# TCA-TWAS Data Simulation with heritability and genetics correlation

Gary Qiurui Ma, Brandon Jew

August 2019

## 1 TCA model

Let $Z_h^i$ be the gene expression level of individual $i \in 1, ...n$ in cell type $h \in 1, ...k$ at gene $j$, and let $\mathbf{X} \in \mathbb{R}^{n \times m}$ be a matrix of $m$ cis-snps, and let $\mathbf{C^1} \in \mathbb{R}^{n \times p_1}$ be a matrix of $p_1$ covariates. We assume:

$$Z_h^i = (x_i)^T \beta_h + c_i^1 \gamma_h + \epsilon_z \qquad (1)$$
$$\epsilon_z \sim \mathcal{N}(0, \sigma_z^2)$$

where $x_i$ is the $i$-th row of $\mathbf{X}$ (corresponding to the m snps of the $i$-th individual), $\beta_h$ is a $m$-th length vector of corresponding effect size for the $m$ snps in the $h$-th cell type, which is also the $h$-th column of $\mathcal{B}$, $c_i^1$ is the $i$-th row of $\mathbf{C^1}$ (corresponding to the $p_1$ covariates of the $i$-th individual), $\gamma_h$ is a $p_1$-th length vector, and $\epsilon_z$ an i.i.d. component of variation

We assume the observed bulk level gene expression are convolved signals from $k$ different cell types. We denote $\mathbf{W} \in \mathbb{R}^{n \times k}$ as a matrix of cell-type proportions of $k$ cell types for each of the $n$ individuals, and $\mathbf{C^2} \in \mathbb{R}^{n \times p_2}$ as a matrix of $p_2$ global covariates potentially affecting the observed bulk level gene expression. TCA model for $G_i$, the observed bulk level gene expression of the $i$-th individual in $j$-th gene, is as follows:

$$G_i = c_i^2 \delta + \sum_{h=1}^{k} w_{hi} z_{hi} + \epsilon_g \qquad (2)$$
$$\epsilon_g \sim \mathcal{N}(0, \sigma_g^2)$$

where $c_i^2$ is the $i$-th columnn of $\mathbf{C^2}$ (corresponding to the $p_2$ snps of the $i$-th individual), $\delta$ is a $p_2$-th length vector of corresponding effect size for the $p_2$ global covariates in the $j$-th gene, and $\epsilon_g$ is a component of i.i.d. variation that models measurement noise.

1

# 2  Data Simulations

## 2.1  SNPs data generation

For each gene $j$, the $d$-th snps cis-snps data for person $i$ is sampled with binomial distribution $x_{id} \sim \text{Binomial}(2, \text{MAF}_d), d \in 1, ...m$. The $x_d$ is centered and scaled to zero mean and unit variance. The $p$-th cell-type specific covariate for person $i$ is sampled from random normal distribution $c_{pi}^1 \sim \mathcal{N}(0,1)$. The $p$-th global covariate for bulk level gene expression for persion $i$ is generated similarly: $c_{pi}^2 \sim \mathcal{N}(0,1)$

For $d$-th cis-snps, sparsity is enforced on the effect size of this snps on $k$ cell-type specific gene expression as follow:

$$\beta_{dh} = Y_1^h \times Y_2 \tag{3}$$
$$Y_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_\beta)$$
$$Y_2 \sim \text{Bernoulli}(1 - pslab)$$

Where $\beta_{dh}$ is the $d$-th entry in $\beta_h$ (corresponding to the effect size of $d$-th cis-snps on gene expression for $h$-th cell type); $Y_1^h$ is the $h$-th entry in a multivariate normal random variable $Y_1$, whose covariance matrix is defined by $\mathbf{\Sigma}_\beta \in \mathbf{R}^{k \times k}$ (the effect size of one snps on $k$ cell types are dependent); $Y_2$ is drawn from a Bernoulli distribution with $pslab$ chance of getting zero.

Covariance of $\beta_d$ could be expressed as a parametric function of $pslab$ and $\mathbf{\Sigma}_\beta$:

$$\begin{aligned}
\text{Var}\,(\beta_{dh}) &= E[(Y_1^h)^2(Y_2)^2] - E[Y_1^h]^2 E[Y_2]^2 \\
&= E[(Y_1^h)^2]E[(Y_2)^2] \\
&= \mathbf{\Sigma}_\beta^{hh} \times (1 - pslab) \\
\text{Cov}\,(\beta_{dh_1}, \beta_{dh_2}) &= E[Y_d^{h_1} Y_d^{h_2} Y_{2h_1} Y_{2h_2}] - E[Y_d^{h_1}]E[Y_d^{h_2}]E[Y_2]^2 \\
&= \mathbf{\Sigma}_\beta^{h_1 h_2} \times (1 - pslab)^2 \quad (h_1 \neq h_2)
\end{aligned}$$

$$\tag{4}$$
$$\tag{5}$$

For $p$-th covariate, the effect size on gene expression of cell-type $k$ is sampled from a normal distribution $\gamma_{ph} \sim \mathcal{N}(0, \sigma_\gamma^2)$. Similar to $\gamma$, $\delta$ is generated in the same way.

Cell type weight for $i$-th person, which is a length-$k$ vector, is sampled from Dirichlet distribution: $w_i \sim \text{Dirichlet}(k, \alpha), \alpha \in \mathbf{R}^k$

## 2.2  Cell Type Specific Heritability

Denoting $h$-th cell-type gene expression level as a R.V. $z_h$:

$$z_h = \sum_{d=1}^m x_d \beta_h^d + \sum_{p=1}^{p_1} c_p^1 \gamma_h^p + \epsilon_z \tag{6}$$

where $d$-th snps being $x_d$, effect size of $d$-th snps for $h$-th cell type being $\beta_h^d$, $p$-th covariates being $c_p^1$, effect size of $p$-th covariate for $h$-th cell type being $\gamma_h^p$, noise being $\epsilon_z$. Then heritability for $h$-th cell type is defined as:

$$h_{snps}^2 = \frac{\mathrm{Var}\left(\sum_{d=1}^m x_d \beta_h^d\right)}{\mathrm{Var}\left(z_h\right)} \tag{7}$$

$$= \frac{\mathrm{Var}\left(\sum_{d=1}^m x_d \beta_h^d\right)}{\mathrm{Var}\left(\sum_{d=1}^m x_d \beta_h^d\right) + \mathrm{Var}\left(\sum_{p=1}^{p_1} c_p^1 \gamma_h^p\right) + \mathrm{Var}\left(\epsilon_z\right)}$$

The second line follows as snps, covariates and noise are assumed to be independent. Examining the nominator gives us:

$$\mathrm{Var}\left(\sum_{d=1}^m x_d \beta_h^d\right) = E[(\sum_{d=1}^m x_d \beta_h^d)^2] - E[\sum_{d=1}^m x_d \beta_h^d]^2$$

$$= E[\sum_{d=1}^m x_d^2 (\beta_h^d)^2 + 2 \sum_{d_1 \neq d_2} x_{d_1} x_{d_2} \beta_h^{d_1} \beta_h^{d_2}]$$

$$= \sum_{d=1}^m E[x_d^2 (\beta_h^d)^2] + 2 \sum_{d_1 \neq d_2} E[x_{d_1} x_{d_2} \beta_h^{d_1} \beta_h^{d_2}]$$

$$= \sum_d m(E[x_d^2] - 0)(E[(\beta_h^d)^2 - 0])$$

$$= m \, \mathrm{Var}\left(\beta_h\right) \tag{8}$$

The second line stems from the fact that $E[x_d \beta_d^h] = 0$, the fifth line is a result of $X$ being centered and scaled to unit variance. Similarly, denominator in (7) is calculated as:

$$\mathrm{Var}\left(\sum_{p=1}^{p_1} c_p^1 \gamma_h^p\right) = p_1 \, \mathrm{Var}\left(\gamma_h\right) \tag{9}$$

Bringing (8) and (9) into (7) brings:

$$h_{snps}^2 = \frac{m \, \mathrm{Var}\left(\beta_h\right)}{m \, \mathrm{Var}\left(\beta_h\right) + p_1 \, \mathrm{Var}\left(\gamma_h\right) + \mathrm{Var}\,\epsilon_z}$$

$$\mathrm{Var}\left(\beta_h\right) = \frac{h_{snps}^2 (p_1 \, \mathrm{Var}\left(\gamma_h\right) + \mathrm{Var}\,\epsilon_z)}{m(1 - h_{snps}^2)} \tag{10}$$

Bring (10) into (4) essentialy illustrates that we therefore set

$$\mathbf{\Sigma}_\beta^{hh} = \frac{h_{snps}^2 (p_1 \sigma_\gamma^2 + \sigma_z^2)}{(1 - pslab)(1 - h_{snps}^2)m} \tag{11}$$

In summary, to specify cell-type specific heritability, we specify $h_{snps}^2$, $\sigma_z$, $pslab$, $\sigma_{\gamma_{hj}}$ and calculate diagonal entries in $\Sigma_\beta$.

## 2.3 Genetics Correlation Simulation

According to Rheenen, et.al [Rhe+19],for gene $j$-th expression $\overrightarrow{Z_{h_1}}, \overrightarrow{Z_{h_2}}$ for two cell types, we have

$$z_{h_1} = \sum_{d=1}^{m} x_d \beta_{h_1}^d + \sum_{p=1}^{p_1} c_p^1 \gamma_{h_1}^p + \epsilon_z$$

$$z_{h_2} = \sum_{d=1}^{m} x_d \beta_{h_2}^d + \sum_{p=1}^{p_1} c_p^1 \gamma_{h_2}^p + \epsilon_z$$

the genetic correlation $\rho$ of the gene expression in two cell-types is defined as

$$\rho^{h_1 h_2} = \frac{\text{Cov} \left( \sum_{d=1}^{m} x_d \beta_{h_1}^d, \sum_{d=1}^{m} x_d \beta_{h_2}^d \right)}{\sqrt{\text{Var} \left( \sum_{d=1}^{m} x_d \beta_{h_1}^d \right) \text{Var} \left( \sum_{d=1}^{m} x_d \beta_{h_2}^d \right)}} \tag{12}$$

With (8) we have

$$\rho^{h_1 h_2} = \frac{\text{Cov} \left( \sum_{d=1}^{m} x_d \beta_{h_1}^d, \sum_{d=1}^{m} x_d \beta_{h_2}^d \right)}{m \sqrt{\text{Var} \, \beta_{h_1} \, \text{Var} \, \beta_{h_2}}}$$

$$= \frac{m \, \text{Cov} \left( \beta_{h_1}, \beta_{h_2} \right)}{m \sqrt{\text{Var} \, \beta_{h_1} \, \text{Var} \, \beta_{h_2}}}$$

$$= \rho_\beta^{\beta_{h_1} \beta_{h_2}}$$

$$\text{Cov} \left( \beta_{h_1}, \beta_{h_2} \right) = \rho^{h_1 h_2} \sqrt{\text{Var} \, \beta_{h_1} \, \text{Var} \, \beta_{h_2}} \tag{13}$$

where $\rho_\beta$ is the correlation of effect-size for a gene. (13) essentially requires that correlation of $\beta$ being the same with the genetic correlation. With $\beta$'s correlation matrix specified here and the diagonal entries in $\beta$'s covariance matrix specified by (10), the $\boldsymbol{\Sigma}_\beta$ in (5) could be calculated as

$$\boldsymbol{\Sigma}_\beta^{h_1 h_2} \times (1 - pslab)^2 = \rho^{h_1 h_2} \sqrt{\text{Var} \, \beta_{h_1} \, \text{Var} \, \beta_{h_2}}$$

$$\boldsymbol{\Sigma}_\beta^{h_1 h_2} = \frac{\rho^{h_1 h_2}}{(1 - pslab)^2} \sqrt{\text{Var} \, \beta_{h_1} \, \text{Var} \, \beta_{h_2}} \quad (h_1 \neq h_2) \tag{14}$$

Combinging (14) with (4) presents

$$\Sigma_\beta^{h_1 h_2} = \begin{cases} \frac{\rho^{h_1 h_2}}{(1 - pslab)^2} \sqrt{\text{Var} \, \beta_{h_1} \, \text{Var} \, \beta_{h_2}} & h_1 \neq h_2 \\ \frac{1}{(1 - pslab)} \text{Var} \left( \beta_{h_1} \right) & h_1 = h_2 \end{cases} \tag{15}$$

Where $\text{Var} \left( \beta_h \right)$ given by (10). Recall the definition for correlation

$$-1 \leq \rho = \frac{\text{Cov} \left( X, Y \right)}{\sqrt{\text{Var} \, X \, \text{Var} \, Y}} \leq 1$$

enforces the following constraint:

$$-\sqrt{\Sigma_\beta^{h_1 h_1} \Sigma_\beta^{h_2 h_2}} \le \Sigma_\beta^{h_1 h_2} \le \sqrt{\Sigma_\beta^{h_1 h_1} \Sigma_\beta^{h_2 h_2}}$$

$$pslab - 1 = \rho^{h_1 h_1} \rho^{h_2 h_2}(pslab - 1) \le \rho^{h_1 h_2} \le \rho^{h_1 h_1} \rho^{h_2 h_2}(1 - pslab) = 1 - pslab$$

The constraint makes sense as when $pslab$ is large, then all entries become zero and correlation disappear. Further notice that this is just a multivariate extension for the cell type specific heritability in section(2.2)

In summary, in order to specify the genetics correlation between effect sizes among different cell types, a correlation matrix $\rho \in R^{K \times K}$ is to be provided. This $\rho$ together with $h_{snps}^2$ (heritability for each cell type) determines the parameter for the distribution of effect size $\beta$. Sampling $\beta$ accordnig to (3),(15) and (10) shall guarantee that both heritablity and genetic correlation is as desired.

## 2.4 Bulk Heritability Estimate

Recall equation(1) and equation(2)

$$Z_h^i = (x_i)^T \beta_h + c_i^1 \gamma_h + \epsilon_z$$

$$\epsilon_z \sim \mathcal{N}(0, \sigma_z^2)$$

$$G_i = c_i^2 \delta + \sum_{h=1}^{k} w_{hi} z_{hi} + \epsilon_g$$

$$\epsilon_g \sim \mathcal{N}(0, \sigma_g^2)$$

Suppose the genetic effects $(X\beta)$ have some covariance structure $\Sigma_{X\beta} \in \mathbb{R}^{k \times k}$, the covariance effects $(C^1\gamma)$ have some covariance structure $\Sigma_{C^1\gamma} \in \mathbb{R}^{k \times k}$. Then the covariance of Z across cell types $\Sigma_Z = \Sigma_{X\beta} + \Sigma_{C^1\gamma} + \text{diag}(\sigma_z^2)$

$$
\begin{aligned}
\text{Var}(G) &= \text{Var}(\sum_{h=1}^{k} w_h z_h + \sum_{p=1}^{p_2} c_p^2 \delta_p + \epsilon_g) \\
&= \text{Var}(\sum_{h=1}^{k} w_h z_h) + \text{Var}(\sum_{p=1}^{p_2} c_p^2 \delta_p) + \sigma_{\epsilon_g}^2 \\
&= \left( \sum_{h=1}^{k} \sum_{l=1}^{k} \text{Cov}(w_h z_h, w_l z_l) \right) + p_2 \sigma_\delta^2 + \sigma_{\epsilon_g}^2 \\
&= \left( \sum_{h=1}^{k} \sum_{l=1}^{k} E[w_h z_h w_l z_l] \right) + p_2 \sigma_\delta^2 + \sigma_{\epsilon_g}^2 \quad \text{(assuming each } z \text{ is centered)} \\
&= \left( \sum_{h=1}^{k} \sum_{l=1}^{k} E[w_h w_l] E[z_h z_l] \right) + p_2 \sigma_\delta^2 + \sigma_{\epsilon_g}^2 \quad (16)
\end{aligned}
$$

Where

$$E[z_h z_l] = \Sigma_{z\{h,l\}} = \Sigma_{X\beta} + \Sigma_{C^1\gamma} + \text{diag}(\sigma_g^2)$$

$$E[w_h w_l] = \begin{cases} \frac{\tilde{\alpha_h}(1-\tilde{\alpha_h})}{\alpha_0+1} + \tilde{\alpha_h}^2 & h = l \\ \tilde{\alpha_h}\tilde{\alpha_l}(1 - \frac{1}{\alpha_0+1}) & h \neq l \end{cases}$$

Define each entry $\{h,l\}$ of $\Sigma_\alpha \in \mathcal{R}^{k \times k}$ to be $E[w_h w_l]$. Then

$$\begin{aligned} \text{var}(G) &= \text{sum}(\Sigma_\alpha \odot \Sigma_Z) + p_2\sigma_\gamma^2 + \sigma_{\epsilon_g}^2 \\ &= \text{sum}(\Sigma_\alpha \odot \Sigma_{X\beta}) + \text{sum}(\Sigma_\alpha \odot \Sigma_{C^1\gamma}) \\ &\quad + \text{sum}(\Sigma_\alpha \odot \text{diag}(\sigma_g^2)) + p_2\sigma_\gamma^2 + \sigma_{\epsilon_g}^2 \end{aligned} \tag{17}$$

where the 'sum' operator is the sum of each entry of the argument matrix. Heritability of bulk expression in this model can then be defined as

$$h_{\text{bulk}}^2 = \frac{\text{sum}(\Sigma_\alpha \odot \Sigma_{X\beta})}{\text{var}(G)} \tag{18}$$

Therefore, by varying $\sigma_\gamma$ and $\sigma_{\epsilon_g}$, bulk level heritability could be modified as well.

# References

[Rhe+19]   Wouter van Rheenen et al. "Genetic correlations of polygenic disease traits: from theory to practice". In: *Nature Reviews Genetics* (2019). ISSN: 1471-0064. DOI: 10.1038/s41576-019-0137-z. URL: https://doi.org/10.1038/s41576-019-0137-z.