

Cross-organism prediction of drug hepatotoxicity by sparse group factor analysis

Tommi Suvitaival

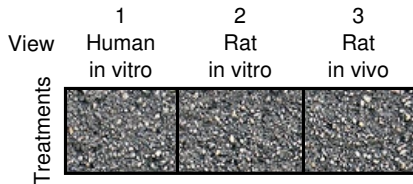
Juuso A. Parkkinen Seppo Virtanen Samuel Kaski



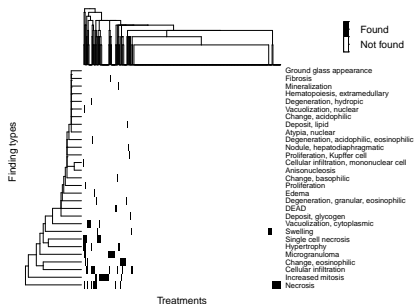
July 19-20, 2013 – CAMDA

Starting point

High-dimensional gene-expression data from 3 types of organisms

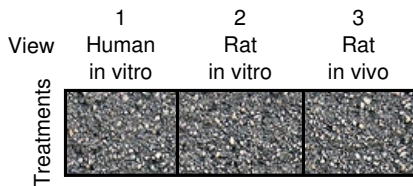


Sparse pathological data of rat *in vivo*

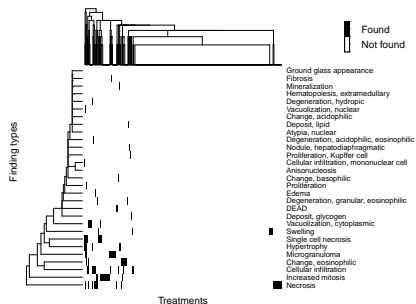


Starting point

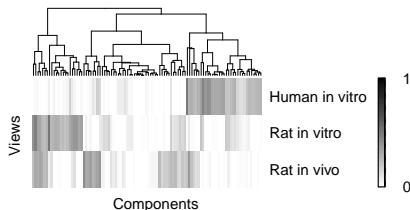
High-dimensional gene-expression data from 3 types of organisms



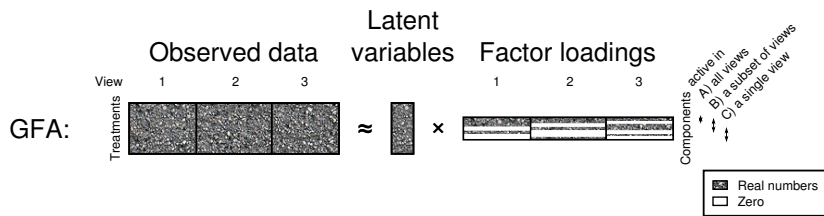
Sparse pathological data of rat *in vivo*



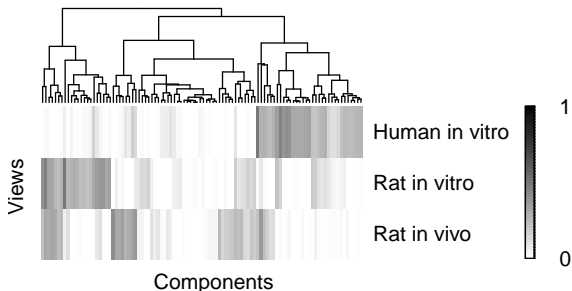
1. Can we replace the animal study with *in vitro* assay?
2. Can we predict the liver injury in humans using toxicogenomics data from animals?



Group factor analysis (GFA)



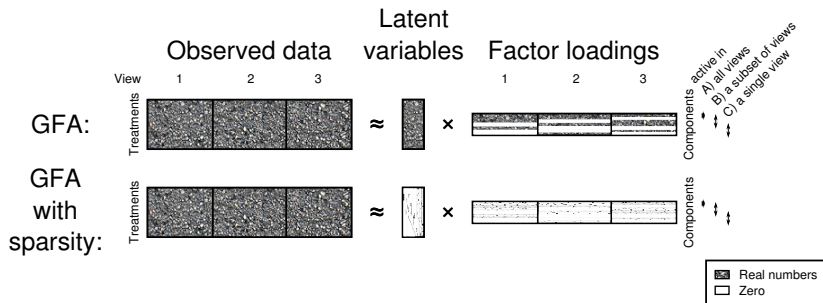
Making generalizations across organisms



Shared components

- ▶ associations between views
- ▶ cross-view prediction

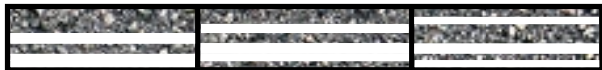
GFA with sparsity (1)



GFA with and without sparsity



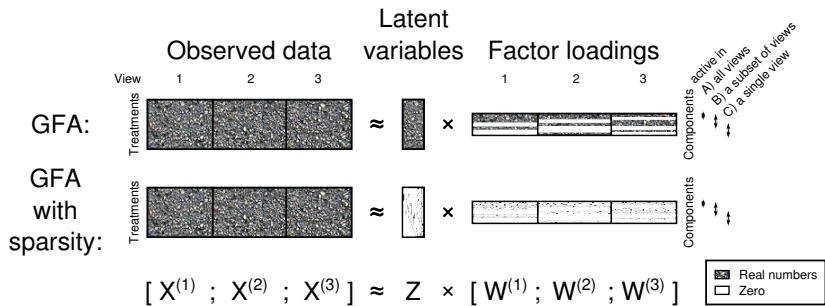
×



×



GFA with sparsity (2)



Sparsity – why

Sparsity in the model is encouraged due to

1. High dimensionality of the gene expression microarray data sets \Rightarrow Sparsity in terms of variables
2. Strong sparsity of the pathology data
3. Treatments heterogeneous by their effects \Rightarrow Sparsity in terms of samples

Sparsity – how

1. Sparsity in terms of variables

⇒

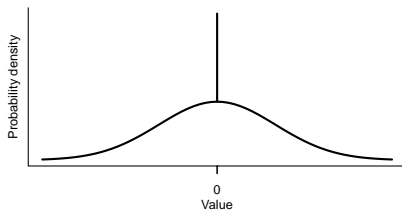
*Spike-and-slab prior** for factor loadings matrix \mathbf{W}

2. Sparsity in terms of samples

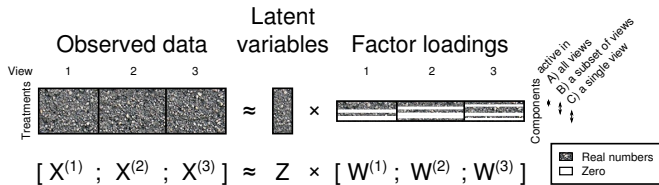
⇒

Spike-and-slab prior for latent variables \mathbf{Z}

*



GFA – model



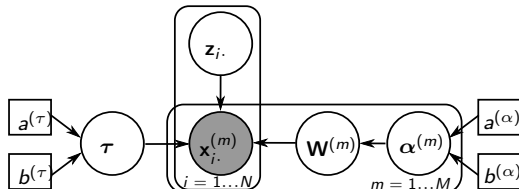
$$x_{i \cdot}^{(m)} \sim \mathcal{N}(z_{i \cdot} \mathbf{W}^{(m)}, \tau_m^{-1} \mathbf{I})$$

$$z_{i \cdot} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbf{w}_{k \cdot}^{(m)} \sim \mathcal{N}\left(\mathbf{0}, \frac{1}{\alpha_k^{(m)}} \mathbf{I}\right)$$

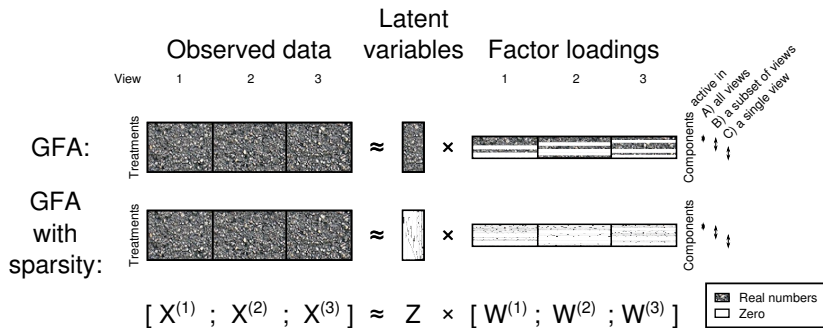
$$\alpha_k^{(m)} \sim \text{Gamma}(a^{(\alpha)}, b^{(\alpha)})$$

$$\tau_m \sim \text{Gamma}(a^{(\tau)}, b^{(\tau)})$$



i : samples, m : views

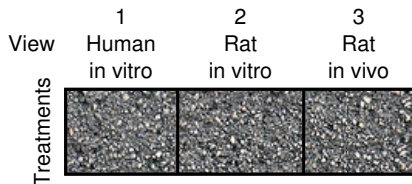
GFA with sparsity – model



GFA	GFA with sparsity
$\mathbf{x}_{i \cdot}^{(m)} \sim \mathcal{N}(\mathbf{z}_i \cdot \mathbf{W}^{(m)}, \tau_m^{-1} \mathbf{I})$	$\mathbf{x}_{i \cdot}^{(m)} \sim \mathcal{N}(\mathbf{z}_i \cdot \mathbf{W}^{(m)}, (\mathbf{\Lambda}^{(m)})^{-1})$
$\mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$	$\mathbf{z}_{ik} \sim \mathbf{H}_k^{(z)} \mathcal{N}\left(0, \frac{1}{\alpha_{ik}^{(z)}}\right) + (1 - \mathbf{H}_k^{(z)}) \delta_0$
$\mathbf{w}_{k \cdot}^{(m)} \sim \mathcal{N}\left(0, \frac{1}{\alpha_k^{(m)}} \mathbf{I}\right)$	$\mathbf{W}_{dk}^{(m)} \sim \mathbf{H}_{dk}^{(m)} \mathcal{N}\left(0, \frac{1}{\alpha_{dk}^{(m)}}\right) + (1 - \mathbf{H}_{dk}^{(m)}) \delta_0$

Data representation – gene expression

- ▶ Treatments that occur in all 3 types of organism:
 - ▶ 119 compounds
 - ▶ dosage levels *middle* & *high*
 - ▶ time points *8/9 h* & *24 h*
- ▶ Average differential expression over the replicates of each treatment
 - ⇒ Treatment = sample for the model
 - ⇒ Matching treatments between the 3 transcriptomic *views*
 $\mathbf{X}_{\text{in vitro}}^{\text{human}}$, $\mathbf{X}_{\text{in vitro}}^{\text{rat}}$ and $\mathbf{X}_{\text{in vivo}}^{\text{rat}}$

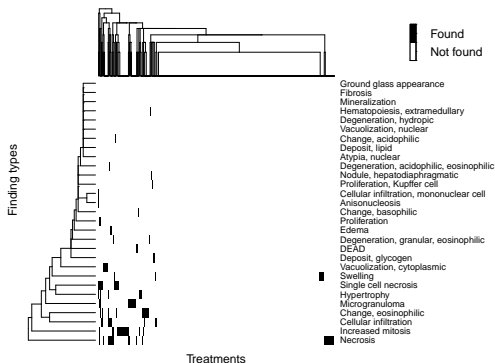


Data representation – histopathology of the liver

Grade-weighted count of each pathological finding type over the replicates of a treatment

⇒ Pathology view

$Y_{\text{in vivo}}^{\text{rat}}$ with matching treatments to the 3 transcriptomic views



Results

Our tasks:

1. Predict liver damage of rats *in vivo* based on cell-level transcriptomic responses in the 3 types of model organisms
2. Test how well the transcriptomic cell-level responses generalize to known effects of the compounds on humans

Analysis: model organisms' generalizability to organ level

Training: Learn associations between the views

- ▶ 3 transcriptomic views $\mathbf{X}_{\text{in vitro}}^{\text{human}}$, $\mathbf{X}_{\text{in vitro}}^{\text{rat}}$ and $\mathbf{X}_{\text{in vivo}}^{\text{rat}}$
- ▶ Pathology view $\mathbf{Y}_{\text{in vivo}}^{\text{rat}}$

Testing: Predict the pathological findings $\mathbf{Y}_{\text{in vivo}}^{\text{rat}}$

- ▶ Given one of the transcriptomic views

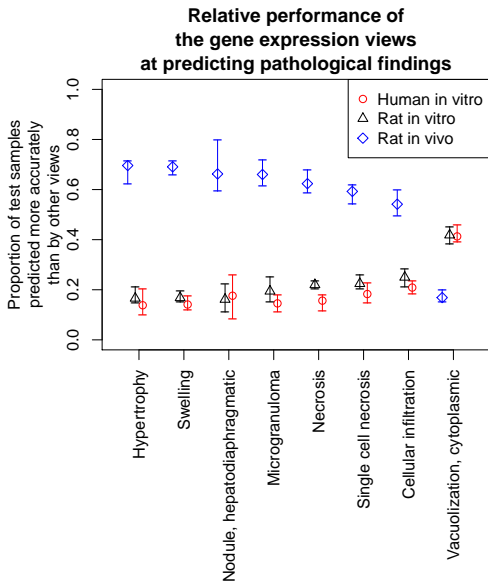
Analysis: model organisms' generalizability to organ level

Training: Learn associations between the views

- ▶ 3 transcriptomic views $\mathbf{X}_{\text{in vitro}}^{\text{human}}$, $\mathbf{X}_{\text{in vitro}}^{\text{rat}}$ and $\mathbf{X}_{\text{in vivo}}^{\text{rat}}$
- ▶ Pathology view $\mathbf{Y}_{\text{in vivo}}^{\text{rat}}$

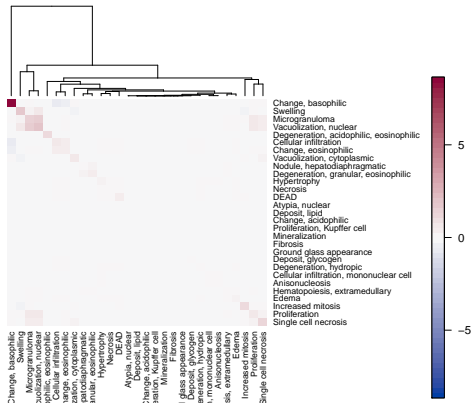
Testing: Predict the pathological findings $\mathbf{Y}_{\text{in vivo}}^{\text{rat}}$

- ▶ Given one of the transcriptomic views



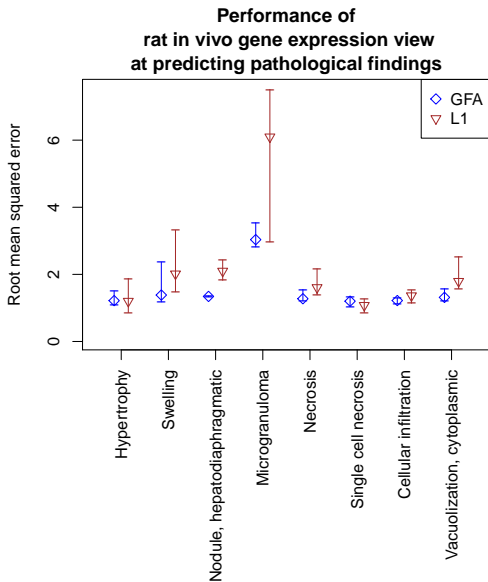
Sparsity in the target view

- ▶ $\mathbf{W}^T\mathbf{W}$ reveals the similarity of component activities between the variables
- ▶ Thanks to sparsity, projections to many variables are 0
- ▶ The model automatically decides which variables to explain by
 - A. coherent components
 - B. noise parameter



Prediction: drug hepatotoxicity based on gene expression

- ▶ Given $\mathbf{X}_{\text{in vivo}}^{\text{rat}}$,
predict $\mathbf{Y}_{\text{in vivo}}^{\text{rat}}$
- ▶ Same prediction task
using ℓ_1 -regularized
linear regression



Translation over model organisms to humans

- ▶ How do the transcriptional changes in model organisms generalize system-level effects in humans?
- ▶ Can the model learn structure relevant to the properties of the compounds in an unsupervised way?

Translation over model organisms to humans (1)

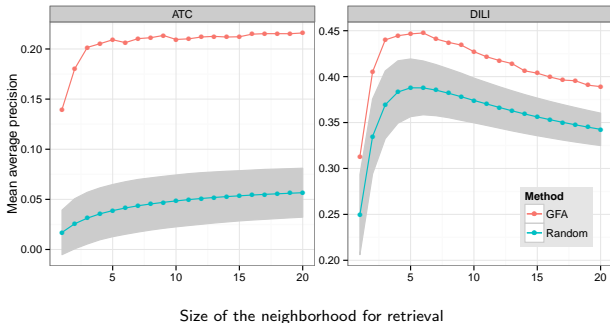
We quantify the success of translation by the retrieval of similar compounds

- ▶ Ground-truth:
 - A. Anatomical Therapeutic Chemical (ATC) Classification System's labels (level 4)
 - B. Drug-induced liver injury (DILI) labels
- ▶ Model: GFA with sparsity for the transcriptomic views of the model organisms, $\mathbf{X}_{\text{in vitro}}^{\text{human}}$, $\mathbf{X}_{\text{in vitro}}^{\text{rat}}$ and $\mathbf{X}_{\text{in vivo}}^{\text{rat}}$
- ▶ Measure: Average precision in the retrieval of similar compounds in the latent space

Translation over model organisms to humans (2)

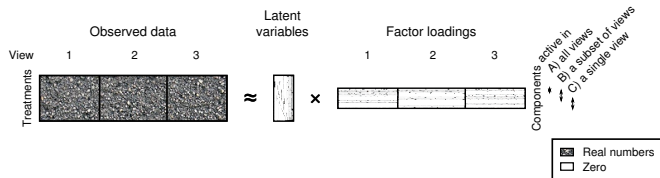
We quantify the success of translation by the retrieval of similar compounds

- ▶ Ground-truth:
 - A. Anatomical Therapeutic Chemical (ATC) Classification System's labels (level 4)
 - B. Drug-induced liver injury (DILI) labels
- ▶ Model: GFA with sparsity for the transcriptomic views of the model organisms, $\mathbf{X}_{\text{in vitro}}^{\text{human}}$, $\mathbf{X}_{\text{in vitro}}^{\text{rat}}$ and $\mathbf{X}_{\text{in vivo}}^{\text{rat}}$
- ▶ Measure: Average precision in the retrieval of similar compounds in the latent space



Conclusions

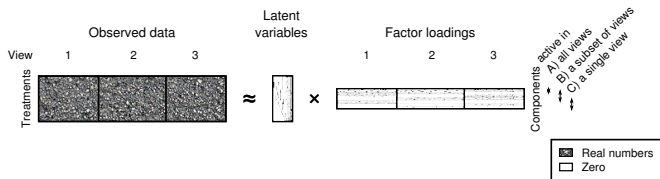
- ▶ GFA reveals associations between the views
- ▶ Associations indicate what generalizes between the views
- ▶ Sparsity helps in this decision
- ▶ Latent representation allows us to explore structure in the data in an unsupervised way



Discussion

We can

- ▶ analyse the similarity of model organisms
- ▶ learn what generalizes from the model organisms to humans



Funding:

- ▶ The Academy of Finland
 - ▶ Finnish Centre of Excellence in Computational Inference Research COIN, 251170
 - ▶ Computational Modeling of the Biological Effects of Chemicals, 140057
- ▶ Finnish Doctoral Programme in Computational Sciences FICS
- ▶ Helsinki Doctoral Programme in Computer Science