

Collaborative Phenotype Inference from Comorbid Substance Use Disorders and Genotypes

BIBM 2017@Kansas City

Presenter: Jin Liu

University of Connecticut

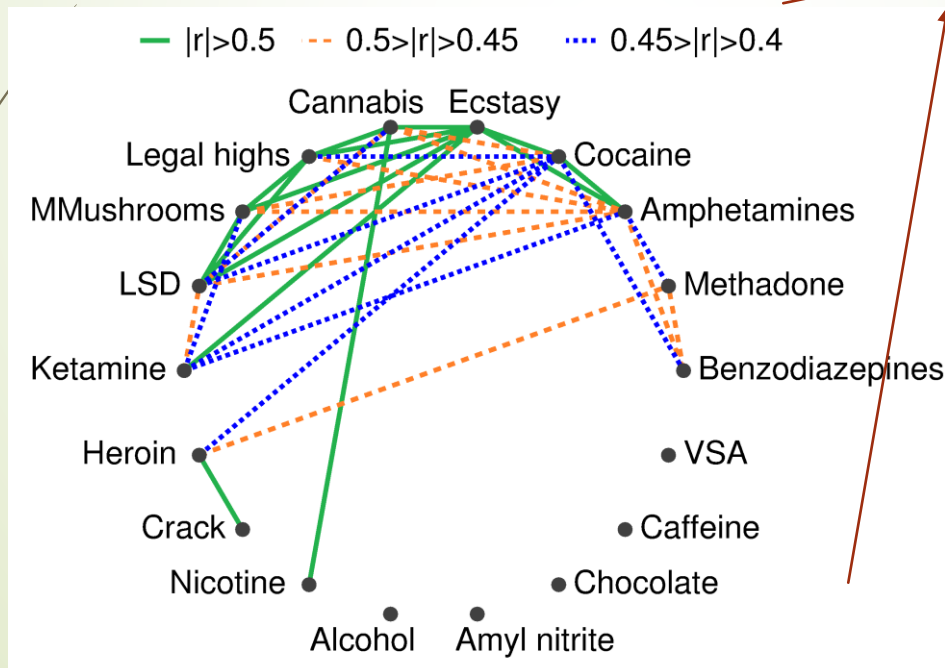
Joint work with Jiangwen Sun, Xinyu Wang, Henry Kranzler, Joel Gelernter and Jinbo Bi

Comorbid substance use disorders (CSUD)

very dangerous

- Heatstroke
- Immunity problems
- Seizures
- Coma
- Heart attack
- Respiratory failure
- Liver damage
- Overdose

validated by strong correlations



USE BOTH
GENOTYPES
RELATED PHENOTYPES
TO INFER
DIAGNOSTIC
CRITERIA FOR
UNREPORTED SUD

Dependence on different substances are correlated both phenomenologically and biologically.

Inferring SUD diagnostic criteria

Our phenotypic imputation problem

$\mathbf{F} =$

	opioid	cocaine
Patient ₁	No observ	
Patient ₂		
...		
		No observ
	No observ	
	No observ	
		No observ
Patient _m	No observ	
	f_1 $f_2 \dots$	f_n

additional useful information:
 associated genetic variants; known similarities
 between comorbid disorders.
 Often referred to as side or auxiliary information in
 matrix completion

Classical low-rank matrix completion problem

$\mathbf{F} =$

	Avatar	Toy Story	Mickey Blue Eyes	Barbie	Matrix	Spider-Man
Tom Hanks	2		4	5	?	
Julia Roberts	5		4			1
Samuel L. Jackson			5	?	2	
Tom Cruise		1		5		4
Brad Pitt			4		2	
George Clooney	4	5		1		

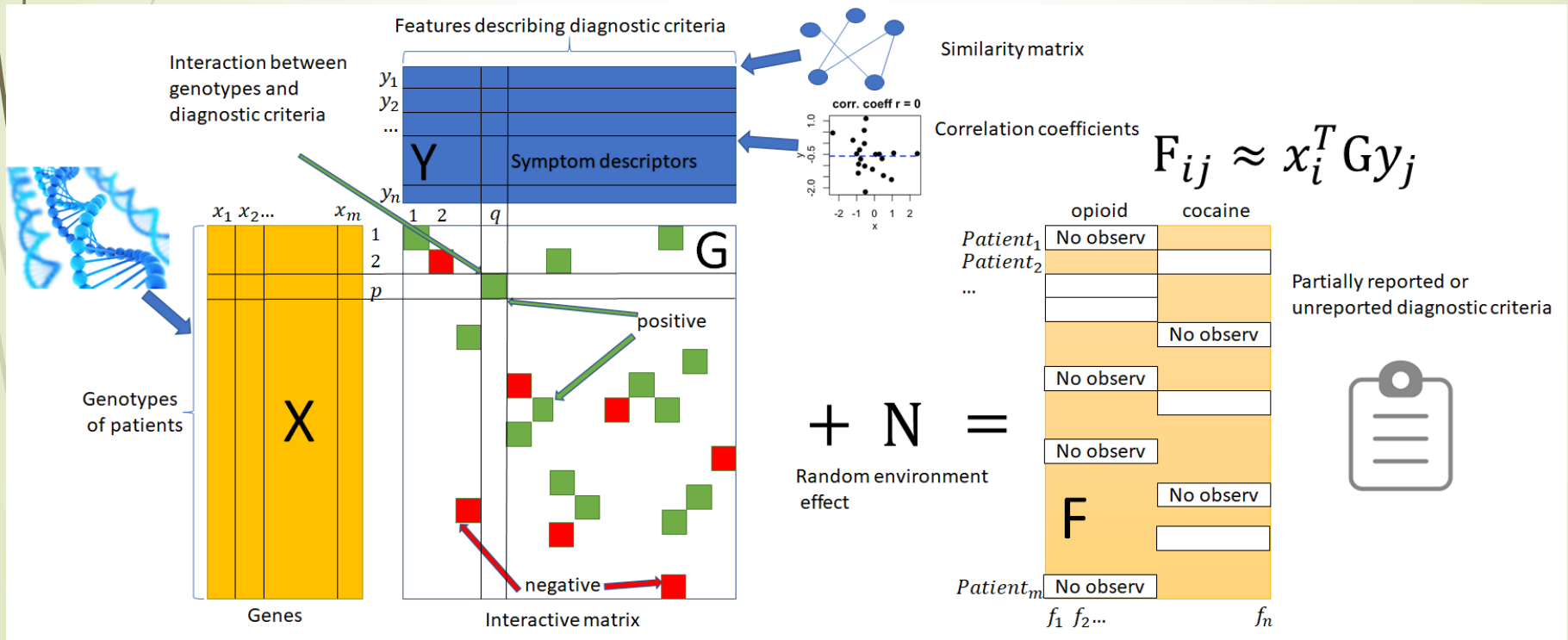
Optimization Problem:

$$\min_{\mathbf{E}} \|\mathbf{E}\|_* \quad \text{subject to} \quad R_{\Omega}(\mathbf{E}) = R_{\Omega}(\mathbf{F})$$

where Ω is the set of indices of **observed entries** in \mathbf{F} and $\|\cdot\|_*$ computes the **nuclear norm (low-rank regularizer)**.

The proposed method

- Our collaborative inference method of CSUD diagnostic criteria using side information.
- Side information of patients: Genotypes
- Side information of diagnostic criteria: Sampled Corr Coef matrix; Similarity matrix



The proposed method

$$\min_{\mathbf{G}} \frac{1}{2} \|\mathbf{X}^T \mathbf{G} \mathbf{Y} - \mathbf{E}\|_F^2 + \lambda_G g(\mathbf{G}) + \lambda_E \|\mathbf{E}\|_*,$$

s. t. $R_\Omega(\mathbf{E}) = R_\Omega(\mathbf{F})$

- The proposed method uses a **low-rank matrix** \mathbf{E} to directly approximate matrix \mathbf{F} and then estimates \mathbf{E} from matrix \mathbf{X} and \mathbf{Y} .

Proposed method

$$\min_{\mathbf{G}} \frac{1}{2} \|\mathbf{X}^T \mathbf{G} \mathbf{Y} - \mathbf{E}\|_F^2 + \lambda_G g(\mathbf{G}) + \lambda_E \|\mathbf{E}\|_*,$$

s. t. $R_\Omega(\mathbf{E}) = R_\Omega(\mathbf{F})$

- The proposed method uses a **low-rank matrix E** to directly approximate matrix F and then estimates E from matrix X and Y.
- The proposed model can identify crucial interactions between specific genotypes and diagnostic criteria **by enforcing the sparsity in G**. ($g(\mathbf{G}) = \|\mathbf{G}\|_1$)

Adaptive LADMM Algorithm

- ▶ We propose a new stochastic Linearized Alternative Direction Method of Multipliers (**StoLADMM**) algorithm
- ▶ by substituting $\mathbf{C} = \mathbf{E} - \mathbf{X}^T \mathbf{G} \mathbf{Y}$.

- ▶ The augmented Lagrangian function is given by

$$\begin{aligned} & \mathcal{L}(\mathbf{E}, \mathbf{G}, \mathbf{C}, \mathbf{M}_1, \mathbf{M}_2, \beta) \\ &= \frac{1}{2} \|\mathbf{C}\|_F^2 + \lambda_G \|\mathbf{G}\|_1 + \lambda_E \|\mathbf{E}\|_* + \langle \mathbf{M}_1, R_\Omega(\mathbf{E} - \mathbf{F}) \rangle \\ &+ \langle \mathbf{M}_2, \mathbf{E} - \mathbf{X}^T \mathbf{G} \mathbf{Y} - \mathbf{C} \rangle + \frac{\beta}{2} \|R_\Omega(\mathbf{E} - \mathbf{F})\|_F^2 \\ &+ \frac{\beta}{2} \|\mathbf{E} - \mathbf{X}^T \mathbf{G} \mathbf{Y} - \mathbf{C}\|_F^2 \end{aligned}$$

- ▶ Solve each variable alternatively.

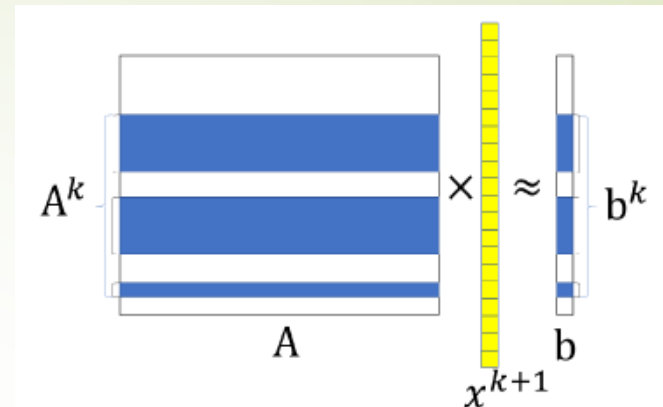
Our efficient stochastic algorithm

► Effectiveness

1. convergence in expectation
2. global optimal solution for our convex optimization problem

► Efficiency

1. Save memory costs
2. Can utilize parallel computing to speed up the algorithm
3. Without sacrificing performance notably.



Algorithm 1 The StoLADMM algorithm to solve $\mathbf{C}^k, \mathbf{G}^k, \mathbf{E}^k, k = 1, \dots, K$

Input: \mathbf{X}, \mathbf{Y} and $R_\Omega(\mathbf{F})$ with parameters $\lambda_G, \lambda_E, \tau_A, \tau_B, \rho$ and β_{max} .

Output: $\mathbf{C}, \mathbf{G}, \mathbf{E}$;

- 1: Initialize $\mathbf{E}^0, \mathbf{G}^0, \mathbf{M}_1^0, \mathbf{M}_2^0$. Compute $\mathbf{A} = \mathbf{Y}^T \otimes \mathbf{X}^T$. $k = 0$,
repeat;
- 2: $\mathbf{C}^{k+1} = \frac{\beta}{\beta+1}(\mathbf{E}^k - \mathbf{X}^T \mathbf{G}^k \mathbf{Y} + \mathbf{M}_2^k / \beta)$;
- 3: $\mathbf{G}^{k+1} = \text{reshape}(\max(|\mathbf{g}^k - f_1^k / \tau_A| - \frac{\lambda_G}{\tau_A \beta}, 0) \odot \text{sgn}(\mathbf{g}^k - f_1^k / \tau_A))$ where f_1^k can be computed by (5);
- 4: $\mathbf{E}^{k+1} = \text{SVT}(\mathbf{E}^k - (f_2^k + f_3^k) / (2\tau_B), \lambda_E / 2(\beta\tau_B))$ where f_2^k and f_3^k can be computed by (6);
- 5: $\mathbf{M}_1^{k+1} = \mathbf{M}_1^k + \beta(R_\Omega(\mathbf{E}^{k+1}) - \mathbf{F})$.
- 6: $\mathbf{M}_2^{k+1} = \mathbf{M}_2^k + \beta(\mathbf{E}^{k+1} - \mathbf{X}^T \mathbf{G}^{k+1} \mathbf{Y} - \mathbf{C}^{k+1})$.
- 7: $k = k + 1$ until convergence;
Return $\mathbf{C}, \mathbf{G}, \mathbf{E}$;

Experimental results

➤ Compared methods:

➤ MAXIDE

M. Xu, R. Jin, and Z. hua Zhou. Speedup matrix completion with side information: Application to multi-label learning. *Advances in Neural Information Processing Systems 26*, pages 2301–2309, 2013

➤ IMC

N. Natarajan and I. S. Dhillon. Inductive matrix completion for predicting gene–disease associations. *Bioinformatics*, 30(12):i60–i68, 2014

➤ DirtyIMC

K.-Y. Chiang, C.-J. Hsieh, and I. S. Dhillon. Matrix completion with noisy side information. *Advances in Neural Information Processing Systems 28*, pages 3429–3437, 2015.

➤ The relative mean squared error (**RMSE**) is used as the performance measurement.

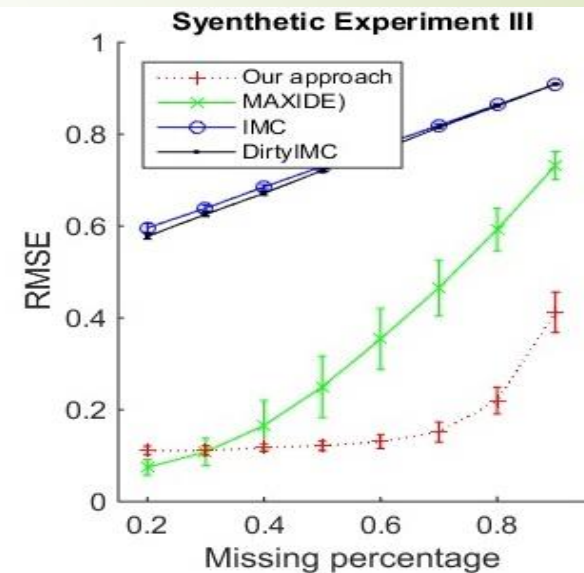
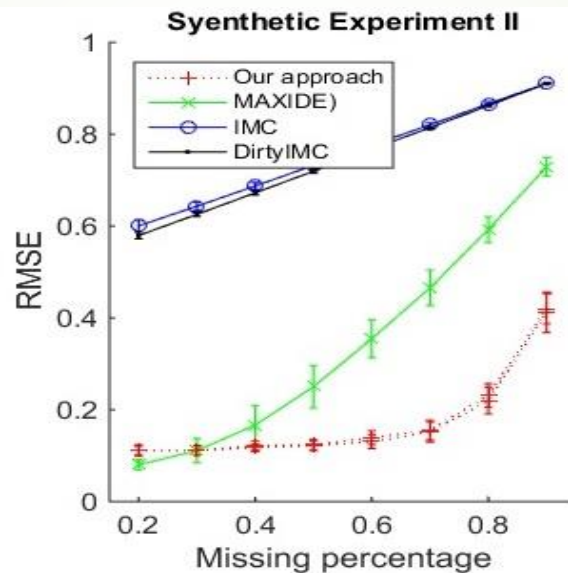
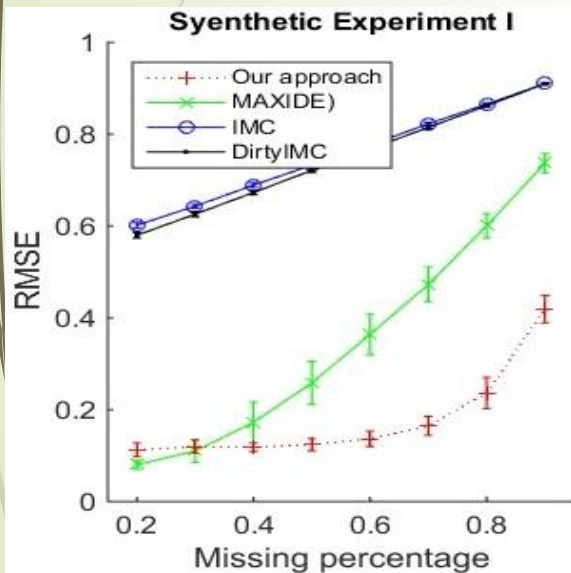
$$\text{RMSE} = \frac{\|R_{\bar{\Omega}}(\mathbf{X}^T \mathbf{G} \mathbf{Y} - \mathbf{F})\|_F^2}{\|R_{\bar{\Omega}}(\mathbf{F})\|_F^2}$$

Experimental results

- Synthetic Datasets:
 - \mathbf{X} and \mathbf{Y} were generated from **Gaussian**, **Poisson** and **Gamma** distributions.
 - \mathbf{G} contains **20%** of non-zero components.
 - $\mathbf{F} = \mathbf{X}^T \mathbf{G} \mathbf{Y} + \mathbf{N}$ where \mathbf{N} represents the noise.
 - Then, the values of \mathbf{F} were dropped by [**20% – 80%**] to test the recovery rate of the methods.

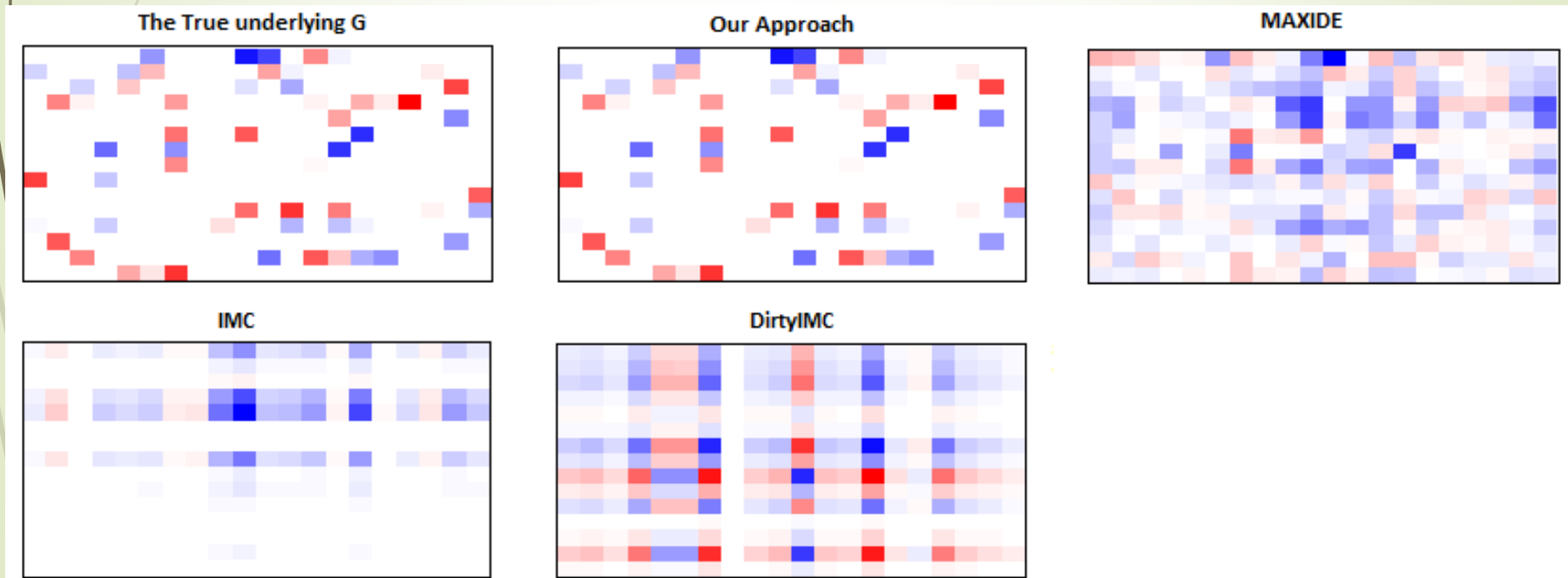
Experimental results

- Synthetic Datasets:
 - RMSE for all compared methods



Experimental results

- Synthetic Datasets:
 - RMSE for all compared methods
 - Recovery of true underlying \mathbf{G} .



Experimental results

➤ Synthetic Datasets

➤ Cormorbid Substance Use Data:

- A total of **7,189 subjects** were aggregated from family and case-control based genetic studies of cocaine use disorder (CUD) and opioid use disorder (OUD).
- The **383 genetic variants** identified in our GWAS were used as side feature matrix X with the size 7189 by 383.
- The correlations between **22 CUD and OUD symptoms** formed a correlation matrix which was used as side features matrix Y with the size 22 by 22.
- We randomly removed the phenotypes of $q\%$ CSUD patients associated with either opioid or cocaine use (not both). Then our partially observed F is the matrix with the size **7189 by 22**, which needs inference.

Experimental results

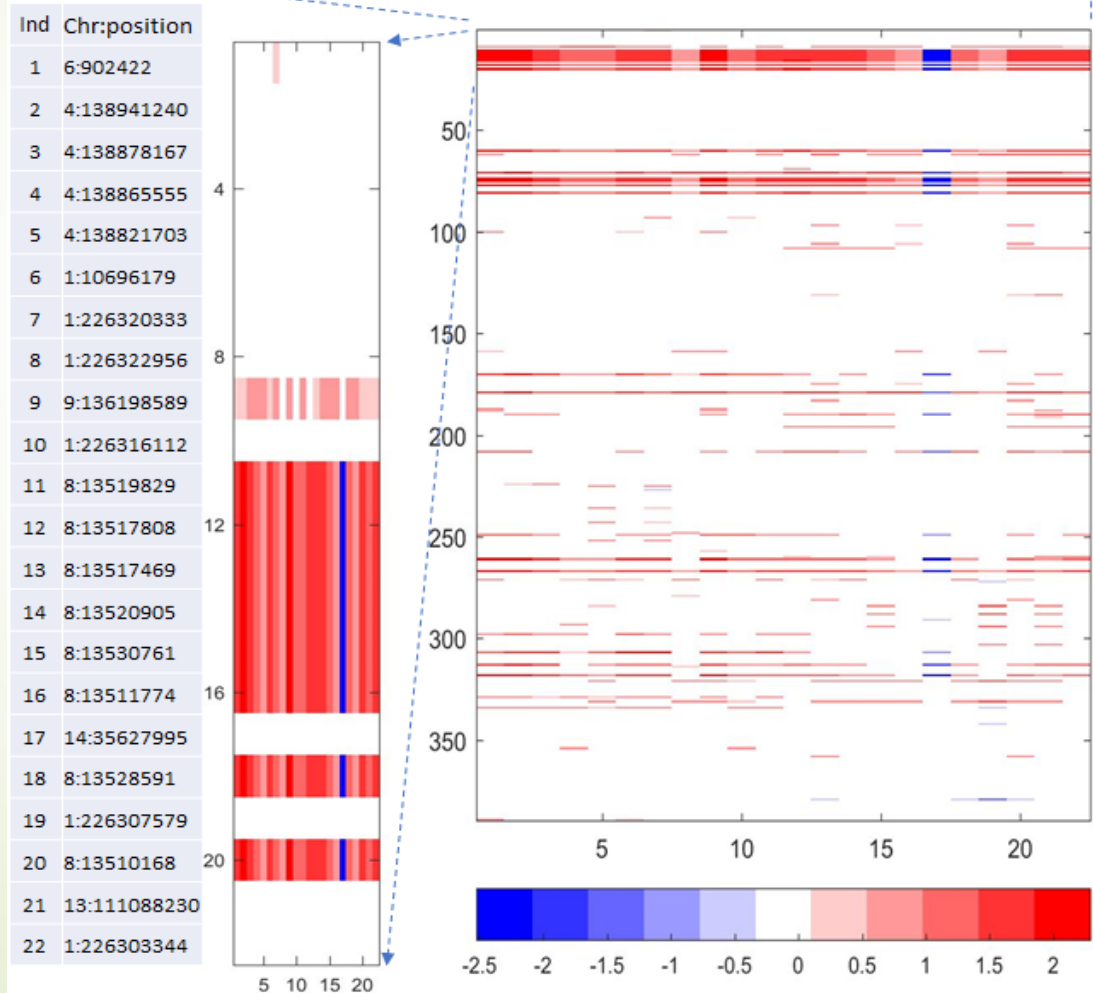
- Synthetic Datasets
- Cormorbid Substance Use Data:

q		StoLADMM	LADMM	DirtyIMC	IMC	MAXIDE	BM
20%	RMSE	0.236	0.231	0.297	0.230	0.235	0.567
	time(s)	30.938	664.515	45.366	21.053	4732.718	NaN
40%	RMSE	0.226	0.234	0.298	0.235	0.236	0.582
	time(s)	29.953	982.212	21.063	20.803	3772.202	NaN
60%	RMSE	0.228	0.236	0.301	0.237	0.235	0.581
	time(s)	28.719	815.841	20.269	36.737	4718.916	NaN
80%	RMSE	0.236	0.237	0.303	0.239	0.241	0.585
	time(s)	30.547	877.886	23.906	32.872	4011.692	NaN
100%	RMSE	0.223	0.239	0.303	0.246	0.242	0.574
	time(s)	30.172	489.770	22.922	24.653	3695.292	NaN

TABLE II: The inference results on the Opioid-Cocaine data.

Interaction Matrix

- | | |
|---|--|
| P1: Larger/longer Coc use than intended | P12: Larger/longer Opi use than intended |
| P2: Failed efforts to stop on Coc | P13: Failed efforts to stop on Opi |
| P3: Much time spent in Coc related activities | P14: Much time spent in Opi related activities |
| P4: Strong desire to use Coc | P15: Strong desire to use Opi |
| P5: Coc-effects interfered with life | P16: Opi-effects interfered with life |
| P6: Coc use despite of its interference | P17: Opi use despite of its interference |
| P7: Major activities reduced by Coc use | P18: Major activities reduced by Opi use |
| P8: Physical hazard caused by Coc use | P19: Physical hazard caused by Opi use |
| P9: Coc use knowing it threatening health | P20: Opi use knowing it threatening health |
| P10: Coc tolerance | P21: Opi tolerance |
| P11: Coc withdrawal syndrome | P22: Opi withdrawal syndrome |



Conclusion

- ▶ We adopted a matrix completion approach to infer SUD criteria using both correlation among criteria of different conditions and genotypes as side information.
- ▶ By imposing sparse prior on the model parameters, the method can find a sparse interactive matrix that connects specific genotypes to diagnostic criteria.
- ▶ We introduced an efficient stochastic LADMM algorithm to solve the optimization problem in this method.
- ▶ The empirical evaluation shows that our method can significantly enhance the running efficiency with minimal adverse effects on the imputation accuracy.



Any Questions?

Thank you.