

# PROBABILISTIC ANALYSIS OF GENETIC ASSOCIATIONS WITH CLINICAL FEATURES IN CANCER

MELANIE F. PRADIER<sup>1,2</sup>, JULIA E. VOGT<sup>2</sup>, STEFAN STARK<sup>2</sup>,  
THEOFANIS KARALETSOS<sup>2</sup>, FERNANDO PEREZ-CRUZ<sup>1</sup> AND GUNNAR RÄTSCH<sup>2</sup>

1. University Carlos III in Madrid, Spain and 2. Memorial Sloan-Kettering Cancer Center, New York, USA

Email to: melanie@tsc.uc3m.es



## MOTIVATION

**Objective:** Finding meaningful relationships between genetic information and clinical features.



**Why does it matter?**

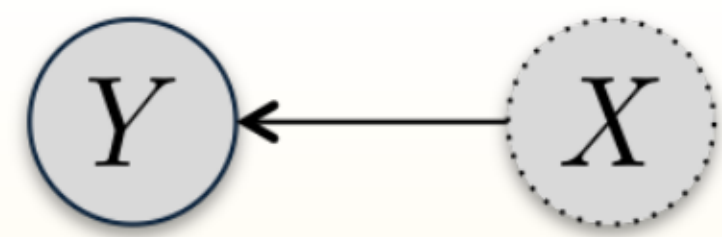
- Improved Diagnosis
- Risk Identification
- Biological Insight

**Our models should be able to deal with...**

- Epistasis
- Pleiotropy
- Confounders



## LINEAR MIXED MODEL (STATE-OF-THE-ART)



$$y = X\beta + u + \epsilon \quad (1)$$

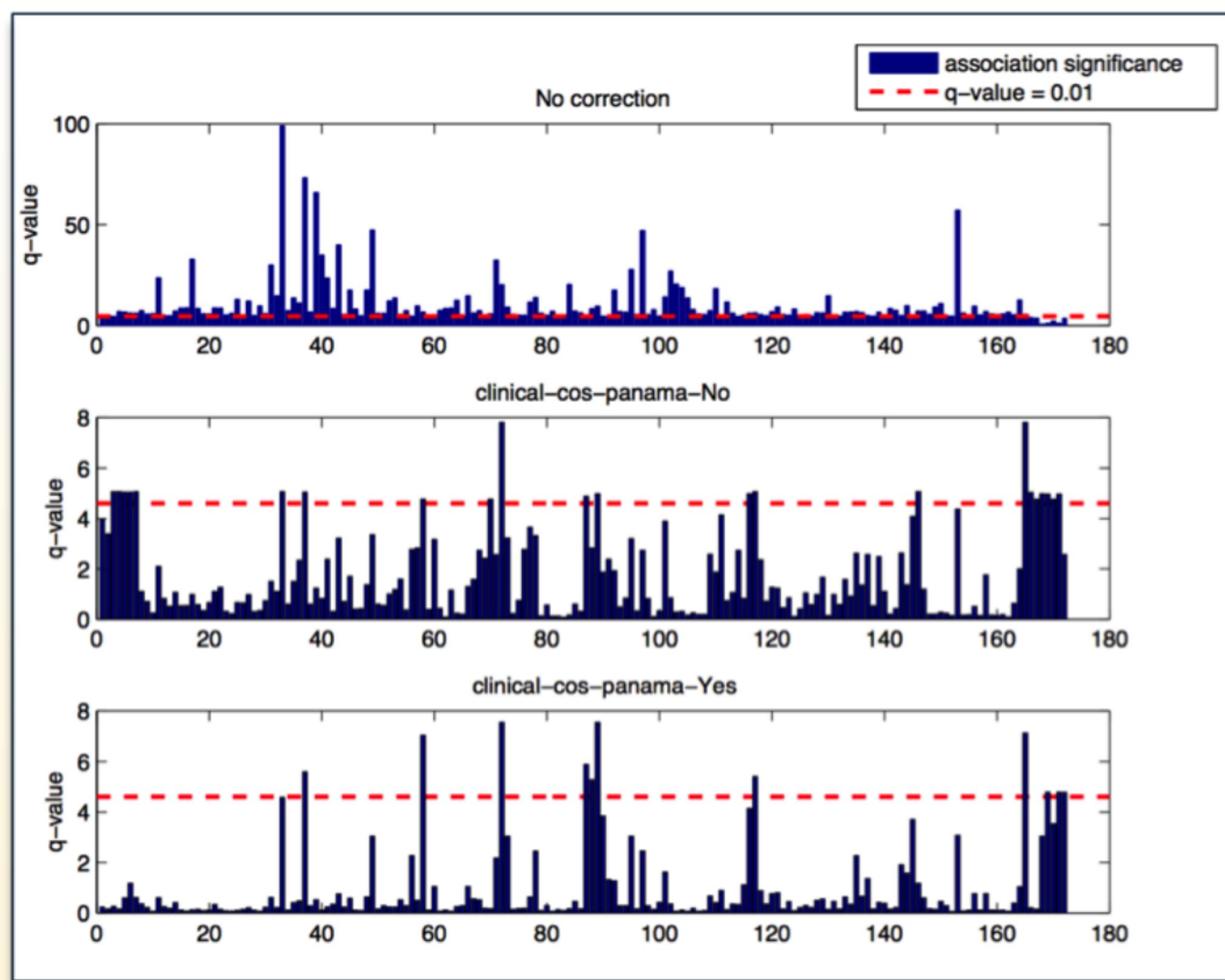
where

$$u \sim \mathcal{N}(0, K) \quad (2)$$

$$\epsilon \sim \mathcal{N}(0, \sigma_e^2 I) \quad (3)$$

- medical observation  $y \in \mathbb{R}^{P \times 1}$
- genetic variation matrix  $X \in \{0, 1\}^{P \times G}$
- fixed effects  $\beta \in \mathbb{R}^{G \times 1}$ , random effects  $u \in \mathbb{R}^{P \times 1}$
- $K$ : covariance matrix (e.g.  $K = CC^T$ )
- Classical model in statistical genetics [1]
- Rank Transformation to make  $Y$  observations Gaussian
- Inference of hidden confounders  $W$  using PANAMA [2], resulting in  $y = X\beta + u + \epsilon + W\gamma$

## RESULTS LMM - MOST CONSERVATIVE



10 significant hits with clinical covariates, 165 without correction. Examples:

Gene	MAF	q-value	$\beta$	Sentence prototype
APC	112	0.0037	0.33	He underwent a colonoscopy which revealed a pedunculated polyp in the ascending colon.
ALK	40	0.0063	0.57	The patient showed a mild decrease in her blood counts.
HNF1A	13	0.0028	0.70	The patient was tearful presented with depressed affect and mood.
TRAF7	11	0.0008	0.59	He has a history of adenoid cystic carcinoma of the salivary gland.

## CONCLUSION AND FUTURE WORK

- LMM is easy to interpret and gives False Discovery Rate estimate
- PFM can deal with epistasis naturally and finds better clinical feature representation
- Sparsity of topics
- Conditional model of  $\theta$  given  $C$
- Non-parametric extension

## FUNDING

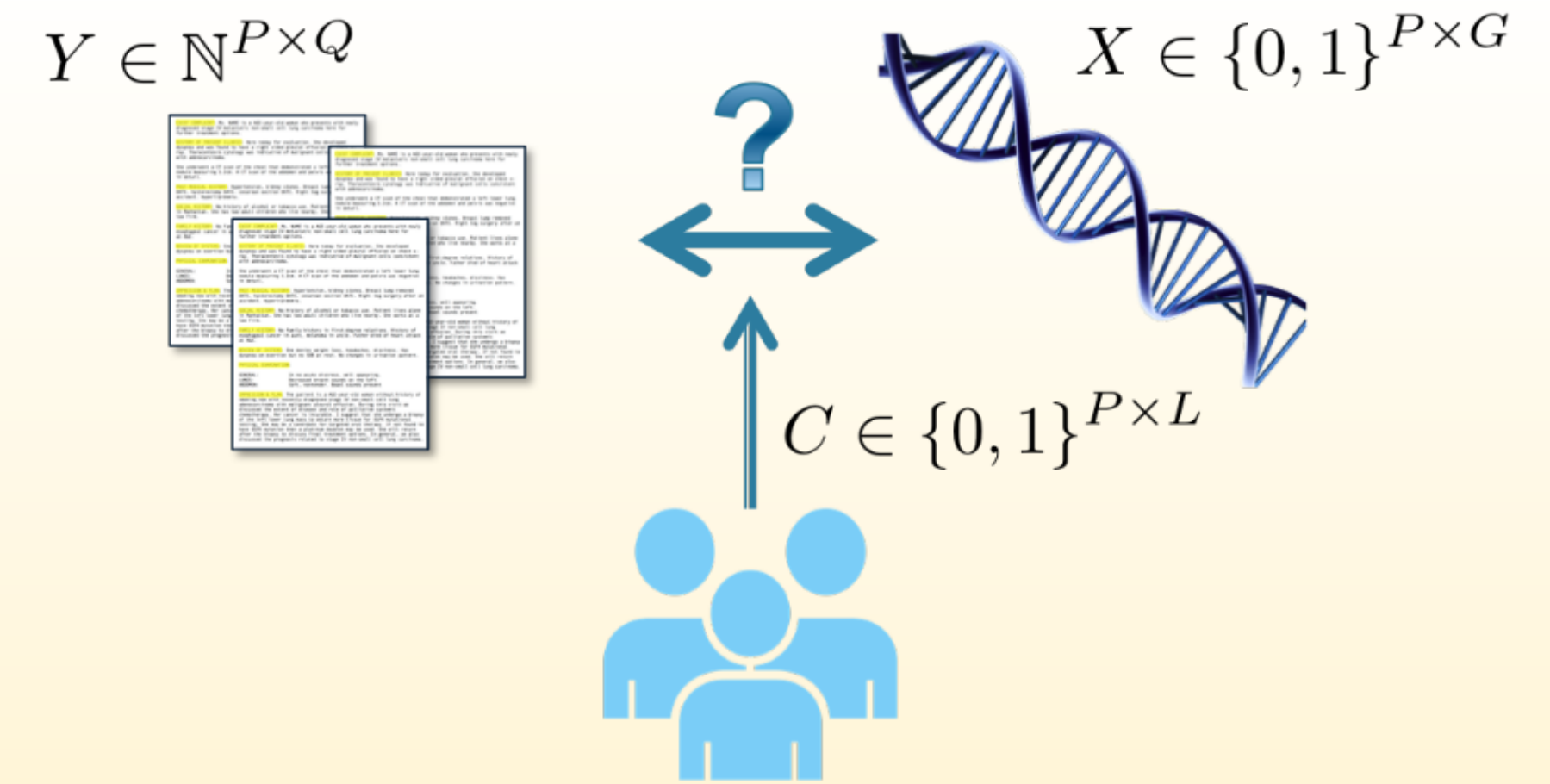
This project has received funding from the European Union's Seventh Framework Programme for research, technological development and demonstration under grant agreement no. 316861. This work is based on unpublished MSKCC data.



## PROBLEM SETTING

### Notation

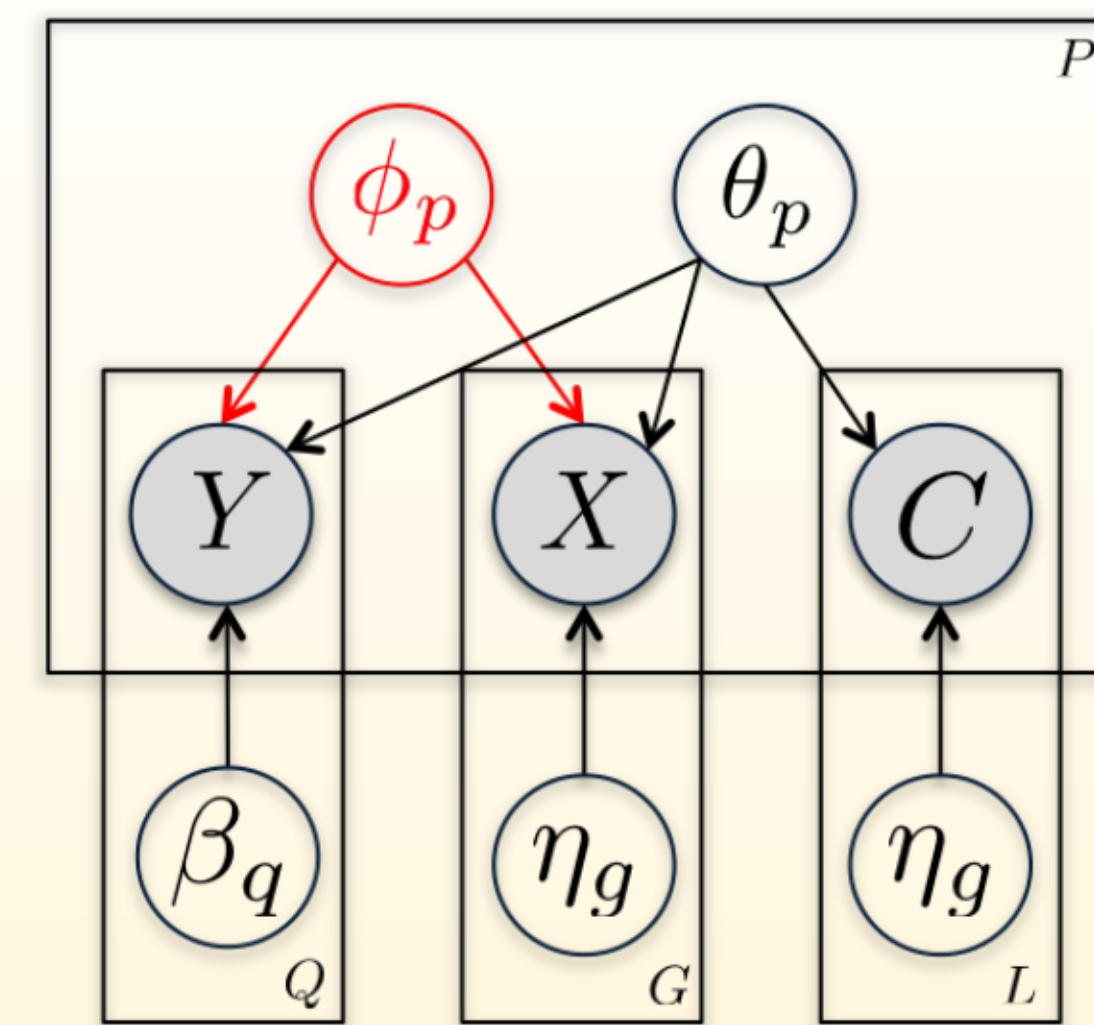
- $P$  patients
- $G$  genes
- $Q$  features
- $L$  covariates



Features are medical observations extracted from the Medical Health Records. This work considers sentence clusters and word topics.

## POISSON FACTORIZATION MODEL

The idea is to build a joint generative model of the clinical text and genetic information.



- Medical factors  $\beta_q$
- Genetic factors  $\eta_g$
- New interesting associations captured by  $\phi$ , confounders by  $\theta$
- Based on [3], variational inference

$$x_{pg} \sim \text{Poisson}(\phi_p \eta'_g + \theta_p \eta_g) \quad (4)$$

$$y_{pq} \sim \text{Poisson}(\phi_p \beta'_q + \theta_p \beta_q) \quad (5)$$

$$c_{pl} \sim \text{Poisson}(\theta_p \zeta_q) \quad (6)$$

$$\text{rest} \sim \text{Gamma}(a, b) \quad (7)$$

## RESULTS PFM - PRELIMINARY

- $P = 1265$  patients,  $\sim 15k$  documents,  $G = 343$  genes,  $K = 20$  factors,  $Q = 1k$  words out of 30k selected with highest it-idf index.
- Associations in grey generated by  $\theta$ , in color by  $\phi$ .
- Associations with genes in bold have a q-value  $< 5.10^{-4}$ .

bladder b hepatitis unilateral cisplatin gemcitabine invasive muscle neoadjuvant cystoscopy	0.068504 0.042724 0.029279 0.028239 0.018153 0.016724 0.016155 0.014433 0.014266 0.013865
female mastectomy pap invasive woman ductal er tamoxifen smear mammogram	0.029603 0.028461 0.022585 0.019685 0.017589 0.017533 0.017472 0.017209 0.017020 0.016214
rectal colorectal adenocarcinoma folfox anal cea rectum colon scorer bevacizumab	0.043272 0.031529 0.030553 0.025718 0.021779 0.013466 0.013129 0.012224 0.011782 0.010074

hepatocellular ajcc sorafenib protocol flow murmur splenomegaly peril touch displaced	0.036975 0.029099 0.026971 0.020808 0.017106 0.016528 0.016478 0.016387 0.015962 0.015881
prostate psa gleason prostatectomy adenocarcinoma androgen protocol lupron urinary radical	0.119886 0.062841 0.030116 0.018958 0.017361 0.012359 0.011871 0.011870 0.011405 0.011151
neuroendocrine pleasant gentleman issues spinal informed arise embolization cva octreotide	0.037745 0.026469 0.023654 0.016134 0.015970 0.015514 0.015472 0.015224 0.014779 0.014425

- A False Discovery Rate estimate is obtained by plugging the inferred word topics to the LMM as new input matrix  $Y$ .

## BASIC REFERENCES

1. C. Lippert, F. P. Casale, B. Rakitsch, and O. Stegle, "LIMIX: genetic analysis of multiple traits", bioRxiv, 2014.
2. N. Fusi, O. Stegle, and N. D. Lawrence, "Joint Modelling of Confounding Factors and Prominent Genetic Regulators Provides Increased Accuracy in Genetical Genomics Studies", PLoS Comput Biol, vol. 8, no. 1, p. e1002330, Jan. 2012.
3. P. Gopalan and D. Blei, "Content-based recommendations with Poisson factorization", presented at the Advances in Neural Information Processing Systems 27, 2014.