

Supplementary Information: Genome-wide association between transcription factor expression and chromatin accessibility reveals regulators of chromatin accessibility.

Supplementary Methods

As mentioned in the main text, we use the following linear mixed model,

$$\mathbf{y} = \mathbf{x}_i \beta^i + \boldsymbol{\delta}^i + \boldsymbol{\epsilon}^i,$$

where \mathbf{y} is a vector of motif accessibility scores across n cell lines, \mathbf{x}_i is the expression vector of gene i , β^i is the effect size of gene i :

$$\boldsymbol{\epsilon} \sim N_n(\mathbf{0}, \sigma_r^2 \mathbf{I}_n),$$

and

$$\boldsymbol{\delta} \sim N_n(\mathbf{0}, \sigma_e^2 \mathbf{C}_e),$$

with

$$\mathbf{C}_e = \frac{1}{p} \sum_{i=1}^p \mathbf{x}_i \mathbf{x}_i^T.$$

The likelihood function of is

$$f(\mathbf{y}) = \frac{1}{(2\pi)^{n/2} |\sigma_e^2 \mathbf{C}_e + \sigma_r^2 \mathbf{I}_n|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\mathbf{y} - \mathbf{x}_i \beta^i)^T (\sigma_e^2 \mathbf{C}_e + \sigma_r^2 \mathbf{I}_n)^{-1} (\mathbf{y} - \mathbf{x}_i \beta^i)\right).$$

We define the spectral decomposition of \mathbf{C}_e as:

$$\mathbf{C}_e = \boldsymbol{\Gamma} \boldsymbol{\Lambda} \boldsymbol{\Gamma}^T.$$

For any values of σ_e^2 and σ_r^2 we have

$$(\sigma_e^2 \mathbf{C}_e + \sigma_r^2 \mathbf{I}_n) = \boldsymbol{\Gamma} (\sigma_e^2 \boldsymbol{\Lambda}_e + \sigma_r^2 \mathbf{I}_n) \boldsymbol{\Gamma}^T,$$

i.e.: the eigenvectors of the mixture matrix are constant w.r.t the mixing parameters. Set

$$\begin{aligned} \mathbf{y}' &= \boldsymbol{\Gamma}^T \mathbf{y}, \\ \mathbf{x}' &= \boldsymbol{\Gamma}^T \mathbf{x}. \end{aligned}$$

Since the likelihood is invariant to rotations, we have

$$f(\mathbf{y}') = f(\mathbf{y})$$

and

$$f(\mathbf{y}') = \frac{1}{(2\pi)^{n/2} |\sigma_e^2 \boldsymbol{\Lambda}_e + \sigma_r^2 \mathbf{I}_n|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\mathbf{y}' - \mathbf{x}'_i \beta^i)^T (\sigma_e^2 \boldsymbol{\Lambda}_e + \sigma_r^2 \mathbf{I}_n)^{-1} (\mathbf{y}' - \mathbf{x}'_i \beta^i)\right).$$

Reparametrizing with

$$\gamma = \sigma_r^2 / \sigma_e^2,$$

the log-likelihood becomes

$$l(\mathbf{y}') = \frac{n}{2} \log(2\pi) - \frac{n}{2} \sum \log(\sigma_e^2(\lambda_i + \gamma)) - \frac{n}{2} \sum \frac{(y'_k - x'_{ki}\beta_i)^2}{(2\sigma_e^2(\lambda_i + \gamma))}$$

Partial derivation shows that the maximum of the log likelihood is reached at

$$\hat{\beta}^i = \sum_k \frac{y'_k x'_{ki}}{(\lambda_k + \hat{\gamma})} / \sum_k \frac{(x'_{ki})^2}{(\lambda_k + \hat{\gamma})},$$

and

$$\hat{\sigma}_e^2 = \sum_k \frac{(y'_k - x'_{ki}\beta^i)^2}{(\lambda_k + \hat{\gamma})} / \sum_k \frac{n}{(\lambda_k + \hat{\gamma})}.$$

Reducing the 3 parameter optimization problem to a one parameter optimization over γ . p -values can be obtained by the likelihood ratio test for null hypothesis that $\beta = 0$.