

# Purity Calculation

## Kullback-Leibler divergence and Purity

To measure how much the differences are between the treatment group and the placebo group, we apply the Kullback-Leibler (KL) divergence, which measures how one probability distribution  $F_1$  is different from another probability distribution  $F_2$ .

$$D_{KL}(F_1||F_2) = \int_{-\infty}^{+\infty} f_1(x) \log\left(\frac{f_1(x)}{f_2(x)}\right) dx \quad (1)$$

where  $f_1$  and  $f_2$  denote the probability density functions (pdf) of  $F_1$  and  $F_2$ , separately. The larger the KL divergence between distributions is, the more "pure" the distributions are. Besides,  $D_{KL}(F_1||F_2) \geq 0$ . Similarly, the  $D_{KL}(F_2||F_1)$  is also always larger than or equals to 0.

Based on the Kullback-Leibler divergence, we define the *purity*, which represent how much the differences between the treatment group distribution  $F_1$  and the placebo group distribution  $F_2$ . We define the purty function of the summation of two Kullback-Leibler divergence as

$$\begin{aligned} \text{purity} &= D_{KL}(F_1||F_2) + D_{KL}(F_2||F_1) \\ &= \int_{-\infty}^{+\infty} f_1(x) \log\left(\frac{f_1(x)}{f_2(x)}\right) dx + \int_{-\infty}^{+\infty} f_2(x) \log\left(\frac{f_2(x)}{f_1(x)}\right) dx \end{aligned} \quad (2)$$

where

$$f_1(x) \sim MVN(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), (\boldsymbol{\mu}_1 : p \times 1, \boldsymbol{\Sigma}_1 : p \times p)$$

$$f_2(x) \sim MVN(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2), (\boldsymbol{\mu}_2 : p \times 1, \boldsymbol{\Sigma}_2 : p \times p)$$

Let's calculate the purity value by calculating  $\int f_1 \log f_1$ ,  $\int f_2 \log f_2$ ,  $\int f_1 \log f_2$ , and  $\int f_2 \log f_1$ .

**Part**  $\int f_1 \log f_1$

$$\begin{aligned} \int f_1 \log f_1 &= E_1 \left\{ -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\boldsymbol{\Sigma}_1|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)' (\boldsymbol{\Sigma}_1)^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) \right\} \\ &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\boldsymbol{\Sigma}_1|) - \frac{1}{2} E_1 [(\mathbf{x} - \boldsymbol{\mu}_1)' (\boldsymbol{\Sigma}_1)^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)] \end{aligned}$$

And

$$\begin{aligned} E_1 [(\mathbf{x} - \boldsymbol{\mu}_1)' (\boldsymbol{\Sigma}_1)^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)] &= E_1 [tr((\mathbf{x} - \boldsymbol{\mu}_1)' (\boldsymbol{\Sigma}_1)^{-1} (\mathbf{x} - \boldsymbol{\mu}_1))] \\ &= E_1 [tr((\boldsymbol{\Sigma}_1)^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) (\mathbf{x} - \boldsymbol{\mu}_1)')] \\ &= tr(E_1 [(\boldsymbol{\Sigma}_1)^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) (\mathbf{x} - \boldsymbol{\mu}_1)']) \\ &= tr(\boldsymbol{\Sigma}_1^{-1} E_1 [(\mathbf{x} - \boldsymbol{\mu}_1) (\mathbf{x} - \boldsymbol{\mu}_1)']) \\ &= tr(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_1) = tr(\mathbf{I}_p) = p \end{aligned}$$

Therefore,

$$\int f_1 \log f_1 = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\boldsymbol{\Sigma}_1|) - \frac{p}{2} \quad (3)$$

Similarly,

$$\int f_2 \log f_2 = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\boldsymbol{\Sigma}_2|) - \frac{p}{2} \quad (4)$$

**Part**  $\int f_1 \log f_2$

$$\begin{aligned} \int f_1 \log f_2 &= E_1\left(-\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_2|) - \frac{1}{2}(\mathbf{x} - \mu_2)' \Sigma_2^{-1}(\mathbf{x} - \mu_2)\right) \\ &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_2|) - \frac{1}{2} E_1[(\mathbf{x} - \mu_2)' \Sigma_2^{-1}(\mathbf{x} - \mu_2)] \end{aligned}$$

And

$$\begin{aligned} &E_1[(\mathbf{x} - \mu_2)' \Sigma_2^{-1}(\mathbf{x} - \mu_2)] \\ &= E_1[(\mathbf{x} - \mu_1 + \mu_1 - \mu_2)' \Sigma_2^{-1}(\mathbf{x} - \mu_1 + \mu_1 - \mu_2)] \\ &= E_1[(\mathbf{x} - \mu_1)' \Sigma_2^{-1}(\mathbf{x} - \mu_1) + (\mu_1 - \mu_2)' \Sigma_2^{-1}(\mathbf{x} - \mu_1) \\ &\quad + (\mathbf{x} - \mu_1)' \Sigma_2^{-1}(\mu_1 - \mu_2) + (\mu_1 - \mu_2)' \Sigma_2^{-1}(\mu_1 - \mu_2)] \\ &= E_1[(\mathbf{x} - \mu_1)' \Sigma_2^{-1}(\mathbf{x} - \mu_1)] + (\mu_1 - \mu_2)' \Sigma_2^{-1} E_1(\mathbf{x} - \mu_1) + \\ &\quad E_1(\mathbf{x} - \mu_1)' \Sigma_2^{-1}(\mu_1 - \mu_2) + (\mu_1 - \mu_2)' \Sigma_2^{-1}(\mu_1 - \mu_2) \\ &= E_1[(\mathbf{x} - \mu_1)' \Sigma_2^{-1}(\mathbf{x} - \mu_1)] + 0 + 0 + (\mu_1 - \mu_2)' \Sigma_2^{-1}(\mu_1 - \mu_2) \\ &= E_1[\text{tr}(\mathbf{x} - \mu_1)' \Sigma_2^{-1}(\mathbf{x} - \mu_1)] + (\mu_1 - \mu_2)' \Sigma_2^{-1}(\mu_1 - \mu_2) \\ &= E_1[\text{tr}(\Sigma_2^{-1}(\mathbf{x} - \mu_1)'(\mathbf{x} - \mu_1))] + (\mu_1 - \mu_2)' \Sigma_2^{-1}(\mu_1 - \mu_2) \\ &= \text{tr}(E_1[\Sigma_2^{-1}(\mathbf{x} - \mu_1)'(\mathbf{x} - \mu_1)]) + (\mu_1 - \mu_2)' \Sigma_2^{-1}(\mu_1 - \mu_2) \\ &= \text{tr}(\Sigma_2^{-1} E_1[(\mathbf{x} - \mu_1)'(\mathbf{x} - \mu_1)]) + (\mu_1 - \mu_2)' \Sigma_2^{-1}(\mu_1 - \mu_2) \\ &= \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_1 - \mu_2)' \Sigma_2^{-1}(\mu_1 - \mu_2) \end{aligned}$$

Therefore,

$$\int f_1 \log f_2 = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_2|) - \frac{1}{2} \{ \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_1 - \mu_2)' \Sigma_2^{-1}(\mu_1 - \mu_2) \} \quad (5)$$

Similarly,

$$\int f_2 \log f_1 = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_1|) - \frac{1}{2} \{ \text{tr}(\Sigma_1^{-1} \Sigma_2) + (\mu_1 - \mu_2)' \Sigma_1^{-1}(\mu_1 - \mu_2) \} \quad (6)$$

Then the purity function is

$$\begin{aligned} \text{purity} &= \int f_1 \log f_1 + \int f_2 \log f_2 - \int f_2 \log f_1 - \int f_1 \log f_2 \\ &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_1|) - \frac{p}{2} \\ &\quad -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_2|) - \frac{p}{2} \\ &\quad -(-\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_2|) - \frac{1}{2} \{ \text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_1 - \mu_2)' \Sigma_2^{-1}(\mu_1 - \mu_2) \}) \\ &\quad -(-\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma_1|) - \frac{1}{2} \{ \text{tr}(\Sigma_1^{-1} \Sigma_2) + (\mu_1 - \mu_2)' \Sigma_1^{-1}(\mu_1 - \mu_2) \}) \\ &= -p + \frac{1}{2} \text{tr}(\Sigma_1^{-1} \Sigma_2) + \frac{1}{2} \text{tr}(\Sigma_2^{-1} \Sigma_1) \\ &\quad + \frac{1}{2} [(\mu_1 - \mu_2)' \Sigma_1^{-1}(\mu_1 - \mu_2)] + \frac{1}{2} [(\mu_1 - \mu_2)' \Sigma_2^{-1}(\mu_1 - \mu_2)] \end{aligned} \quad (7)$$

In our settings, we fit the outcome variable with a linear mixed model

$$\mathbf{Y}_1 = \mathbf{S}(\beta_1 + \mathbf{b}_1 + \mathbf{\Gamma}_1(\alpha' \mathbf{x})) + \epsilon$$

$$\mathbf{Y}_2 = \mathbf{S}(\beta_2 + \mathbf{b}_2 + \mathbf{\Gamma}_2(\alpha' \mathbf{x})) + \epsilon$$

and treat the coefficient  $\mathbf{z}_i = \beta_i + \mathbf{b}_i + \mathbf{\Gamma}_i(\alpha' \mathbf{x})$ ,  $i \in \{1, 2\}$  as a MVN, that is,

$$\mathbf{z}_1 | \mathbf{x} \sim \text{MVN}(\beta_1 + \mathbf{\Gamma}_1(\alpha' \mathbf{x}), \mathbf{D}_1)$$

$$\mathbf{z}_2|x \sim MVN(\beta_2 + \Gamma_2(\alpha'x), D_2)$$

Therefore,  $\mu_1 = \beta_1 + \Gamma_1(\alpha'x)$ ,  $\mu_2 = \beta_2 + \Gamma_2(\alpha'x)$ ,  $\Sigma_1 = D_1$ ,  $\Sigma_2 = D_2$

Therefore

$$\begin{aligned} & (\mu_1 - \mu_2)'(D_1^{-1} + D_2^{-1})(\mu_1 - \mu_2) \\ &= (\beta_1 - \beta_2 + (\Gamma_1 - \Gamma_2)\alpha'x)'(D_1^{-1} + D_2^{-1})(\beta_1 - \beta_2 + (\Gamma_1 - \Gamma_2)\alpha'x) \\ &= (\beta_1 - \beta_2)'(D_1^{-1} + D_2^{-1})(\beta_1 - \beta_2) \\ &+ 2[(\beta_1 - \beta_2)'(D_1^{-1} + D_2^{-1})(\Gamma_1 - \Gamma_2)x'\alpha \\ &+ \alpha'xx'\alpha((\Gamma_1 - \Gamma_2))'(D_1^{-1} + D_2^{-1})((\Gamma_1 - \Gamma_2))] \end{aligned}$$

And purity function in terms of  $\beta_i, \Gamma_i, D_i$  is,

$$\begin{aligned} \text{purity} = g(\alpha'x) = & -p + \frac{1}{2}\text{tr}(D_2^{-1}D_1) + \frac{1}{2}\text{tr}(D_1^{-1}D_2) \\ & + \frac{1}{2}\{(\beta_1 - \beta_2)'(D_1^{-1} + D_2^{-1})(\beta_1 - \beta_2) \\ & + 2[(\beta_1 - \beta_2)'(D_1^{-1} + D_2^{-1})(\Gamma_1 - \Gamma_2)x'\alpha \\ & + \alpha'xx'\alpha((\Gamma_1 - \Gamma_2))'(D_1^{-1} + D_2^{-1})((\Gamma_1 - \Gamma_2))]\} \end{aligned}$$

The expectation of the purity function is

$$\begin{aligned} G(\alpha) = E(g(\alpha)) = & A_0 + A_1\mu'_x\alpha + \frac{A_2}{2}E[\alpha'xx'\alpha] \\ = & A_0 + A_1\mu'_x\alpha + \frac{A_2}{2}[\alpha'(\Sigma_x + \mu'_x\mu_x)\alpha] \end{aligned}$$

where

- $E(x) = \mu_x$
- $\text{Var}(x) = \Sigma_x$
- $A_0 = -p + \frac{1}{2}\text{tr}(D_2^{-1}D_1) + \frac{1}{2}\text{tr}(D_1^{-1}D_2) + \frac{1}{2}(\beta_1 - \beta_2)'(D_1^{-1} + D_2^{-1})(\beta_1 - \beta_2)$
- $A_1 = (\beta_1 - \beta_2)'(D_1^{-1} + D_2^{-1})(\Gamma_1 - \Gamma_2)$
- $A_2 = (\Gamma_1 - \Gamma_2)'(\hat{D}_1^{-1} + D_2^{-1})(\Gamma_1 - \Gamma_2)$