

# Tradict - mathematical details

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

## 3 Contents

4	1 Preliminaries	1
5	2 Model	1
6	3 Training	2
7	3.1 Inference of $z_m$ given $\mu^{(m)}$ and $\Sigma^{(m)}$	3
8	3.2 Complete inference of $\mu^{(m)}$ , $\Sigma^{(m)}$ , and $z_m$	4
9	4 Prediction	5
10	4.1 MAP estimation of gene and program abundances	5
11	4.2 Posterior density estimation of gene and program abundances	6
12	4.2.1 Sampling $z_m$ via MCMC	6
13	4.2.2 Sampling program and gene abundances	7
14	5 References	8

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.

## 19 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  $q$ .

## 22 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 [1, 2, 3, 4]. 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:

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

30 where  $C_\lambda$  is a normalization constant. The mean of this distribution is given by  $\lambda$ , just as the Poisson.

31 We begin by building a predictive model of gene expression, and thereafter discuss a predictive model  
 32 for the expression of transcriptional programs. Let  $z_j$  denote the log-latent abundance of gene  $j$ , such that  
 33  $\exp(z_j)$  is the latent abundance of that gene (in TPM) whose measured abundance is given by  $t_j$ . Let  
 34  $T_j = t_j o$  be the measured total number of transcripts of gene  $j$ . Here  $o$  is the sequencing depth in millions  
 35 of reads of the sample under consideration. We assume then,

$$z \sim \mathcal{N}(\mu, \Sigma)$$

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

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

40 Note that we could model the TPM measurements directly in the second layer by assuming  $t_j \sim$   
 41  $\text{Continuous-Poisson}(\exp(z_j))$ ; however, this formulation does not consider sequencing depth, which can be a  
 42 valuable source of information when inferring latent abundances for rare/poorly sampled genes [5].

43 During prediction, we are interested in building a predictive model between markers and all genes in the  
 44 transcriptome. Therefore, we need to consider a conditional model of the transcriptome given the log-latent  
 45 abundances of the markers. Let  $m$  be the set of indices for the given panel of selected markers, which are the  
 46 subset of genes Tradict selects as representative of the transcriptome. To perform prediction we therefore  
 47 need  $p(z_{-m}|z_m)$ . We have,

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

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

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

48 Here,  $\mu^{(m)}$  and  $\Sigma^{(m)}$  refer to mean vector and covariance matrix of  $z_m$ . Given these, the conditional  
 49 mean of the log-latent abundances for all non-marker genes can be obtained through Gaussian conditioning.  
 50 Specifically, for two normally distributed row-vector variables  $a$  and  $b$  the conditional mean of  $b$  given  $a$  is  
 51 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  $\sigma_{ab}$  is the cross-covariance between  
 52  $a$  and  $b$ , and  $\Sigma_a$  and  $\Sigma_b$  are the covariance matrices of  $a$  and  $b$ , respectively.

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

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

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

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

## 60 3 Training

61 As described in the Supplemental Information, given an estimate of  $z_m$ ,  $\hat{z}_m$ , inference of  $\mu_s$ ,  $\mu_{z_{-m}}$ ,  $\sigma_{z_m,s}$  and  
 62  $\sigma_{z_m,z_{-m}}$  is straightforward. In lag transforming the entire training TPM expression matrix,  $t \in \mathbb{R}^{\text{samples} \times \text{genes}}$ ,  
 63 we have an estimate of  $z$ ,  $\hat{z} = \text{lag}(t)$  [5]. Thus, an estimate of  $\mu_{z_{-m}}$  is given by the usual column-wise sample  
 64 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  $\mu_s, \hat{\mu}_s$ , 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.

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

Suppose Tradict has estimates of  $\mu^{(m)}$  and  $\Sigma^{(m)}$  given by  $\hat{\mu}^{(m)}$  and  $\hat{\Sigma}^{(m)}$ , and let  $T_m = t_m(o \times \mathbf{1}_{1 \times \text{markers}})$  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 of sample sequencing depths in millions of reads. Given these, we would like to calculate the maximum *a posteriori* (MAP) estimate of  $\hat{z}_m = \text{argmax}_{z_m} p(z_m | o, T_m, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)})$ .

The posterior distribution over  $z_m$  is given by

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

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

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

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

respectively. Taken together, this suggests an iterative conditional modes algorithm [6] in which we maximize the posterior one column of  $z$  at a time, while conditioning on all others.

Let  $\hat{z}_m$  denote our current estimate of  $z_m$ . Let  $m(j)$  denote the index of the  $j^{\text{th}}$  marker and let  $m(-j)$

91 denote the indices of all markers but the  $j^{th}$  one. The above sub-objective is given by,

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

92 Differentiating we get,

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

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

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

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

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

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

## 100 3.2 Complete inference of $\mu^{(m)}$ , $\Sigma^{(m)}$ , and $z_m$

101 For complete inference we use the following iterative conditional modes algorithm [6]:

- 102 • Initialize  $T_m = t_m(o \times \mathbf{1}_{1 \times \text{markers}})$ ,  $\hat{z}_m = \text{lag}(t_m)$ .
- 103 • 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:
- 104     – Update  $\hat{\mu}^{(m)}$  and  $\hat{\Sigma}^{(m)}$ :

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

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

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

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

## 4 Prediction

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

### 4.1 MAP estimation of gene and program abundances

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

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

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

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

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

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

Similarly, for the entire transcriptome we have,

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

where,

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

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

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

137 The Moment Generating Function of a Normal random variable  $X$  with mean  $\mu$  and variance  $\sigma^2$  is given by  
 138  $M(t) = \mathbb{E}[\exp(tX)] = \exp(\mu t + \sigma^2 t^2/2)$ . Therefore we have,

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

139 where,

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

140 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  
 141 between  $z_m$  and  $z_j$ . These were learned from the training data during encoding.

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

## 144 4.2 Posterior density estimation of gene and program abundances

145 The above predictions represent point estimates. Ideally, we would like to know the uncertainty around these  
 146 estimates. Given measurements of the representative markers, we can estimate the posterior distribution of  
 147 expression values for transcriptional programs and the non-markers, and therein calculate any point estimates  
 148 and/or measures of uncertainty. Recall that for transcriptional programs:

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

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

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

150 Given  $z_m$ , the distribution of expression values are simple normal distributions with analytically available  
 151 means and covariances. However, because  $z_m$  is unknown, we must factor into our estimate its distribution,  
 152 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  
 153 strategy to estimate the posterior density of programs and non-markers will therefore be to sample from  
 154 the posterior of  $z_m$ , and then given these draws, sample from the conditional Normal distribution of each  
 155 program and non-marker gene.

### 156 4.2.1 Sampling $z_m$ via MCMC

157 To sample  $z_m$  we use Metropolis-Hastings Markov Chain Monte Carlo (MCMC) sampling [7], using the  
 158 following posterior density function:

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

159 Note that we do not require the marginal distribution for Metropolis-Hastings sampling.  
 160 As our proposal distribution we use:

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

161 Here  $z^{(i)}$  is the  $i^{th}$  draw from the sampler, and  $\gamma$  represents the width (variance) of the proposal distribution.  
 162 To choose this width, we examine, for a schedule of proposal widths (50 logarithmically spaced widths between  
 163  $10^{3.5}$  and  $10^{-1}$ ), which width gives an acceptance rate closest to 0.234 – the ideal rate for a high dimensional  
 164 parameter [7]. Using this width, we sample 20,100 times from the sampler. We burn-in the first 100 samples  
 165 and keep every  $100^{th}$  sample thereafter (to offset the effects of the chain’s auto-correlation) as our draws  
 166 from the distribution. Note that we initialize the chain at the MAP estimate of  $z_m$ . This ensures the chain  
 167 is stationary from the beginning.

#### 168 4.2.2 Sampling program and gene abundances

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

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

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

173 where

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

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

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

176 where

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

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

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

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

183 From the  $M$  samples we have from the conditional posterior distribution of each program and gene,  
 184 we can estimate properties of the posterior distribution. As point estimates for expression we can use the  
 185 posterior mean or mode. As confidence estimates for expression we can build credible intervals.

## 186 5 References

- 187 [1] C H Ho and Hong Kong. The multivariate Poisson-log normal distribution. 2, 1989.
- 188 [2] L. Madsen and D. Dalthorp. Simulating correlated count data. *Environmental and Ecological Statistics*,  
189 14(2):129–148, March 2007.
- 190 [3] Surojit Biswas, Meredith Mcdonald, Derek S Lundberg, Jeffery L Dangl, and Vladimir Jojic. Learning  
191 Microbial Interaction Networks from Metagenomic Count Data. In *Research in Computational Molecular*  
192 *Biology*, volume 1, pages 32–43, 2015.
- 193 [4] Hao Wu, Xinwei Deng, and Naren Ramakrishnan. Sparse Estimation of Multivariate Poisson Log-Normal  
194 Models from Count Data. *arXiv*, 2016.
- 195 [5] Surojit Biswas. The latent logarithm. *arXiv*, pages 1–11, 2016.
- 196 [6] Julian Besag. On the Statistical Analysis of Dirty Pictures. *Journal of the Royal Statistical Society*,  
197 48(3):259–302, 1986.
- 198 [7] Andrew Gelman, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin.  
199 *Bayesian Data Analysis*. Chapman & Hall, 3rd edition, 2013.