

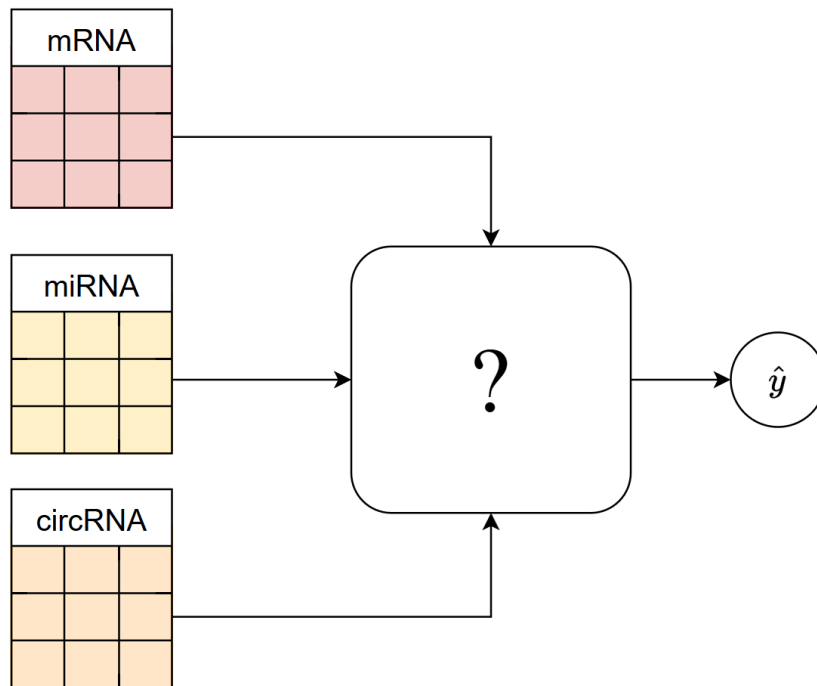
Graph Neural Networks for Identification of Robust Biomarkers from Multi-Omics Data

Bc. Jan Lubojacký

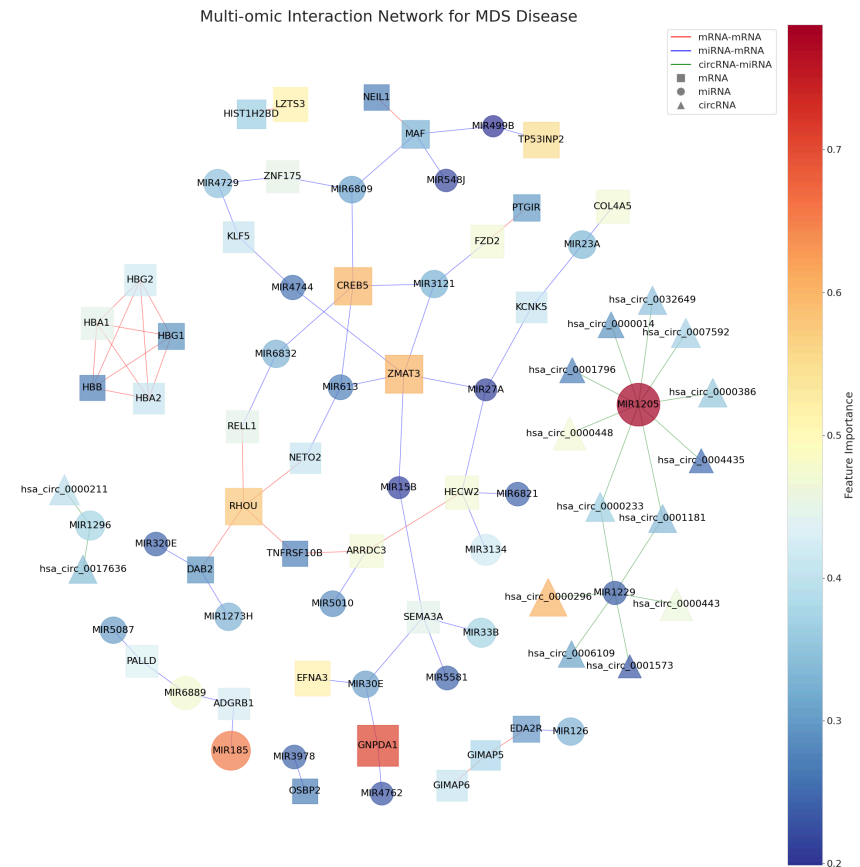
Supervisor: Ing. Jiří Kléma, Ph.D.
Project reviewer: Ing. Jan Drchal, Ph.D.

Medical Electronics and Bioinformatics
Faculty of Electrical Engineering
Department of Cybernetics

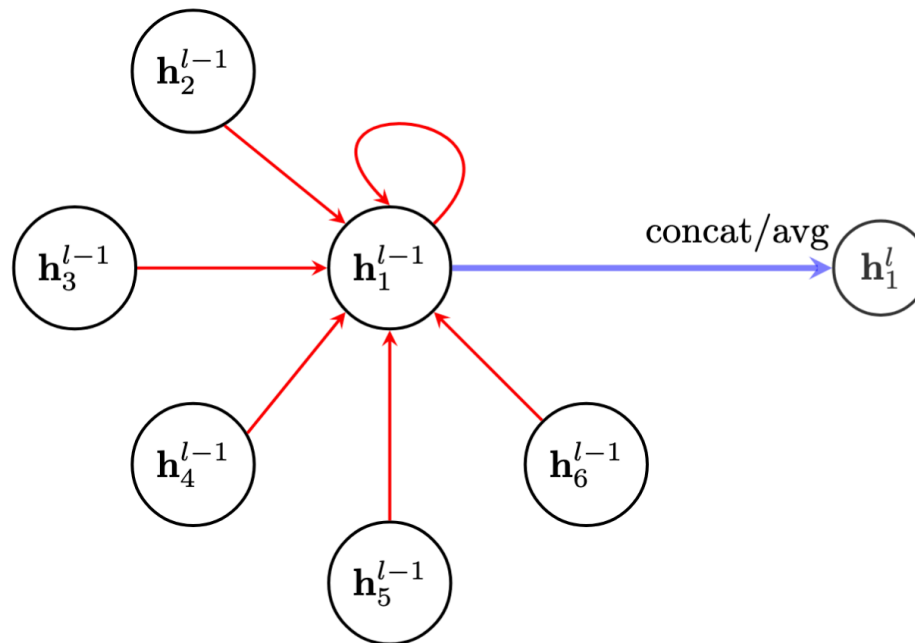
Making predictions



Extracting biomarkers



Graph Neural Networks



$$h_u^l = \text{UPDATE}(h_u^{l-1}, \text{AGGREGATE}(h_v^{l-1}, \forall v \in \mathcal{N}(u)))$$

$$h_u^l = \text{UPDATE}(h_u^{l-1}, \text{AGGREGATE}(h_v^{l-1}, \forall v \in \mathcal{N}(u)))$$

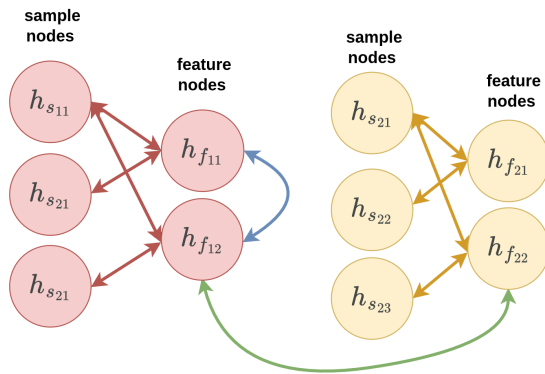
- **Graph Convolutional Networks**

$$h_u^l = \sigma \left(\mathbf{w}_{\text{self}}^l h_u^{l-1} + \mathbf{w}_{\text{neigh}}^l \sum_{v \in \mathcal{N}(u)} h_v^{l-1} + \mathbf{b}^l \right)$$

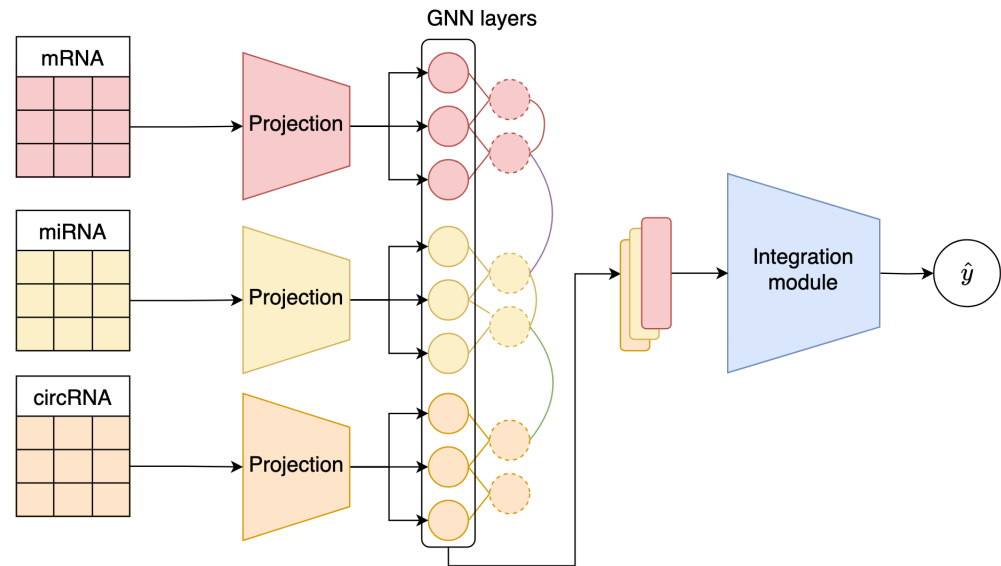
- **Graph Attention Networks**

$$h_u^l = \sigma \left(\sum_{v \in \mathcal{N}(u)} \alpha(v, u) \mathbf{w}^l h_v^{l-1} \right)$$

Bipartite GNN



Input graph structure



Model architecture

Integration modules

- **linear integrator**

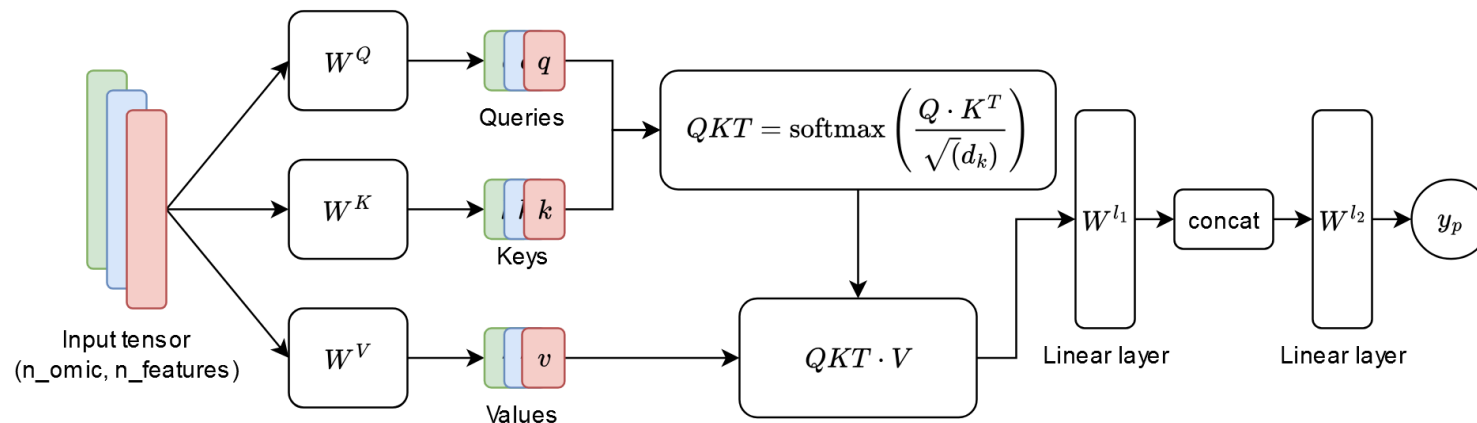
- $\hat{y} = \sigma(W \cdot \parallel_{i=1}^m x^{(i)} + b)$

- **view correlation discovery network**

- $C_{a_1, \dots, a_m} = \prod_{i=1}^m x_{a_i}^{(i)}; a = 1, \dots, m; C \in \mathbb{R}^{|x_1| \times \dots \times |x_m|}$

- $\hat{y} = \text{VCDN}(C)$

- **attention integrator**



Models

- KNN
- Linear SVMs
- XGBoost
- Linear NN
- MOGONET
- BipartiteGNN

Models

- KNN
- Linear SVMs
- XGBoost
- Linear NN
- MOGONET
- BipartiteGNN

Datasets

- TCGA-BRCA
 - 483 samples
 - 4 classes
 - mRNA, miRNA, CNV, DNA methylation
- MDS
 - 3 tasks
 - Disease (74 samples)
 - Risk (53 samples)
 - Mutation (26 samples)
 - 2 classes

Models

- KNN
- Linear SVMs
- XGBoost
- Linear NN
- MOGONET
- BipartiteGNN

Datasets

- TCGA-BRCA
 - 483 samples
 - 4 classes
 - mRNA, miRNA, CNV, DNA methylation
- MDS
 - 3 tasks
 - Disease (74 samples)
 - Risk (53 samples)
 - Mutation (26 samples)
 - 2 classes

Evaluation

for each trial in 1 to N:

1. sample hyperparameters

for each CV split:

2. split data

3. feature selection

4. normalize features

5. fit model

6. test model

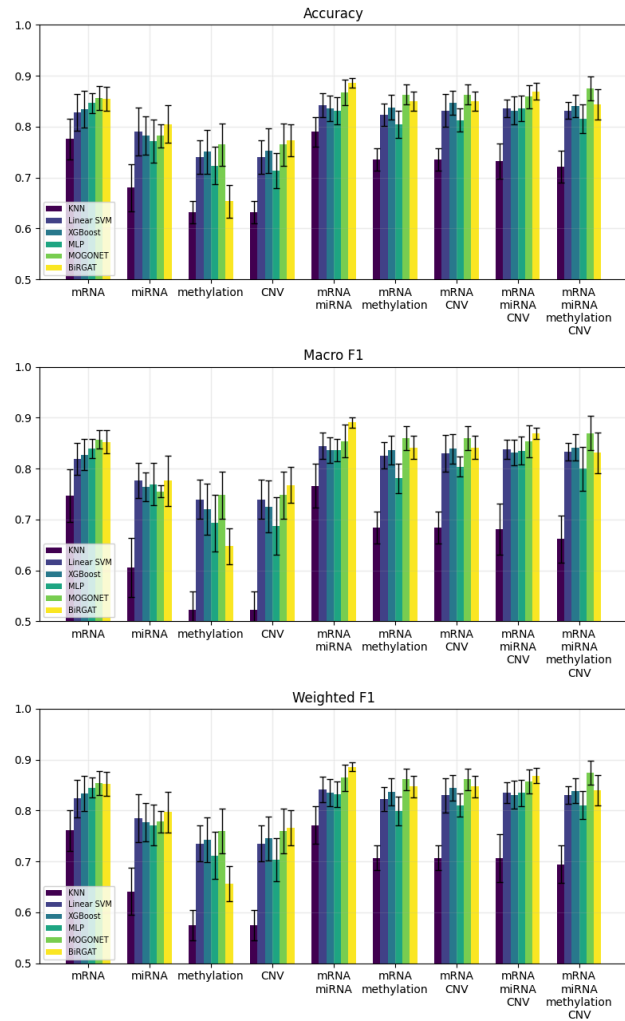
7. record performance P

8. Average P across folds

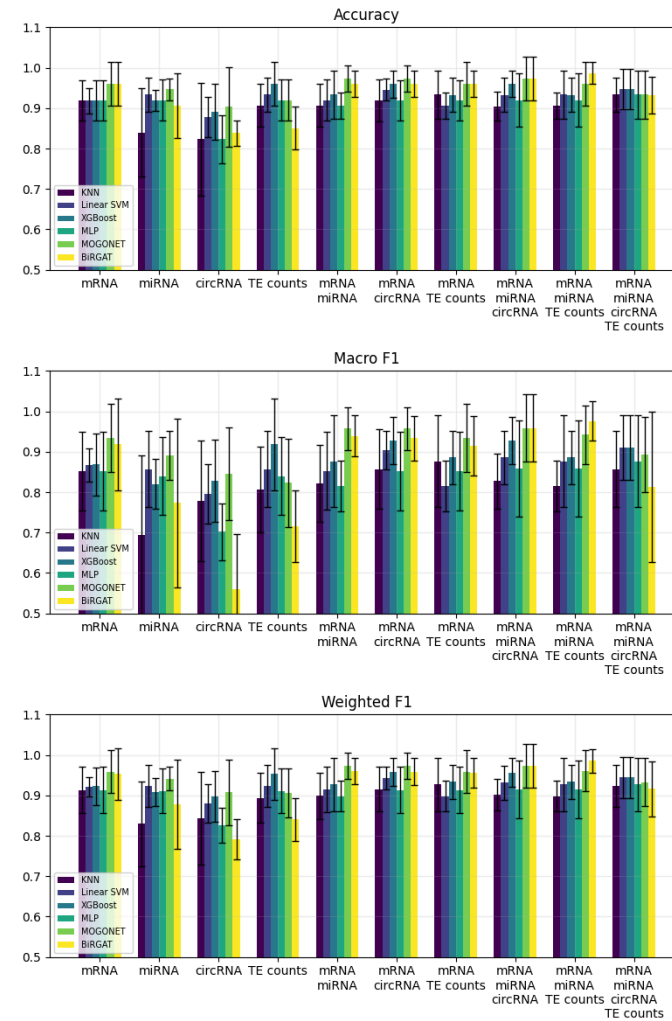
9. If $P > P_{best}$: $P_{best}=P$

return P_{best}

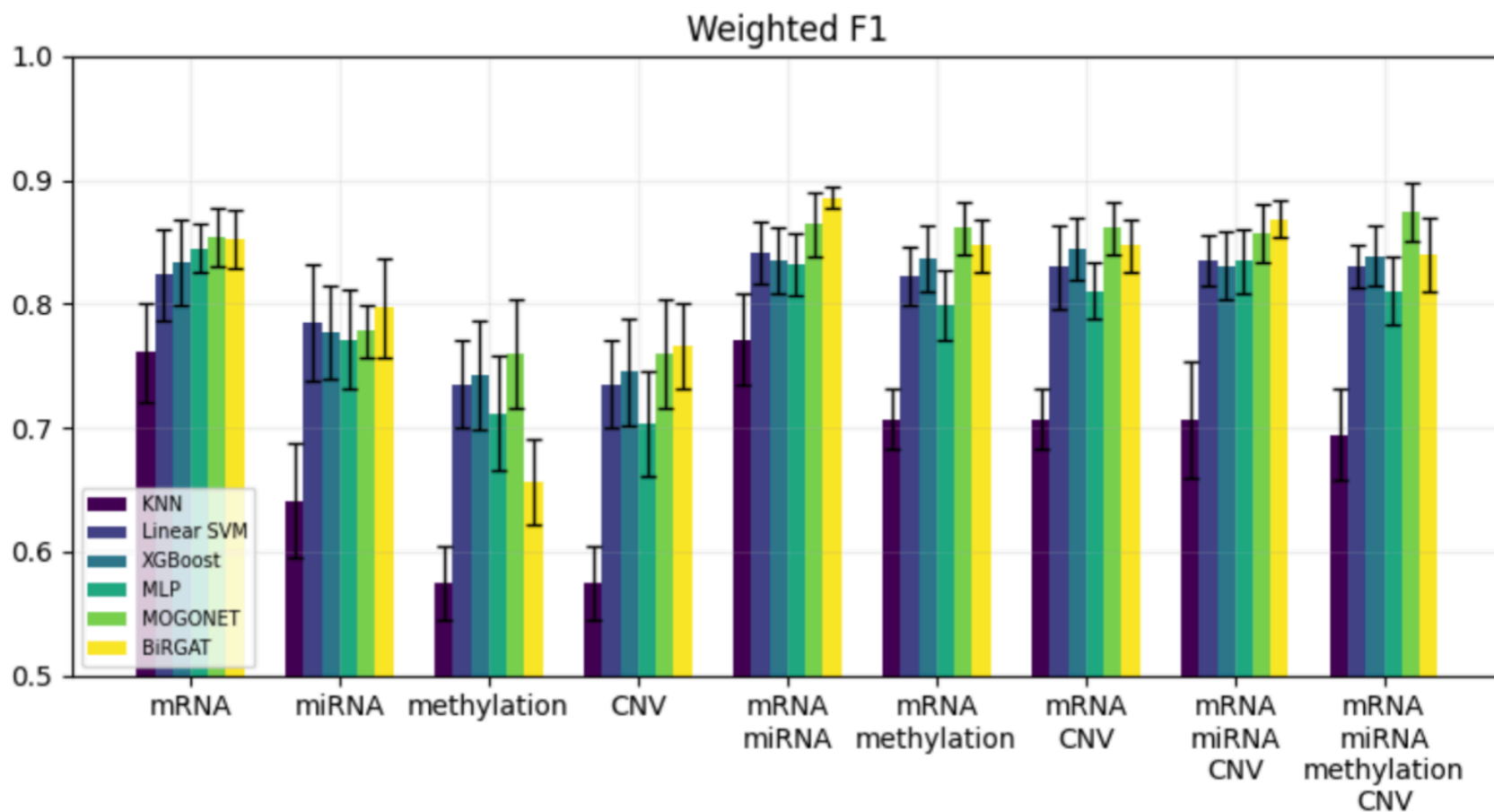
Results



• TCGA-BRCA results



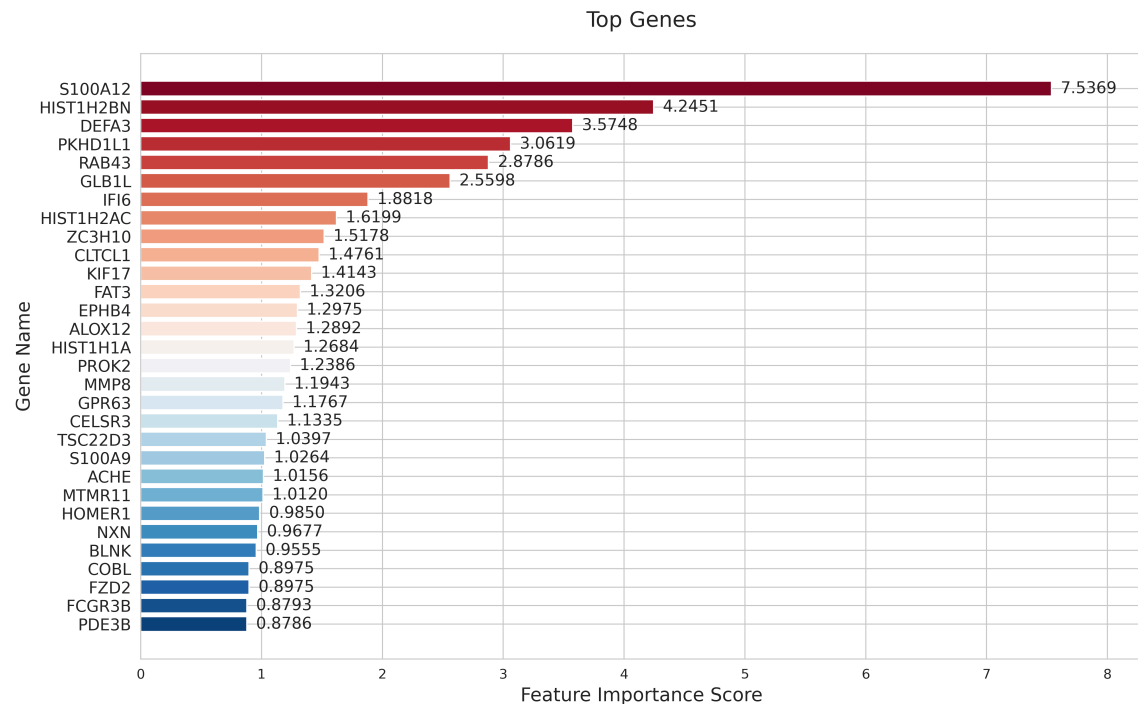
• MDS disease results



- Model performances in terms of Weighted F1 score on the TCGA-BRCA dataset

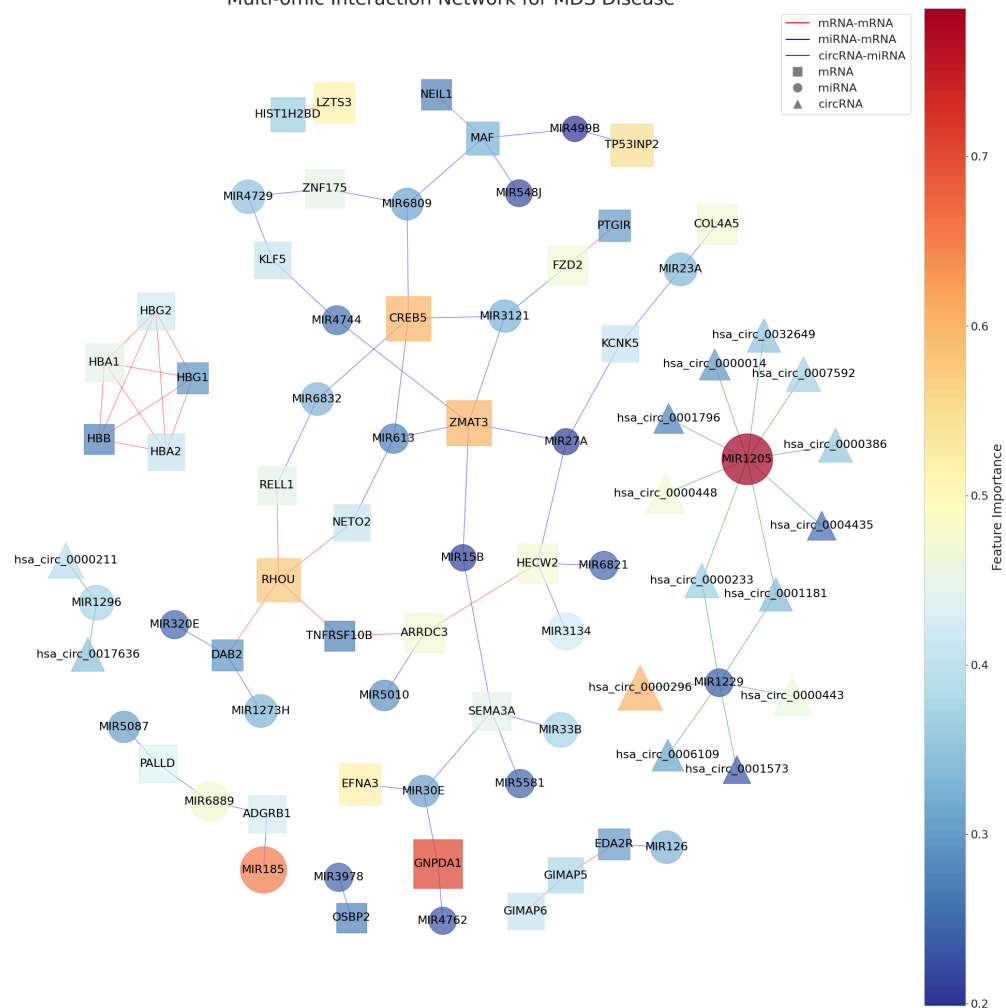
Extracting Biomarkers

- XGBoost linear booster weights
- Feature permutation for GNNs
- Comparison against traditional methods



- feature importances (MDS risk, XGBoost)

Multi-omic Interaction Network for MDS Disease



Thank You
for your attention

Reviewer question

- U sloupcových grafů v obrázcích 8.1, 8.2 a 8.3 mi není zřejmý přesný význam černých čar posazených na vrchní část každého sloupce. Předpokládám, že se jedná o konfidenční intervaly. Zejména u obr. 8.3 jsou tyto intervaly značně široké a pro jednotlivé metody se výrazně překrývají. Můžete na základě grafů říci něco o statistické významnosti výsledků?

