

# Bayesian Hybrid Matrix Factorisation for Data Integration

## Supplementary Materials

20th International Conference on Artificial Intelligence and Statistics (AISTATS 2017).

### Contents

<b>1</b>	<b>Models</b>	<b>2</b>
1.1	Matrix factorisation models . . . . .	2
1.2	Matrix factorisation with ARD and importance values . . . . .	8
1.3	Hybrid matrix factorisation model . . . . .	10
1.3.1	Model definition . . . . .	10
1.3.2	Gibbs sampler . . . . .	11
<b>2</b>	<b>Model discussion</b>	<b>17</b>
2.1	Software . . . . .	17
2.2	Complexity . . . . .	17
2.3	Missing values and predictions . . . . .	17
2.4	Initialisation . . . . .	17
2.5	Relation to tensor decomposition . . . . .	18
<b>3</b>	<b>Datasets and preprocessing</b>	<b>21</b>
3.1	Drug sensitivity datasets . . . . .	21
3.2	Preprocessing drug sensitivity values . . . . .	22
3.3	Features . . . . .	22
3.4	Methylation and gene expression datasets . . . . .	25
<b>4</b>	<b>Additional experiments</b>	<b>27</b>
4.1	Run-time comparison . . . . .	27
4.2	Model selection . . . . .	28
4.3	Initialisation . . . . .	28
4.4	Importance value . . . . .	31
4.5	Negativity constraints . . . . .	35
4.6	Factorisation types . . . . .	37

# 1 Models

In the model for Hybrid Matrix Factorisation (HMF), we combine models for Bayesian matrix factorisation, tri-factorisation, and tri-factorisation of similarity kernels. For each model, there are three versions: nonnegative, semi-nonnegative, and real-valued. In this section, we give the details for each of these nine models, as well as the Gibbs sampling algorithms.

In addition, for the semi-nonnegative and real-valued versions we can either use a univariate Gaussian (resulting in individual draws, but those can be drawn in parallel per column), or a multivariate Gaussian as the posterior (resulting in row-wise draws of new values, but each row can be drawn in parallel at the same time). We give Gibbs sampling algorithms for both options. For the nonnegative models, where the posterior is a truncated normal, this is technically also possible. However, we did not succeed in finding an efficiently implemented library for multivariate truncated normal draws, and therefore do not support it.

We first introduce each of the models separately in Section 1.1, extending them with Automatic Relevance Determination in Section 1.2, and finally explain how we combine them into one to form the HMF model in Section 1.3.

## 1.1 Matrix factorisation models

The matrix factorisation models are:

- Bayesian matrix factorisation (BMF), nonnegative matrix factorisation (BNMF), semi-nonnegative matrix factorisation (BSNMF). We decompose  $\mathbf{D} \approx \mathbf{F} \cdot \mathbf{G}^T$ . The Gibbs samplers are given in Table 2.
- Bayesian matrix tri-factorisation (BMTF), nonnegative matrix tri-factorisation (BNMTF), and semi-nonnegative matrix tri-factorisation (BSNMTF). We decompose  $\mathbf{R} \approx \mathbf{F} \cdot \mathbf{S} \cdot \mathbf{G}^T$ . The Gibbs samplers are given Table 3.
- Bayesian similarity matrix tri-factorisation (BSMTF), nonnegative similarity matrix tri-factorisation (BNSMTF), and semi-nonnegative similarity matrix tri-factorisation (BSNSMTF). We decompose  $\mathbf{C} \approx \mathbf{F} \cdot \mathbf{S} \cdot \mathbf{F}^T$ . The Gibbs sampler are given in Table 4.

The model priors and Gibbs sampling posteriors are given Table 1. Here,

- $\mathbf{R}, \mathbf{D}, \mathbf{C} \in \mathbb{R}^{I \times J}$ , with  $i = 1..I, j = 1..J$ .
- $\mathbf{F} \in \mathbb{R}^{I \times K}$ , with  $i = 1..I, k = 1..K$ .
- $\mathbf{S} \in \mathbb{R}^{K \times L}$  (matrix tri-factorisation), or  $\mathbf{S} \in \mathbb{R}^{K \times K}$  (similarity matrix tri-factorisation), with  $k = 1..K, l = 1..L$ .
- $\mathbf{G} \in \mathbb{R}^{J \times K}$  (matrix factorisation), or  $\mathbf{G} \in \mathbb{R}^{J \times L}$  (matrix tri-factorisation), with  $j = 1..J, k = 1..K, l = 1..L$ .
- $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is a multivariate Gaussian distribution with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ .

Table 1: Model definitions, and Gibbs sampling posteriors.

Model name	Likelihood	Priors	Posteriors	Table
BMF	$D_{ij} \sim \mathcal{N}(D_{ij}   \mathbf{F}_i \cdot \mathbf{G}_j, \tau^{-1})$	$F_{ik} \sim \mathcal{N}(F_{ik}   0, \lambda_F^{-1})$ $G_{jk} \sim \mathcal{N}(G_{jk}   0, \lambda_G^{-1})$ $\tau \sim \mathcal{G}(\tau   \alpha, \beta)$	$\mathcal{N}(F_{ik}   \mu_{ik}^F, (\tau_{ik}^F)^{-1})$ $\mathcal{N}(G_{jk}   \mu_{jk}^G, (\tau_{jk}^G)^{-1})$ $\mathcal{G}(\tau   \alpha^*, \beta^*)$	2
BMF (multivariate)			$\mathcal{N}(\mathbf{F}_i   \boldsymbol{\mu}_i^F, \boldsymbol{\Sigma}_i^F)$ $\mathcal{N}(\mathbf{G}_j   \boldsymbol{\mu}_j^G, \boldsymbol{\Sigma}_j^G)$	2
BNMF	$D_{ij} \sim \mathcal{N}(D_{ij}   \mathbf{F}_i \cdot \mathbf{G}_j, \tau^{-1})$	$F_{ik} \sim \mathcal{E}(F_{ik}   \lambda_F)$ $G_{jk} \sim \mathcal{E}(G_{jk}   \lambda_G)$ $\tau \sim \mathcal{G}(\tau   \alpha, \beta)$	$\mathcal{TN}(F_{ik}   \mu_{ik}^F, \tau_{ik}^F)$ $\mathcal{TN}(G_{jk}   \mu_{jk}^G, \tau_{jk}^G)$ $\mathcal{G}(\tau   \alpha^*, \beta^*)$	2
BSNMF	$D_{ij} \sim \mathcal{N}(D_{ij}   \mathbf{F}_i \cdot \mathbf{G}_j, \tau^{-1})$	$F_{ik} \sim \mathcal{E}(F_{ik}   \lambda_F)$ $G_{jk} \sim \mathcal{N}(G_{jk}   0, \lambda_G^{-1})$ $\tau \sim \mathcal{G}(\tau   \alpha, \beta)$	$\mathcal{TN}(F_{ik}   \mu_{ik}^F, \tau_{ik}^F)$ $\mathcal{N}(G_{jk}   \mu_{jk}^G, (\tau_{jk}^G)^{-1})$ $\mathcal{G}(\tau   \alpha^*, \beta^*)$	2
BSNMF (multivariate)			$\mathcal{N}(\mathbf{G}_j   \boldsymbol{\mu}_j^G, \boldsymbol{\Sigma}_j^G)$	2
BMTF	$R_{ij} \sim \mathcal{N}(R_{ij}   \mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{G}_j, \tau^{-1})$	$F_{ik} \sim \mathcal{N}(F_{ik}   0, \lambda_F^{-1})$ $S_{kl} \sim \mathcal{N}(S_{kl}   0, \lambda_S^{-1})$ $G_{jl} \sim \mathcal{N}(G_{jl}   0, \lambda_G^{-1})$ $\tau \sim \mathcal{G}(\tau   \alpha, \beta)$	$\mathcal{N}(F_{ik}   \mu_{ik}^F, (\tau_{ik}^F)^{-1})$ $\mathcal{N}(S_{kl}   \mu_{kl}^S, (\tau_{kl}^S)^{-1})$ $\mathcal{N}(G_{jl}   \mu_{jl}^G, (\tau_{jl}^G)^{-1})$ $\mathcal{G}(\tau   \alpha^*, \beta^*)$	3
BMTF (multivariate)			$\mathcal{N}(\mathbf{F}_i   \boldsymbol{\mu}_i^F, \boldsymbol{\Sigma}_i^F)$ $\mathcal{N}(\mathbf{S}_k   \boldsymbol{\mu}_k^S, \boldsymbol{\Sigma}_k^S)$ $\mathcal{N}(\mathbf{G}_j   \boldsymbol{\mu}_j^G, \boldsymbol{\Sigma}_j^G)$	3
BNMTF	$R_{ij} \sim \mathcal{N}(R_{ij}   \mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{G}_j, \tau^{-1})$	$F_{ik} \sim \mathcal{E}(F_{ik}   \lambda_F)$ $S_{kl} \sim \mathcal{E}(S_{kl}   \lambda_S)$ $G_{jl} \sim \mathcal{E}(G_{jl}   \lambda_G)$ $\tau \sim \mathcal{G}(\tau   \alpha, \beta)$	$\mathcal{TN}(F_{ik}   \mu_{ik}^F, \tau_{ik}^F)$ $\mathcal{TN}(S_{kl}   \mu_{kl}^S, \tau_{kl}^S)$ $\mathcal{TN}(G_{jl}   \mu_{jl}^G, \tau_{jl}^G)$ $\mathcal{G}(\tau   \alpha^*, \beta^*)$	3
BSNMTF	$R_{ij} \sim \mathcal{N}(R_{ij}   \mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{G}_j, \tau^{-1})$	$F_{ik} \sim \mathcal{E}(F_{ik}   \lambda_F)$ $S_{kl} \sim \mathcal{N}(S_{kl}   0, \lambda_S^{-1})$ $G_{jl} \sim \mathcal{E}(G_{jl}   \lambda_G)$ $\tau \sim \mathcal{G}(\tau   \alpha, \beta)$	$\mathcal{TN}(F_{ik}   \mu_{ik}^F, \tau_{ik}^F)$ $\mathcal{N}(S_{kl}   \mu_{kl}^S, (\tau_{kl}^S)^{-1})$ $\mathcal{TN}(G_{jl}   \mu_{jl}^G, \tau_{jl}^G)$ $\mathcal{G}(\tau   \alpha^*, \beta^*)$	3
BSNMTF (multivariate)			$\mathcal{N}(\mathbf{S}_k   \boldsymbol{\mu}_k^S, \boldsymbol{\Sigma}_k^S)$	3
BSMTF	$C_{ij} \sim \mathcal{N}(C_{ij}   \mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{F}_j, \tau^{-1})$	$F_{ik} \sim \mathcal{N}(F_{ik}   0, \lambda_F^{-1})$ $S_{kl} \sim \mathcal{N}(S_{kl}   0, \lambda_S^{-1})$ $\tau \sim \mathcal{G}(\tau   \alpha, \beta)$	$\mathcal{N}(F_{ik}   \mu_{ik}^F, (\tau_{ik}^F)^{-1})$ $\mathcal{N}(S_{kl}   \mu_{kl}^S, (\tau_{kl}^S)^{-1})$ $\mathcal{G}(\tau   \alpha^*, \beta^*)$	4
BSMTF (multivariate)			$\mathcal{N}(\mathbf{F}_i   \boldsymbol{\mu}_i^F, \boldsymbol{\Sigma}_i^F)$ $\mathcal{N}(\mathbf{S}_k   \boldsymbol{\mu}_k^S, \boldsymbol{\Sigma}_k^S)$	4
BNSMTF	$C_{ij} \sim \mathcal{N}(C_{ij}   \mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{F}_j, \tau^{-1})$	$F_{ik} \sim \mathcal{E}(F_{ik}   \lambda_F)$ $S_{kl} \sim \mathcal{E}(S_{kl}   \lambda_S)$ $\tau \sim \mathcal{G}(\tau   \alpha, \beta)$	$\mathcal{TN}(F_{ik}   \mu_{ik}^F, \tau_{ik}^F)$ $\mathcal{TN}(S_{kl}   \mu_{kl}^S, \tau_{kl}^S)$ $\mathcal{G}(\tau   \alpha^*, \beta^*)$	4
BSNSMTF	$C_{ij} \sim \mathcal{N}(C_{ij}   \mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{F}_j, \tau^{-1})$	$F_{ik} \sim \mathcal{E}(F_{ik}   \lambda_F)$ $S_{kl} \sim \mathcal{N}(S_{kl}   0, \lambda_S^{-1})$ $\tau \sim \mathcal{G}(\tau   \alpha, \beta)$	$\mathcal{TN}(F_{ik}   \mu_{ik}^F, \tau_{ik}^F)$ $\mathcal{N}(S_{kl}   \mu_{kl}^S, (\tau_{kl}^S)^{-1})$ $\mathcal{G}(\tau   \alpha^*, \beta^*)$	4
BSNSMTF (multivariate)			$\mathcal{N}(\mathbf{S}_k   \boldsymbol{\mu}_k^S, \boldsymbol{\Sigma}_k^S)$	4

The Gibbs sampling posteriors can be obtained using Bayes' theorem, for example for BN-MTF:

$$\begin{aligned}
& p(F_{ik}|\tau, \mathbf{F}_{-ik}, \mathbf{S}, \mathbf{G}, \mathbf{R}, h) \\
& \propto p(\mathbf{R}|\tau, \mathbf{F}, \mathbf{S}, \mathbf{G}) \times p(F_{ik}|\lambda_F) \\
& \propto \prod_{j \in \Omega_i^1} \mathcal{N}(R_{ij}|\mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{G}_j, \tau^{-1}) \times \mathcal{E}(F_{ik}|\lambda_F) \\
& \propto \exp \left\{ \frac{\tau}{2} \sum_{j \in \Omega_i^1} (R_{ij} - \mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{G}_j)^2 \right\} \times \exp \{-\lambda_F F_{ik}\} \times u(x) \\
& \propto \exp \left\{ \frac{F_{ik}^2}{2} \left[ \tau \sum_{j \in \Omega_i^1} (\mathbf{S}_k \cdot \mathbf{G}_j)^2 \right] \right. \\
& \quad \left. + F_{ik} \left[ -\lambda_F + \tau \sum_{j \in \Omega_i^1} (R_{ij} - \sum_{k' \neq k} F_{ik'} (\mathbf{S}_{k'} \cdot \mathbf{G}_j)) (\mathbf{S}_k \cdot \mathbf{G}_j) \right] \right\} \times u(x) \\
& \propto \exp \left\{ \frac{\tau_{ik}^F}{2} (F_{ik} - \mu_{ik}^F)^2 \right\} \times u(x) \\
& \propto \mathcal{TN}(F_{ik}|\mu_{ik}^F, \tau_{ik}^F)
\end{aligned}$$

where  $h = \{\lambda_F, \lambda_G, \lambda_S, \alpha, \beta\}$  are the hyperparameters to the model,  $u(x)$  is the unit step function, and  $\Omega_i^1 = \{j \mid (i, j) \in \Omega\}$ ,  $\Omega_j^2 = \{i \mid (i, j) \in \Omega\}$  indicate the observed entries per row and column, respectively.

All the parameter values can be found in Tables 2-4. We use  $\otimes$  to denote the outer product,  $\mathbf{I}$  for the identity matrix, and  $\mathbf{S}_{\cdot l}$  for the  $l$ th column of matrix  $\mathbf{S}$ .

For the similarity matrix factorisation we could have also decided to decompose  $\mathbf{C} = \mathbf{F}\mathbf{F}^T + \mathbf{E}$ , without the intermediate matrix  $\mathbf{S}$ . Ding et al. [2005] includes a good discussion of the benefits to our approach. For this decomposition we do not consider diagonal entries of  $\mathbf{C}$  (in other words,  $(i, i) \notin \Omega$ ,  $i = 1..I$ ) as this leads to third and fourth order terms in the posteriors and makes Gibbs sampling impossible. See Zhang and Yeung [2012] for a non-probabilistic approach that does consider these elements, leading to a very complicated optimisation problem.

Table 2: Gibbs Samplers for Bayesian Matrix Factorisation (BMF, BNMF, BSNMF).

PARAM	UPDATE (GAUSSIAN PRIOR)	UPDATE (EXPONENTIAL PRIOR)
$\alpha^*$	$\alpha + \frac{ \Omega }{2}$	$\alpha + \frac{ \Omega }{2}$
$\beta^*$	$\beta + \frac{1}{2} \sum_{(i,j) \in \Omega} (D_{ij} - \mathbf{F}_i \cdot \mathbf{G}_j)^2$	$\beta + \frac{1}{2} \sum_{(i,j) \in \Omega} (D_{ij} - \mathbf{F}_i \cdot \mathbf{G}_j)^2$
$\tau_{ik}^F$	$\lambda_F + \tau \sum_{j \in \Omega_i^1} G_{jk}^2$	$\tau \sum_{j \in \Omega_i^1} G_{jk}^2$
$\mu_{ik}^F$	$\frac{1}{\tau_{ik}^F} \left( \tau \sum_{j \in \Omega_i^1} (D_{ij} - \sum_{k' \neq k} F_{ik'} G_{jk'}) G_{jk} \right)$	$\frac{1}{\tau_{ik}^F} \left( -\lambda_F + \tau \sum_{j \in \Omega_i^1} (D_{ij} - \sum_{k' \neq k} F_{ik'} G_{jk'}) G_{jk} \right)$
$\tau_{jk}^G$	$\lambda_G + \tau \sum_{i \in \Omega_j^2} F_{ik}^2$	$\tau \sum_{i \in \Omega_j^2} F_{ik}^2$
$\mu_{jl}^G$	$\frac{1}{\tau_{jl}^G} \left( \tau \sum_{i \in \Omega_j^2} (D_{ij} - \sum_{k' \neq k} F_{ik'} G_{jk'}) F_{ik} \right)$	$\frac{1}{\tau_{jl}^G} \left( -\lambda_G + \tau \sum_{i \in \Omega_j^2} (D_{ij} - \sum_{k' \neq k} F_{ik'} G_{jk'}) F_{ik} \right)$
$\Sigma_i^F$	$\left( \lambda^F \mathbf{I} + \tau \sum_{j \in \Omega_i^1} \mathbf{G}_j \otimes \mathbf{G}_j \right)^{-1}$	-
$\mu_i^F$	$\Sigma_i^F \cdot \left( \tau \sum_{j \in \Omega_i^1} D_{ij} \mathbf{G}_j \right)$	-
$\Sigma_j^G$	$\left( \lambda^G \mathbf{I} + \tau \sum_{i \in \Omega_j^2} \mathbf{F}_i \otimes \mathbf{F}_i \right)^{-1}$	-
$\mu_j^G$	$\Sigma_j^G \cdot \left( \tau \sum_{i \in \Omega_j^2} D_{ij} \mathbf{F}_i \right)$	-

Table 3: Gibbs Samplers for Bayesian Matrix Tri-Factorisation (BMTF).

PARAM	UPDATE (GAUSSIAN PRIOR)	UPDATE (EXPONENTIAL PRIOR)
$\alpha^*$	$\alpha + \frac{ \Omega }{2}$	$\alpha + \frac{ \Omega }{2}$
$\beta^*$	$\beta + \frac{1}{2} \sum_{(i,j) \in \Omega} (R_{ij} - \mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{G}_j)^2$	$\beta + \frac{1}{2} \sum_{(i,j) \in \Omega} (R_{ij} - \mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{G}_j)^2$
$\tau_{ik}^F$	$\lambda_F + \tau \sum_{j \in \Omega_i^1} (\mathbf{S}_k \cdot \mathbf{G}_j)^2$	$\tau \sum_{j \in \Omega_i^1} (\mathbf{S}_k \cdot \mathbf{G}_j)^2$
$\mu_{ik}^F$	$\frac{1}{\tau F_{ik}} \left( \tau \sum_{j \in \Omega_i^1} (R_{ij} - \sum_{k' \neq k} F_{ik'} (\mathbf{S}_{k'} \cdot \mathbf{G}_j)) (\mathbf{S}_k \cdot \mathbf{G}_j) \right)$	$\frac{1}{\tau F_{ik}} \left( -\lambda_F + \tau \sum_{j \in \Omega_i^1} (R_{ij} - \sum_{k' \neq k} F_{ik'} (\mathbf{S}_{k'} \cdot \mathbf{G}_j)) (\mathbf{S}_k \cdot \mathbf{G}_j) \right)$
$\tau_{kl}^S$	$\lambda_S + \tau \sum_{(i,j) \in \Omega} F_{ik}^2 G_{jl}^2$	$\tau \sum_{(i,j) \in \Omega} F_{ik}^2 G_{jl}^2$
$\mu_{kl}^S$	$\frac{1}{\tau S_{kl}} \left( \tau \sum_{(i,j) \in \Omega} (R_{ij} - \sum_{(k',l') \neq (k,l)} F_{ik'} S_{k'l'} G_{jl'}) F_{ik} G_{jl} \right)$	$\frac{1}{\tau S_{kl}} \left( -\lambda_S + \tau \sum_{(i,j) \in \Omega} (R_{ij} - \sum_{(k',l') \neq (k,l)} F_{ik'} S_{k'l'} G_{jl'}) F_{ik} G_{jl} \right)$
$\tau_{jl}^G$	$\lambda_G + \tau \sum_{i \in \Omega_j^2} (\mathbf{F}_i \cdot \mathbf{S}_{\cdot,l})^2$	$\tau \sum_{i \in \Omega_j^2} (\mathbf{F}_i \cdot \mathbf{S}_{\cdot,l})^2$
$\mu_{jl}^G$	$\frac{1}{\tau G_{jl}} \left( \tau \sum_{i \in \Omega_j^2} (R_{ij} - \sum_{l' \neq l} G_{jl'} (\mathbf{F}_i \cdot \mathbf{S}_{\cdot,l'})) (\mathbf{F}_i \cdot \mathbf{S}_{\cdot,l}) \right)$	$\frac{1}{\tau G_{jl}} \left( -\lambda_G + \tau \sum_{i \in \Omega_j^2} (R_{ij} - \sum_{l' \neq l} G_{jl'} (\mathbf{F}_i \cdot \mathbf{S}_{\cdot,l'})) (\mathbf{F}_i \cdot \mathbf{S}_{\cdot,l}) \right)$
$\Sigma_i^F$	$\left( \lambda^F \mathbf{I} + \tau \sum_{j \in \Omega_i^1} (\mathbf{S} \cdot \mathbf{G}_j) \otimes (\mathbf{S} \cdot \mathbf{G}_j) \right)^{-1}$	-
$\mu_i^F$	$\Sigma_i^F \cdot \left( \tau \sum_{j \in \Omega_i^1} R_{ij} (\mathbf{S} \cdot \mathbf{G}_j) \right)$	-
$\Sigma_k^S$	$\left( \lambda^S \mathbf{I} + \tau \sum_{(i,j) \in \Omega} F_{ik} (\mathbf{G}_j \otimes \mathbf{G}_j) \right)^{-1}$	-
$\mu_k^S$	$\Sigma_k^S \cdot \left( \tau \sum_{(i,j) \in \Omega} (R_{ij} - \sum_{k' \neq k} F_{ik'} (\mathbf{S}_{k'} \cdot \mathbf{G}_j)) F_{ik} \mathbf{G}_j \right)$	-
$\Sigma_j^G$	$\left( \lambda^G \mathbf{I} + \tau \sum_{i \in \Omega_j^2} (\mathbf{F}_i \cdot \mathbf{S}) \otimes (\mathbf{F}_i \cdot \mathbf{S}) \right)^{-1}$	-
$\mu_j^G$	$\Sigma_j^G \cdot \left( \tau \sum_{i \in \Omega_j^2} R_{ij} (\mathbf{F}_i \cdot \mathbf{S}) \right)$	-

Table 4: Gibbs Samplers for Bayesian Similarity Matrix Tri-Factorisation (BSMTF).

PARAM	UPDATE (GAUSSIAN PRIOR)	UPDATE (EXPONENTIAL PRIOR)
$\alpha^*$	$\alpha + \frac{ \Omega }{2}$	$\alpha + \frac{ \Omega }{2}$
$\beta^*$	$\beta + \frac{1}{2} \sum_{(i,j) \in \Omega} (C_{ij} - \mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{F}_j)^2$	$\beta + \frac{1}{2} \sum_{(i,j) \in \Omega} (C_{ij} - \mathbf{F}_i \cdot \mathbf{S} \cdot \mathbf{F}_j)^2$
$\tau_{ik}^F$	$\lambda_F + \tau \left[ \sum_{j \in \Omega_i^+} (\mathbf{S}_k \cdot \mathbf{F}_j)^2 + \sum_{i' \in \Omega_i^2} (\mathbf{F}_{i'} \cdot \mathbf{S}_{\cdot,k})^2 \right]$	$\tau \left[ \sum_{j \in \Omega_i^+} (\mathbf{S}_k \cdot \mathbf{F}_j)^2 + \sum_{i' \in \Omega_i^2} (\mathbf{F}_{i'} \cdot \mathbf{S}_{\cdot,k})^2 \right]$
$\mu_{ik}^F$	$\frac{1}{\tau_{ik}^F} \left( \tau \sum_{j \in \Omega_i^+} (C_{ij} - \sum_{k' \neq k} F_{ik'}) (\mathbf{S}_{k'} \cdot \mathbf{F}_j) \right. \\ \left. + \tau \sum_{i' \in \Omega_i^2} (C_{i'i} - \sum_{l \neq k} F_{il}) (\mathbf{F}_{i'} \cdot \mathbf{S}_{\cdot,k}) \right)$	$\frac{1}{\tau_{ik}^F} \left( -\lambda_F + \tau \sum_{j \in \Omega_i^+} (C_{ij} - \sum_{k' \neq k} F_{ik'}) (\mathbf{S}_{k'} \cdot \mathbf{F}_j) \right. \\ \left. + \tau \sum_{i' \in \Omega_i^2} (C_{i'i} - \sum_{l \neq k} F_{il}) (\mathbf{F}_{i'} \cdot \mathbf{S}_{\cdot,k}) \right)$
$\tau_{kl}^S$	$\lambda_S + \tau \sum_{(i,j) \in \Omega} F_{ik}^2 F_{jl}^2$	$\tau \sum_{(i,j) \in \Omega} F_{ik}^2 F_{jl}^2$
$\mu_{kl}^S$	$\frac{1}{\tau_{kl}^S} \left( \tau \sum_{(i,j) \in \Omega} (C_{ij} - \sum_{(k',l') \neq (k,l)} F_{ik'} S_{k'l'} F_{j'l'}) F_{ik} F_{jl} \right)$	$\frac{1}{\tau_{kl}^S} \left( -\lambda_S + \tau \sum_{(i,j) \in \Omega} (C_{ij} - \sum_{(k',l') \neq (k,l)} F_{ik'} S_{k'l'} F_{j'l'}) F_{ik} F_{jl} \right)$
$\Sigma_i^F$	$\left( \lambda^F \mathbf{I} + \tau \left[ \sum_{j \in \Omega_i^+} (\mathbf{S} \cdot \mathbf{F}_j) \otimes (\mathbf{S} \cdot \mathbf{F}_j) + \sum_{i' \in \Omega_i^2} (\mathbf{F}_{i'} \cdot \mathbf{S}) \otimes (\mathbf{F}_{i'} \cdot \mathbf{S}) \right] \right)^{-1}$	-
$\mu_i^F$	$\Sigma_i^F \cdot \left( \tau \sum_{j \in \Omega_i^+} C_{ij} (\mathbf{S} \cdot \mathbf{F}_j) + \tau \sum_{i' \in \Omega_i^2} C_{i'i} (\mathbf{F}_{i'} \cdot \mathbf{S}) \right)$	-
$\Sigma_k^S$	$\left( \lambda^S \mathbf{I} + \tau \sum_{(i,j) \in \Omega} F_{ik} (\mathbf{F}_j \otimes \mathbf{F}_j) \right)^{-1}$	-
$\mu_k^S$	$\Sigma_k^S \cdot \left( \tau \sum_{(i,j) \in \Omega} (C_{ij} - \sum_{k' \neq k} F_{ik'}) (\mathbf{S}_{k'} \cdot \mathbf{F}_j) \right) F_{ik} \mathbf{F}_j$	-

## 1.2 Matrix factorisation with ARD and importance values

We will now explain how to extend matrix factorisation with Automatic Relevance Determination (ARD) and importance values.

**ARD** We change the model definition to the following.

$$\begin{aligned} D_{ij} &\sim \mathcal{N}(D_{ij} | \mathbf{F}_i \cdot \mathbf{G}_j, \tau^{-1}) \\ F_{ik} &\sim \mathcal{N}(F_{ik} | 0, (\lambda_k)^{-1}) \quad \text{or} \quad \mathcal{E}(F_{ik} | \lambda_k) \\ G_{jk} &\sim \mathcal{N}(G_{jk} | 0, (\lambda_k)^{-1}) \quad \text{or} \quad \mathcal{E}(G_{jk} | \lambda_k) \\ \tau &\sim \mathcal{G}(\tau | \alpha, \beta) \end{aligned}$$

Notice that the main difference is replacing  $\lambda_F$  and  $\lambda_G$  by the parameter  $\lambda_k$ , for  $k = 1..K$ . We place a Gamma prior over these variables,

$$\lambda_k \sim \mathcal{G}(\lambda_k | \alpha_0, \beta_0).$$

This can similarly be done for matrix tri-factorisation, by placing one ARD over  $\mathbf{F}$  (using  $\lambda_k^F$ ) and another over  $\mathbf{G}$  (using  $\lambda_k^G$ ).

The Gibbs sampling algorithms remain largely the same, simply replacing  $\lambda_F, \lambda_G$  in the updates by  $\lambda_k$ . For the multivariate posteriors, we replace  $\lambda^F \mathbf{I}$  by  $\text{diag}(\boldsymbol{\lambda})$ , a diagonal matrix where the  $k$ th diagonal element is given by  $\lambda_k$ . The posterior Gibbs sampling distribution for  $\lambda_k$  itself can be derived to be another Gamma distribution,

$$p(\lambda_k | \mathbf{D}, \mathbf{F}, \mathbf{G}, \tau) = \mathcal{G}(\lambda_k | \alpha_0^*, \beta_0^*)$$

where

$$\alpha_0^* = \alpha_0 + \frac{I}{2} + \frac{J}{2} \quad \beta_0^* = \beta_0 + \frac{1}{2} \sum_{i=1}^I F_{ik}^2 + \frac{1}{2} \sum_{j=1}^J G_{jk}^2$$

for the real-valued (Gaussian prior) model, and

$$\alpha_0^* = \alpha_0 + I + J \quad \beta_0^* = \beta_0 + \sum_{i=1}^I F_{ik} + \sum_{j=1}^J G_{jk}$$

for the nonnegative (exponential prior) model. Note that each  $\lambda_k$  therefore depends only on the values in the  $k$ th column of  $\mathbf{F}$  and  $\mathbf{G}$ , and those values in turn depend on the value of  $\lambda_k$ : if most values in that column are high,  $\lambda_k$  gets a small value (indicating that the factor is active); and if  $\lambda_k$  has a high value, the  $k$ th column will get a Gibbs sampling posterior around 0, resulting in low values (pushing the other values for this factor down, since it is inactive).



**Importance value** As discussed in the paper, we modify the likelihood with an importances values  $\alpha^n, \alpha^l, \alpha^m$  as follows.

$$\begin{aligned}
p(\boldsymbol{\theta}|\mathbf{R}, \mathbf{D}, \mathbf{C}) &\propto \prod_{t=1}^T p(\mathbf{F}^t|\boldsymbol{\lambda}^t) \prod_{n=1}^N p(\tau^n) \prod_{l=1}^L p(\tau^l) \prod_{m=1}^M p(\tau^m) \\
&\times \prod_{n=1}^N p(\mathbf{R}^n|\mathbf{F}^{t_n}, \mathbf{S}^n, \mathbf{F}^{u_n}, \tau^n)^{\alpha^n} \\
&\times \prod_{l=1}^L p(\mathbf{D}^l|\mathbf{F}^{t_l}, \mathbf{G}^l, \tau^l)^{\alpha^l} \\
&\times \prod_{m=1}^M p(\mathbf{C}^m|\mathbf{F}^{t_m}, \mathbf{S}^m, \tau^m)^{\alpha^m}
\end{aligned}$$

where  $\boldsymbol{\theta}$  is the set of model parameters  $(\mathbf{F}^t, \mathbf{S}^n, \mathbf{S}^m, \mathbf{G}^l, \boldsymbol{\lambda}^t, \tau^n, \tau^l, \tau^m)$ . We effectively repeat the dataset  $\mathbf{R}^n$   $\alpha^n$  times, and similarly for  $\mathbf{D}^l, \mathbf{C}^m$ , requiring a better fit to the data. The Gibbs samplers remain largely the same, with the addition of several  $\alpha$  values in the updates. This is illustrated below for a single dataset  $\mathbf{D} \approx \mathbf{F} \cdot \mathbf{G}^T$  with nonnegative factors, importance value  $\alpha$ , and without ARD.

$$\begin{aligned}
\tau &\sim \mathcal{G}(\tau|\alpha_*, \beta_*) & \alpha_* &= \alpha_\tau + \alpha \frac{|\Omega|}{2} \\
& & \beta_* &= \beta_\tau + \alpha \frac{1}{2} \sum_{(i,j) \in \Omega^t} (D_{ij} - \mathbf{F}_i \cdot \mathbf{G}_j)^2 \\
F_{ik} &\sim \mathcal{TN}(F_{ik}|\mu_{ik}, \tau_{ik}) & \mu_{ik} &= \frac{1}{\tau_{ik}} \left( -\lambda_F + \alpha\tau \sum_{j \in \Omega_i^1} (D_{ij} - \sum_{k' \neq k} F_{ik'} G_{jk'}) G_{jk} \right) \\
& & \tau_{ik} &= \alpha\tau \sum_{j \in \Omega_i^1} G_{jk}^2 \\
G_{jk} &\sim \mathcal{TN}(G_{jk}|\mu_{jk}, \tau_{jk}) & \mu_{jk} &= \frac{1}{\tau_{jk}} \left( -\lambda_G + \alpha\tau \sum_{i \in \Omega_j^2} (D_{ij} - \sum_{k' \neq k} F_{ik'} G_{jk'}) F_{ik} \right) \\
& & \tau_{jk} &= \alpha\tau \sum_{i \in \Omega_j^2} (F_{ik})^2
\end{aligned}$$

Similar derivations can be done for matrix tri-factorisation, and real-valued versions.

### 1.3 Hybrid matrix factorisation model

The hybrid matrix factorisation (HMF) model combines all the ideas presented in the previous two sections. Recall we are given three types of datasets:

1. Main datasets  $\mathbf{R} = \{\mathbf{R}^1, \dots, \mathbf{R}^N\}$ , relating two different entity types. Each dataset  $\mathbf{R}^n \in \mathbb{R}^{I_{t_n} \times I_{u_n}}$  relates entity types  $E_{t_n}, E_{u_n}$ . We use matrix tri-factorisation to decompose it into two entity type factor matrices  $\mathbf{F}^{t_n}, \mathbf{F}^{u_n}$ , and a dataset-specific matrix  $\mathbf{S}^n \in \mathbb{R}^{K_{t_n} \times K_{u_n}}$ .

$$\mathbf{R}^n = \mathbf{F}^{t_n} \mathbf{S}^n (\mathbf{F}^{u_n})^T + \mathbf{E}^n.$$

2. Feature datasets  $\mathbf{D} = \{\mathbf{D}^1, \dots, \mathbf{D}^L\}$ , representing features for an entity type. Each dataset  $\mathbf{D}^l \in \mathbb{R}^{I_{t_l} \times J_l}$  relates an entity type  $E_{t_l}$  to  $J_l$  features. We use matrix factorisation to decompose it into one entity type factor matrix  $\mathbf{F}^{t_l}$ , and a dataset-specific matrix  $\mathbf{G}^l \in \mathbb{R}^{J_l \times K_{t_l}}$ .

$$\mathbf{D}^l = \mathbf{F}^{t_l} (\mathbf{G}^l)^T + \mathbf{E}^l.$$

3. Similarity datasets  $\mathbf{C} = \{\mathbf{C}^1, \dots, \mathbf{C}^M\}$ , giving similarities between entities of the same entity type. Each dataset  $\mathbf{C}^m \in \mathbb{R}^{I_{t_m} \times I_{t_m}}$  relates an entity type  $E_{t_m}$  to itself. We use matrix tri-factorisation to decompose it into a entity type factor matrix  $\mathbf{F}^{t_m}$ , a dataset-specific matrix  $\mathbf{S}^m \in \mathbb{R}^{K_{t_m} \times K_{t_m}}$ , and  $\mathbf{F}^{t_m}$  again.

$$\mathbf{C}^m = \mathbf{F}^{t_m} \mathbf{S}^m (\mathbf{F}^{t_m})^T + \mathbf{E}^m.$$

Observed entries are given by the sets  $\Omega^n = \{(i, j) \mid R_{ij}^n \text{ observed}\}$ ,  $\Omega^l = \{(i, j) \mid D_{ij}^l \text{ observed}\}$ ,  $\Omega^m = \{(i, j) \mid C_{ij}^m \text{ observed}\}$ , respectively.

#### 1.3.1 Model definition

The model likelihood functions are

$$\begin{aligned} R_{ij}^n &\sim \mathcal{N}(R_{ij}^n \mid \mathbf{F}_i^{t_n} \cdot \mathbf{S}^n \cdot \mathbf{F}_j^{s_n}, (\tau^n)^{-1}) \\ D_{ij}^l &\sim \mathcal{N}(D_{ij}^l \mid \mathbf{F}_i^{t_l} \cdot \mathbf{G}_j^l, (\tau^l)^{-1}) \\ C_{ij}^m &\sim \mathcal{N}(C_{ij}^m \mid \mathbf{F}_i^{t_m} \cdot \mathbf{S}^m \cdot \mathbf{F}_j^{t_m}, (\tau^m)^{-1}), \end{aligned}$$

with Bayesian priors

$$\begin{aligned} \tau^n, \tau^l, \tau^m &\sim \mathcal{G}(\tau^* \mid \alpha_\tau, \beta_\tau) \\ F_{ik}^t &\sim \mathcal{E}(F_{ik}^t \mid \lambda_k^t) & \text{or} & \quad F_{ik}^t \sim \mathcal{N}(F_{ik}^t \mid 0, (\lambda_k^t)^{-1}) \\ G_{jk}^l &\sim \mathcal{E}(G_{jk}^l \mid \lambda_k^{t_l}) & \text{or} & \quad G_{jk}^l \sim \mathcal{N}(G_{jk}^l \mid 0, (\lambda_k^{t_l})^{-1}) \\ S_{kl}^n &\sim \mathcal{E}(S_{kl}^n \mid \lambda_S^n) & \text{or} & \quad S_{kl}^n \sim \mathcal{N}(S_{kl}^n \mid 0, (\lambda_S^n)^{-1}) \\ S_{kl}^m &\sim \mathcal{E}(S_{kl}^m \mid \lambda_S^m) & \text{or} & \quad S_{kl}^m \sim \mathcal{N}(S_{kl}^m \mid 0, (\lambda_S^m)^{-1}). \\ \lambda_k^t &\sim \mathcal{G}(\lambda_k^t \mid \alpha_0, \beta_0). \end{aligned}$$

Finally, we add an importance value for each of the  $\mathbf{R}^n, \mathbf{D}^l, \mathbf{C}^m$  datasets, respectively  $\alpha_n, \alpha_l, \alpha_m$ .

### 1.3.2 Gibbs sampler

The Gibbs sampling algorithm has updates that combine the parameter values given in Tables 2, 3, and 4, for the single-dataset matrix factorisations, matrix tri-factorisations, and similarity matrix tri-factorisations. Because there are so many different parts of the models involved, careful notational definition is essential.

The datasets relating a given entity type  $E_t$  are indicated by the following sets,

$$\begin{aligned} U_1^t &= \{n \mid \mathbf{R}^n \in \mathbf{R} \wedge t_n = t\} \\ U_2^t &= \{n \mid \mathbf{R}^n \in \mathbf{R} \wedge u_n = t\} \\ V^t &= \{l \mid \mathbf{D}^l \in \mathbf{D} \wedge t_l = t\} \\ W^t &= \{m \mid \mathbf{C}^m \in \mathbf{C} \wedge t_m = t\}. \end{aligned}$$

Since the updates for the ARD can be different if a feature dataset is decomposed using negative factors or real-valued factors for  $\mathbf{G}^l$ , we also introduce the sets

$$V_+^t = \{l \in V^t \mid \mathbf{G}^l \text{ is nonnegative}\} \quad V_-^t = \{l \in V^t \mid \mathbf{G}^l \text{ is real-valued}\}.$$

Observed entries per row  $i$  and column  $j$  are given by

$$\begin{aligned} \Omega_i^{n1} &= \{j \mid (i, j) \in \Omega^n\} & \Omega_j^{n2} &= \{i \mid (i, j) \in \Omega^n\} \\ \Omega_i^{l1} &= \{j \mid (i, j) \in \Omega^l\} & \Omega_j^{l2} &= \{i \mid (i, j) \in \Omega^l\} \\ \Omega_i^{m1} &= \{j \mid (i, j) \in \Omega^m\} & \Omega_j^{m2} &= \{i \mid (i, j) \in \Omega^m\}. \end{aligned}$$

We obtain the following posterior distributions and parameter values:

#### Noise parameters

$$\begin{aligned} \tau^n &\sim \mathcal{G}(\tau^n \mid \alpha_*^n, \beta_*^n) & \alpha_*^n &= \alpha_\tau + \alpha^n \frac{|\Omega^n|}{2} & \beta_*^n &= \beta_\tau + \alpha^n \frac{1}{2} \sum_{(i,j) \in \Omega^n} (R_{ij}^n - \mathbf{F}_i^{t_n} \cdot \mathbf{S}^n \cdot \mathbf{F}_j^{u_n})^2 \\ \tau^l &\sim \mathcal{G}(\tau^l \mid \alpha_*^l, \beta_*^l) & \alpha_*^l &= \alpha_\tau + \alpha^l \frac{|\Omega^l|}{2} & \beta_*^l &= \beta_\tau + \alpha^l \frac{1}{2} \sum_{(i,j) \in \Omega^l} (D_{ij}^l - \mathbf{F}_i^{t_l} \cdot \mathbf{G}_j^l)^2 \\ \tau^m &\sim \mathcal{G}(\tau^m \mid \alpha_*^m, \beta_*^m) & \alpha_*^m &= \alpha_\tau + \alpha^m \frac{|\Omega^m|}{2} & \beta_*^m &= \beta_\tau + \alpha^m \frac{1}{2} \sum_{(i,j) \in \Omega^m} (C_{ij}^m - \mathbf{F}_i^{t_m} \cdot \mathbf{S}^m \cdot \mathbf{F}_j^{t_m})^2 \end{aligned}$$

**ARD** If  $\mathbf{F}^t$  contains nonnegative factors:

$$\begin{aligned} \lambda_k^t &\sim \mathcal{G}(\lambda_k^t \mid \alpha_k^t, \beta_k^t) & \alpha_k^t &= \alpha_0 + I_t + \sum_{l \in V_+^t} I_{t_l} + \sum_{l \in V_-^t} \frac{I_{t_l}}{2} \\ & & \beta_k^t &= \beta_0 + \sum_{i=1}^{I_t} F_{ik} + \sum_{l \in V_+^t} \sum_{j=1}^{J_l} G_{jk} + \sum_{l \in V_-^t} \frac{1}{2} \sum_{j=1}^{J_l} G_{jk}^2. \end{aligned}$$

If  $\mathbf{F}^t$  contains real-valued factors:

$$\lambda_k^t \sim \mathcal{G}(\lambda_k^t | \alpha_k^t, \beta_k^t) \quad \alpha_k^t = \alpha_0 + \frac{I_t}{2} + \sum_{l \in V_+^t} I_{t_l} + \sum_{l \in V_-^t} \frac{I_{t_l}}{2}$$

$$\beta_k^t = \beta_0 + \frac{1}{2} \sum_{i=1}^{I_t} F_{ik}^2 + \sum_{l \in V_+^t} \sum_{j=1}^{J_l} G_{jk} + \sum_{l \in V_-^t} \frac{1}{2} \sum_{j=1}^{J_l} G_{jk}^2.$$

**Dataset-specific factor matrices** If  $\mathbf{G}^l, \mathbf{S}^n, \mathbf{S}^m$  contain nonnegative factors:

$$G_{jk}^l \sim \mathcal{TN}(G_{jk}^l | \mu_{jk}^l, \tau_{jk}^l) \quad \mu_{jk}^l = \frac{1}{\tau_{jk}^l} \left( -\lambda_k^{t_l} + \tau^l \alpha^l \sum_{i \in \Omega_j^{l^2}} (D_{ij}^l - \sum_{k' \neq k} F_{ik'}^{t_l} G_{jk'}^l) F_{ik}^{t_l} \right)$$

$$\tau_{jk}^l = \tau^l \alpha^l \sum_{i \in \Omega_j^{l^2}} (F_{ik}^{t_l})^2$$

$$S_{kl}^n \sim \mathcal{TN}(S_{kl}^n | \mu_{kl}^n, \tau_{kl}^n) \quad \mu_{kl}^n = \frac{1}{\tau_{kl}^n} \left( -\lambda_S^n + \tau^n \alpha^n \sum_{(i,j) \in \Omega^n} (R_{ij}^n - \sum_{(k',l') \neq (k,l)} F_{ik'}^{t_n} S_{k'l'}^n F_{jl'}^{u_n}) F_{ik}^{t_n} F_{jl}^{u_n} \right)$$

$$\tau_{kl}^n = \tau^n \alpha^n \sum_{(i,j) \in \Omega^n} (F_{ik}^{t_n})^2 (F_{jl}^{u_n})^2$$

$$S_{kl}^m \sim \mathcal{TN}(S_{kl}^m | \mu_{kl}^m, \tau_{kl}^m) \quad \mu_{kl}^m = \frac{1}{\tau_{kl}^m} \left( -\lambda_S^m + \tau^m \alpha^m \sum_{(i,j) \in \Omega^m} (C_{ij}^m - \sum_{(k',l') \neq (k,l)} F_{ik'}^{t_m} S_{k'l'}^m F_{jl'}^{t_m}) F_{ik}^{t_m} F_{jl}^{t_m} \right)$$

$$\tau_{kl}^m = \tau^m \alpha^m \sum_{(i,j) \in \Omega^m} (F_{ik}^{t_m})^2 (F_{jl}^{t_m})^2$$

If  $\mathbf{G}^l, \mathbf{S}^n, \mathbf{S}^m$  contain real-valued factors:

$$G_{jk}^l \sim \mathcal{N}(G_{jk}^l | \mu_{jk}^l, (\tau_{jk}^l)^{-1}) \quad \mu_{jk}^l = \frac{1}{\tau_{jk}^l} \left( \tau^l \alpha^l \sum_{i \in \Omega_j^{l2}} (D_{ij}^l - \sum_{k' \neq k} F_{ik'}^{t_l} G_{jk'}^l) F_{ik}^{t_l} \right)$$

$$\tau_{jk}^l = \lambda_k^{t_l} + \tau^l \alpha^l \sum_{i \in \Omega_j^{l2}} (F_{ik}^{t_l})^2$$

$$\mathbf{G}_j^l \sim \mathcal{N}(\mathbf{G}_j^l | \boldsymbol{\mu}_j^l, \boldsymbol{\Sigma}_j^l) \quad \boldsymbol{\mu}_j^l = \boldsymbol{\Sigma}_j^l \cdot \left( \tau^l \alpha^l \sum_{i \in \Omega_j^{l2}} D_{ij}^l \mathbf{F}_i^{t_l} \right)$$

$$\boldsymbol{\Sigma}_j^l = \left( \text{diag}(\boldsymbol{\lambda}^t) + \tau^l \alpha^l \sum_{i \in \Omega_j^{l2}} \mathbf{F}_i^{t_l} \otimes \mathbf{F}_i^{t_l} \right)^{-1}$$

$$S_{kl}^n \sim \mathcal{N}(S_{kl}^n | \mu_{kl}^n, (\tau_{kl}^n)^{-1}) \quad \mu_{kl}^n = \frac{1}{\tau_{kl}^n} \left( \tau^n \alpha^n \sum_{(i,j) \in \Omega^n} (R_{ij}^n - \sum_{(k',l') \neq (k,l)} F_{ik'}^{t_n} S_{k'l'}^n F_{jl'}^{u_n}) F_{ik}^{t_n} F_{jl}^{u_n} \right)$$

$$\tau_{jk}^n = \lambda_S^n + \tau^n \alpha^n \sum_{(i,j) \in \Omega^n} (F_{ik}^{t_n})^2 (F_{jl}^{u_n})^2$$

$$\mathbf{S}_k^n \sim \mathcal{N}(\mathbf{S}_k^n | \boldsymbol{\mu}_k^n, \boldsymbol{\Sigma}_k^n) \quad \boldsymbol{\mu}_k^n = \boldsymbol{\Sigma}_k^n \cdot \left( \tau^n \alpha^n \sum_{(i,j) \in \Omega^n} (R_{ij}^n - \sum_{k' \neq k} F_{ik'}^{t_n} (\mathbf{S}_{k'}^n \cdot \mathbf{F}_j^{u_n})) F_{ik}^{t_n} \mathbf{F}_j^{u_n} \right)$$

$$\boldsymbol{\Sigma}_k^n = \left( \lambda_S^n \mathbf{I} + \tau^n \alpha^n \sum_{(i,j) \in \Omega^n} F_{ik}^{t_n} (\mathbf{F}_j^{u_n} \otimes \mathbf{F}_j^{u_n}) \right)^{-1}$$

$$S_{kl}^m \sim \mathcal{N}(S_{kl}^m | \mu_{kl}^m, (\tau_{kl}^m)^{-1}) \quad \mu_{kl}^m = \frac{1}{\tau_{kl}^m} \left( \tau^m \alpha^m \sum_{(i,j) \in \Omega^m} (C_{ij}^m - \sum_{(k',l') \neq (k,l)} F_{ik'}^{t_m} S_{k'l'}^m F_{jl'}^{t_m}) F_{ik}^{t_m} F_{jl}^{t_m} \right)$$

$$\tau_{kl}^m = \lambda_S^m + \tau^m \alpha^m \sum_{(i,j) \in \Omega^m} (F_{ik}^{t_m})^2 (F_{jl}^{t_m})^2$$

$$\mathbf{S}_k^m \sim \mathcal{N}(\mathbf{S}_k^m | \boldsymbol{\mu}_k^m, \boldsymbol{\Sigma}_k^m) \quad \boldsymbol{\mu}_k^m = \boldsymbol{\Sigma}_k^m \cdot \left( \tau^m \alpha^m \sum_{(i,j) \in \Omega^m} (C_{ij}^m - \sum_{k' \neq k} F_{ik'}^{t_m} (\mathbf{S}_{k'}^m \cdot \mathbf{F}_j^{t_m})) F_{ik}^{t_m} \mathbf{F}_j^{t_m} \right)$$

$$\boldsymbol{\Sigma}_k^m = \left( \lambda_S^m \mathbf{I} + \tau^m \alpha^m \sum_{(i,j) \in \Omega^m} F_{ik}^{t_m} (\mathbf{F}_j^{t_m} \otimes \mathbf{F}_j^{t_m}) \right)^{-1}$$

Shared factor matrices If  $\mathbf{F}^t$  contains nonnegative factors:

$$\begin{aligned}
F_{ik}^t &\sim \mathcal{TN}(F_{ik}^t | \mu_{ik}^t, \tau_{ik}^t) & \mu_{ik}^t &= \frac{1}{\tau_{ik}^t} \left( -\lambda_k^t + \sum_{n \in U_1^t} \tau^n \alpha^n \sum_{j \in \Omega_i^{n1}} (R_{ij}^n - \sum_{k' \neq k} F_{ik'}^t (\mathbf{S}_{k'}^n \cdot \mathbf{F}_j^{u_n})) (\mathbf{S}_k^n \cdot \mathbf{F}_j^{u_n}) \right. \\
& & & + \sum_{n \in U_2^t} \tau^n \alpha^n \sum_{i' \in \Omega_i^{n2}} (R_{i'i}^n - \sum_{l \neq k} F_{il}^t (\mathbf{F}_{i'}^{t_n} \cdot \mathbf{S}_{.,l}^n)) (\mathbf{F}_{i'}^{t_n} \cdot \mathbf{S}_{.,k}^n) \\
& & & + \sum_{l \in V^t} \tau^l \alpha^l \sum_{j \in \Omega_i^{l1}} (D_{ij}^l - \sum_{k' \neq k} F_{ik'}^l G_{jk'}^l) G_{jk}^l \left. \right) \\
& & & + \sum_{m \in W^t} \tau^m \alpha^m \left[ \sum_{j \in \Omega_i^{m1}} (C_{ij}^m - \sum_{k' \neq k} F_{ik'}^m (\mathbf{S}_{k'}^m \cdot \mathbf{F}_j^t)) (\mathbf{S}_k^m \cdot \mathbf{F}_j^t) \right. \\
& & & \left. + \sum_{i' \in \Omega_i^{m2}} (C_{i'i}^m - \sum_{l \neq k} F_{il}^m (\mathbf{F}_{i'}^t \cdot \mathbf{S}_{.,l}^m)) (\mathbf{F}_{i'}^t \cdot \mathbf{S}_{.,k}^m) \right] \\
\tau_{ik}^t &= \sum_{n \in U_1^t} \tau^n \alpha^n \sum_{j \in \Omega_i^{n1}} (\mathbf{S}_k^n \cdot \mathbf{F}_j^{u_n})^2 + \sum_{n \in U_2^t} \tau^n \alpha^n \sum_{i' \in \Omega_i^{n2}} (\mathbf{F}_{i'}^{t_n} \cdot \mathbf{S}_{.,k}^n)^2 \\
& + \sum_{l \in V^t} \tau^l \alpha^l \sum_{j \in \Omega_i^{l1}} (G_{jk}^l)^2 \\
& + \sum_{m \in W^t} \tau^m \alpha^m \left[ \sum_{j \in \Omega_i^{m1}} (\mathbf{S}_k^m \cdot \mathbf{F}_j^t)^2 + \sum_{i' \in \Omega_i^{m2}} (\mathbf{F}_{i'}^t \cdot \mathbf{S}_{.,k}^m)^2 \right]
\end{aligned}$$

If  $\mathbf{F}^t$  contains real-valued factors:

$$\begin{aligned}
F_{ik}^t \sim \mathcal{N}(F_{ik}^t | \mu_{ik}^t, (\tau_{ik}^t)^{-1}) \quad \mu_{ik}^t = & \frac{1}{\tau_{ik}^t} \left( \sum_{n \in U_1^t} \tau^n \alpha^n \sum_{j \in \Omega_i^{n1}} (R_{ij}^n - \sum_{k' \neq k} F_{ik'}^t (\mathbf{S}_{k'}^n \cdot \mathbf{F}_j^{u_n})) (\mathbf{S}_k^n \cdot \mathbf{F}_j^{u_n}) \right. \\
& + \sum_{n \in U_2^t} \tau^n \alpha^n \sum_{i' \in \Omega_i^{n2}} (R_{i'i}^n - \sum_{l \neq k} F_{i'l}^t (\mathbf{F}_{i'}^{t_n} \cdot \mathbf{S}_{.,l}^n)) (\mathbf{F}_{i'}^{t_n} \cdot \mathbf{S}_{.,k}^n) \\
& + \sum_{l \in V^t} \tau^l \alpha^l \sum_{j \in \Omega_i^{l1}} (D_{ij}^l - \sum_{k' \neq k} F_{ik'}^l G_{jk'}^l) G_{jk}^l \\
& \left. + \sum_{m \in W^t} \tau^m \alpha^m \left[ \sum_{j \in \Omega_i^{m1}} (C_{ij}^m - \sum_{k' \neq k} F_{ik'}^m (\mathbf{S}_{k'}^m \cdot \mathbf{F}_j^t)) (\mathbf{S}_k^m \cdot \mathbf{F}_j^t) \right. \right. \\
& \left. \left. + \sum_{i' \in \Omega_i^{m2}} (C_{i'i}^m - \sum_{l \neq k} F_{i'l}^m (\mathbf{F}_{i'}^t \cdot \mathbf{S}_{.,l}^m)) (\mathbf{F}_{i'}^t \cdot \mathbf{S}_{.,k}^m) \right] \right)
\end{aligned}$$

$$\begin{aligned}
\tau_{ik}^t = & \lambda_k^t + \sum_{n \in U_1^t} \tau^n \alpha^n \sum_{j \in \Omega_i^{n1}} (\mathbf{S}_k^n \cdot \mathbf{F}_j^{u_n})^2 + \sum_{n \in U_2^t} \tau^n \alpha^n \sum_{i' \in \Omega_i^{n2}} (\mathbf{F}_{i'}^{t_n} \cdot \mathbf{S}_{.,k}^n)^2 \\
& + \sum_{l \in V^t} \tau^l \alpha^l \sum_{j \in \Omega_i^{l1}} (G_{jk}^l)^2 \\
& + \sum_{m \in W^t} \tau^m \alpha^m \left[ \sum_{j \in \Omega_i^{m1}} (\mathbf{S}_k^m \cdot \mathbf{F}_j^t)^2 + \sum_{i' \in \Omega_i^{m2}} (\mathbf{F}_{i'}^t \cdot \mathbf{S}_{.,k}^m)^2 \right]
\end{aligned}$$

$$\begin{aligned}
\mathbf{F}_i^t \sim \mathcal{N}(\mathbf{F}_i^t | \boldsymbol{\mu}_i^t, \boldsymbol{\Sigma}_i^t) \quad & \boldsymbol{\mu}_i^t = \boldsymbol{\Sigma}_i^t \cdot \left( \sum_{n \in U_1^t} \tau^n \alpha^n \sum_{j \in \Omega_i^{n1}} R_{ij}^n (\mathbf{S}^n \cdot \mathbf{F}_j^{u_n}) + \sum_{n \in U_2^t} \tau^n \alpha^n \sum_{i' \in \Omega_i^{n2}} R_{i'i}^n (\mathbf{F}_{i'}^{t_n} \cdot \mathbf{S}^n) \right. \\
& + \sum_{l \in V^t} \tau^l \alpha^l \sum_{j \in \Omega_i^{l1}} D_{ij}^l \mathbf{G}_j^l \\
& \left. + \sum_{m \in W^t} \tau^m \alpha^m \left[ \sum_{j \in \Omega_i^{m1}} C_{ij}^m (\mathbf{S}^m \cdot \mathbf{F}_j^t) + \sum_{i' \in \Omega_i^{m2}} C_{i'i}^m (\mathbf{F}_{i'}^t \cdot \mathbf{S}^m) \right] \right) \\
\boldsymbol{\Sigma}_i^t = & \left( \text{diag}(\boldsymbol{\lambda}_k^t) + \sum_{n \in U_1^t} \tau^n \alpha^n \sum_{j \in \Omega_i^{n1}} (\mathbf{S}^n \cdot \mathbf{F}_j^{u_n}) \otimes (\mathbf{S}^n \cdot \mathbf{F}_j^{u_n}) \right. \\
& + \sum_{n \in U_2^t} \tau^n \alpha^n \sum_{i' \in \Omega_i^{n2}} (\mathbf{F}_{i'}^{t_n} \cdot \mathbf{S}^n) \otimes (\mathbf{F}_{i'}^{t_n} \cdot \mathbf{S}^n) \\
& + \sum_{l \in V^t} \tau^l \alpha^l \sum_{j \in \Omega_i^{l1}} (\mathbf{G}_j^l \otimes \mathbf{G}_j^l) \\
& + \sum_{m \in W^t} \tau^m \alpha^m \left[ \sum_{j \in \Omega_i^{m1}} (\mathbf{S}^m \cdot \mathbf{F}_j^t) \otimes (\mathbf{S}^m \cdot \mathbf{F}_j^t) \right. \\
& \left. \left. + \sum_{i' \in \Omega_i^{m2}} (\mathbf{F}_{i'}^t \cdot \mathbf{S}^m) \otimes (\mathbf{F}_{i'}^t \cdot \mathbf{S}^m) \right] \right)^{-1}
\end{aligned}$$



## 2 Model discussion

### 2.1 Software

We have provided an open-source Python implementation of all models discussed in the paper, available at <https://github.com/ThomasBrouwer/HMF>. We furthermore provide all datasets, preprocessing scripts, and Python code for the experiments. Please refer to the README in the Github project.

### 2.2 Complexity

The updates for the Gibbs sampler for Bayesian matrix tri-factorisation have time complexity  $\mathcal{O}(IJK^2L)$ , compared to  $\mathcal{O}(IJK^2)$  for Bayesian matrix factorisation. For HMF the complexity becomes  $\mathcal{O}((N + M + L)I^2K^3)$  where  $I = \max_t I_t$  and  $K = \max_t K_t$ . Notice that **our model scales linearly in the number of observed datasets**. Furthermore, the random draws for columns of the factor matrices are independent of each other, and therefore the parameter updates can be formulated as efficient joint matrix operations and new values drawn in parallel. Alternatively, the draws can be done per row by using a multivariate posterior, and then all these row-wise draws can be done in parallel as well.

### 2.3 Missing values and predictions

Missing values can be indicated to the model through the mask sets  $\Omega^n, \Omega^l, \Omega^m$ . Note that this also means that if specific feature values are missing for one of the entities, these features can still be included for the other entities, simply by marking them as unobserved when we do not know their value. This is much better than imputing those values, for example using the row or column average, as the model will still fit to those imputed values.

The missing values can then be predicted, by using the posterior draws of the Gibbs sampler (after burn-in and thinning) to estimate the posteriors of the factor matrices. For example, if we wish to predict missing values in the matrix  $\mathbf{D}^l$ , we estimate  $\mathbf{F}^{t_i}$  and  $\mathbf{G}^l$ , and predictions for the missing entries are given by  $\mathbf{F}^{t_i} \cdot (\mathbf{G}^l)^T$ .

### 2.4 Initialisation

Gibbs sampling can easily get stuck in a local minimum of posterior likelihood, and therefore initialisation of the random variables is essential to obtain a good solution. There are two obvious ways to do this. Since the user specifies the values of the hyperparameters,  $\alpha_0, \beta_0, \alpha_\tau, \beta_\tau, \lambda_S^n, \lambda_S^m$ , we can use the model definition to initialise the variables  $\mathbf{F}^t, \mathbf{G}^l, \mathbf{S}^n, \mathbf{S}^m, \tau^n, \tau^l, \tau^m, \lambda_k^t$  either using the expectation of the prior model distribution, or by randomly drawing their value according to that distribution.

Alternatively, we can initialise the entity type factor matrices  $\mathbf{F}^t$  using  $K$ -means clustering, as suggested by Ding et al. [2006], and initialise the dataset-specific matrices  $\mathbf{S}^n, \mathbf{S}^m, \mathbf{F}^l$  using least squares. This can be done using the Moore-Penrose pseudo-inverse (+), as long as the dataset-specific matrices ( $\mathbf{S}^n, \mathbf{S}^m, \mathbf{G}^l$ ) are real-valued. For example,

$$\mathbf{S}^n = (\mathbf{F}^{t_n})^+ \cdot \mathbf{R}^n \cdot ((\mathbf{F}^{u_n})^T)^+.$$

If the datasets are not real-valued, we can still initialise  $\mathbf{S}^n$  or the other factor matrices in this way, but then set all values below zero to zero. We measure the effectiveness of the different initialisation methods in Section 4.3, which shows that this combination of  $K$ -means and least squares initialisation generally gives the fastest convergence.

## 2.5 Relation to tensor decomposition

Multiple matrix factorisation and tri-factorisation methods are closely linked with tensor decomposition. Here, we explore some of these connections. In particular, we show that the CANDECOP/PARAFAC (CP, Harshman [1970]) method is a less general version of the multiple matrix tri-factorisation (MMTF) part of our HMF model; and furthermore that the Tucker Decomposition (TD, Tucker [1966]), without its orthogonality constraints, is equivalent to MMTF. All three decompositions are illustrated in Figure 1, and we define them mathematically below.

The CP method decomposes a given tensor  $\mathbf{R} \in \mathbb{R}^{I \times J \times N}$  into the sum of  $K$  rank-1 tensors. This is effectively a generalisation of matrix factorisation to three (rather than two) dimensions, with each dimension getting its own factor matrix:  $\mathbf{F}^1 \in \mathbb{R}^{I \times K}$ ,  $\mathbf{F}^2 \in \mathbb{R}^{J \times K}$ ,  $\mathbf{S} \in \mathbb{R}^{N \times K}$ . Overall, we perform the factorisation  $\mathbf{R} = \mathbf{F}^1 \otimes \mathbf{F}^2 \otimes \mathbf{S}$ , where  $\otimes$  denotes the matrix outer product. Each individual entry in  $\mathbf{R}$  is decomposed as follows:

$$R_{ijn} = \sum_{k=1}^K F_{ik}^1 \cdot F_{jk}^2 \cdot S_{nk}. \quad (1)$$

The Tucker decomposition is defined similarly, but in addition to the three factor matrices, we also get a core tensor  $\mathbf{G} \in \mathbb{R}^{K \times L \times Q}$ , and the factor matrices have its own number of latent factors  $K, L, Q$ ;  $\mathbf{F}^1 \in \mathbb{R}^{I \times K}$ ,  $\mathbf{S} \in \mathbb{R}^{J \times L}$ ,  $\mathbf{F}^2 \in \mathbb{R}^{N \times Q}$ . We now factorise  $\mathbf{R} = \mathbf{G} \cdot_1 \mathbf{F}^1 \cdot_2 \mathbf{F}^2 \cdot_3 \mathbf{S}$ , where  $\cdot_i$  denotes the matrix dot product using the  $i$ th dimension of tensor  $\mathbf{G}$ . Individual entries in  $\mathbf{R}$  are decomposed as:

$$R_{ijn} = \sum_{k=1}^K \sum_{l=1}^L \sum_{q=1}^Q F_{ik}^1 \cdot F_{jl}^2 \cdot S_{nq} \cdot G_{klq}. \quad (2)$$

Now consider the multiple matrix tri-factorisation part of our HMF model. Say we are given  $N$  datasets  $\mathbf{R}^n$ , all spanning the same two entity types  $E_1, E_2$ , with  $I$  rows,  $J$  columns,  $K$  row factors (for entity type  $E_1$ ), and  $L$  row factors (for entity type  $E_2$ ). Performing MMTF on these datasets can be seen as concatenating the  $N$  matrices into one big tensor. Each

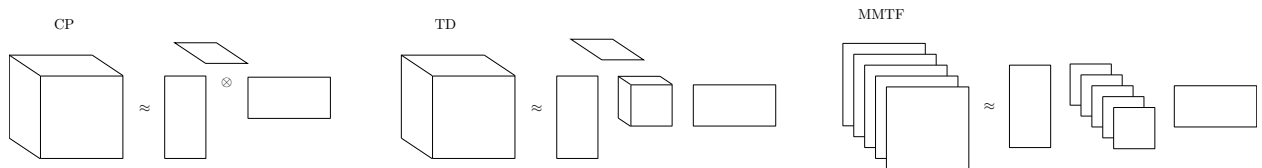


Figure 1: Overview of the CANDECOP/PARAFAC (CP, left), Tucker Decomposition (TD, middle), and multiple matrix tri-factorisation (MMTF, right) methods. CP uses the outer product ( $\otimes$ ), whereas TD and MMTF use the matrix product.

entry is decomposed as:

$$R_{ij}^n = \sum_{k=1}^K \sum_{l=1}^L F_{ik}^1 \cdot F_{jl}^2 \cdot S_{kl}^n. \quad (3)$$

Firstly, compare this with the TD formulation. If we define a new matrix  $H_{kl}^n = \sum_{q=1}^Q S_{nq} \cdot G_{klq} = \mathbf{S}_n \cdot \mathbf{G}_{kl}$ , we can rewrite the TD expression as:

$$R_{ijn} = \sum_{k=1}^K \sum_{l=1}^L F_{ik}^1 \cdot F_{jl}^2 \cdot H_{kl}^n. \quad (4)$$

Note that this is now equivalent to our MMTF expression in Equation 3, with  $F_{ik}^1 \leftrightarrow F_{ik}^1, F_{jl}^2 \leftrightarrow F_{jl}^2, S_{kl}^n \leftrightarrow H_{kl}^n$ . Since both the  $\mathbf{S}$  and  $\mathbf{G}$  matrices have to be inferred by our model, we can in fact merge them into one – as long as no further constraints are placed on them individually. Often in the TD method the three factor matrices ( $\mathbf{F}^1, \mathbf{F}^2, \mathbf{S}$ ) have orthogonality constraints placed on them. If these constraints are dropped, TD is equivalent to MMTF.

Moving on to CP, consider constraining our MMTF model to have only diagonal entries in the  $\mathbf{S}^n$  matrices (so that  $S_{kl}^n = 0$  for  $k \neq l$ ). Equation 3 then becomes:

$$R_{ij}^n = \sum_{k=1}^K F_{ik}^1 \cdot F_{jk}^2 \cdot S_{kk}^n. \quad (5)$$

This is equivalent to the CP formulation in Equation 1, with  $F_{ik}^1 \leftrightarrow F_{ik}^1, F_{jk}^2 \leftrightarrow F_{jk}^2, S_{kk}^n \leftrightarrow S_{nk}$ , showing that CP is a constrained version of MMTF, where the middle factor matrices  $S_{kl}^n$  are constrained to be diagonal.

Finally, we wanted to validate that the more general formulation offered by our HMF model is necessary to obtain good predictive performances, by comparing it with the CP method. We constrained our HMF model to have diagonal  $\mathbf{S}^n$  matrices (we call this method HMF CP), but otherwise the exact same Bayesian priors and settings. The results are given in Tables 5 and 6. We can see that in the in-matrix prediction setting CP still does fairly well, although our unconstrained MMTF model (HMF D-MTF) outperforms it on all four datasets. However, in the out-of-matrix prediction setting CP does not manage to give sensible predictions, barely doing better than the gene average baseline.

Table 5: Mean squared error (MSE) of 10-fold in-matrix cross-validation results on the drug sensitivity datasets. The best performances are highlighted in bold.

Method	GDSC $IC_{50}$	CTRP $EC_{50}$	CCLE $IC_{50}$	CCLE $EC_{50}$
HMF D-MF	0.0775	0.0919	0.0592	<b>0.1062</b>
HMF D-MTF	<b>0.0768</b>	<b>0.0908</b>	<b>0.0558</b>	0.1073
HMF CP	0.0796	0.0913	0.0560	0.1104

Table 6: Mean squared error (MSE) of 10-fold out-of-matrix cross-validation results on the methylation datasets. The best performances are highlighted in bold.

Method	GM, PM to GE	GE, GM to PM	GE, PM to GM
Gene average	1.009	1.008	1.009
HMF D-MF	<b>0.788</b>	<b>0.735</b>	<b>0.602</b>
HMF D-MTF	0.850	0.798	0.640
HMF S-MF	0.820	0.794	0.672
HMF CP	1.006	0.972	0.968

### 3 Datasets and preprocessing

#### 3.1 Drug sensitivity datasets

We will now describe the preprocessing steps undertaken for the drug sensitivity datasets used in the paper. We used four different datasets:

- Genomics of Drug Sensitivity in Cancer (GDSC v5.0, Yang et al. [2013]) – giving the natural log of  $IC_{50}$  values for 139 drugs across 707 cell lines, with 80% observed entries.
- Cancer Therapeutics Response Portal (CTRP v2, Seashore-Ludlow et al. [2015]) – giving  $EC_{50}$  values for 545 drugs across 887 cell lines, with 80% observed entries.
- Cancer Cell Line Encyclopedia (CCLE, Barretina et al. [2012]) – giving both  $IC_{50}$  and  $EC_{50}$  values for 24 drugs across 504 cell lines, with 96% and 63% observed entries respectively.

$IC_{50}$  values indicate the required drug concentration needed to reduce the activity of a given cell line (cancer type in a tissue) by half. We thus measure when an undesired effect has been inhibited by half. With  $EC_{50}$  values we measure the maximal (desired) effect a drug can have on a cell line, and then measure the concentration of the drug where we achieve half of this value. In both cases, a lower value is better.

In this paper we are most interested in enhancing predictive power by integrating different datasets. Therefore we focus on drugs and cell lines for which at least two of the four datasets have values available, giving 52 drugs and 630 cell lines. Venn diagrams displaying the overlaps between drugs and cell lines are given in Figures 2a and 2b, respectively. The CTRP dataset contains a large number of small molecule probes (311) causing very little intersection with the other datasets. We also filtered out cell lines with no features available, as discussed in Subsection 3.3.

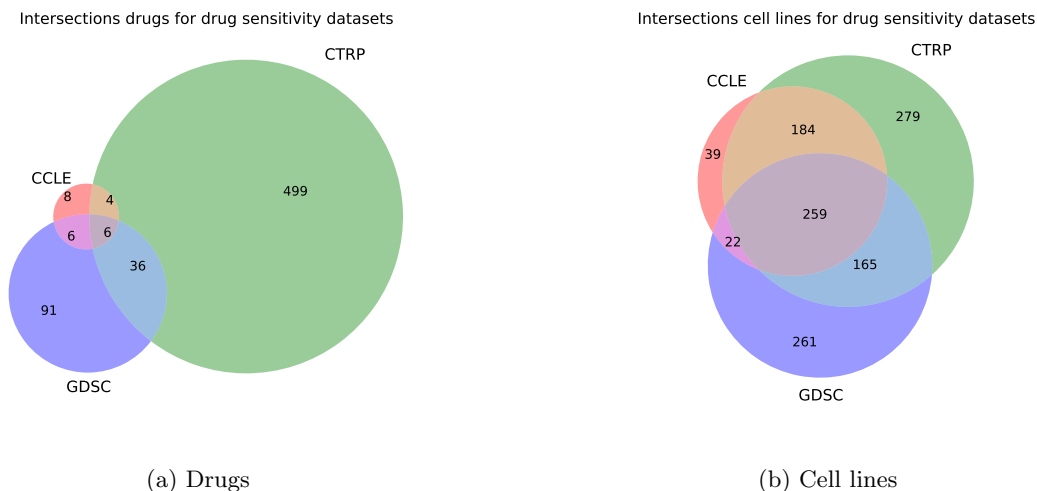


Figure 2: Venn diagrams of the drugs and cell lines in the four datasets.

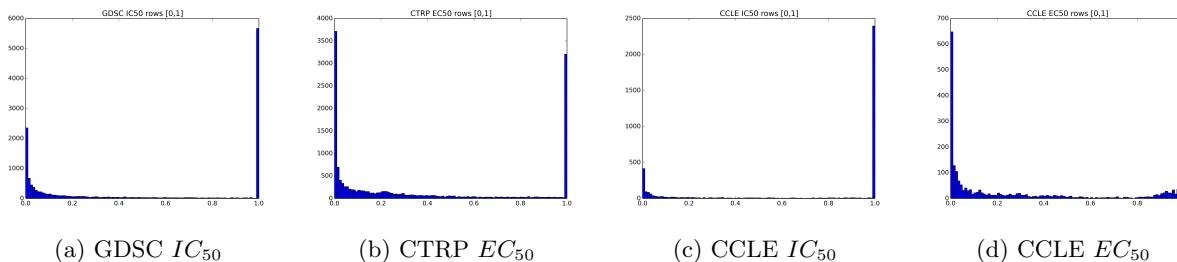


Figure 3: Plots of the distribution of values in the drug sensitivity datasets, after capping the extremely high values in the CTRP  $EC_{50}$  and GDSC  $IC_{50}$  datasets to 20.

### 3.2 Preprocessing drug sensitivity values

The CCLE and CTRP datasets all give the drug concentration levels, but the GDSC dataset gives the natural log transform of these values. We undo this transform by taking the exponent of each value. The drug sensitivity values for the CCLE  $IC_{50}$  and  $EC_{50}$  datasets lie in the range  $[0,8]$  and  $[0,10]$ , but the other two datasets sporadically have extremely large values. This is a result of the curve fitting procedure used to approximate  $IC_{50}$  and  $EC_{50}$ , and in those cases it indicates an inefficient drug for the cell line. We cap all values above 20 to 20 to resolve this issue, and as a result obtain a similar shape of distribution of values to the CCLE datasets. Finally, we map the values in each row (per cell line) to the range  $[0,1]$ . This is shown in Figure 3, where we see that the data tends to be bimodal.

### 3.3 Features

We also want to incorporate feature information, which is readily provided by the GDSC dataset for 399 of the 630 cell lines. This gave us gene expression information (13321 genes, positive values), copy number variations (CNV; 426 features, count data), and mutations (82, binary data). We filtered out cell lines without this information.

For the drugs we used the PubChem Identifier Exchange Service (<https://pubchem.ncbi.nlm.nih.gov/identexchange/identexchange.cgi>) to obtain the PubChem identifiers for all the drugs. Where there were multiple, we used the one in the GDSC database, or otherwise the first one. We then used the PaDeL-Descriptor software (<http://www.yapcwsoft.com/dd/padeldescriptor/>) to extract 1D and 2D descriptors, as well as Pubchem fingerprints of functional groups in the drugs. The GDSC dataset has drug target information available for 48 out of the 52 drugs, which we extracted from their website and encoded as a binary dataset. For the four remaining targets we mark the entries as unknown using the mask matrix in the feature dataset and kernel. We removed features with the same value across all drugs.

For each of the feature datasets we constructed a similarity kernel. For binary data we used a Jaccard kernel, and for real-valued data we first standardised each feature to have zero mean and unit variance, and then used a Gaussian kernel to compute similarities, with as variance parameter the number of features. The resulting distributions of kernel similarity values are given in Figure 4. We found that adding the similarity kernels in our HMF model

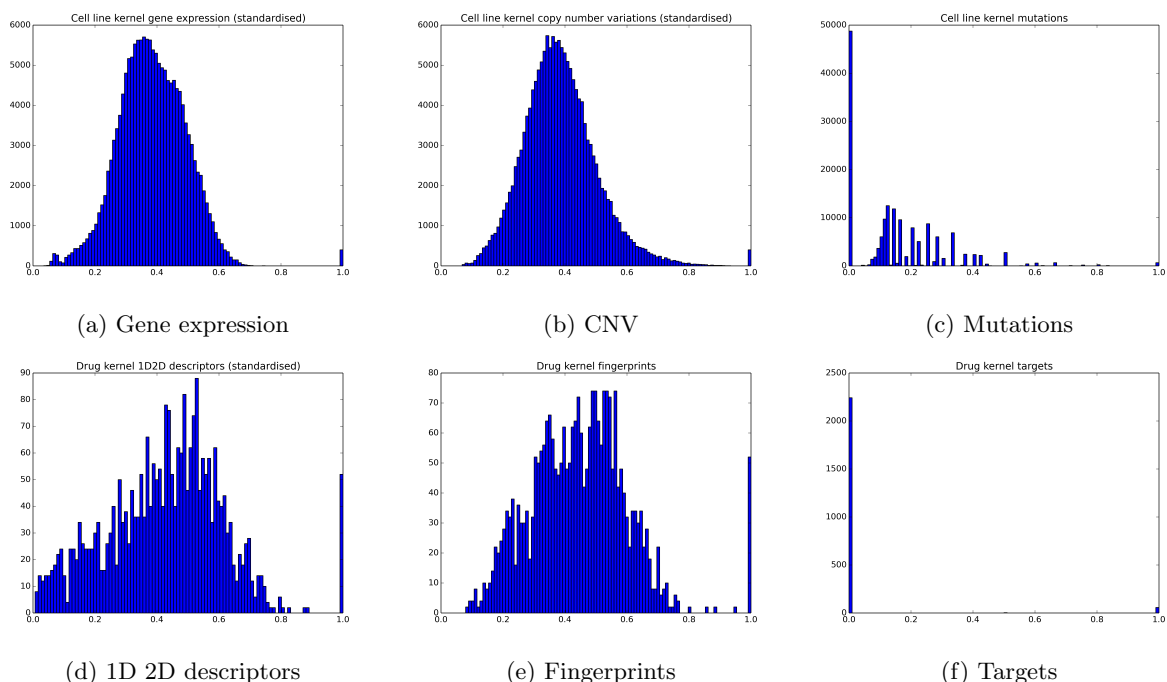


Figure 4: Plots of the distribution of similarity values in the kernel matrices based on the feature datasets. The similarity scores based on drug fingerprints, drug targets, and cell line mutations were constructed using a Jaccard kernel. The scores for gene expression data, copy number variations, and 1D and 2D descriptors were computed using a Gaussian kernel after standardising the values per feature (zero mean, unit variance) and using the number of features as the variance parameter.

did not improve the predictions. This is because adding dissimilar datasets makes it very hard (if not impossible) to find a good solution, and instead converges to a bad one.

The datasets are summarised in Table 7, and represented graphically in Figure 5. Of the remaining datasets (52 drugs by 399 cell lines), 95.1% of the entries have a value in at least one of the four datasets, and 62.9% have an entry in two or more. GDSC  $IC_{50}$  has 67.9% observed entries, CTRP  $EC_{50}$  has 72.3%, CCLE  $IC_{50}$  has 18.8%, and CCLE  $EC_{50}$  has 11.4%. This is also shown in Figure 5. Individually, the GDSC dataset contains entries for 48 drugs and 399 cell lines (73.6% observed), CTRP spans 46 drugs and 379 cell lines (86.0% observed), and finally CCLE  $IC_{50}$  and  $EC_{50}$  have 16 drugs and 253 cell lines (96.4% and 58.6% observed).

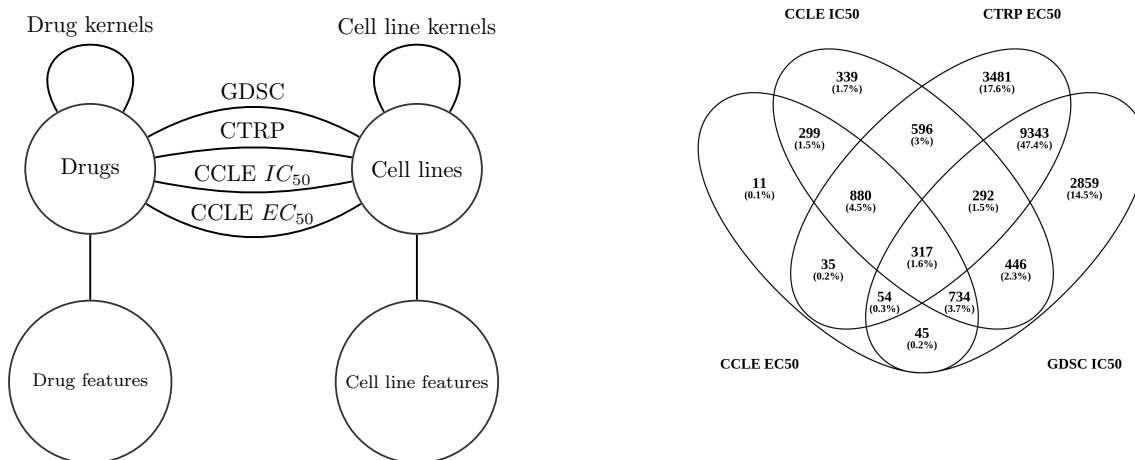


Figure 5: The datasets used for drug sensitivity prediction on the left; and the overlap of values between the four datasets on the right.

Table 7: Drug sensitivity datasets and feature datasets, summarising the entities they relate, the sizes of the datasets, and the fraction of observed entries.

Dataset	Entities	Size	Observed
GDSC $IC_{50}$	Cell lines, drugs	399, 52	67.9%
CTRP $EC_{50}$	Cell lines, drugs	399, 52	72.3%
CCLE $IC_{50}$	Cell lines, drugs	399, 52	18.8%
CCLE $EC_{50}$	Cell lines, drugs	399, 52	11.4%
Targets	Drugs, targets	52, 53	92.3%
1D2D	Drugs, 1D 2D descriptors	52, 1160	100%
FP	Drugs, fingerprints	52, 495	100%
Targets kernel	Drugs, drugs	52, 52	85.2%
1D2D kernel	Drugs, drugs	52, 52	100%
FP kernel	Drugs, drugs	52, 52	100%
CNV	Cell lines, CNVs	399, 426	100%
Mutations	Cell lines, mutations	52, 82	100%
Gene expression kernel	Cell lines, cell lines	399, 399	100%
CNV kernel	Cell lines, cell lines	399, 399	100%
Mutations kernel	Cell lines, cell lines	399, 399	100%



### 3.4 Methylation and gene expression datasets

The preprocessing for the methylation and gene expression datasets is much simpler. We obtained three datasets from the The Cancer Genome Atlas (TCGA, Koboldt et al. [2012]): promoter-region methylation (PM), gene body methylation (GM), and gene expression (GE) profiles for 254 breast cancer patients. This dataset originally spanned 13966 genes, but we focused on 160 breast cancer driver genes, given by the IntOGen database (Gonzalez-Perez et al. [2013]). We standardised each of the datasets to have zero mean and unit standard deviation per gene. Plots of the datasets containing the 160 genes, before and after standardising, can be found in Figure 6.

To construct the similarity kernels for the HMF S-MF model, giving the similarity of samples, we used a Gaussian kernel with  $\sigma^2 = \text{no. genes}$ . The kernel value distributions are plotted in Figure 7. Notice the Gaussian-like distribution of values between  $[0,1]$ . If we generate values randomly from our HMF model’s probabilistic definition (for example using value 1 for all hyperparameters, and then randomly sampling values from the prior distributions), we obtain a similar distribution of values.  $\sigma^2$  was chosen in such a way to match the model definition more closely.

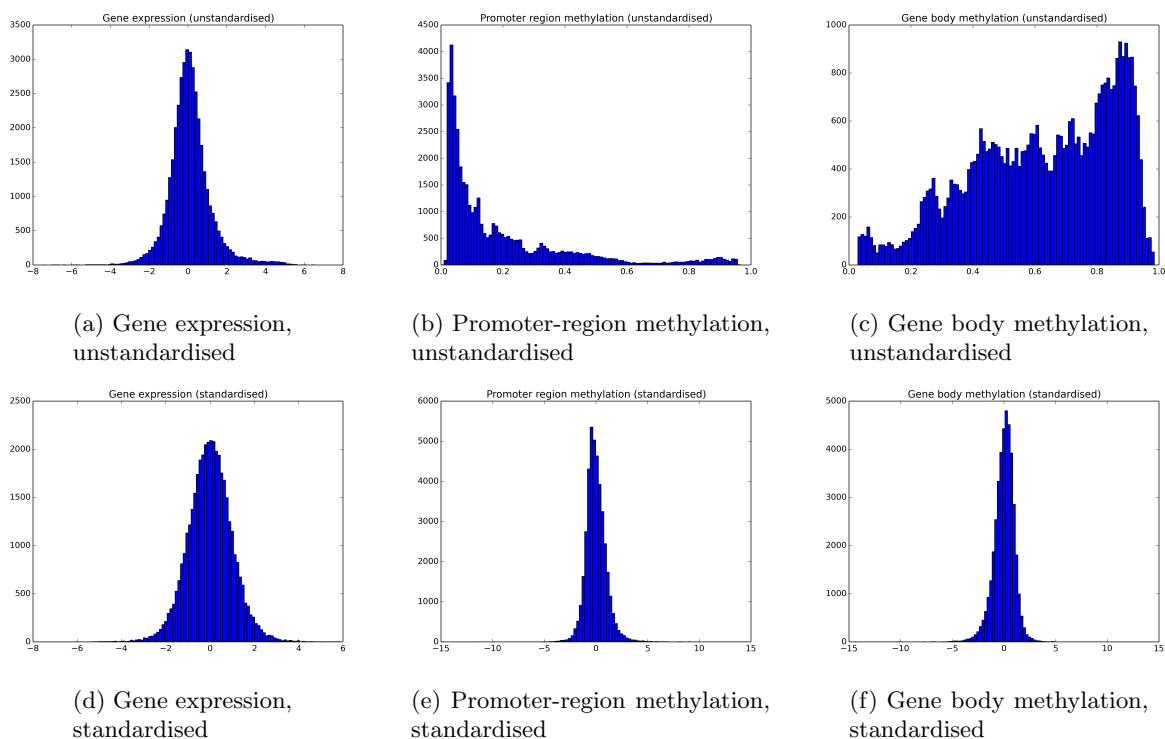
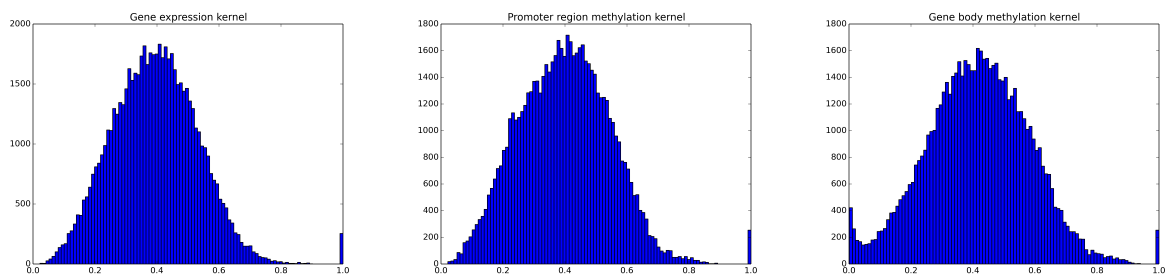


Figure 6: Plots of the distribution of the methylation datasets, before standardising (top row) and after (bottom row).



(a) Gene expression kernel      (b) Promoter-region methylation kernel      (c) Gene body methylation kernel

Figure 7: Plots of the distribution of the methylation kernels.

## 4 Additional experiments

We performed several additional experiments on the drug sensitivity and methylation datasets, to explore the advantages and limits of our models. In particular, we looked at: the run-time performance of the models; the effectiveness of automatic model selection using the Automatic Relevance Determination (ARD) Bayesian prior; the best initialisation approach; the best values to use for the dataset importances  $\alpha$ ; and the trade-offs of our hybrid choices: effects of factorisation types and nonnegativity constraints on predictive performance.

### 4.1 Run-time comparison

In this section we give a rough indication of the run-time performances of the models we compared in our experiments. We compare the time it takes for most of the matrix factorisation methods (HMF D-MF, HMF D-MTF, BNMF, BNMTF, NMF, NMTF, and Multiple NMF) on the four drug sensitivity datasets. In particular, we give the number of iterations we used in our experiments for each model, the total run-time it took to train each model, and the time per iteration of the model updates. We use  $K_t = 10, K = 10, L = 10$  for all models.

The run-time performances on the drug sensitivity datasets are given in Table 8. Often the matrix factorisation variants are faster than the matrix tri-factorisation versions (up to six times, in the case of NMF and NMTF). The non-probabilistic methods give the fastest run-time (both per iteration, and total), but our experiments show that their predictive performance is worse and they are more prone to overfitting. Our HMF models incur a slightly higher run-time than the other Bayesian models (BNMF, BNMTF), due to having to consider all four datasets at the same time, but only by a factor of roughly two (despite there being four datasets in total).

Note that when we use cross-validation to measure the predictive performances, our HMF models do not need to run nested cross-validation to choose the dimensionality ( $K_t$ ), whereas the other matrix factorisation models (BNMF, BNMTF, NMF, NMTF, Multiple NMF) do. This actually makes our models faster in cross-validation.

Table 8: Run-time performances of the matrix factorisation models on each of the four drug sensitivity datasets, giving the total number of iterations (It) used to train each model in our experiments (the same on each dataset), and then for each dataset the average number of seconds per iteration (s / it), and the total time in seconds to train a single model (Total).

Model	Iterations	GDSC $IC_{50}$		CTRP $EC_{50}$		CCLE $IC_{50}$		CCLE $EC_{50}$	
		s / it	Total	s / it	Total	s / it	Total	s / it	Total
HMF D-MF	200	0.120	23.9	0.119	23.8	0.058	11.6	0.052	10.5
HMF D-MTF	200	0.303	60.6	0.293	58.6	0.148	29.5	0.148	29.6
BNMF	1000	0.061	61.2	0.059	58.8	0.034	34.3	0.037	37.0
BNMTF	500	0.112	56.1	0.107	53.5	0.052	26.1	0.056	28.1
NMF	1000	0.007	6.9	0.007	6.6	0.002	2.3	0.002	2.0
NMTF	1000	0.037	36.6	0.033	33.4	0.009	9.3	0.010	9.8
Multiple NMF	1000	0.017	16.9	0.017	17.0	0.012	11.6	0.011	11.0

All run-time experiments were conducted on a MacBook Pro laptop, with 2.2 GHz Intel Core i7 processor, 16 GB 1600 MHz DDR3 memory, and an Intel Iris Pro 1536 MB Graphics card.

## 4.2 Model selection

Our model employs the ARD prior to perform automatic model selection. In this experiment, we verify how effective this is. We repeat the in-matrix cross-validation experiments on the four drug sensitivity experiments, for the HMF D-MF (multiple matrix factorisation) and HMF D-MTF (multiple matrix tri-factorisation) models. We vary the values for  $K_t$ , using the values [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 12, 14, 16, 18, 20, 25, 30].

Usually when we add more factors to a matrix factorisation models, it gives the model more freedom to fit well to the data, eventually leading to overfitting. Hence, the average cross-validation performance should initially go down, and then go up again as it starts overfitting. If the ARD works as desired, and we add more factors ( $K_t$  increases), the ARD will turn them off and use a similar number of factors to before. This should result in less overfitting than the equivalent model with no ARD, resulting in a flatter curve going back up as the number of factors increase.

We compare the effects of adding ARD to the HMF model in Figure 8. Here, we clearly see that adding ARD to the model consistently reduces overfitting on all four drug sensitivity datasets. ARD is not perfect, and we can still see that the curve goes up as the values for  $K_t$  increase, but this overfitting is significantly less severe than the models without ARD.

**Usage recommendations** Always use ARD to reduce overfitting in the model. Even though ARD does not always entirely eliminate the need for model selection (there is still some overfitting as  $K_t$  becomes very large), it generally makes it much less critical to try a large range of dimensionalities to find the best one. Instead, trying one or a couple will prove just as effective.

## 4.3 Initialisation

The initialisation method can have a huge impact on performance. Initialise too well, and it leads to overfitting very quickly. Initialise poorly, and your model may not converge to a good solution.

As discussed in Section 2.4, there are several ways to initialise the Gibbs sampling parameter values. Here, we measure the convergence of the HMF D-MF and HMF D-MTF models (nonnegative shared factors, real-valued private factors) on the drug sensitivity datasets. We try the following initialisation approaches:

1. **Exp:** All parameters initialised using expectation.
2. **Random:** ARD initialised using expectation, all other parameters using random draws.

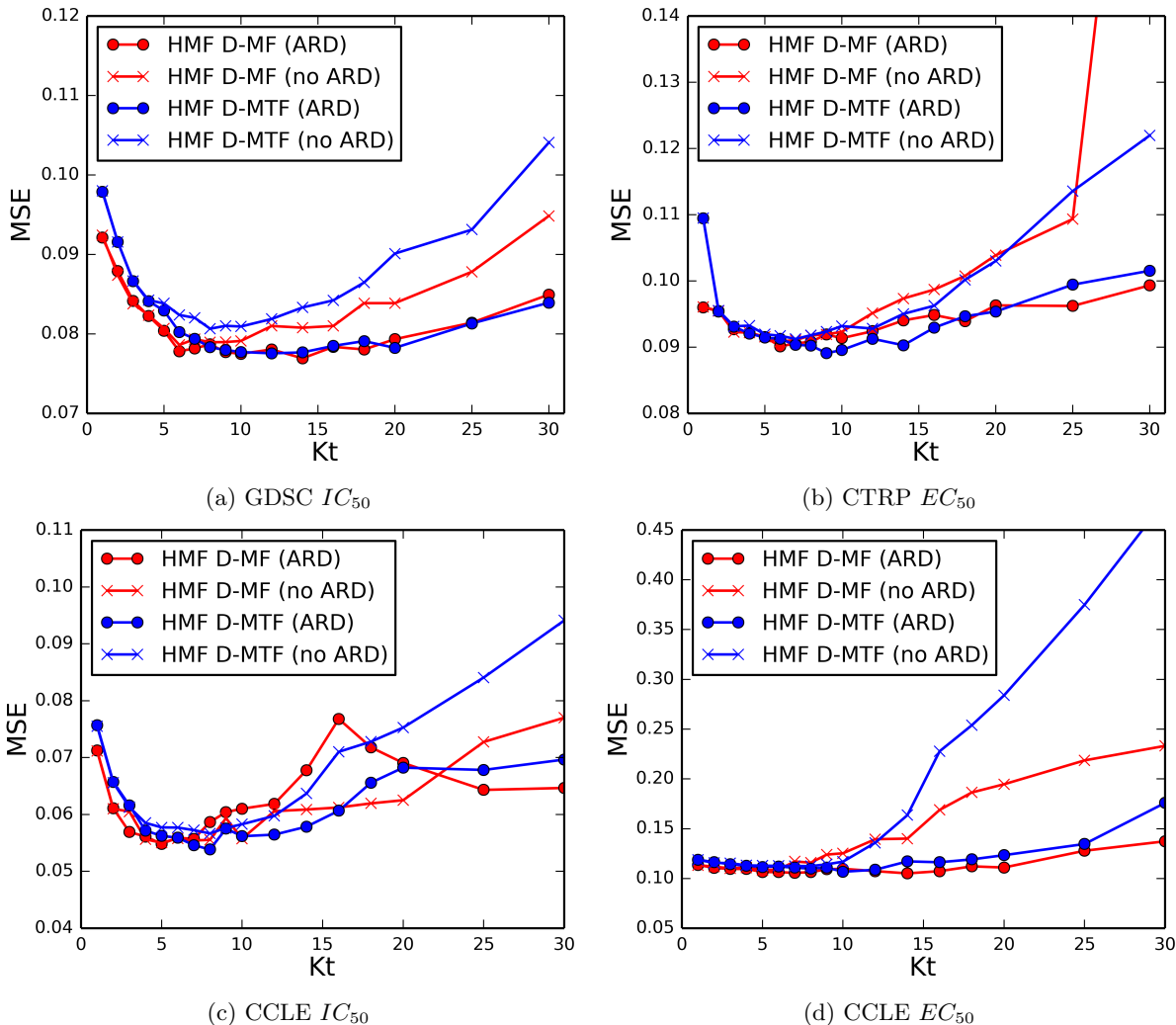


Figure 8: Graphs showing the cross-validation performance of in-matrix predictions on the drug sensitivity datasets, where we vary the dimensionality  $K_t$  for our HMF models (HMF D-MF in red, HMF D-MTF in blue), both for HMF with ARD (o), and without (x). Adding ARD clearly reduces overfitting as  $K_t$  increases.

3.  **$K$ -means, exp:** Entity type factor matrices  $F^t$  initialised using  $K$ -means, all other parameters using expectation.
4.  **$K$ -means, random:** Entity type factor matrices  $F^t$  initialised using  $K$ -means, ARD initialised using expectation, and all other factor matrices using random draws.
5. **Exp, least squares:** Entity type factor matrices  $F^t$  and ARD initialised using expectation, and all other factor matrices using least squares.
6. **Random, least squares:** Entity type factor matrices  $F^t$  initialised using random draws, ARD initialised using expectation, and all other factor matrices using least squares.
7.  **$K$ -means, least squares:** Entity type factor matrices  $F^t$  initialised using  $K$ -means, ARD initialised using expectation, and all other factor matrices using least squares.

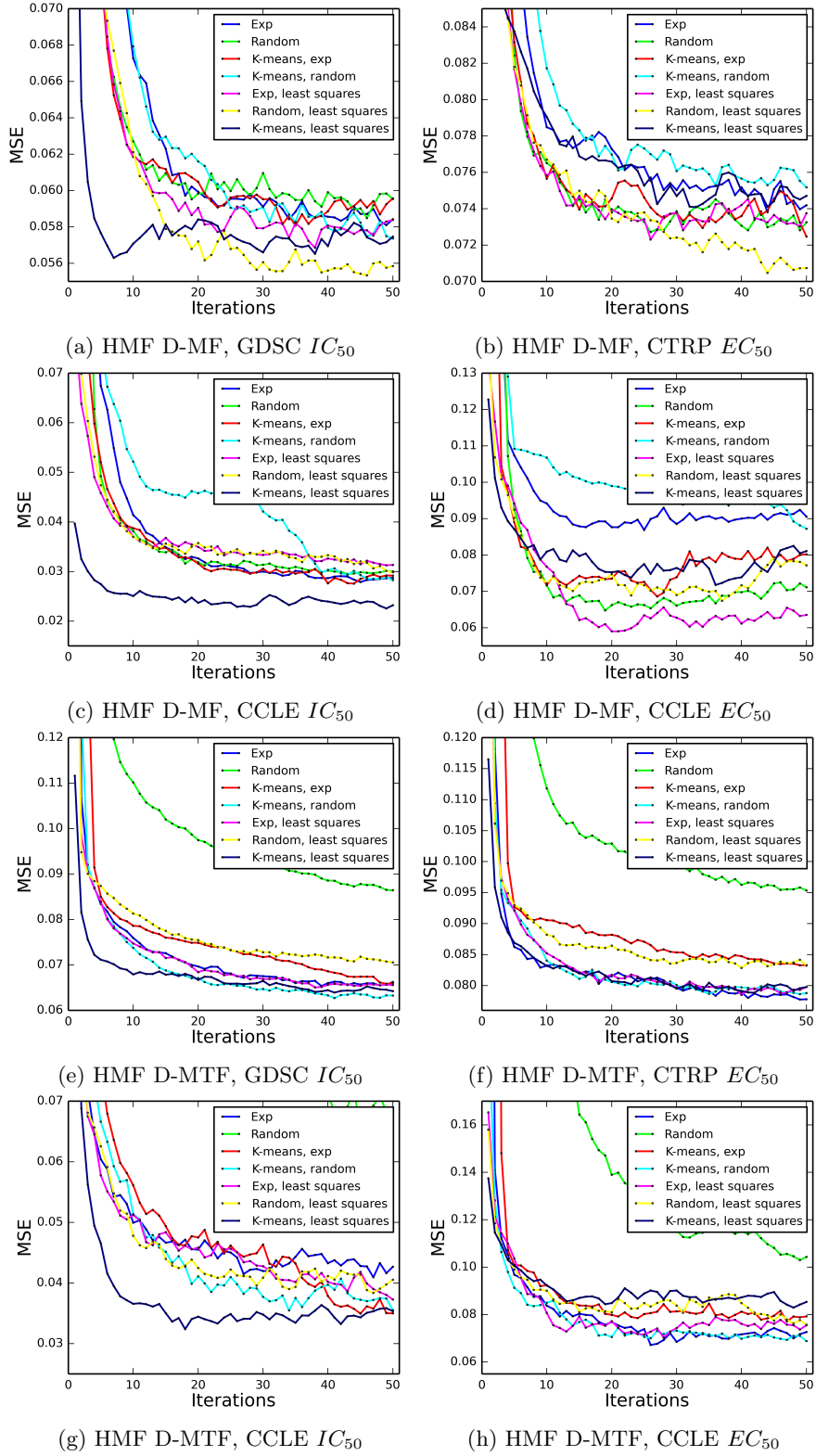


Figure 9: Graphs showing the convergence of the HMF D-MF (top two rows) and HMF D-MTF (bottom two rows) models on the four drug sensitivity datasets, for the seven different initialisation approaches.

The plots of convergence for HMF D-MF and HMF D-MTF are given in Figure 9. We can see that the  $K$ -means with least squares initialisation strategy (dark blue) provides the fastest convergence on half of the datasets, oftentimes significantly faster. Random for  $\mathbf{F}^t$  and least squares also performs well. The other strategies sometimes provide faster convergence, but none of them do so consistently.

**Recommendation** The fastest convergence is generally provided by combining  $K$ -means initialisation for  $\mathbf{F}^t$  with least squares for the other factor matrices. Random initialisation for  $\mathbf{F}^t$  with least squares also works well.

#### 4.4 Importance value

We experimented with different values for the importance values  $\alpha$ , for the out-of-matrix prediction setting. Specifically, we consider the case where we predict one dataset (either gene expression, promoter-region methylation, or gene body methylation) using the other two datasets as additional datasets.

We vary the value of  $\alpha$  for the dataset we are trying to predict ( $\alpha_0$ ) as well as for the two other datasets we are learning from ( $\alpha_1, \alpha_2$ ), using the values [0.25, 0.5, 1.0, 1.5, 2.0]. We use  $K_t = 10$  for all experiments, non-negative factors for the shared factor matrices, and real-valued for the private ones. For initialisation we use  $K$ -means for the shared factor matrices, and least squares for the private (real-valued) ones.

We perform 10-fold cross-validation, taking out 10% of the samples each time for the dataset we are trying to predict, and then measuring the mean squared error (MSE) of predictions. The average performances can be found in Table 9, for both HMF D-MF (multiple matrix factorisation) and HMF D-MTF (multiple matrix tri-factorisation).

For the HMF D-MF model (left column) we see that datasets with low predictivity (such as GE and PM – note the high MSE) have the best parameter values when the importance of the dataset to be predicted ( $\alpha_0$ ) is low, and the importance of the datasets to learn from ( $\alpha_1, \alpha_2$ ) is high. This is presumably because higher importance values lead to a better fit to the data, and if there is low predictivity, we should not fit to the data too much (otherwise we might overfit). In contrast, when the predictivity is high (such as GM), the importance value for all datasets should not be too low, because this results in a poor fit to the data and hence poor predictions.

For the HMF D-MTF model (right column) we see a similar effect, in that if the predictivity is better, the best values for the importance increase. However, for this approach all of the importance values should generally be set low if the datasets are different.

We conducted the same experiment for the approach based on similarity kernels, where we use the same ones as used in the out-of-matrix predictions from the paper. For the dataset we are trying to predict we decompose it using matrix factorisation (HMF S-MF). Matrix tri-factorisation could have also been chosen, but because the third dataset is not shared in this case, it is equivalent to matrix factorisation (in fact, we get very similar tables as the

ones shown for HMF S-MF). As before, we use nonnegative factors for the shared matrices, real-valued factors for the private ones, and  $K$ -means and least squares for initialisation. The results can be found in Table 10. Here, we see that the importance value is much less important, as long as it is not lower than 1.0.

**Recommendation** If the datasets are very dissimilar and have low predictivity, use a low importance value for the main dataset for which we are trying to predict values (like 0.5). When using HMF D-MF, use a higher importance value for the other datasets (like 1.5), but when using HMF D-MTF, use a low importance value as well (like 0.5). If the datasets are more similar, use normal importance values (1.0) for all datasets. Finally, when using similarity kernels (HMF S-MF), use the normal importance value for the kernels (1.0).



Table 9: Performances of out-of-matrix cross-validation results for HMF D-MF (left column) and D-MTF (right column), where we vary the importance value for the dataset we are trying to predict ( $\alpha_0$ ), and for the other two datasets we are learning from ( $\alpha_1, \alpha_2$ ). We have three different datasets (gene expression, GE; gene body methylation, GM; and promoter region methylation, PM). We therefore have three different prediction settings. We have highlighted the most promising parameter value areas in green, and the least promising in red.

HMF D-MF, GM + PM $\rightarrow$ GE						HMF D-MTF, GM + PM $\rightarrow$ GE					
GE ( $\alpha_0$ )	GM ( $\alpha_1$ ), PM ( $\alpha_2$ )					GE ( $\alpha_0$ )	GM ( $\alpha_1$ ), PM ( $\alpha_2$ )				
	0.25	0.5	1.0	1.5	2.0		0.25	0.5	1.0	1.5	2.0
0.25	0.878	0.843	0.835	0.831	0.832	0.25	0.866	0.906	0.941	0.950	0.954
0.5	0.845	0.849	0.829	0.832	0.829	0.5	0.874	0.870	0.914	0.941	0.945
1.0	0.848	0.916	1.195	1.322	0.831	1.0	0.933	0.985	0.940	0.942	0.961
1.5	0.871	0.948	1.288	1.407	1.470	1.5	0.959	0.958	1.065	1.026	1.000
2.0	0.896	0.966	1.392	1.711	1.782	2.0	0.957	1.080	1.211	1.145	1.147

HMF D-MF, GE + GM $\rightarrow$ PM						HMF D-MTF, GE + GM $\rightarrow$ PM					
PM ( $\alpha_0$ )	GE ( $\alpha_1$ ), GM ( $\alpha_2$ )					PM ( $\alpha_0$ )	GE ( $\alpha_1$ ), GM ( $\alpha_2$ )				
	0.25	0.5	1.0	1.5	2.0		0.25	0.5	1.0	1.5	2.0
0.25	0.859	0.799	0.795	0.799	0.798	0.25	0.862	0.893	0.943	0.946	0.947
0.5	0.811	0.812	0.783	0.783	0.784	0.5	0.841	0.885	0.898	0.943	0.944
1.0	0.789	0.814	0.987	0.788	0.783	1.0	0.852	0.848	1.012	0.885	0.904
1.5	0.784	0.809	1.070	1.247	0.870	1.5	0.876	0.866	0.900	1.080	0.904
2.0	0.798	0.835	1.111	1.280	1.255	2.0	0.884	0.881	0.900	1.096	1.099

HMF D-MF, GE + PM $\rightarrow$ GM						HMF D-MTF, GE + PM $\rightarrow$ GM					
GM ( $\alpha_0$ )	GE ( $\alpha_1$ ), PM ( $\alpha_2$ )					GM ( $\alpha_0$ )	GE ( $\alpha_1$ ), PM ( $\alpha_2$ )				
	0.25	0.5	1.0	1.5	2.0		0.25	0.5	1.0	1.5	2.0
0.25	0.790	0.708	0.703	0.697	0.701	0.25	0.773	0.775	0.851	0.854	0.861
0.5	0.724	0.670	0.687	0.688	0.685	0.5	0.761	0.746	0.774	0.818	0.850
1.0	0.698	0.657	0.659	0.670	0.685	1.0	0.782	0.757	0.754	0.786	0.779
1.5	0.698	0.655	0.667	0.657	0.666	1.5	0.832	0.752	0.745	0.783	0.795
2.0	0.689	0.668	0.671	0.676	0.665	2.0	0.836	0.811	0.802	0.791	0.805

Table 10: Performances of out-of-matrix cross-validation results for HMF S-MF, where we vary the importance value for the dataset we are trying to predict ( $\alpha_0$ ), and for the other two datasets we are learning from ( $\alpha_1, \alpha_2$ ). We have three different datasets (gene expression, GE; gene body methylation, GM; and promoter region methylation, PM). We therefore have three different prediction settings. We have highlighted the most promising parameter value areas in green, and the least promising in red.

HMF S-MF, GM + PM $\rightarrow$ GE						HMF S-MF, GE + GM $\rightarrow$ PM					
GE ( $\alpha_0$ )	GM ( $\alpha_1$ ), PM ( $\alpha_2$ )					GE ( $\alpha_0$ )	GM ( $\alpha_1$ ), PM ( $\alpha_2$ )				
	0.25	0.5	1.0	1.5	2.0		0.25	0.5	1.0	1.5	2.0
0.25	0.873	0.870	0.871	0.880	0.879	0.25	0.858	0.858	0.866	0.859	0.854
0.5	0.852	0.850	0.853	0.852	0.849	0.5	0.829	0.830	0.823	0.832	0.828
1.0	0.839	0.835	0.846	0.843	0.842	1.0	0.824	0.813	0.815	0.814	0.814
1.5	0.872	0.839	0.837	0.842	0.840	1.5	0.853	0.804	0.806	0.812	0.808
2.0	0.945	0.849	0.834	0.845	0.833	2.0	0.898	0.858	0.810	0.807	0.809

HMF S-MF, GE + PM $\rightarrow$ GM					
GE ( $\alpha_0$ )	GM ( $\alpha_1$ ), PM ( $\alpha_2$ )				
	0.25	0.5	1.0	1.5	2.0
0.25	0.762	0.771	0.782	0.796	0.797
0.5	0.722	0.724	0.740	0.742	0.753
1.0	0.701	0.706	0.707	0.718	0.721
1.5	0.756	0.703	0.702	0.709	0.710
2.0	0.757	0.710	0.702	0.697	0.701

## 4.5 Negativity constraints

Negativity constraints have two advantages. Firstly, they make it easier to analyse the factor values after performing matrix factorisation, and for example identify clusters in the data. With real-valued factors this can be a lot more complicated. Secondly, the nonnegativity constraints can reduce overfitting on sparse or noisy datasets.

We measured the effects of nonnegativity constraints on our HMF models. On the drug sensitivity datasets we tried three variants: nonnegative (all factor matrices are nonnegative), real-valued (all factor matrices are real-valued), and semi-nonnegative (shared factor matrices are nonnegative, dataset-specific ones are real-valued). On the gene expression and methylation datasets we tried the real-valued and semi-nonnegative versions. For the variants containing real-valued factor matrices, we also see whether there is a difference between row-wise and column-wise posterior draws.

We ran 10-fold cross-validation for in-matrix predictions of the drug sensitivity datasets. We use  $K_t = 10$ ,  $K$ -means initialisation for the  $\mathbf{F}^t$  matrices, and least squares initialisation for the other matrices. Results are given in Table 11, where we see that the entirely real-valued models consistently outperform all other versions, for both HMF D-MF and D-MTF. In addition, row-wise draws perform better than column-wise ones.

Similarly, we ran 10-fold cross-validation for out-of-matrix predictions of the drug sensitivity datasets. We ran the HMF D-MF, HMF D-MTF, and HMF S-MF models, again with  $K_t = 10$ , and  $K$ -means and least squares initialisation. For HMF D-MF we used importance value 0.5 for all three datasets. For HMF D-MTF, 1.5 for the main dataset we are trying to predict, and 0.5 for the others. For HMF S-MF we used 1.0 for all datasets. These results are given in Table 12, showing again that the real-valued version performs better for HMF D-MF, D-MTF, and S-MF. Column draws sometime do best, but generally row-wise draws are the best options.

Although nonnegativity can reduce the chance for overfitting, it comes at the cost of worse fitting to the data, as the nonnegativity makes convergence harder. This probably explains the lower predictive performance in this experiment for the nonnegative and semi-nonnegative models. The Bayesian nature of the models already reduces overfitting, potentially making the nonnegativity redundant. However, if the data is very sparse, and hence overfitting is more likely, the nonnegativity could be a great option.

To explore the advantages of nonnegativity in sparse settings, we repeat the experiments for sparse data predictions that were performed in the main paper on the CTRP drug sensitivity dataset, comparing the five different versions of HMF D-MF and D-MTF from the previous section. We vary the fraction of observed data, splitting the data randomly into train and test 10 times, and taking the average performance of predictions. This is shown for both methods in Figure 10. We do see that when the sparsity increases to very high levels like 90%, the nonnegative model (in dark blue) outperforms the real-valued model with row draws (in light blue). This is especially true for the HMF D-MTF model.

Table 11: Performances of in-matrix cross-validation results for HMF D-MF (top table) and D-MTF (bottom table) on the drug sensitivity datasets, where we vary the nonnegativity of the matrices. We try nonnegative, semi-nonnegative, and real-valued variants, and for the real-valued variants try both row-wise draws and column-wise draws. The best performances are highlighted in bold.

HMF D-MF		MSE			
$F_t$	$G_l$	GDSC $IC_{50}$	CTRP $EC_{50}$	CCLE $IC_{50}$	CCLE $EC_{50}$
Nonnegative	Nonnegative	0.0783	0.0907	0.0544	0.1083
Nonnegative	Real-valued (row)	0.0764	0.0899	0.0541	0.1069
Nonnegative	Real-valued (column)	0.0760	0.0903	0.0572	0.1050
Real-valued (row)	Real-valued (row)	<b>0.0758</b>	<b>0.0878</b>	<b>0.0519</b>	<b>0.1028</b>
Real-valued (column)	Real-valued (column)	0.0761	0.0889	0.0521	0.1029

HMF D-MTF		MSE			
$F_t$	$S_n$	GDSC $IC_{50}$	CTRP $EC_{50}$	CCLE $IC_{50}$	CCLE $EC_{50}$
Nonnegative	Nonnegative	0.0802	0.0916	0.0573	0.1133
Nonnegative	Real-valued (row)	0.0773	0.0899	0.0558	0.1110
Nonnegative	Real-valued (column)	0.0780	0.0904	0.0549	0.1116
Real-valued (row)	Real-valued (row)	<b>0.0770</b>	<b>0.0883</b>	<b>0.0535</b>	<b>0.1011</b>
Real-valued (column)	Real-valued (column)	0.0799	0.0904	0.0572	0.1101

Table 12: Performances of out-of-matrix cross-validation results for HMF D-MF (top table), D-MTF (middle table), and S-MF (bottom table), on the gene expression and methylation datasets, where we vary the nonnegativity of the matrices. We try semi-nonnegative and real-valued variants, and also row-wise draws and column-wise draws. The best performances are highlighted in bold.

HMF D-MF		MSE		
$F_t$	$G_l$	GM, PM to GE	GE, GM to PM	GE, PM to GM
Nonnegative	Real-valued (row)	0.880	0.792	0.670
Nonnegative	Real-valued (column)	0.851	0.803	0.672
Real-valued (row)	Real-valued (row)	<b>0.834</b>	0.775	0.651
Real-valued (column)	Real-valued (column)	0.839	<b>0.769</b>	<b>0.647</b>

HMF D-MTF		MSE		
$F_t$	$S_n$	GM, PM to GE	GE, GM to PM	GE, PM to GM
Nonnegative	Real-valued (row)	0.959	0.864	0.755
Nonnegative	Real-valued (column)	0.985	0.869	0.777
Real-valued (row)	Real-valued (row)	<b>0.893</b>	<b>0.822</b>	0.756
Real-valued (column)	Real-valued (column)	0.909	0.837	<b>0.738</b>

HMF S-MF		MSE		
$F_t$	$S_n$	GM, PM to GE	GE, GM to PM	GE, PM to GM
Nonnegative	Real-valued (row)	0.846	0.818	0.713
Nonnegative	Real-valued (column)	0.836	0.811	0.723
Real-valued (row)	Real-valued (row)	<b>0.815</b>	<b>0.805</b>	<b>0.697</b>
Real-valued (column)	Real-valued (column)	0.849	0.823	0.721

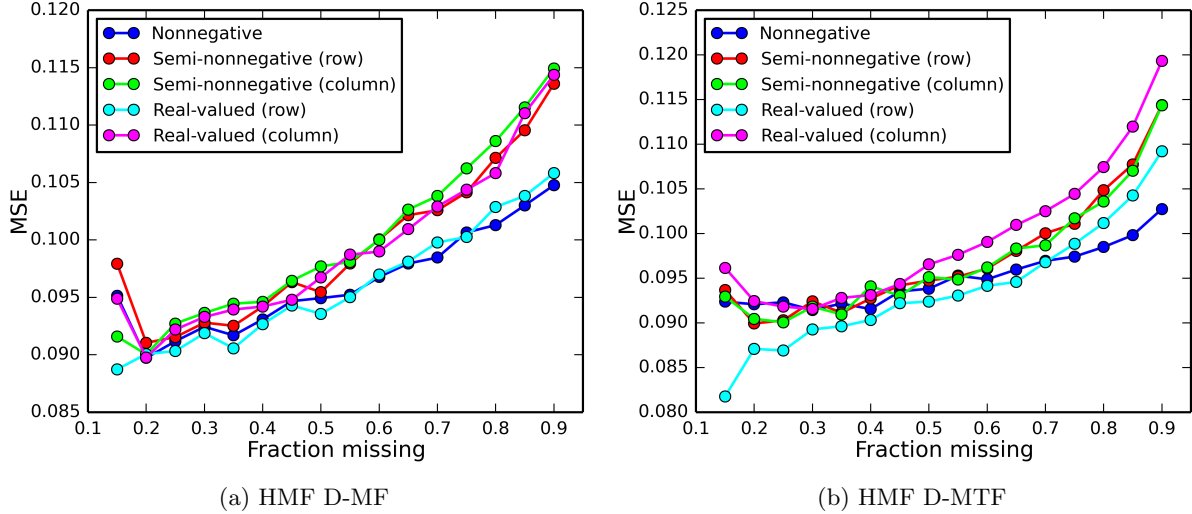


Figure 10: Graphs showing the performance for the different negativity options for the HMF D-MF (left) and HMF D-MTF (right) models, as the sparsity of the CTRP dataset increases. We plot the average mean squared error (MSE) across 10 random training and data splits.

**Recommendation** It is generally best use Gaussian priors for all factor matrices, resulting in entirely real-valued models. Row-wise draws most often give better performances than column-wise draws. In the case of very sparse matrices, nonnegative models can reduce overfitting.

#### 4.6 Factorisation types

Finally, we experiment with the choice of factorisation types. Recall that each dataset can be factorised either using matrix factorisation ( $\mathbf{D}_l$ ) or matrix tri-factorisation ( $\mathbf{R}_n$ ). For the drug sensitivity dataset, there are therefore a number of hybrid factorisation possibilities: using matrix factorisation for all (as used in the main paper; HMF D-MF), using matrix tri-factorisation for all (HMF D-MTF), using matrix factorisation on one dataset and tri-factorisation for the other three, or using matrix factorisation on two datasets and tri-factorisation on two as well. Note that applying matrix tri-factorisation on only one dataset is equivalent to using matrix factorisation on all four (since the second factor matrix is not shared with any other dataset). In this section we explore some of these choices.

##### Methylation data

We firstly consider the methylation datasets (GE, GM, and PM). We computed the Spearman correlation of values in each pair of these datasets, as given in Table 13. Here we see that GM and PM are (weakly) positively correlated, and GE and PM are (weakly negatively) correlated. We then performed 10-fold out-of-matrix cross-validation experiments (as before), varying the factorisation types on the three datasets. For simplicity we used  $K_t = 10$  and  $\alpha = 0.5$  for all datasets,  $K$ -means initialisation for the shared factor matrices ( $\mathbf{F}_t$ ) and least squares for the private ones ( $\mathbf{G}_l, \mathbf{S}_n$ ).

The results are given in Table 14, where the left three columns indicate the factorisation

Table 13: Spearman correlation between values of each of the dataset pairs.

	GE	GM	PM
GE	-	-0.07	-0.12
GM	-0.07	-	0.14
PM	-0.12	0.14	-

Table 14: Performances of out-of-matrix cross-validation results on the methylation datasets, where we vary the factorisation types on each of the three datasets. The best performances are highlighted in bold.

Factorisation type			Performance		
GE	GM	PM	GE	GM	PM
<i>R</i>	<i>R</i>	<i>R</i>	0.876	0.744	0.864
<i>D</i>	<i>R</i>	<i>R</i>	0.877	0.717	0.949
<i>R</i>	<i>D</i>	<i>R</i>	1.128	0.698	0.960
<i>R</i>	<i>R</i>	<i>D</i>	<b>0.860</b>	0.692	0.834
<i>D</i>	<i>D</i>	<i>D</i>	0.869	<b>0.663</b>	<b>0.799</b>

types, and the right three give the predictive performances on each of the datasets (each column corresponds to an experiment where we predict missing rows in that dataset). As can be seen, using multiple matrix factorisation (***D*** for all matrices) gives the best performance most of the time, which is unsurprising since the three datasets are so weakly correlated. The best performance for GE (*R*, *R*, *D*) is most likely due to random variations of performance in the cross-validation procedure (the splitting of data into train and test sets is done randomly each time).

### Drug sensitivity data

Similarly, we do this experiment for the four drug sensitivity datasets, for in-matrix predictions. These four datasets have much higher correlation (as they are repeated experiments – the same experiment conducted by different biological labs), as shown in Table 15, giving the Spearman correlation between the overlapping observed entries in each pair of the datasets. You can also find the overlaps between the datasets in Table 16, from the main paper. Here we see that:

- CCLE  $IC_{50}$  and CCLE  $EC_{50}$  are highly correlated, and have a big overlap.
- GDSC  $IC_{50}$  and CCLE  $IC_{50}$  are highly correlated, but few GDSC entries are also in CCLE. In contrast, many CCLE entries are in GDSC.
- GDSC  $IC_{50}$  and CCLE  $EC_{50}$  are (relatively) weakly correlated, but have a very small overlap. Overlap wise the same applies as for CCLE  $IC_{50}$ .
- Very few CTRP  $EC_{50}$  entries are in CCLE  $IC_{50}$  or  $EC_{50}$ , but many are in GDSC  $IC_{50}$ .

For the experiments we used  $K_t = 10$  and  $\alpha = 1.0$  for all datasets,  $K$ -means initialisation for the shared factor matrices ( $\mathbf{F}_t$ ) and least squares for the private ones ( $\mathbf{G}_l, \mathbf{S}_n$ , and if we used matrix factorisation we shared the row factors (corresponding to drugs). The results are given in Table 17. There are a lot of results, so we will consider them one column at a time.

Table 15: Spearman correlation between values of each of the dataset pairs.

	GDSC $IC_{50}$	CTRP $EC_{50}$	CCLE $IC_{50}$	CCLE $EC_{50}$
GDSC $IC_{50}$	-	0.47	0.59	0.39
CTRP $EC_{50}$	0.47	-	0.44	0.45
CCLE $IC_{50}$	0.59	0.44	-	0.65
CCLE $EC_{50}$	0.39	0.45	0.65	-

Table 16: Overview of the four drug sensitivity dataset after preprocessing and filtering.

Dataset	Number cell lines	Number drugs	Fraction observed	Overlap with other datasets			
				GDSC $IC_{50}$	CTRP $EC_{50}$	CCLE $IC_{50}$	CCLE $EC_{50}$
GDSC $IC_{50}$	399	48	73.57%	-	52.25%	9.34%	6.00%
CTRP $EC_{50}$	379	46	86.03%	57.39%	-	11.96%	7.37%
CCLE $IC_{50}$	253	16	96.42%	44.19%	51.51%	-	55.06%
CCLE $EC_{50}$	252	16	58.88%	28.52%	31.87%	55.28%	-

Table 17: Performances of in-matrix cross-validation results on the drug sensitivity datasets, where we vary the factorisation types on each of the four datasets. The best performances are highlighted in bold.

Factorisation type				Performance			
GDSC $IC_{50}$	CTRP $EC_{50}$	CCLE $IC_{50}$	CCLE $EC_{50}$	GDSC $IC_{50}$	CTRP $EC_{50}$	CCLE $IC_{50}$	CCLE $EC_{50}$
<i>R</i>	<i>R</i>	<i>R</i>	<i>R</i>	0.0767	<b>0.0889</b>	<b>0.0537</b>	0.1071
<i>D</i>	<i>R</i>	<i>R</i>	<i>R</i>	0.0765	0.0901	0.0546	0.1098
<i>R</i>	<i>D</i>	<i>R</i>	<i>R</i>	<b>0.0758</b>	0.0909	0.0543	0.1079
<i>R</i>	<i>R</i>	<i>D</i>	<i>R</i>	0.0765	0.0890	0.0584	<b>0.1055</b>
<i>R</i>	<i>R</i>	<i>R</i>	<i>D</i>	0.0768	0.0893	<b>0.0537</b>	0.1078
<i>D</i>	<i>D</i>	<i>R</i>	<i>R</i>	0.0763	0.0899	0.0569	0.1079
<i>D</i>	<i>R</i>	<i>D</i>	<i>R</i>	0.0766	0.0901	0.0553	0.1090
<i>D</i>	<i>R</i>	<i>R</i>	<i>D</i>	0.0769	0.0898	0.0554	0.1064
<i>R</i>	<i>D</i>	<i>D</i>	<i>R</i>	0.0765	0.0901	0.0566	0.1060
<i>R</i>	<i>D</i>	<i>R</i>	<i>D</i>	0.0772	0.0906	0.0543	0.1082
<i>R</i>	<i>R</i>	<i>D</i>	<i>D</i>	0.0771	0.0892	0.0542	0.1076
<i>D</i>	<i>D</i>	<i>D</i>	<i>D</i>	0.0776	0.0910	0.0562	0.1064

- **GDSC  $IC_{50}$**  (first column) – multiple matrix tri-factorisation (first row) achieves better performance than multiple matrix factorisation (last row), but a slight improvement can be achieved using a hybrid approach.
- **CTRP  $EC_{50}$**  (second column) – the best performance is achieved when using multiple matrix tri-factorisation on all datasets. Note that the CTRP dataset has a very similar correlation to all three other datasets.
- **CCLE  $IC_{50}$**  (third column) – the best performance is achieved when using multiple matrix tri-factorisation on all datasets, or when only CCLE  $EC_{50}$  is not decomposed as an *R* matrix, but using matrix factorisation (*D*) instead.
- **CCLE  $EC_{50}$**  (last column) – the best performances are achieved when the CCLE  $IC_{50}$  and  $EC_{50}$  are not both decomposed as an *R* matrix, but instead one or either is decomposed as a *D* matrix.

The last two items seem to imply that having a large overlap of entries and high correlation (such as CCLE  $IC_{50}$  and  $EC_{50}$ ), does not necessarily mean we should use matrix tri-factorisation on both datasets. The best hybrid combination of factorisations is not so obvious, but by trying out multiple candidates a suitable hybridity can be found.

**Recommendation** For dissimilar datasets, with low correlation (like the methylation data), it is best to use multiple matrix factorisation. When the datasets are very similar, with high correlation (like the drug sensitivity data), matrix tri-factorisation can give better results. These two models will generally give very good performance already. One of the hybrid combinations of matrix factorisation and tri-factorisation can sometimes lead to even better results. Nested cross-validation can be used to find the best hybrid combination.



## Bibliography

- J. Barretina, G. Caponigro, N. Stransky, K. Venkatesan, A. A. Margolin, S. Kim, C. J. Wilson, et al. The Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug sensitivity. *Nature*, 483(7391):603–7, Mar. 2012.
- C. Ding, X. He, and H. D. Simon. On the Equivalence of Nonnegative Matrix Factorization and Spectral Clustering. In *Proceedings of the fifth SIAM International Conference on Data Mining (SDM)*, pages 606–610, 2005.
- C. Ding, T. Li, W. Peng, and H. Park. Orthogonal nonnegative matrix t-factorizations for clustering. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, page 126, New York, New York, USA, Aug. 2006. ACM Press.
- A. Gonzalez-Perez, C. Perez-Llamas, J. Deu-Pons, D. Tamborero, M. P. Schroeder, A. Jene-Sanz, A. Santos, and N. Lopez-Bigas. IntOGen-mutations identifies cancer drivers across tumor types. *Nature Methods*, 10(11):1081–1082, Sept. 2013. ISSN 1548-7091. doi: 10.1038/nmeth.2642.
- R. A. Harshman. Foundations of the PARAFAC procedure: Models and conditions for an explanatory multimodal factor analysis. *UCLA Working Papers in Phonetics*, 16(10):1–84, 1970.
- D. C. Koboldt, R. S. Fulton, M. D. McLellan, H. Schmidt, J. Kalicki-Veizer, J. F. McMichael, L. L. Fulton, et al. Comprehensive molecular portraits of human breast tumours. *Nature*, 490(7418):61–70, Sept. 2012.
- B. Seashore-Ludlow, M. G. Rees, J. H. Cheah, M. Cokol, E. V. Price, M. E. Coletti, V. Jones, et al. Harnessing Connectivity in a Large-Scale Small-Molecule Sensitivity Dataset. *Cancer discovery*, 5(11):1210–23, Nov. 2015.
- L. R. Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):279–311, Sept. 1966.
- W. Yang, J. Soares, P. Greninger, E. J. Edelman, H. Lightfoot, S. Forbes, N. Bindal, et al. Genomics of Drug Sensitivity in Cancer (GDSC): a resource for therapeutic biomarker discovery in cancer cells. *Nucleic Acids Research*, 41(Database issue):D955–61, Jan. 2013.
- Y. Zhang and D.-Y. Yeung. Overlapping community detection via bounded nonnegative matrix tri-factorization. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, page 606, New York, New York, USA, Aug. 2012. ACM Press.