# BAYESIAN POISSON FACTORIZATION FOR GENETIC ASSOCIATIONS WITH CLINICAL FEATURES IN CANCER
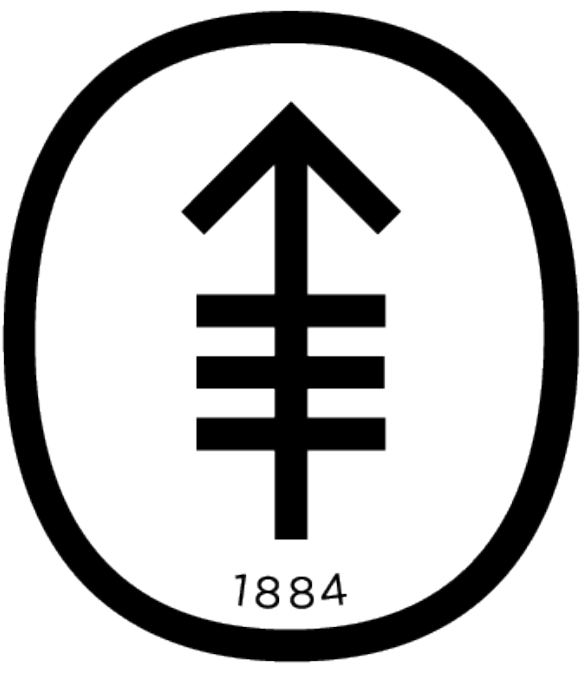
M. F. PRADIER[1], T. KARALETSOS[3], S. STARK[3], J.E. VOGT[3], F. PEREZ-CRUZ[1,2], AND G. RÄTSCH[3]
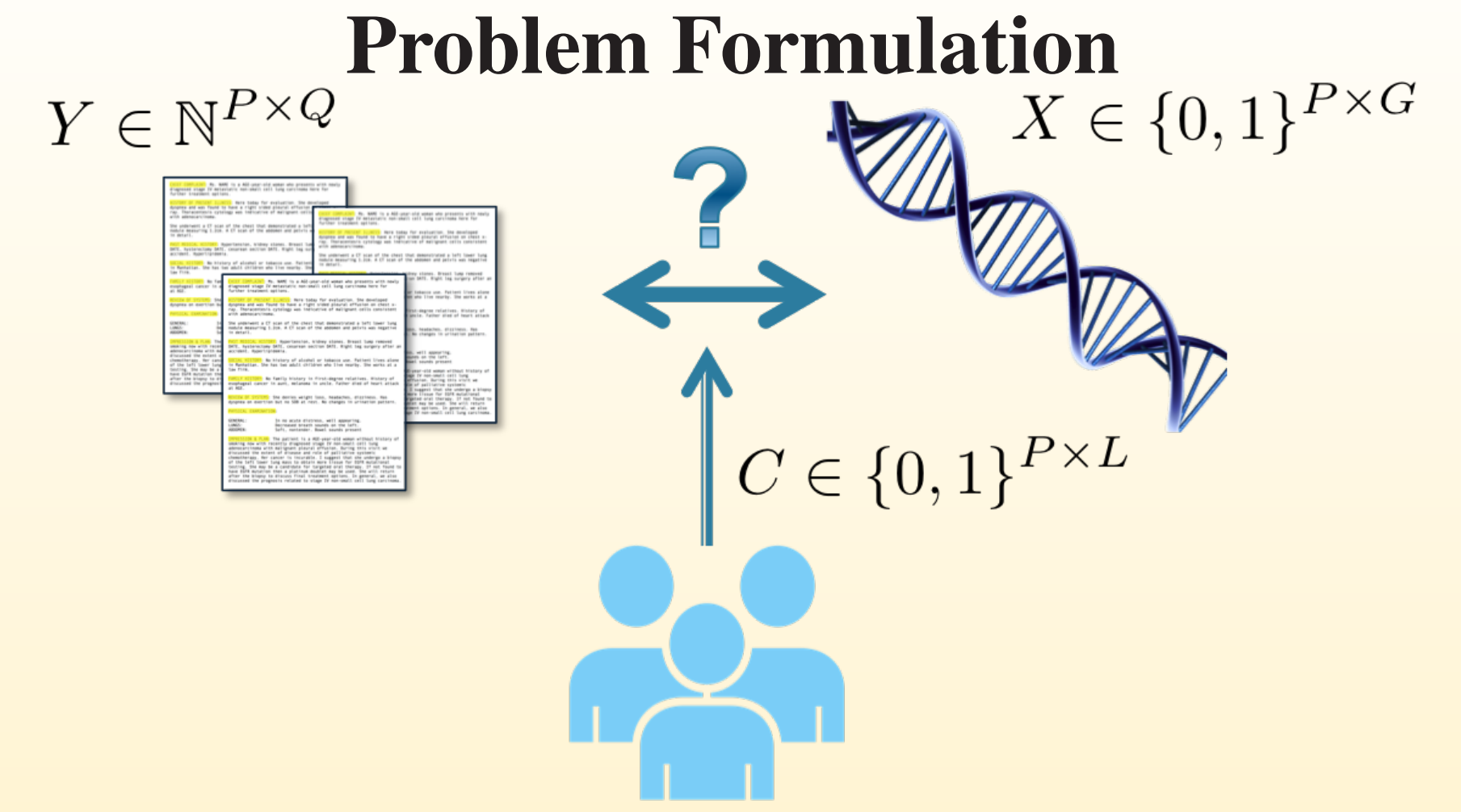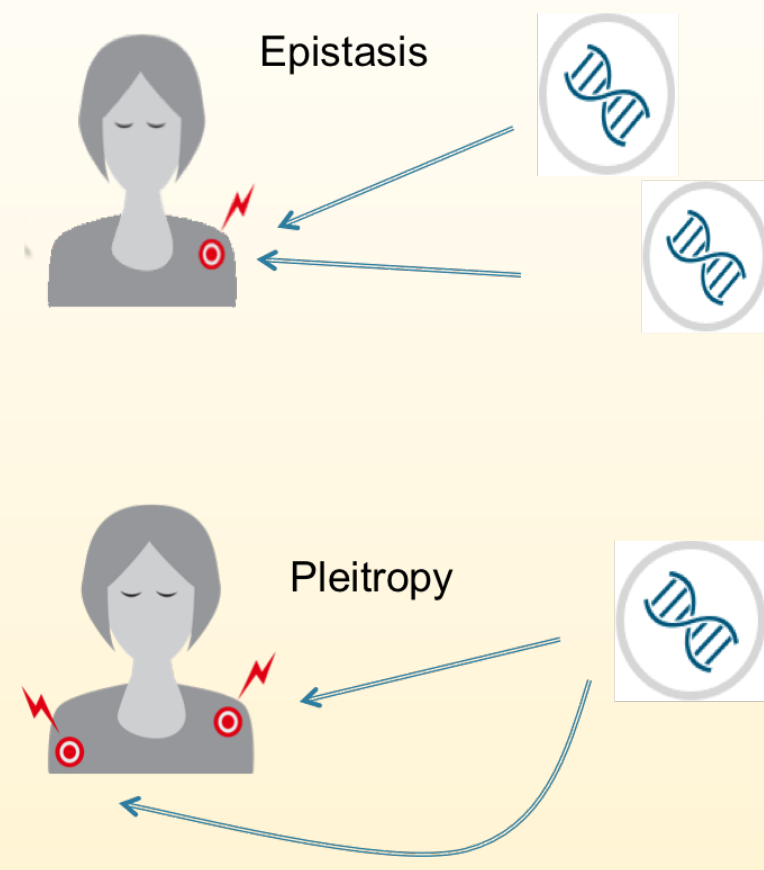
1. University Carlos III in Madrid, Spain, 2. Bell Labs, Alcatel-Lucent, New Jersey, USA and 3. Memorial Sloan-Kettering Cancer Center, New York, USA

{melanie,fernando}@tsc.uc3m.es    {karaletsos,starks,vogt,raetsch}@cbio.mskcc.org

## INTRODUCTION

### Motivation

- Cancer ≡ set of complex genetic diseases not very well understood yet.

- Our aim: Exploratory Analysis. Get meaninful genotype-phenotype associations by looking at Electronic Health Records (EHR) and genomic data.

- We would like our model to account for:
  - Confounders
  - Pleiotropy
  - Epistasis

### Problem Formulation

$Y \in \mathbb{N}^{P \times Q}$  $X \in \{0,1\}^{P \times G}$

$C \in \{0,1\}^{P \times L}$



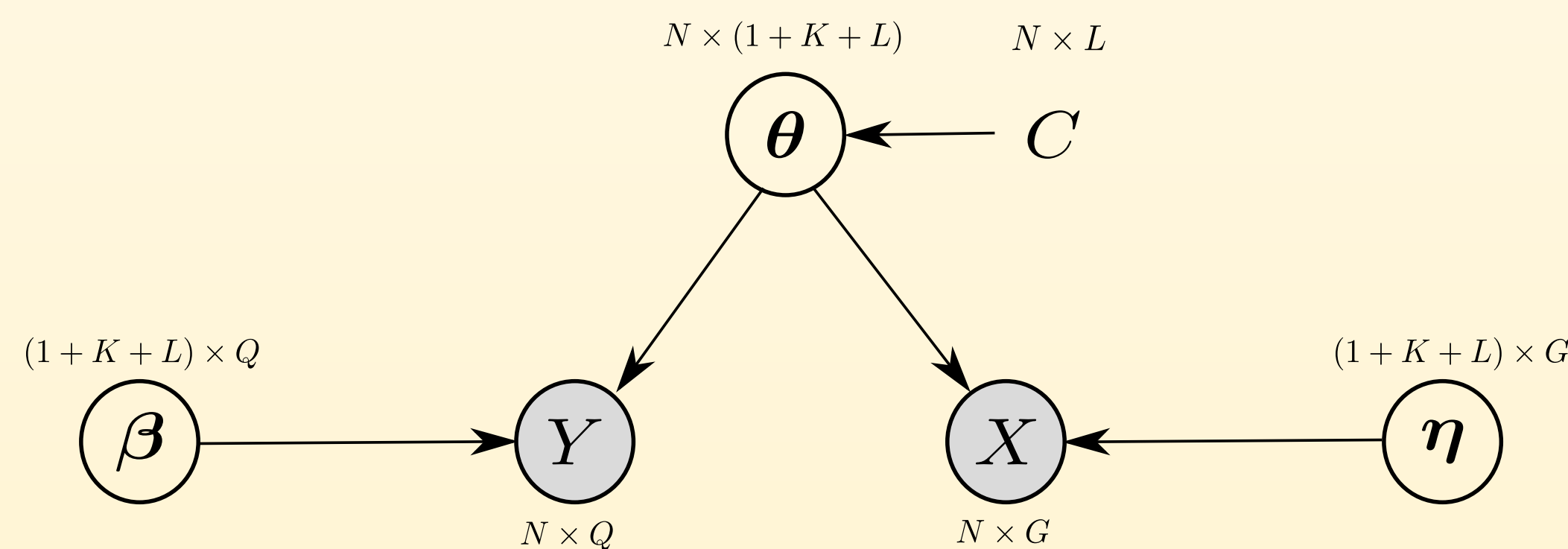## CONFOUNDER-CORRECTED POISSON FACTORIZATION MODEL

- We propose a model that directly finds associations between somatic mutations and clinical features through hidden factors. We call it Bayesian Confounder-corrected Poisson Factorization model (CC-PFM).

- Based on a PFM model for recommendation systems [1].

- Notation:
  - $N$ patients, $G$ mutations, $Q$ words, $L$ cancer types and $K$ latent topics.
  - $\beta$ and $\eta$ are the clinical and genetic factors respectively.
  - $\theta$ is the weight matrix: it captures the presence/activation of each factor per patient.

$$y_{nq}|\boldsymbol{\theta_{n.}}, \boldsymbol{\beta_{.q}} \sim \text{Poisson}(\beta_{0q} + \boldsymbol{\theta'}_{n.}\boldsymbol{\beta'_{.q}} + \boldsymbol{\theta''}_{n.}\boldsymbol{\beta''_{.q}}) \quad (1)$$
$$x_{ng}|\boldsymbol{\theta_{n.}}, \boldsymbol{\eta_{.g}} \sim \text{Poisson}(\eta_{0g} + \boldsymbol{\theta'}_{n.}\boldsymbol{\eta'_{.g}} + \boldsymbol{\theta''}_{n.}\boldsymbol{\eta''_{.g}}) \quad (2)$$
$$\theta_{nr} \sim \text{Gamma}(a,b), \quad \beta_{rq} \sim \text{Gamma}(c,d), \quad \eta_{rg} \sim \text{Gamma}(e,f), \quad (3)$$

where $\boldsymbol{\theta} = [\mathbf{1}_N, \boldsymbol{\theta'}_{N\times K}, \boldsymbol{\theta''}_{N\times L} \odot C_{N\times L}]$ is an $N \times (1+K+L)$ matrix, $\mathbf{1}_N$ is a column vector of ones for the bias term, $\odot$ is the Hadamard product, $\boldsymbol{\beta} = [\beta_0; \boldsymbol{\beta'}_{K\times Q}; \boldsymbol{\beta''}_{L\times Q}]$ is a $(1+K+L) \times Q$ matrix, and $\boldsymbol{\eta} = [\eta_0; \boldsymbol{\eta'}_{K\times G}; \boldsymbol{\eta''}_{L\times G}]$ is a $(1+K+L) \times Q$ matrix of genetic factors.



We can also write the likelihood as:

$$y_{nq}|\boldsymbol{\theta_{n.}}, \boldsymbol{\beta_{.q}} \sim \text{Poisson}(\boldsymbol{\theta_{n.}}\boldsymbol{\beta_{.q}}) \quad (4)$$
$$x_{ng}|\boldsymbol{\theta_{n.}}, \boldsymbol{\eta_{.g}} \sim \text{Poisson}(\boldsymbol{\theta_{n.}}\boldsymbol{\eta_{.g}}) \quad (5)$$

- We use mean-field variational inference for learning.
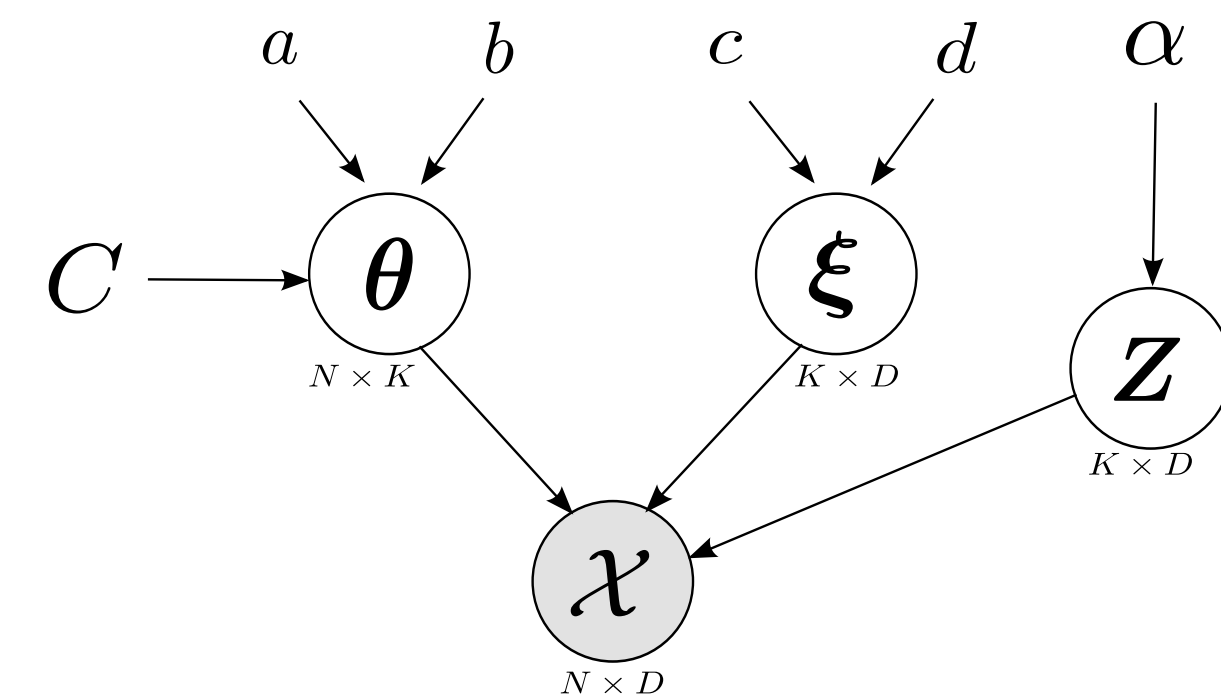
## NON-PARAMETRIC SPARSE CC-PFM

- We introduce 2 binary matrices $Z_\beta$ and $Z_\eta$ that work as a mask on the factor matrices. In other words, we replace the Gamma priors by Spike-and-slab priors. The likelihoods are now:
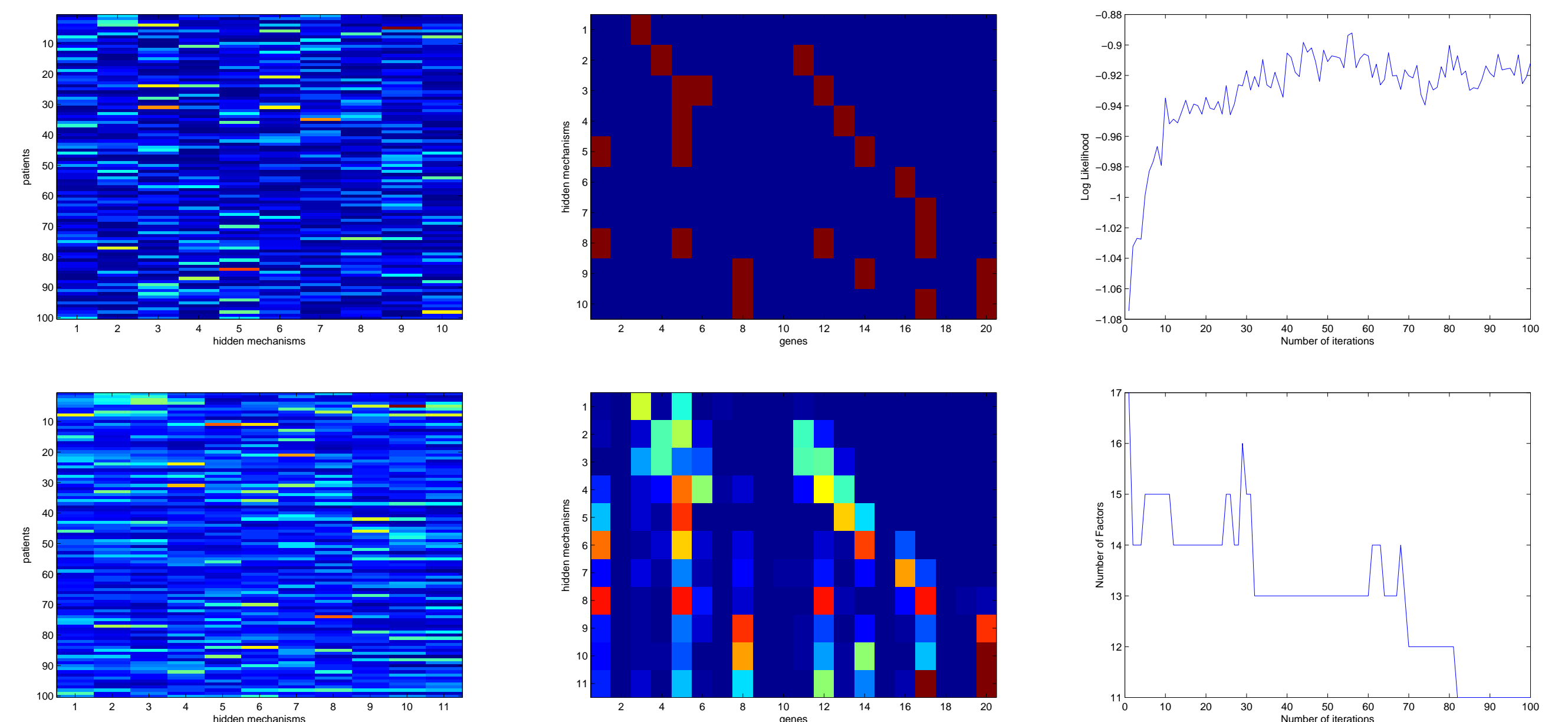
$$y_{nq}|\boldsymbol{\theta_{n.}}, \boldsymbol{\beta_{.q}} \sim \text{Poisson}\big(\boldsymbol{\theta_{n.}}(Z_\beta \odot \boldsymbol{\beta_{.q}})\big) \quad (6)$$
$$x_{ng}|\boldsymbol{\theta_{n.}}, \boldsymbol{\eta_{.g}} \sim \text{Poisson}\big(\boldsymbol{\theta_{n.}}(Z_\eta \odot \boldsymbol{\eta_{.g}})\big) \quad (7)$$

- We use an Indian Buffet Process prior on the matrices $Z_\beta$ and $Z_\eta$ to make $K \to \infty$. That is, $Z \sim \text{IBP}(\alpha)$ where $Z = [Z_\eta, Z_\beta]^\text{T}$. This model extends the model in [2].

- Simpler notation: $\boldsymbol{\xi} = [\boldsymbol{\beta}, \boldsymbol{\eta}]$ and $\mathcal{X} = [Y, X]$.

- Inference: Gibbs sampling + slice sampling for matrix $Z$ [3].



**The Idea:** We show $\theta$ and $Z$, ground truth (top) and inferred (bottom).



## RESULTS

| | Textual topics $\beta'_{k.}$ | Genetic topics $\eta'_{k.}$ |
|---|---|---|
| Factor 0 (Bias) | demonstrated, oncologist, suv, died, involvement | TP53 |
| Factor 1 | adenocarcinoma, pleural, woman, smoker | PIK3CA, RB1 |
| Factor 2 | pelvic, female, woman, endometrial, vaginal | NRAS |
| Factor 3 | cisplatin, squamous, icterus, kg, exertion | SPEN |
| Factor 4 | m, icterus, colon, fluid, ascites, cavity, hepatomegaly | KRAS |
| Factor 5 | folfox, colorectal, anc, colon, oxaliplatin | APC, KRAS, CIC |
| Factor 6 | brain, hemangiopericytoma, female, parietal | FUBP1, AXL |
| Factor 7 | breast, woman, adjuvant, female, mastectomy | PIK3CA, CDH1 |

*Free Associations*

| | Textual topics $\beta''_{l.}$ | Genetic topics $\eta''_{l.}$ |
|---|---|---|
| Appendiceal | mucinous, debulking, intraperitoneal, appendectomy | KRAS, GNAS |
| Bladder | bladder, urothelial, gemcitabine, invasive, cisplatin | TERT, KDM6A |
| Breast Carcinoma | breast, mastectomy, invasive, husband, female | PIK3CA, GATA3 |
| Melanoma | melanoma, ipilimumab, database, toe, temozolomide | MYCN, RAD51 |
| Small Cell Lung | etoposide, smoker, cisplatin, reassessment, irinotecan | RB1 |
| Soft Tissue | sarcoma, gentleman, adjuvant, ifosfamide, c, adriamycin | MYOD1 |

*Cancer Specific*

### Future work

- *Sparsity*: IBP ties the total number and per-row number of hidden features. This property is undesired a priori, and should be relaxed.

- *Validation*: What is the statistical significance of the inferred factors? Direct testing using statistical tests for stratified categorical data [4].

- *Scalability*: Inference using Stochastic Variational Inference.

- *Flexibility*: Introduce better conditional dependence on the confounders.

## FUNDING

## BASIC REFERENCES

1  P. Gopalan and D. Blei, "Content-based recommendations with Poisson factorization", presented at the Advances in Neural Information Processing Systems 27, 2014.

2  S. K. Gupta, D. Phung, and S. Venkatesh, "A nonparametric Bayesian Poisson gamma model for count data," in 2012 21st International Conference on Pattern Recognition (ICPR), 2012, pp. 1815–1818.

3  Y. W. Teh, "Stick-breaking construction for the Indian buffet process," in In Proceedings of the International Conference on Artificial Intelligence and Statistics, p. 2007.

4  K. Fukumizu, A. Gretton, X. Sun, and B. Schölkopf, "Kernel measures of conditional dependence," in In Adv. NIPS, 2008.