

Tradict - mathematical details

Surojit Biswas, Konstantin Kerner, Paulo José Pereira Lima Texeira,
Jeffery L. Dangl, Vladimir Jojic, Philip A. Wigge

3 Contents

This document describes the full mathematical details for the concepts presented in the “Tradict algorithm” section, “Building a predictive Multivariate Normal Continuous-Poisson hierarchical model” subsection of the Materials and Methods in the Supplemental Information. Specifically, we present exactly how Tradict uses a selected set of markers to 1) complete training, and 2) to perform prediction.

8 1 Preliminaries

For a matrix A , $A_{:,i}$ and $A_{i,:}$ index the i^{th} column and row, respectively. For a set of indices, q , we use $-q$ to refer to all indices not specified by q .

11 2 Model

Tradict uses a Continuous-Poisson Multivariate Normal (CP-MVN) hierarchical model to model the expression of transcriptional programs and all genes in the transcriptome. Multivariate Normal hierarchies have been explored in the past as a means of modeling correlation structure among count based random variables [REF]. However, given we will be working with abundances as transcripts per million (TPM), which are non-negative (can equal zero) and fractional, we relax the integral assumption of the Poisson so it is continuous on $[0, \infty)$. Specifically, we define the continuous relaxation of the Poisson distribution (hereafter, Continuous-Poisson) to have the following density function:

$$f(x|\lambda) = C_\lambda \frac{e^{-\lambda} \lambda^x}{\Gamma(x+1)}$$

where C_λ is a normalization constant [REF]. The mean of this distribution is given by λ , just as the Poisson. We begin by building a predictive model of gene expression, and thereafter discuss a predictive model for the expression of transcriptional programs. Let z_j denote the log-latent abundance of gene j , such that $\exp(z_j)$ is the latent abundance of that gene (in TPM) whose measured abundance is given by t_j . Let $T_j = t_j o$ be the measured total number of transcripts of gene j . Here o is the sequencing depth in millions of reads of the sample under consideration. We assume then,

$$z \sim \mathcal{N}(\mu, \Sigma)$$
$$T_j \sim \text{Continuous-Poisson}(\exp(z_j) o)$$

where μ and Σ are of dimension $1 \times \#$ -genes and $\#$ -genes \times $\#$ -genes, respectively. In effect, we are assuming that the measured number of transcripts for gene j is a noisy realization of a latent abundance $\exp(z_j)$ times the sequencing depth, o . The dependencies between log-latent abundances (the z_j 's) are then encoded by the covariance matrix of the Multivariate Normal layer of the model.

Note that we could model the TPM measurements directly in the second layer by assuming $t_j \sim \text{Continuous-Poisson}(\exp(z_j))$; however, this formulation does not consider sequencing depth, which can be a valuable source of information when inferring latent abundances for rare/poorly sampled genes ?).

During decoding, we are interested in building a predictive model between markers and all genes in the transcriptome. Therefore, we need to consider a conditional model of the transcriptome given the log-latent abundances of the markers. Let m be the set of indices for the given panel of selected markers, which are the subset of genes Tradict selects as representative of the transcriptome. To perform prediction we therefore need $p(z_{-m}|z_m)$, and given this we would like to ultimately compute as our estimate of the abundance of all genes in the transcriptome $\hat{T} = \text{argmax}_T p(T|z_m)$. We have,

$$\begin{aligned} z_m &\sim \mathcal{N}(\mu^{(m)}, \Sigma^{(m)}) \\ z_{-m}|z_m &\sim \mathcal{N}(\mu_{z_{-m}|z_m}, \Sigma_{z_{-m}|z_m}) \\ T_j &\sim \text{Continuous-Poisson}(\exp(z_j)o) \end{aligned}$$

Here, $\mu^{(m)}$ and $\Sigma^{(m)}$ refer to mean vector and covariance matrix of z_m . Given these, the conditional mean of the log-latent abundances for all non-marker genes can be obtained through Gaussian conditioning. Specifically, for two normally distributed row-vector variables a and b the conditional mean of b given a is given by $\mu_{b|a} = \mu_b + (a - \mu_a)\Sigma_a^{-1}\sigma_{ab}$ and $\Sigma_{b|a} = \Sigma_b - \sigma_{ab}^T\Sigma_a^{-1}\sigma_{ab}$, where σ_{ab} is the cross-covariance between a and b , and Σ_a and Σ_b are the covariance matrices of a and b , respectively.

Given the expression of a transcriptional program is a linear combination of the latent abundances of its constituent genes, they will be normally distributed given 1) Central Limit Theorem, and 2) the latent abundances themselves are normally distributed (convolutions of normals are normals). Let s be the expression of all transcriptional programs. We posit the following model,

$$\begin{aligned} z_m &\sim \mathcal{N}(\mu^{(m)}, \Sigma^{(m)}) \\ s|z_m &\sim \mathcal{N}(\mu_{s|z_m}, \Sigma_{s|z_m}) \end{aligned}$$

To use these models for prediction, we must learn their parameters from training data. This would complete the process of encoding described in the Supplemental Information. Specifically, we need to learn $\mu^{(m)}$, $\Sigma^{(m)}$, μ_s , $\mu_{z_{-m}}$, $\sigma_{z_m,s}$ and $\sigma_{z_m,z_{-m}}$.

3 Training

As described in the Supplemental Information, given an estimate of z_m , \hat{z}_m , inference of μ_s , $\mu_{z_{-m}}$, $\sigma_{z_m,s}$ and $\sigma_{z_m,z_{-m}}$ is straightforward. In lag transforming the entire training TPM expression matrix, $t \in \mathbb{R}^{\text{samples} \times \text{genes}}$, we have an estimate of z , $\hat{z} = \text{lag}(t)$?). Thus, an estimate of $\mu_{z_{-m}}$ is given by the usual column-wise sample mean of \hat{z}_{-m} .

Let $\Lambda \in \mathbb{R}^{\text{genes} \times \text{transcriptional programs}}$ be a matrix of principal component 1 coefficients over genes for each transcriptional program. Note, that $\Lambda_{ij} = 0$ if gene i is not in transcriptional program j . An estimate of s is given by $\hat{s} = \hat{z}\Lambda$, and so an estimate for μ_s , $\hat{\mu}_s$, is given by the usual column-wise mean of \hat{s} .

Given \hat{z}_m the cross-covariances, $\sigma_{z_m,s}$ and $\sigma_{z_m,z_{-m}}$, are given by the usual sample cross-covariance between \hat{z}_m and \hat{s} and between \hat{z}_m and \hat{z}_{-m} , respectively.

Now, though we could use the lag-transformed values of t_m as our estimate for z_m , we have an opportunity to improve this estimate by virtue of having to estimate $\mu^{(m)}$ and $\Sigma^{(m)}$. More specifically, given z_m , estimates of $\mu^{(m)}$ and $\Sigma^{(m)}$ are given by – up to some regularization – the usual sample mean and covariance of z_m . Furthermore, given $\mu^{(m)}$ and $\Sigma^{(m)}$, we can update our estimate of z_m to the maximum of its posterior distribution. This suggests an alternating iterative procedure in which we iterate 1) estimation of $\mu^{(m)}$ and $\Sigma^{(m)}$, and 2) maximum *a posteriori* inference of z_m until convergence of their joint likelihood. It is the \hat{z}_m that we obtain from this procedure that we use in the cross-covariance calculations above. The following section details this procedure.

68 3.1 Inference of z_m given $\mu^{(m)}$ and $\Sigma^{(m)}$

69 Suppose Tradict has estimates of $\mu^{(m)}$ and $\Sigma^{(m)}$ given by $\hat{\mu}^{(m)}$ and $\hat{\Sigma}^{(m)}$, and let $T_m = t_m(o \times \mathbf{1}_{1 \times \text{markers}})$
70 be a matrix of the total measured number of transcripts for each marker. Here $o \in \mathbb{R}^{\text{samples} \times 1}$ is a vector
71 of sample sequencing depths in millions of reads. Given these, we would like to calculate the maximum *a*
72 *posteriori* (MAP) estimate of $\hat{z}_m = \text{argmax}_{z_m} p(z_m | o, T_m, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})$.
73 The posterior distribution over z_m is given by

$$\begin{aligned} p(z_m | o, T_m, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) &= \frac{p(T_m | o, z_m, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) p(z_m | \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})}{\int_k p(T_m | o, k, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) p(k | \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) dk} \\ &\propto \prod_{i=1}^n p(T_{im} | o, z_{im}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) p(z_{im} | \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) \\ &= \prod_{i=1}^n \left[\prod_{j=1}^{|m|} C_{[\exp(z_{ij})o_i]} [\exp(z_{ij})o_i]^{T_{ij}} e^{-[\exp(z_{ij})o_i]} / \Gamma(T_{ij} + 1) \right] \\ &\quad \times \frac{1}{\sqrt{2\pi|\hat{\Sigma}^{(m)}|}^{|m|}} \exp \left(-\frac{1}{2} (z_{i:} - \hat{\mu}^{(m)}) \text{inv} \left(\hat{\Sigma}^{(m)} \right) (z_{i:} - \hat{\mu}^{(m)})^T \right) \end{aligned}$$

74 where for notational clarity we have used $\text{inv}(\cdot)$ to represent matrix inverse.

75 Given z is a matrix parameter, this may be difficult to solve directly. However, note that given z_{ij} , T_{ij}
76 is conditionally independent of $T_{i,-j}$. Additionally, given $z_{i,-j}$, z_{ij} is normally distributed with mean and
77 covariance

$$\begin{aligned} a_{ij} &= \mu_j^{(m)} + (z_{i,-j} - \mu_{-j}^{(m)}) \text{inv} \left(\Sigma_{-j,-j}^{(m)} \right) \Sigma_{-j,j}^{(m)} \\ \sigma_{m(j)} &= \Sigma_{j,j}^{(m)} - \Sigma_{j,-j}^{(m)} \text{inv} \left(\Sigma_{-j,-j}^{(m)} \right) \Sigma_{-j,j}^{(m)} \end{aligned}$$

78 respectively. Taken together, this suggests an iterative conditional modes algorithm ?) in which we maximize
79 the posterior one column of z at a time, while conditioning on all others.

80 Let \hat{z}_m denote our current estimate of z_m . Let $m(j)$ denote the index of the j^{th} marker and let $m(-j)$
81 denote the indices of all markers but the j^{th} one. The above sub-objective is given by,

$$\begin{aligned} \hat{z}_{im(j)} &= \text{argmax}_{z_{im(j)} | z_{im(-j)}} \log p(z_{im(j)} | T_{im(j)}, o_i, \hat{z}_{im(-j)}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) \\ &= \text{argmax}_{z_{im(j)} | z_{im(-j)}} \log p(T_{im(j)} | z_{im(j)}, o_i, \hat{z}_{im(-j)}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) p(z_{im(j)} | \hat{z}_{im(-j)}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) \\ &= \text{argmax}_{z_{im(j)} | z_{im(-j)}} \log p(T_{im(j)} | z_{im(j)}, o_i) p(z_{im(j)} | \hat{z}_{im(-j)}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) \\ &= \text{argmax}_{z_{im(j)} | z_{im(-j)}} \log \left[[\exp(z_{im(j)})o_i]^{T_{im(j)}} e^{-[\exp(z_{im(j)})o_i]} \exp \left(-\frac{1}{2\sigma_{m(j)}} (z_{im(j)} - a_{im(j)})^2 \right) \right] \\ &= \text{argmax}_{z_{im(j)} | z_{im(-j)}} T_{im(j)} \exp(z_{im(j)})o_i - \exp(z_{im(j)})o_i - \frac{1}{2\sigma_{m(j)}} (z_{im(j)} - a_{im(j)})^2 \end{aligned}$$

82 Differentiating we get,

$$\begin{aligned} \frac{\partial}{\partial z_{im(j)}} T_{im(j)} z_{im(j)} o_i - \exp(z_{im(j)}) o_i - \frac{1}{2\sigma_{m(j)}} (z_{im(j)} - a_{im(j)})^2 \\ = T_{im(j)} o_i - \exp(z_{im(j)}) o_i - \frac{1}{\sigma_{m(j)}} (z_{im(j)} - a_{im(j)}) \end{aligned}$$

83 Because $z_{im(j)}$ appears as a linear and exponential term, we cannot solve this gradient analytically. We
 84 therefore utilize Newton-Raphson optimization. For this we also require the Hessian, which is given by,

$$\begin{aligned} \frac{\partial}{\partial z_{im(j)}} T_{im(j)} o_i - \exp(z_{im(j)}) o_i - \frac{1}{\sigma_{m(j)}} (z_{im(j)} - a_{im(j)}) \\ = -\exp(z_{im(j)}) o_i - \frac{1}{\sigma_{m(j)}} < 0 \end{aligned}$$

85 Notice the Hessian is always negative-definite, which implies each update has a single, unique optimum.

86 In practice, the Newton-Raphson updates can be performed in vectorized fashion iteratively for each
 87 column of z . We generally find that this optimization takes 5-15 iterations (full passes over all columns
 88 of z) and less than a minute to converge. We refer to the program that performs these calculations as

89 $\hat{z}_m = \text{MAP_Z} \left(t, o, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)} \right)$.

90 3.2 Complete inference of $\mu^{(m)}$, $\Sigma^{(m)}$, and z_m

91 For complete inference we use the following iterative conditional modes algorithm ?):

- 92 • Initialize $T_m = t_m(o \times \mathbf{1}_{1 \times \text{markers}})$, $\hat{z}_m = \text{lag}(t_m)$.
- 93 • Until convergence of $\log p(T_m | o, \hat{z}_m, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) + \log p(\hat{z}_m | \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})$, iterate:
- 94 – Update $\hat{\mu}^{(m)}$ and $\hat{\Sigma}^{(m)}$:

$$\begin{aligned} \hat{\mu}^{(m)} &= \frac{1}{\# \text{samples}} \sum_i \hat{z}_{im} \\ \hat{\Sigma}^{(m)} &= \frac{1}{\# \text{samples} - 1} \sum_i (\hat{z}_{im} - \hat{\mu}^{(m)})^T (\hat{z}_{im} - \hat{\mu}^{(m)}) + \lambda \text{diag} \left[\text{cov} \left(\hat{z}_m^{(\text{init})} \right) \right] \end{aligned}$$

- 95 – Update $\hat{z}_m = \text{MAP_Z} \left(t, o, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)} \right)$.

96 Here $\text{diag}(x)$ of the square matrix x returns an equivalently sized matrix with only the diagonal of x preserved
 97 and 0's for the off-diagonal terms. $\text{cov}(\cdot)$ denotes the usual sample covariance matrix.

98 Note that in this algorithm we have added a regularization to the estimate of the covariance matrix.
 99 This is done in order to ensure stability and to avoid infinite-data-likelihood singularities that arise from
 100 singular covariance matrices. This is most often happens when a gene's TPM abundance is mostly zero (i.e.
 101 there is little data for the gene), giving the multivariate normal layer an opportunity to increase the data
 102 likelihood (via the determinant of the covariance matrix) by tightly coupling this gene's latent abundance
 103 to that of another gene, thereby producing a singularity. This regularization is probabilistically equivalent
 104 to adding an Inverse-Wishart prior over $\Sigma^{(m)}$. The parameter λ controls the strength of the regularization.
 105 In practice, we find $\lambda = 0.1$ leads to good predictive performance, stable (non-singular) covariance matrices,
 106 and reasonably quick convergence.

107 4 Prediction

108 During prediction we are given new measured TPM measurements for our markers, $t_m^* \in \mathbb{R}^{\text{query samples} \times |m|}$,
 109 and we must make predictions about the expression of all transcriptional programs and the remaining non-
 110 marker genes. We have two options available to us: 1) Calculate a point (MAP) estimate or 2) calculate the
 111 complete posterior distribution over each non-marker gene and transcriptional program in a fully Bayesian
 112 manner. The former option is faster, but the second gives more information on the uncertainty of the
 113 prediction. We therefore implement both options in Tradict and detail their derivation below. Note that
 114 knowing the entire posterior distribution allows one to derive whatever estimator they would like, and so
 115 option 2, informationally speaking, supersedes option 1.

116 4.1 MAP estimation of gene and program abundances

117 We first need an estimate of the log-latent abundances \hat{z}_m^* associated with t_m^* . Given the estimates $\hat{\mu}^{(m)}$ and
 118 $\hat{\Sigma}^{(m)}$ obtained from the training data, we obtain these estimates as

$$\hat{z}_m^* = \text{MAP_Z} \left(t_m^*, \mathbf{1}_{\text{query samples} \times 1}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)} \right)$$

119 Given the inferred marker latent abundances, we let our estimates of s^* and t_m^* be the maximizers of
 120 their probability distribution. In other words, $\hat{s}^* = \text{argmax}_{s^*} p(s^* | \hat{z}_m^*)$ and $\hat{t}_m^* = \text{argmax}_{t_m^*} p(t_m^* | \hat{z}_m^*)$.

121 Our estimate for the expression of all transcriptional programs is given by

$$\text{argmax}_{s^*} p(s^* | \hat{z}_m^*) = \mathbb{E}[s^* | \hat{z}_m^*] = \mu_{s^* | \hat{z}_m^*} = \hat{\mu}_s + \left(\hat{z}_m^* - \hat{\mu}^{(m)} \right) \text{inv} \left(\hat{\Sigma}^{(m)} \right) \hat{\sigma}_{z_m, s}.$$

122 Here, $\hat{\mu}_s$ and $\hat{\sigma}_{z_m, s}$ represent estimates of the unconditional mean of s and the cross-covariance matrix
 123 between z_m and s previously learned during encoding.

124 Similarly, for the entire transcriptome we have,

$$\hat{t}_{ij}^* = \text{argmax}_t p(t | \hat{z}_{im}^*) = \exp \left(\mu_{z_{ij} | \hat{z}_{im}^*} \right).$$

125 where,

$$\mu_{z_{ij} | \hat{z}_{im}^*} = \hat{\mu}_j + \left(\hat{z}_{im}^* - \hat{\mu}^{(m)} \right) \text{inv} \left(\hat{\Sigma}^{(m)} \right) \hat{\sigma}_{z_m, z_j}$$

126 We could also use the expected value of t as our estimate.

$$\begin{aligned} \mathbb{E}[t_{ij}^* | \hat{z}_{im}^*] &= \int_{-\infty}^{\infty} \mathbb{E}[t_{ij}^* | z_{ij}] p(z_{ij} | \hat{z}_{im}^*) dz_{ij} \\ &= \int_{-\infty}^{\infty} \exp(z_{ij}) \mathcal{N}(z_{ij} | \mu_{z_{ij} | \hat{z}_{im}^*}, \Sigma_{z_{ij} | \hat{z}_{im}^*}) dz_{ij} \\ &= \mathbb{E}_{\mathcal{N}}[\exp(z_{ij}) | \hat{z}_{im}^*] \end{aligned}$$

127 The Moment Generating Function of a Normal random variable X with mean μ and variance σ^2 is given by

128 $M(t) = \mathbb{E}[\exp(tX)] = \exp(\mu t + \sigma^2 t^2 / 2)$. Therefore we have,

$$\mathbb{E}[t_{ij}^* | \hat{z}_{im}^*] = \mathbb{E}_{\mathcal{N}}[\exp(z_{ij}) | \hat{z}_{im}^*] = M(1) = \exp \left(\mu_{z_{ij} | \hat{z}_{im}^*} + \frac{1}{2} \Sigma_{z_{ij} | \hat{z}_{im}^*} \right)$$

129 where,

$$\begin{aligned} \mu_{z_{ij} | \hat{z}_{im}^*} &= \hat{\mu}_j + \left(\hat{z}_{im}^* - \hat{\mu}^{(m)} \right) \text{inv} \left(\hat{\Sigma}^{(m)} \right) \hat{\sigma}_{z_m, z_j} \\ \Sigma_{z_{ij} | \hat{z}_{im}^*} &= \hat{\sigma}_{jj} - \hat{\sigma}_{z_m, z_j}^T \text{inv} \left(\hat{\Sigma}^{(m)} \right) \hat{\sigma}_{z_m, z_j} \end{aligned}$$

130 Here, $\hat{\mu}_j$ and $\hat{\sigma}_{z_m, z_j}$ represent estimates of the unconditional mean of z_j and the cross-covariance matrix
 131 between z_m and z_j . These were learned from the training data during encoding.

132 Though this predictor is unbiased, it does not produce a good prediction for most samples. This is due
 133 to the right-skew of the Poisson, which drags its mean away from the most likely values.

134 4.2 Posterior density estimation of gene and program abundances

135 The above predictions represent point estimates. Ideally, we would like to know the uncertainty around these
 136 estimates. Given measurements of the representative markers, we can estimate the posterior distribution of

137 expression values for transcriptional programs and the non-markers, and therein calculate any point estimates
 138 and/or measures of uncertainty. Recall that for transcriptional programs:

$$z_m \sim \mathcal{N}(\mu^{(m)}, \Sigma^{(m)})$$

$$s|z_m \sim \mathcal{N}(\mu_{s|z_m}, \Sigma_{s|z_m})$$

139 And similarly for genes (among which the marker genes are included) we have:

$$z_m \sim \mathcal{N}(\mu^{(m)}, \Sigma^{(m)})$$

$$z_{-m}|z_m \sim \mathcal{N}(\mu_{z_{-m}|z_m}, \Sigma_{z_{-m}|z_m})$$

$$T_j \sim \text{Continuous-Poisson}(\exp(z_j)o)$$

140 Given z_m , the distribution of expression values are simple normal distributions with analytically available
 141 means and covariances. However, because z_m is unknown, we must factor into our estimate its distribution,
 142 which is both a function of observed data (t_m, o) and prior information (in the form of $\hat{\mu}^{(m)}$ and $\hat{\Sigma}^{(m)}$). Our
 143 strategy to estimate the posterior density of programs and non-markers will therefore be to sample from
 144 the posterior of z_m , and then given these draws, sample from the conditional Normal distribution of each
 145 program and non-marker gene.

146 4.2.1 Sampling z_m via MCMC

147 To sample z_m we use Metropolis-Hastings Markov Chain Monte Carlo (MCMC) sampling [REF], using the
 148 following posterior density function:

$$p(z_m|o, T_m, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}) = \frac{p(T_m|o, z_m, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})p(z_m|\hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})}{\int_k p(T_m|o, k, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})p(k|\hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})dk}$$

$$\propto \prod_{i=1}^n p(T_{im}|o, z_{im}, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})p(z_{im}|\hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})$$

$$= \prod_{i=1}^n \left[\prod_{j=1}^{|m|} C_{[\exp(z_{ij})o_i]} [\exp(z_{ij})o_i]^{T_{ij}} e^{-[\exp(z_{ij})o_i]} / \Gamma(T_{ij} + 1) \right]$$

149 Note that we do not require the marginal distribution for Metropolis-Hastings sampling.

150 As our proposal distribution we use:

$$z_m^{(i+1)} = \mathcal{N}(z_m^{(i)}, \gamma \mathbb{I}_{|m| \times |m|}).$$

151 Here $z^{(i)}$ is the i^{th} draw from the sampler, and γ represents the width (variance) of the proposal distribution.
 152 To choose this width, we examine, for a schedule of proposal widths (50 logarithmically spaced widths between
 153 $10^{3.5}$ and 10^{-1}), which width gives an acceptance rate closest to 0.4 [REF]. Using this width, we sample
 154 20,000 times from the sampler. We burn-in the first 100 samples and keep every 100th sample thereafter (to
 155 offset the effects of the chain's auto-correlation) as our draws from the distribution. Note that we initialize
 156 the chain at the MAP estimate of z_m . This ensures the chain is stationary from the beginning.

157 4.2.2 Sampling program and gene abundances

158 Given our $M = 200$ draws, $\left[z_m^{(i)} \right]_{i=1}^M$, we can sample from the conditional distribution of each program and
 159 gene.

160 Our $i^{(th)}$ draw from the posterior distribution over all programs is obtained from sampling the following
 161 Multivariate-Normal,

$$s^{(i)}|z_m^{(i)} \sim \mathcal{N}\left(\mu_{s|z_m^{(i)}}, \Sigma_{s|z_m^{(i)}}\right)$$

162 where

$$\begin{aligned}\mu_{s|z_m^{(i)}} &= \hat{\mu}_s + \left(z_m^{(i)} - \hat{\mu}^{(m)}\right) \text{inv}\left(\hat{\Sigma}^{(m)}\right) \hat{\sigma}_{z_m, s} \\ \Sigma_{s|z_m^{(i)}} &= \hat{\Sigma}_s - \hat{\sigma}_{z_m, s}^T \text{inv}\left(\hat{\Sigma}^{(m)}\right) \hat{\sigma}_{z_m, s}\end{aligned}$$

163 Similarly, our $i^{(th)}$ draw from the posterior distribution over all genes *could be* obtained from sampling
 164 the following Multivariate-Normal,

$$z_{-m}^{(i)}|z_m^{(i)} \sim \mathcal{N}\left(\mu_{z_{-m}|z_m^{(i)}}, \Sigma_{z_{-m}|z_m^{(i)}}\right)$$

165 where

$$\begin{aligned}\mu_{z_{-m}|z_m^{(i)}} &= \hat{\mu}_{z_{-m}} + \left(z_m^{(i)} - \hat{\mu}^{(m)}\right) \text{inv}\left(\hat{\Sigma}^{(m)}\right) \hat{\sigma}_{z_m, z_{-m}} \\ \Sigma_{z_{-m}|z_m^{(i)}} &= \hat{\Sigma}_{z_{-m}} - \hat{\sigma}_{z_m, z_{-m}}^T \text{inv}\left(\hat{\Sigma}^{(m)}\right) \hat{\sigma}_{z_m, z_{-m}}\end{aligned}$$

166 However, given the size of $\Sigma_{z_{-m}|z_m^{(i)}}$ (approximately 21000×21000), this is not easily doable. Recall, though,
 167 that one of our basic assumptions is that the conditional mean abundance of all genes given the abundance
 168 of our markers has the covariance structure of all genes sufficiently built in. Thus, we assume

$$\mathcal{N}\left(\mu_{z_{-m}|z_m^{(i)}}, \Sigma_{z_{-m}|z_m^{(i)}}\right) \approx \mathcal{N}\left(\mu_{z_{-m}|z_m^{(i)}}, \text{diag}\left(\Sigma_{z_{-m}|z_m^{(i)}}\right)\right)$$

169 Here $\text{diag}(\cdot)$ replaces all off-diagonal entries with zeros. Consequently, we only need to compute the diagonal
 170 entries of the conditional covariance matrix, and furthermore, given the conditional mean of each, we can
 171 sample it's abundance in parallel and independently of all others.

172 5 References

- 173 Surojit Biswas. The latent logarithm. *arXiv*, pages 1–11, 2016.
- 174 Julian Besag. On the Statistical Analysis of Dirty Pictures. *Journal of the Royal Statistical Society*, 48(3):259–
 175 302, 1986.