

Genetics and population analysis

Sparse group factor analysis for biclustering of multiple data sources

Kerstin Bunte^{*,†}, Eemeli Leppäaho, Inka Saarinen and Samuel Kaski*

Helsinki Institute for Information Technology HIIT, Department of Computer Science, Aalto University, Finland

*To whom correspondence should be addressed.

†Present address: School of Computer Science, University of Birmingham, Edgbaston B15 2TT, UK

Associate Editor: Oliver Stegle

Received on December 28, 2015; revised on March 14, 2016; accepted on April 10, 2016

Abstract

Motivation: Modelling methods that find structure in data are necessary with the current large volumes of genomic data, and there have been various efforts to find subsets of genes exhibiting consistent patterns over subsets of treatments. These biclustering techniques have focused on one data source, often gene expression data. We present a Bayesian approach for joint biclustering of multiple data sources, extending a recent method Group Factor Analysis to have a biclustering interpretation with additional sparsity assumptions. The resulting method enables data-driven detection of linear structure present in parts of the data sources.

Results: Our simulation studies show that the proposed method reliably infers biclusters from heterogeneous data sources. We tested the method on data from the NCI-DREAM drug sensitivity prediction challenge, resulting in an excellent prediction accuracy. Moreover, the predictions are based on several biclusters which provide insight into the data sources, in this case on gene expression, DNA methylation, protein abundance, exome sequence, functional connectivity fingerprints and drug sensitivity.

Availability and Implementation: <http://research.cs.aalto.fi/pml/software/GFAsparse/>

Contacts: kerstin.bunte@googlemail.com or samuel.kaski@aalto.fi

1 Introduction

Numerous clustering approaches have advanced to extract knowledge from sets of e.g. gene expression experiments, when conditions of the samples are either not known or researchers are interested in dependencies within or across experiments. Conditions or treatments can affect the expression levels of certain genes only, and similarly, many genes are likely to be co-regulated under certain conditions only. For this purpose, biclustering techniques have been developed (Cheng and Church, 2000; Hartigan, 1972; Lazzeroni *et al.*, 2002; Morgan and Sonquist, 1963). Biclustering is traditionally defined as simultaneously clustering both rows and columns in a data matrix. Depending on the metric and the data, different approaches have emerged, aiming to cluster genes based on their expression levels being the same, differing by a constant, or being linearly dependent, with respect to different conditions (Madeira and

Oliveira, 2004). Hochreiter *et al.* (2010) introduced a generative approach called Factor Analysis for Bicluster Acquisition (FABIA), accounting for linear dependencies between gene expression and conditions. The biclusters are factors of the measurement matrix, and hence can be overlapping in both genes and conditions, whereas many approaches are limited to distinct clusters. Each bicluster can also include oppositely regulated genes (up- and down-regulated) across conditions. Similar approaches have been proposed by Carvalho *et al.* (2008) and Gao *et al.* (2014), with the latter one additionally focusing on the inference of gene co-expression networks. FABIA has also been extended to better suit genotype data (Hochreiter, 2013).

Waltman *et al.* (2010) proposed an algorithm for simultaneous biclustering of heterogeneous multiple species data collections. They investigate the identification of conserved co-regulated gene groups

(modules) by comparing genome-wide datasets for closely related organisms and the evolution of gene regulatory networks. Most genes are unlikely to be co-regulated under every possible condition, and exclusive gene clusters cannot capture the complexity of transcriptional regulatory networks. Their proposed approach aims to identify meaningful condition-dependent conserved modules, integrating data across the same genes present in multiple species.

Inferring bicluster structure jointly from multiple data sources is potentially more accurate than analysis of a single set, and the discovered relationships between the sources may offer new insights. In this paper, we extend a recent generative Bayesian modelling approach, group factor analysis (GFA) ([Klami et al., 2015; Suvitaival et al., 2014; Virtanen et al., 2012](#)). GFA was developed for exploratory analysis of multiple data sources (views), resulting in an interpretable group-sparse factorization of the data collection. When the factors are additionally variable-wise sparse, as a result of introducing suitable priors, they are interpretable as biclusters of multiple co-occurring data sources that need not share the same features (genes; as opposed to [Waltman et al., 2010](#)). As a factor model GFA further shares the favourable properties of FABIA. We demonstrate its use in a multi-view drug sensitivity prediction task that the previous bicluster methods could not handle naturally: the approach shows superior prediction performance, and is able to infer meaningful structure present in subsets of the data.

2 Methods

2.1 Factor analysis

Factor Analysis for Bicluster Acquisition ([Hochreiter et al., 2010](#)) assumes preprocessed and filtered gene expression data $\mathbf{Y} \in \mathbb{R}^{N \times D}$. Every row represents a sample and every column a gene. Therefore the value $y_{i,j}$ corresponds to the expression level of the j th gene in the i th sample. A bicluster is defined as a set of rows that are similar for a set of columns, and vice versa. The model for K biclusters is

$$\mathbf{Y} = \sum_{k=1}^K \mathbf{x}_{:,k} \mathbf{w}_{:,k}^\top + \epsilon, \quad (1)$$

where each factor k is defined by an outer product of the k th columns of the factor matrix $\mathbf{X} \in \mathbb{R}^{N \times K}$ and the loading matrix $\mathbf{W} \in \mathbb{R}^{D \times K}$, and $\epsilon \in \mathbb{R}^{N \times D}$ is normally distributed noise: $\epsilon_{n,d} \sim \mathcal{N}(0, \sigma_d^2)$. The factors are the biclusters, with $|\mathbf{w}_{d,k}|$ indicating the (soft) membership of gene d in bicluster k , and $|\mathbf{x}_{n,k}|$ likewise for sample n . Both the \mathbf{W} and \mathbf{X} are given sparse priors, more specifically component-wise independent Laplace distributions:

$$p(\mathbf{W}) = \left(\frac{1}{\sqrt{2}} \right)^{DK} \prod_{k=1}^K \prod_{d=1}^D e^{-\sqrt{2}|\mathbf{w}_{d,k}|} \quad (2)$$

$$p(\mathbf{X}) = \left(\frac{1}{\sqrt{2}} \right)^{NK} \prod_{k=1}^K \prod_{n=1}^N e^{-\sqrt{2}|\mathbf{x}_{n,k}|}. \quad (3)$$

The parameters are inferred with variational expectation maximization, see [Hochreiter et al. \(2010\)](#) for details.

2.2 Sparse group factor analysis for biclustering

Group factor analysis has been proposed as an extension to factor analysis for finding factors capturing joint variability between datasets instead of individual variables ([Klami et al., 2015; Suvitaival et al., 2014; Virtanen et al., 2012](#)). It is designed to deal with several data sources $\mathbf{Y}^{(1)} \in \mathbb{R}^{N \times D_1}, \dots, \mathbf{Y}^{(M)} \in \mathbb{R}^{N \times D_M}$ (called views) of

dimensionality D_m with N co-occurring observations. GFA models the n th sample of the m th data view as

$$\mathbf{y}_n^{(m)} \sim \mathcal{N}\left(\mathbf{W}^{(m)} \mathbf{x}_n, \tau_m^{-1} \mathbf{I}_{D_m}\right), \quad (4)$$

where τ_m is the noise precision of view m . The loading matrix $\mathbf{W}^{(m)}$ is given a sparse prior that allows omitting any component k from affecting drug data view $\mathbf{Y}^{(m)}$. This enables a group-sparse factorization, where components may be (i) specific to a single data view, (ii) shared between all the data views or (iii) shared between any subset of the data views. As we are interested in finding biclusters, we introduce a prior that is additionally variable-wise sparse, that is, across the elements of the matrices $\mathbf{W}^{(m)}$ and \mathbf{X} . This is done similarly to how [Khan et al. \(2014\)](#) produced variable-wise sparsity, but now for both variables and samples to produce biclusters. Namely we use the following spike and slab priors ([Suvitaival et al. \(2014\)](#) included sparsity for samples and mentioned the connection to biclustering, but the interpretation was not explored further.):

$$x_{n,k} \sim b_{n,k}^{(x)} N\left(0, \left(\alpha_k^{(x)}\right)^{-1}\right) + (1 - b_{n,k}^{(x)}) \delta_0 \quad (5)$$

$$w_{d,k}^{(m)} \sim b_{d,k}^{(m)} N\left(0, \left(\alpha_k^{(m)}\right)^{-1}\right) + (1 - b_{d,k}^{(m)}) \delta_0 \quad (6)$$

$$b_{n,k}^{(:)} \sim \text{Bernoulli}(\pi_k^{(:)}) \quad \pi_k^{(:)} \sim \text{Beta}(a^\pi, b^\pi) \quad \alpha_k^{(:)} \sim \text{Gamma}(a^\alpha, b^\alpha) \quad (7)$$

where the binary $b_{d,k}^{(m)}$ determines whether the component (bicluster) k is active in the d th feature of $\mathbf{Y}^{(m)}$ (for all non-zero n in $b_{n,k}^{(x)}$), $\alpha_k^{(m)}$ the probability of $b_{d,k}^{(m)} = 1$. The prior is analogous for samples (i.e. the rows of the data matrices) through $b_{n,k}^{(x)}$. Effectively, the spike and slab prior will set weights that affect the data (likelihood) only little to 0, which allows direct biclustering interpretations without a need for arbitrary thresholding afterwards. The model is completed with a gamma prior for the noise precision parameters τ_m and uninformative hyperpriors ($a^\pi, b^\pi, a^\alpha, b^\alpha, a^\tau, b^\tau = 1$).

In this formulation, the data source information (feature grouping) is used in three ways: (i) the noise precision (τ_m) is the same for all the features in a view, (ii) the binary vector $b_{:,k}^{(m)}$ has a common probability ($\pi_k^{(m)}$) of being active and (iii) the scale of a component ($\alpha_k^{(m)}$) is shared within a view. With an uninformative Gamma prior, often called Automatic Relevance Determination prior, this third property implements the group sparsity. The second property implies that a specific feature d (in view m) is more likely to be active in a bicluster, if many of the features in view m belong to the said bicluster, and vice versa. This allows explaining variance that is not present in all the data sources, but dense in some of them, more robustly. Given data with significant (source specific) structured variation respective components can help to detect biclusters more accurately.

The formulation above assumes that all the data views have co-occurring samples. We also extend GFA for joint modelling of datasets that are paired in two modes (see [Fig. 1](#)), i.e. $\{\mathbf{Y}^{(1,1)}, \dots, \mathbf{Y}^{(M_1,1)}, \mathbf{Y}^{(1,2)}, \dots, \mathbf{Y}^{(M_2,2)}\}$, where $\mathbf{Y}^{(m,2)} \in \mathbb{R}^{D_1, N_2}$ is paired with the features of $\mathbf{Y}^{(1,1)}$. Both the modes will have a set of components identical to the ones presented above with one exception, and hence we will not repeat the details of the priors here. The exception is that the view paired in both the modes is generated from the components of both the modes, as

$$y_{i,j}^{(1,1)} \sim \mathcal{N}\left(\mathbf{w}_j^{(1,1)} \mathbf{x}_i^{(1)} + \mathbf{w}_j^{(1,2)} \mathbf{x}_i^{(2)}, \tau_{1,1}^{-1}\right). \quad (8)$$

As the priors remain conjugate, the model can be inferred using Gibbs sampling, resulting in linear complexity in both N and D , but

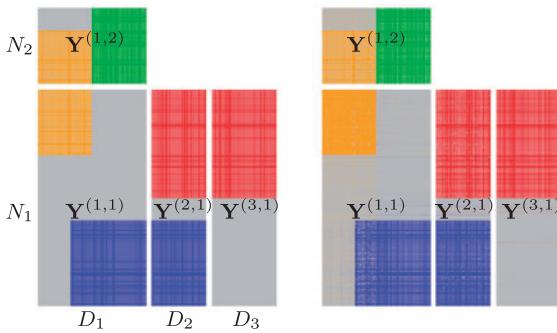


Fig. 1. Left: four non-overlapping biclusters (coloured blocks) used in the multi-view data (gray area). Right: The biclusters inferred by GFA

cubic w.r.t. the number of components K . The parameters shown in this paper, and used in predictive tasks, will be the posterior means. In the following sections we will test FABIA and GFA in several simulation studies and a drug sensitivity analysis, incorporating genetic data available from the DREAM project.

3 Simulation study

We show an illustrative example of bicluster inference, generating a collection of datasets, $\{\mathbf{Y}^{(1,1)}, \mathbf{Y}^{(2,1)}, \mathbf{Y}^{(3,1)}, \mathbf{Y}^{(1,2)}\}$, with 200 samples and dimensions (100, 50, 60) for $\mathbf{Y}^{(:,1)}$, and 100 samples and dimension (70) for $\mathbf{Y}^{(1,2)}$. The data collection was generated with four biclusters and additional noise with variance 1. The non-zero parts of \mathbf{x} and \mathbf{w} for the biclusters were drawn from $\mathcal{N}(0, 1)$, but truncated between absolute values 1 and 2 for illustrative purposes. The data structure is shown in Figure 1 (left); for clarity the biclusters are non-overlapping blocks. We inferred the component structure of these data using GFA; the posterior mean of the biclusters is visualized in Figure 1 (right). GFA can clearly infer this kind of component structure very accurately.

GFA has been designed for joint modelling of multiple datasets. However, when the data consist of one set only, FABIA and GFA are essentially the same model; in the current implementations there is the technical difference that FABIA has a continuous valued sparsity prior for \mathbf{X} and \mathbf{W} , whereas GFA implements a discrete choice with the spike-and-slap. We first investigate the effect of this technical difference by comparing GFA with FABIA on single-view data (FABIA1), and then investigate how much multi-view data helps, by comparing GFA against FABIA for which data are concatenated into a single matrix (FABIA2).

For the simulation studies, we construct data from the generative model Eq. (4), with matrices \mathbf{X} and \mathbf{W} generated such that each element is either zero or sampled randomly from the normal distribution $\mathcal{N}(0, 1)$ to build the bicluster(s). The resulting data matrices \mathbf{Y} are given to the methods, which then return the bicluster estimates $\mathbf{x}'_{:,k}\mathbf{w}_{:,k}^\top$. They are compared to the true biclusters to analyze the models' performance. FABIA is run with the correct number of biclusters K , and the results are reported for a range of thresholds. GFA learns the cluster number by driving unnecessary ones to zero, and we used a component number 5 above the correct K . The final biclustering is based on 101 posterior samples (2000 burn-in samples, 20 thinning): if the majority of $(\mathbf{x}_{:,k}\mathbf{w}_{:,k}^{(m)\top})_{ij}$ are non-zero in the posterior samples, then $\mathbf{Y}_{ij}^{(m)}$ is assigned to the k th bicluster, otherwise not. All the simulation studies are repeated

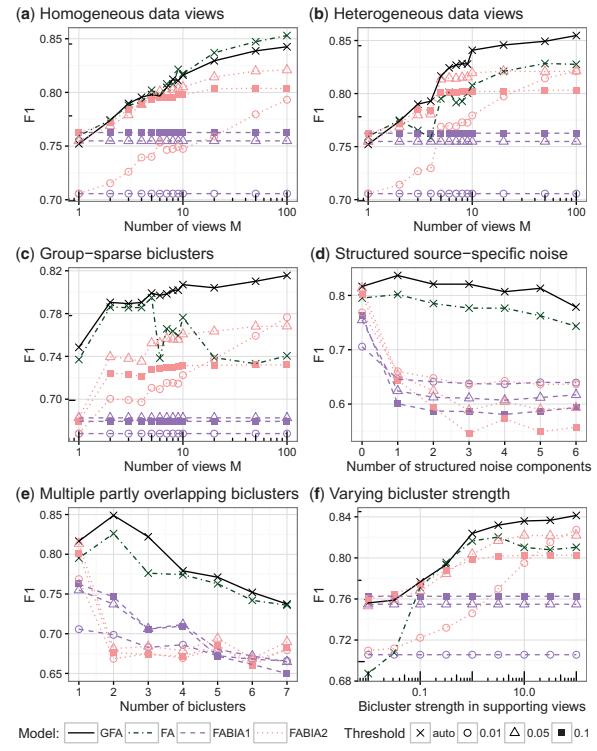


Fig. 2. Simulated experiments comparing the abilities of GFA, FA and FABIA to detect data-generating biclustering. (a) and (b) report the F_1 scores over a varying number of data views (M) present, for homogeneous and heterogeneous data collections, respectively. The biclusters are further assumed group-sparse in (c). In (d), the problem is made more challenging by adding structured noise on top of the signal, whereas the number of biclusters is varied in (e). In (f), the strength of the bicluster is varied in the supporting views. FABIA1 uses only the data matrix of interest, whereas FABIA2 and FA have side information concatenated with it; both FABIAS are reported for thresholds (0.01, 0.05, 0.10) determining the biclusters

10 times and we report the average F_1 score for detecting the true bicluster structure:

$$F_1 = \frac{2TP}{2TP + FN + FP} \quad (9)$$

where TP , FN and FP denote the number *true positives*, *false negatives* and *false positives*, respectively, summed over all the elements of the data matrix. To be able to evaluate the effects of different sparsity priors and inference techniques, we also compare to FA with similar priors as GFA, using the full concatenated data, and inferred with Gibbs sampling. This is done by changing Eq. (5) to Eq. (7) to allow only single $\alpha_k^{(1)}$ and $\pi_k^{(1)}$ across the data features, and by giving each feature independent noise precision τ_d .

By default we use $M=5$ data views, $N=50$ samples, $D_m=100$ features per data view and one bicluster active in 70% of the samples and the features. Bicluster and noise variance in the view of interest are set to 1. To make the data views heterogeneous, the other 4 data views are given bicluster and noise variances of $((0.2, 0.2), (5, 5), (0.2, 5), (5, 0.2))$. We report the mean performance of the methods in Figure 2, in six different experimental settings:

- All data views are set to be homogeneous having variance of 1 in the bicluster and the noise residual, and the number of data views is varied. Due to the homogeneity of the data views, GFA has no advantage over FA. FA outperforms FABIA on high-dimensional large number of views.

- (b) Similar to (a), but with heterogeneous data with respect to the bicluster and noise residual strength. The multi-view approach of GFA is superior, while FA and FABIA have similar performance.
- (c) Similar to (b), but with sparse biclusters w.r.t. the samples and group sparse w.r.t. the features (present in all but every 3rd view). This matches GFA's assumptions leading to superior results.
- (d) Additional view-specific noise components (0–6) were added on top of the bicluster signal (as $\mathbf{x}_{\text{noise}} \mathbf{w}_{\text{noise}}^T$, each vector element sampled from $\mathcal{N}(0, 1)$). The GFA-type of priors and inference are clearly more robust against the structured noise.
- (e) We evaluated the accuracy of detecting 1–7 partly overlapping biclusters. GFA outperforms FA, whereas FABIA does not seem very robust when the number of biclusters increases; using the additional data can even decrease its performance.
- (f) The strength of the biclusters is varied for the additional views (precision α ranging from 10^{-2} to 10^2). GFA is more accurate when the additional views have strong biclusters, and has a similar performance with FABIA when weaker. FA suffers compared to FABIA when the additional views get less relevant.

Across all the studies discussed above, GFA was able to detect the correct number of biclusters (and additional noise components) exactly in 90.4% of the runs, and overestimated it by 1, 2 or 3 in 9%, 0.4% and 0.1% of the runs, respectively. The extra components were modelling artificial structure detected in the residual noise and did not resemble the bicluster structure. To confirm that they did not give an advantage in the comparison, we also ran FABIA with five extra components, resulting in consistently worse performance compared to the reported runs, where FABIA was given the true component amount. The standard deviations of the mean $F1$ scores, averaged over the 10 independent repetitions, ranged from 0.003 to 0.01. Inferring a single model took on average 22 s, 45 s, 0.04 s and 0.5 s for GFA, FA, FABIA1 and FABIA2, respectively, demonstrating the efficiency of the EM-algorithm in FABIA.

Our FA implementation is generally, but not consistently, more robust than FABIA in bicluster detection with additional data sources. The advantages of the multi-view setup of GFA (vs concatenation in FA and FABIA2) are most significant when (i) there are plenty of heterogeneous data views, (ii) the biclusters are group-sparse or (iii) the data views are highly heterogeneous. These conditions are realistic in real-life applications.

4 Drug response study

The NCI-DREAM drug sensitivity prediction challenge (Costello *et al.*, 2014) provided publicly available data consisting of gene expression (GE), RNA, DNA methylation (MET), copy number variation (CNV), protein abundance (RPPA) and exome sequence (EX) measurements for 53 human breast cancer cell lines. In the challenge, expression data were based on Affymetrix Genome-Wide Human SNP6.0 Array and Affymetrix GeneChip Human Gene 1.0 ST microarrays. RNA sequencing libraries were prepared using the TruSeq RNA Sample Preparation Kit and whole transcriptome shotgun sequencing was performed. The Mutation status was obtained from exome-capture sequencing and GenomeStudio Methylation Module v1.0 was used to express the methylation for each genome-wide detected CpG locus resulting in values between 0 (completely unmethylated) and 1 (completely methylated) proportional to the degree of methylation at any particular locus. More

details on the preparation of the genomic data for the challenge are provided by Costello *et al.* (2014). Each cell line was exposed to 31 therapeutic compounds and the dose-response values of growth inhibition were collected. The drug response data was revealed only for 35 of the cell lines, and the challenge was to predict the response of the remaining 18 cell lines, ranking them from the most sensitive to the most resistant.

As the drug response prediction problem is extremely challenging, we performed the following steps, learning from Costello *et al.* (2014), to increase the signal-to-noise ratio: We reduced the dimensionality to the 500 genes with the highest average variance over the data views, including the overlapping set of 14 genes appearing in the RPPA dataset. Furthermore, the most predictive data sources for further analysis were chosen by 7-fold cross validation on the 35 training samples with known drug response values. To compare the sources the root mean squared error (RMSE), as well as Pearson and Spearman correlations of the predicted drug responses, were computed averaged over 10 repetitions of the experiments, each with different random splits. We inferred the GFA model for this multi-view data by ignoring the missing drug response data in the likelihood, after which the missing values can be predicted from $\mathbf{X}\mathbf{W}^{(m)}$, where \mathbf{X} for the missing cell lines is inferred based on the other data sources only. The performance of GFA trained with different combinations of data views is shown in Table 1.

The most promising views finally chosen for the bicluster analysis were gene expression, methylation, exome sequence and RPPA measurements, leaving out the copy number variation and RNA. Finally, we ran GFA for the full data (handling the test drug responses as missing values) and reconstructed the missing data averaged over the posterior samples of 50 sampling chains. We gave the model a mildly informative prior assuming signal-to-noise ratio of 0.5. All the sampler chains were initialized with $K=60$, allowing data-driven inference of model complexity (resulting in 48–56 components). A total of 100 sampled parameters were stored for each chain (every 20th sample stored after 10 000 burn-in iterations), resulting in an average runtime of 84 min per chain. The performance was quantified using the same score as in the challenge, that is, the weighted averaged probabilistic concordance index. We achieved a score of 0.592, which would have been placed the first in the

Table 1. Averaged 7-fold cross validation results for GFA on the training set of the DREAM7 drug sensitivity prediction challenge to identify the views showing best prediction performance (**bolded**) for further analysis

Views used	RMSE	Pearson	Spearman
All	1.9	0.031	0.079
GE, MET, CNV, RNA, RPPA	2.3	0.016	0.088
GE, CNV, RNA, RPPA, EX	2.0	0.031	0.078
GE, MET, CNV, RPPA, EX	1.5	0.040	0.085
GE, MET, CNV, RNA, EX	1.8	0.012	0.078
MET, CNV, RNA, RPPA, EX	1.6	0.018	0.058
GE, MET, RNA, RPPA, EX	1.9	0.040	0.089
GE, MET, CNV, RPPA	1.8	0.028	0.071
GE, CNV, RPPA, EX	1.8	0.018	0.074
GE, MET, CNV, EX	1.5	0.024	0.090
MET, CNV, RPPA, EX	2.1	0.020	0.061
GE, MET, RPPA, EX	1.4	0.046	0.087
GE, MET, RPPA	1.9	0.024	0.072
GE, RPPA, EX	1.6	0.016	0.059
GE, MET, EX	1.5	0.042	0.084
MET, RPPA, EX	1.8	0.011	0.075

challenge (winner model reaching 0.583), indicating excellent prediction performance of GFA on this data, possibly stemming from the biclustering nature of the model. Furthermore, we ran GFA utilizing data sources paired in two modes in a similar way with additional functional connectivity fingerprints describing the drugs (FCFP; calculated with PaDEL-Descriptor, Yap, 2011), allowing joint modelling of biological and chemical effects in the measured data. The additional chemical view resulted in a slight increase in the target score, to 0.599. The structure of the latter model is interpreted in the following sections, motivated by the excellent predictive performance.

4.1 Robust components

For interpretation purposes we next sought representative point solutions to describe the posterior distributions. Due to the extremely challenging nature of the problem the total variance explained by individual components is small. To minimize the risk of analyzing patterns occurred by chance, we searched for components that occur consistently across the different sampling chains, making the assumption (which was checked manually) that component indices are reasonably stable within a chain, but can naturally be arbitrarily permuted between chains. To find the matches between chains we averaged the components over the posterior samples within their chain, and compared using cosine similarity. If the similarity of the best match exceeded the threshold 0.80, we considered the components to be the same. Furthermore, we chose to further study components found in at least half of all chains, deemed robust in this procedure. Out of the average 52.6 components inferred by the sampling chains, 25 were on average chosen to the set of robust components. Ideally we would infer the model parameters with a single well-mixing sampling chain, but as the posterior is multimodal (and we do not want to constrain it artificially) the inference problem is extremely challenging, and we resort to the described computational simplification.

We observed that some of the components are very sparse, only containing one or two cell lines and hence most probably explaining outliers in the data. Therefore, we will focus the interpretations on the more dense biclusters only; 3 out of the total 27 biclusters found predicted 1.26%, 0.06% and 0.11% of the total variance in the test data (2.89%, 0.3% and 0.98% in the set of active drugs). There was 1 additional bicluster shared with the drug descriptors, but it had no significant effect to the test data. With the drug sensitivity prediction of these four components only, we received a target score of 0.591.

4.2 Interpretations of the biclusters

For interpretation purposes, we collected the descriptions of the drugs and cell lines used in the challenge (Costello *et al.*, 2014). Some groups of drugs can be identified, which we will abbreviate as: autophagy (au), cell cycle (cc), metabolism (me), regulation (re) and signalling growth (gr) drugs, as well as nuclear factor (nf), protease (pr) and receptor tyrosine kinase (rtk) inhibitors. Furthermore, most of the cell lines represent a subtype of cancer which can be categorized as basal or luminal.

The bicluster structure of the activity patterns for the first component in the drug sensitivity (DS) view, consisting of measurements of sensitivity of cell lines to drugs, is depicted in Figure 3a. Component one distinguishes basal and luminal cell types, without that information being used in the training. The response for all five cell cycle and all four metabolism drugs is positive or above average for most of the basal cell lines, whereas luminal cells show negative activation. Luminal cells respond strongly to regulation drugs,

where the response of basal cells is negative. Component 2 shows high activity patterns for proteasome and cell cycle drugs as depicted in Figure 3b. The other components have relatively small biclusters with only a few active cells and drugs. Component 3, for example, shows cells that are (un)responsive to rtk inhibitors and otherwise mixed groups of cells and drugs (see Fig. 3c). The remaining robust component, in the second mode, was associated with most of the drugs and drug descriptors, and weakly with approximately half of the cell lines (strongly with T47DKBLUC).

Due to the large number of genes in the other views, we show summaries of enrichment of known cancer genes in the components. We performed hypergeometric tests comparing a varying number of the most active genes (i.e. genes with the highest mean absolute values in W corresponding to RPPA and GE views) and random sets of equal size, for the occurrence of known cancer genes (extracted from Stephens *et al.*, 2012) in the most predictive components and all views. Low P-values indicate that the approach is able to detect a significant amount of known (breast) cancer genes in the top active genes of the components when compared to random subsets of genes in the views. Figure 4 shows the results of the hypergeometric test for two robust components. For component 1 (Fig. 4a) we observe highly significant cancer activity already in the 10 most active genes in every view, except in the RPPA data. In the Exome sequence data we observe significant cancer gene activity independently of the size of the subset. For component 2 (Fig. 4b) we observe less significant activity of cancer genes in the gene expression and exome data than in component 1. However, we find high cancer gene activity in the methylation view and very high activity of breast cancer genes in exome sequencing data.

Besides the statistical tests we report the results on the level of individual genes. We condensed the analysis showing the 50 most active genes (in terms of their absolute values) in component 1 in the gene expression (Fig. 3d) and the RPPA view in Figure 3e. Component 1 contains proportional and anti-proportional co-regulated genes as indicated by the intensity of the biclusters depicted. The genes which are known as cancer or breast cancer genes are marked by a black and gray squares, respectively. Figure 3f summarizes the mean participation in biclusters of the cells in seven different components on the top 10 active genes in each of the views. The left side contains the list of all cells sorted by their mean absolute values in the seven most active components, which are depicted in the middle row. The right side contains the list of top genes clustered by their mean activity in these components, accompanied by a shortcut for the view they were taken from and [C] in case they are known as cancer gene in the literature. Component one (coloured red) contains the biggest biclusters with comparably high mean absolute values depicted by thicker connecting lines in lots of cells and genes throughout the different views except the exome sequence data. Even only selecting 10 most active genes from each view delivers at least one known cancer gene. Although components overlap they also depict relationships of different cells in different views.

Furthermore, we performed a Gene Ontology (GO) enrichment analysis on the most active gene sets in the components and gene related views. In GO the genes or gene products are hierarchically classified and grouped into three categories: *molecular function* (mf) describing the molecular activity of a gene, *biological process* (bp) denoting the larger cellular role and *cellular component* (cc) depicting where the function is executed in the cell. The enrichment analysis was performed directly in the GO website (<http://geneontology.org/>), which connects the PANTHER (Mi *et al.*, 2013) classification system with GO annotations. From each of the gene-related views

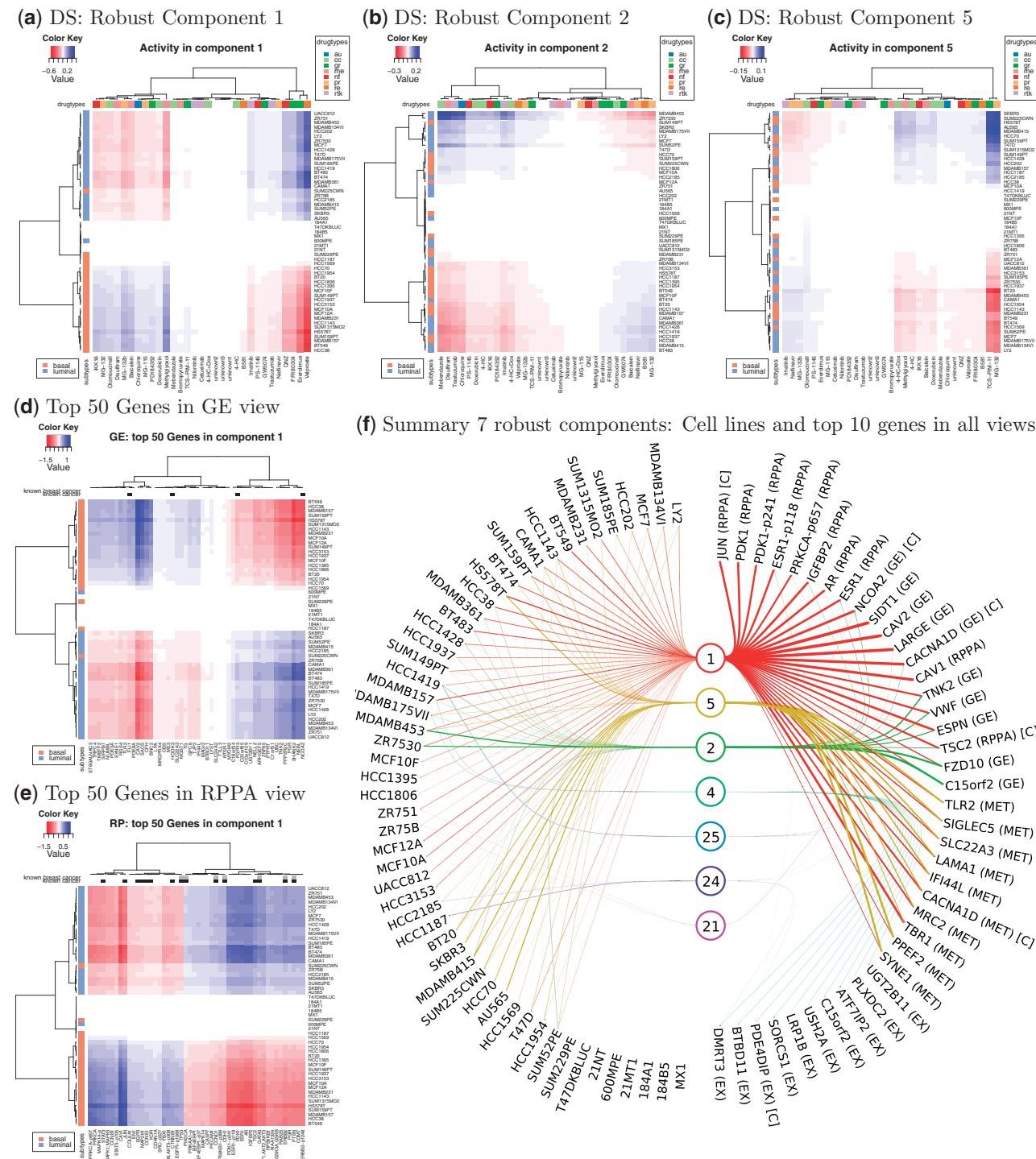


Fig. 3. Biocluster activity patterns of robust components. (a)–(c) The intensity values of cells (in rows) and drugs (in columns) in the drug sensitivity (DS) view of three different components. Component one mainly distinguishes basal and luminal cell lines, while (d) and (e) show the corresponding biocluster activity pattern of the 50 genes with highest mean absolute values in the RPPA and GE views, marking known (breast) cancer genes by a (gray) black square. (f) Biocluster participation of all cells (left) in 7 selected robust components (middle) with respect to their mean absolute intensity values (represented in the thickness of the lines) for the top 10 genes in each of the views (GE, MET, RPPA and EX). Known cancer genes are marked by [C] (Color version of this figure is available at *Bioinformatics* online.)

GE, MET, RPPA and EX we selected a list of the 50 most active genes from the dense robust components. For each of such gene sets we calculated the enrichment for all categories. The result table contains a list of shared GO terms for each gene set together with information about the background and sample frequency, fold enrichment and the *P*-value determined by a binomial statistic.

A *P*-value close to zero indicates the significance of the GO term associated with the provided group of genes.

More than one thousand shared GO terms are returned for the most active gene sets. We condensed the results showing only the most repeating and most significant ones by using a threshold for the *p*-value and showing only GO terms below 10^{-6} , which appear

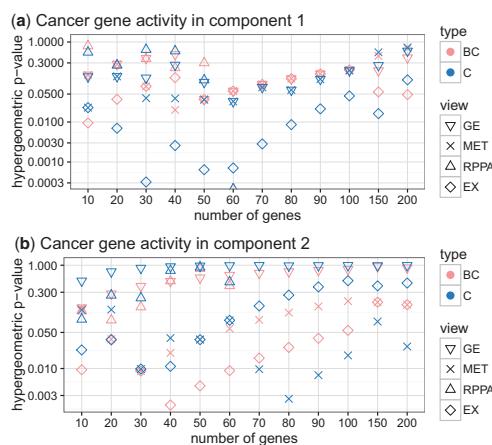


Fig. 4. Hypergeometric test of the activity of known breast cancer (BC) and all cancer (C) genes in the two most predictive components and all views. Low P -values indicate a high number of cancer genes in the top n active genes in comparison with randomly picked sets. See the text for details

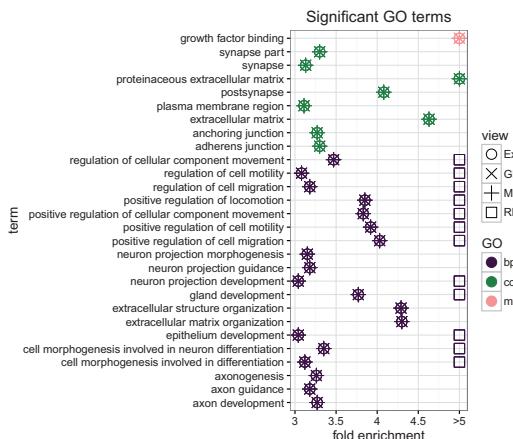


Fig. 5. Enrichment of the most significant GO terms, which occur more than 3 times in all the gene related views and the three dense robust components (Color version of this figure is available at *Bioinformatics* online.)

throughout the views and components more than 3 times. Figure 5 shows the reduced list of significant GO terms with fold enrichment value bigger than 3. This value indicates the magnitude of fold enrichment for the observed set of genes over the expected, thus with values bigger than one the category implying over-representation. For biological process we found most of the repeating significant GO terms in all of the views in nearly all cases. These GO terms are related to cell motility and its regulations.

5 Discussion

We presented sparse group factor analysis as a way of inferring biclusters from heterogeneous multi-source data. The method is able to detect predictive and interpretable structure present in any subset of the data sources, and sparse within the sources. It proved to be robust in this task, as witnessed by the simulation studies and the outstanding performance in the NCI-DREAM drug sensitivity prediction challenge. The biclusters of the joint data identified cancer cell subtypes, grouped drugs by their functional mechanisms, and associated known cancer genes with the drug sensitivity data, all in a data-driven fashion. The shown approach is suitable for exploratory analysis of multiple data sources, giving condensed and interpretable information with respect to the data collection.

In this paper we focused on formulating a model that implements the novel multi-data-source biclustering, and on evaluating the accuracy of the results. Two important questions we did not yet fully address are: (i) could some of the alternative ways of implementing sparsity, substituting the spike-and-slabs of this paper, result in computationally more efficient and still as accurate solutions. (ii) Computational speed. The EM point estimates the single-data-source FABIA algorithm uses would naturally be faster for multiple data sources as well, but the ability to handle uncertainty due to highly noisy and high-dimensional small sample-size data would suffer. Variational approximations would be attractive as they would also help avoid the matchings between the different sampling chains, but deriving variants of the algorithm would be more difficult and variational approximations are known to produce a biased estimate of the uncertainty of the solutions. For large data parallelized sampling solutions would be particularly attractive ways of speeding up computation.

Funding

We thank the Academy of Finland (Finnish Centre of Excellence in Computational Inference Research COIN) for funding.

Conflict of Interest: none declared.

References

- Carvalho,C.M. *et al.* (2008) High-dimensional sparse factor modeling: applications in gene expression genomics. *J. Am. Stat. Assoc.*, **103**, 1438–1456.
- Cheng,Y. and Church,G.M. (2000). Biclustering of expression data. In: *Proc. of the 8th International Conference on Intelligent Systems for Molecular Biology*, pp. 93–103. AAAI Press.
- Costello,J.C. *et al.* (2014) A community effort to assess and improve drug sensitivity prediction algorithms. *Nat. Biotechnol.*, **32**, 1202–1212.
- Gao,C. *et al.* (2014). Differential gene co-expression networks via Bayesian biclustering models. *arXiv preprint arXiv:1411.1997*.
- Hartigan,J.A. (1972) Direct clustering of a data matrix. *J. Am. Stat. Assoc.*, **67**, 123–129.
- Hochreiter,S. (2013) HapFABIA: identification of very short segments of identity by descent characterized by rare variants in large sequencing data. *Nucleic Acids Res.*, **41**, e202.
- Hochreiter,S. *et al.* (2010) FABIA: factor analysis for bicluster acquisition. *Bioinformatics*, **26**, 1520–1527.
- Khan,S.A. *et al.* (2014) Identification of structural features in chemicals associated with cancer drug response: a systematic data-driven analysis. *Bioinformatics*, **30**, i497–i504.
- Klami,A. *et al.* (2015) Group factor analysis. *IEEE Trans. Neural Netw. Learn. Syst.*, **26**, 2136–2147.
- Lazzeroni,L. *et al.* (2002) Plaid models for gene expression data. *Stat. Sin.*, **12**, 61–86.
- Madeira,S.C. and Oliveira,A.L. (2004) Biclustering algorithms for biological data analysis: A survey. *IEEE/ACM Trans. Comput. Biol. Bioinf.*, **1**, 24–45.
- Mi,H. *et al.* (2013) Large-scale gene function analysis with the PANTHER classification system. *Nat. Protoc.*, **8**, 1551–1566.
- Morgan,J.N. and Sonquist,J.A. (1963) Problems in the analysis of survey data, and a proposal. *J. Am. Stat. Assoc.*, **58**, 415–434.
- Stephens,P.J. *et al.* (2012) The landscape of cancer genes and mutational processes in breast cancer. *Nature*, **486**, 400–404.
- Suvitalo,T. *et al.* (2014) Cross-organism toxicogenomics with group factor analysis. *Syst. Biomed.*, **2**, 71–80.
- Virtanen,S. *et al.* (2012). Bayesian group factor analysis. In: Lawrence,N. and Girolami,M. (eds), *Proc. of the 15th International Conference on Artificial Intelligence and Statistics*, pp. 1269–1277.
- Waltman,P. *et al.* (2010) Multi-species integrative biclustering. *Genome Biol.*, **11**, R96.
- Yap,C.W. (2011) PaDEL-descriptor: an open source software to calculate molecular descriptors and fingerprints. *J. Comput. Chem.*, **32**, 1466–1474.