P_{\alpha} \\ &= \arg \max_{\alpha} \int_{\alpha'x} P_{\alpha'x} d\alpha'x \\ &= \arg \max_{\alpha} \int_{\alpha'x} \frac{(f_1(z|\alpha'x) - f_2(z|\alpha'x))^2}{f_1(z|\alpha'x) + f_2(z|\alpha'x)} d\alpha'x \end{aligned}$$