

# Kernelized Bayesian Matrix Factorization (KBMF)

Mehmet Gönen<sup>1</sup>, Muhammad Ammad-ud-din<sup>1</sup>, Suleiman A. Khan<sup>1</sup>, Samuel Kaski<sup>1,2</sup>

<sup>1</sup> Helsinki Institute for Information Technology HIIT, Department of Information and Computer Science, Aalto University, Finland

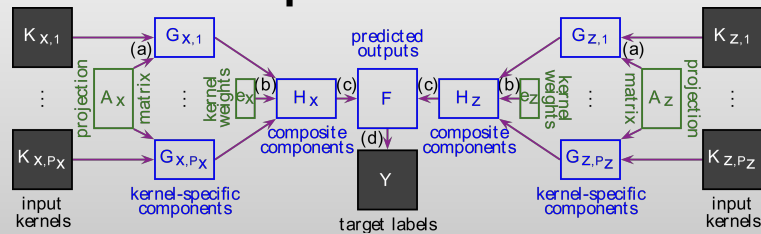
<sup>2</sup> Helsinki Institute for Information Technology HIIT, Department of Computer Science, University of Helsinki, Finland

<http://research.ics.aalto.fi/mi/>

## Abstract

- We extend **kernelized matrix factorization**
  - > with a fully Bayesian treatment,
  - > with an ability to work with multiple side information sources.
- Side information is necessary for making out-of-matrix predictions (e.g., cold-start predictions in recommender systems).
- We mainly discuss bipartite graph inference, where the output matrix is binary.
- We show the performance of our method
  - > by predicting drug–protein interactions on two data sets.

## Proposed Method



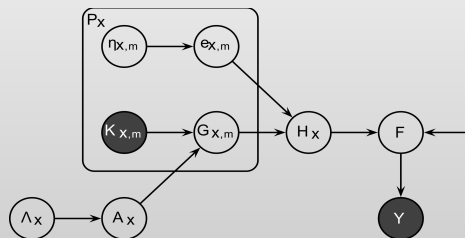
(a) Kernel based non-linear dimensionality reduction

(b) Multiple kernel learning

(c) Matrix factorization

(d) Binary classification (if data is binary)

## Probabilistic Model



(a) kernel-based nonlinear dimensionality reduction

$$\begin{aligned} \lambda_{x,s}^i &\sim \mathcal{G}(\lambda_{x,s}^i; \alpha_x, \beta_x) & \forall(i, s) \\ a_{x,s}^i | \lambda_{x,s}^i &\sim \mathcal{N}(a_{x,s}^i; 0, (\lambda_{x,s}^i)^{-1}) & \forall(i, s) \\ g_{x,m,i}^s | a_{x,s}, k_{x,m,i} &\sim \mathcal{N}(g_{x,m,i}^s; a_{x,s}^T k_{x,m,i}, \sigma_g^2) & \forall(m, s, i) \end{aligned}$$

(b) multiple kernel learning

$$\begin{aligned} \eta_{k,m} &\sim \mathcal{G}(\eta_{k,m}; \alpha_\eta, \beta_\eta) & \forall m \\ \epsilon_{x,m} \eta_{k,m} &\sim \mathcal{N}(\epsilon_{x,m}; 0, \eta_{k,m}^{-1}) & \forall m \\ h_{x,i}^s | \{\epsilon_{x,m}, g_{x,m,i}^s\}_{m=1}^{P_x} &\sim \mathcal{N}\left(h_{x,i}^s; \sum_{m=1}^{P_x} \epsilon_{x,m} g_{x,m,i}^s, \sigma_h^2\right) & \forall(s, i) \end{aligned}$$

(c) matrix factorization

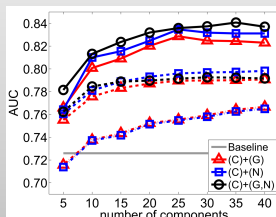
$$f_j^i | h_{x,i}, h_{z,j} \sim \mathcal{N}(f_j^i; h_{x,i}^T h_{z,j}, 1) \quad \forall(i, j)$$

(d) binary classification

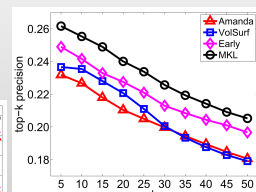
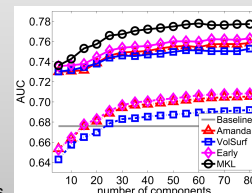
$$y_j^i | f_j^i \sim \delta(f_j^i y_j^i > \nu) \quad \forall(i, j)$$

## Drug–Protein Interaction Data Sets

- A drug–protein network by Yamanishi et al. (2008)
  - > 445 drugs, 664 proteins and 2926 validated interactions
  - C: chemical similarity for drugs
  - G: genomic similarity for proteins
  - N: network similarity for proteins
- 5 replication of 5-fold CV over drugs



- Another drug–protein interaction network by Khan et al. (2012)
  - > 855 drugs, 800 proteins, and 4659 validated Interactions
- Two standard 3D chemical structure descriptors for drugs:
  - > Amanda (Duran et al., 2008) and VolSurf (Cruciani et al., 2000)
- Gaussian kernel whose width is selected as  $\sqrt{D}$
4. 5 replications of 5-fold cross validation over drugs
5. An extra task of finding or retrieving drugs with similar functions



KBMF is statistically significantly better than KPMF of Zhou et al. (2012) according to paired t-test ( $p < 0.01$ ) on both data sets.

(1) M. Gönen, S. A. Khan, and S. Kaski. Kernelized Bayesian Matrix Factorization. In Proceedings of ICML 2013, the 30th International Conference on Machine Learning, volume 28 of JMLR W&CP, pages 864–872. JMLR, 2013. Implementations in Matlab are available at <http://research.ics.aalto.fi/mi/> (2) S. A. Khan, A. Faisal, J. P. Mpindi, J. A. Parkkinen, T. Kallioikoski, A. Poso, O. P. Kallioniemi, K. Wennerberg, and S. Kaski. Comprehensive data-driven analysis of the impact of chemoinformatic structure on the genome-wide biological response profiles of cancer cells to 1159 drugs. BMC Bioinformatics, 13(112), 2012.