BayReL: Bayesian Relational Learning for Multi-omics Data Integration

Introduction

BayReL is a novel Bayesian relational learning method:

- Integrating high-dimensional multi-omics data
- Taking advantage of a priori known relationships modeled as a graph at each corresponding view
- Inferring the relational interactions as a bipartite graph without any preknown interactions across views
- Exploiting non-linear and deep transformations of data
- Enabling Bayesian interpretation

Method

Embedding nodes to the latent space

Independently parametrize the distribution over the adjacency matrix of each view

$$\int p_{\theta}(\mathcal{G}, \mathcal{U}) \, \mathrm{d}\mathcal{U} = \prod_{v=1}^{V} \int p_{\theta}(\mathbf{A}^{v}, \mathbf{U}_{v}) \, \mathrm{d}\mathbf{U}_{v} = \prod_{v=1}^{V} \int p_{\theta}(\mathbf{A}^{v} \,|\, \mathbf{U}_{v}) \, p(\mathbf{U}_{v}) \, \mathrm{d}\mathbf{U}_{v}$$

where the prior distribution of U_v 's are standard Gaussian. We further approximate the distribution of \mathcal{U} with a factorized posterior distribution

$$q(\mathcal{U} \mid \mathcal{X}, \mathcal{G}) = \prod_{v=1}^{V} q(\mathbf{U}_{v} \mid \mathbf{X}_{v}, G_{v}) = \prod_{v=1}^{V} \prod_{i=1}^{N_{v}} q(\mathbf{u}_{i,v} \mid \mathbf{X}_{v}, G_{v})$$
$$q(\mathbf{u}_{i,v} \mid \mathbf{X}_{v}, G_{v}) = \mathcal{N} \left(\mathbf{u}_{i,v} \mid \boldsymbol{\mu}_{i,v}, \operatorname{diag}(\boldsymbol{\sigma}_{i,v}^{2})\right),$$
$$\boldsymbol{\mu}_{v} = \varphi_{v}^{\operatorname{emb}, \mu}(\mathbf{X}_{v}, G_{v}), \quad \log(\boldsymbol{\sigma}_{v}) = \varphi_{v}^{\operatorname{emb}, \sigma}(\mathbf{X}_{v}, G_{v})$$

 ϕ_v^{emb} : variants of graph neural networks including GCN, GIN, GraphSAGE



Schematic illustration of BayReL

Constructing relational multi-partite graph

The distribution of bi-adjacency matrices are defined as follows

$$p(\mathbf{A}^{vv'} | \mathbf{U}_{v}, \mathbf{U}_{v'}) = \prod_{i=1}^{N_{v}} \prod_{j=1}^{N_{v'}} \text{Bernoulli}\left(\mathbf{A}_{ij}^{vv'} | \varphi^{\text{sim}}(\mathbf{u}_{i,v}, \mathbf{u}_{j,v'})\right),$$

where $\varphi_{sim} = \sigma(\boldsymbol{u}_{i,v} \ \boldsymbol{u}_{j,v}^T)$ is a score function measuring the similarity between the latent representations of nodes.

Inferring view-specific latent variables

Parametrize the distributions for node attributes at each view independently

$$\int p_{\theta}(\mathcal{X}, \mathcal{Z} | \mathcal{G}, \mathcal{A}, \mathcal{U}) \, \mathrm{d}\mathcal{Z} = \prod_{v=1}^{V} \prod_{i=1}^{N_{v}} \int p_{\theta}\left(\mathbf{z}_{i,v} | \mathcal{G}, \mathcal{A}, \mathcal{U}\right) \, p_{\theta}(\mathbf{x}_{i,v} | \mathbf{z}_{i,v}) \, \mathrm{d}\mathbf{z}_{i,v}$$

Prior construction

$$p_{\theta}(\mathcal{Z} \mid \mathcal{G}, \mathcal{A}, \mathcal{U}) = \prod_{v=1}^{V} \prod_{i=1}^{N_{v}} p_{\theta}(\mathbf{z}_{i,v} \mid \mathcal{G}, \mathcal{A}, \mathcal{U}), \qquad p_{\theta}(\mathbf{z}_{i,v} \mid \mathcal{G}, \mathcal{A}, \mathcal{U}) = \mathcal{N}(\boldsymbol{\mu}_{i,v}^{\text{prior}}, \boldsymbol{\sigma}_{i,v}^{\text{prior}}),$$
$$\boldsymbol{\mu}^{\text{prior}} = \varphi^{\text{prior}, \boldsymbol{\mu}}(\mathcal{A}, \mathcal{U}), \quad \boldsymbol{\sigma} = \varphi^{\text{prior}, \boldsymbol{\sigma}}(\mathcal{A}, \mathcal{U})$$

Approximate the posterior

$$q(\mathcal{Z} \mid \mathcal{X}, \mathcal{G}) = \prod_{v=1}^{V} \prod_{i=1}^{N_v} q(\mathbf{z}_{i,v} \mid \mathbf{X}_v, G_v), \qquad q(\mathbf{z}_{i,v} \mid \mathbf{X}_v, G_v) = \mathcal{N}(\boldsymbol{\mu}_{i,v}^{\text{post}}, \boldsymbol{\sigma}_{i,v}^{\text{post}})$$
$$\boldsymbol{\mu}^{\text{post}} = \{\varphi_v^{\text{post}, \boldsymbol{\mu}}(\mathbf{X}_v, G_v)\}_{v=1}^{V}, \quad \boldsymbol{\sigma}^{\text{post}} = \{\varphi_v^{\text{post}, \boldsymbol{\sigma}}(\mathbf{X}_v, G_v)\}_{v=1}^{V}$$

Overall likelihood and learning

Marginal likelihood

$$p_{\theta}(\mathcal{X}, \mathcal{G}) = \int \prod_{v=1}^{V} p_{\theta}(\mathbf{X}_{v} | \mathbf{Z}_{v}) p_{\theta}(\mathbf{Z}_{v} | \mathcal{G}, \mathcal{A}, \mathcal{U}) p(\mathcal{A} | \mathcal{U}) p(\mathcal{G} | \mathcal{U}) p(\mathcal{U}) d\mathbf{Z}_{1} \dots d\mathbf{Z}_{V} d\mathcal{A} d\mathcal{U}.$$

Evidence Lower Bound (ELBO)

$$\mathcal{L} = \sum_{v=1}^{V} \left[\mathbb{E}_{q_{\phi}(\mathbf{Z}_{v} \mid \mathcal{G}, \mathcal{X})} \log p_{\theta}(\mathbf{X}_{v} \mid \mathbf{Z}_{v}) + \mathbb{E}_{q_{\phi}(\mathbf{Z}_{v}, \mathcal{U} \mid \mathcal{G}, \mathcal{X})} \log p_{\theta}(\mathbf{Z}_{v} \mid \mathcal{G}, \mathcal{A}, \mathcal{U}) - \mathbb{E}_{q_{\phi}(\mathbf{Z}_{v} \mid \mathcal{G}, \mathcal{X})} q_{\phi}(\mathbf{Z}_{v} \mid \mathcal{G}, \mathcal{X}) \right] - \mathrm{KL} \left(q_{\phi}(\mathcal{U} \mid \mathcal{G}, \mathcal{X}) \mid \mid p(\mathcal{U}) \right),$$

Experiments

Microbiome-metabolome interactions in cystic fibrosis

Data description

- 172 patients with CF
- 138 unique microbial taxa
- 462 metabolite features
- Graph density of input networks:
- Microbiome network: 0.102
- Metabolite network: 0.011

Evaluation metrics

- Positive accuracy: accuracy to identify the validated molecules interacting with *P. aeruginosa*
- Negative accuracy: accuracy of not detecting common targets between anaerobic microbes and notable pathogen





Left: Positive vs negative accuracy in CF dataset. Right: A sub-network of dependency graph consisting of *P. aeruginosa*.

miRNA-mRNA interactions in breast cancer

Data description

- 1156 patients with BRCA
- 11872 genes
- 432 miRNA

Networks:

- Gene regulatory networks (GRN): GENIE3 R package
- miRNA-miRNA: MISIM v2.0
- miRNA-mRNA validation: miRNet

Table 1: Prediction sensitivity (in %) in TCGA for different percentage of training samples.

		BCCA			BayReL	
Avg. degree	289 (25%)	# of training samples 578 (50%)	1156 (100%)	289 (25%)	# of training samples 578 (50%)	1156 (100%)
0.20 0.30 0.40	$\begin{vmatrix} 17.4 \pm 0.8 \\ 26.0 \pm 0.8 \\ 35.4 \pm 0.8 \end{vmatrix}$	17.6 ± 1.0 26.4 ± 1.0 35.5 ± 0.7	21.0 ± 0.0 31.1 ± 0.7 41.1 ± 0.2	$\begin{vmatrix} 31.9 \pm 3.0 \\ 45.8 \pm 3.1 \\ 57.6 \pm 4.4 \end{vmatrix}$	32.1 ± 1.0 45.9 ± 1.5 58.7 ± 1.3	34.0 ± 2.5 47.4 ± 2.6 59.5 ± 3.0

Precision medicine in acute myeloid leukemia

Data description

- 30 AML patients
- 53 drugs
- 9071 genes

Networks:

- GRN: 14 AML cell lines
- drug-drug: action mechanisms
- drug-gene validation: DGIdb

Evaluation metrics

• Prediction sensitivity of identifying reported drug-gene interactions

	Table 2: Comp	parison of predic	tion sensitivity (in %) in AML da	taset for differe	nt graph densities.	
Avg. degree	0.10	0.15	0.20	0.25	0.30	0.40	0.50

Avg. degree	0.10	0.15	0.20	0.25	0.30	0.40	0.50
SRCA	8.03	12.00	17.15	20.70	26.85	34.93	45.79
BCCA	9.65 ± 0.75	14.34 ± 0.06	18.96 ± 0.42	23.29 ± 0.52	28.22 ± 0.66	38.02 ± 2.15	46.88 ± 1.88
BayReL	15.56 ± 0.75	21.70 ± 0.65	27.20 ± 0.17	32.43 ± 1.02	37.76 ± 0.85	47.90 ± 0.43	56.76 ± 0.50

 Consistency of significant gene-drug interactions in two different AML datasets with 30 patients and 14 cell lines: KL divergence between two inferred bipartite graphs are 0.38 and 0.66 for BayReL and BCCA, respectively.

References

Arto Klami et al. (2013). "Bayesian Canonical Correlation Analysis." In Journal of Machine Learning Research.

James T. Morton et al. (2019). "Learning representations of microbe–metabolite interactions." In Nature methods.

Su-In Lee et al. (2018). "A machine learning approach to integrate big data for precision medicine in acute myeloid leukemia." In Nature communications.

