



ČVUT

ČESKÉ VYSOKÉ
UČENÍ TECHNICKÉ
V PRAZE

Directional association inference challenged by severe dropout in single-cell RNAseq data

Eliška Dvořáková

dvorael1@fel.cvut.cz

supervisor: Prof. Joe Song

Motivation

- **Single-cell RNA sequencing (scRNA) allows to gain RNA-seq from small amount of initial material. (single cell)**
- **New field of study focus on data from single cell.**
- **scRNA suffers from 0.9 dropout.**
- **Methods for Bulk RNA-seq performs close to random**
- **New methods FunChisq [1].**



ČVUT

ČESKÉ VYSOKÉ
UČENÍ TECHNICKÉ
V PRAZE

Bulk RNA-seq

- **Input material: multiple cell**
- **No dropout**
- **1000 samples**

Bulk			
Row	105	117	116
	139	102	117
	99	110	95
		Column	

Single-cell RNA-seq

- **Input material: single cell**
- **0.9 dropout**
- **1000 samples**

		scRNA	
Row		35.1	26.1
	36.1	2.1	1.1
	27.1	4.1	0.1
		Column	



ČVUT

ČESKÉ VYSOKÉ
UČENÍ TECHNICKÉ
V PRAZE

Methods

- 1. Pearson's Correlation test**
- 2. Mutual Information**
- 3. Conditional Entropy (directional)**
- 4. Pearson's χ^2 -test**
- 5. Functional χ^2 -test (directional) [1]**

Simulated dataset

- **Size: 200 tables**
- **Dimension : 3x3**
- **Samples: 1000**
- **Noise: 0.0, 0.2**
- **Dropout: 0.2, 0.8, 0.9, 0.99**
- **Configurations:**
 - 1. Detection of relationship**
 - 2. Detection of relationship direction**

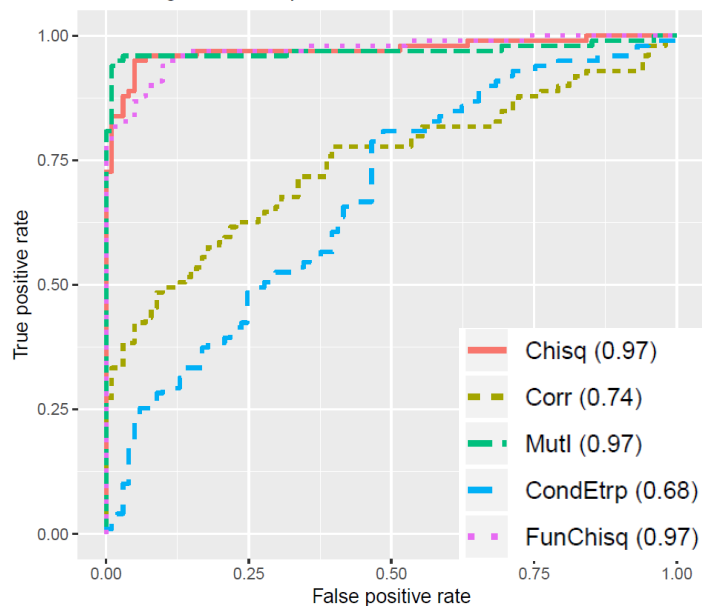


ČVUT

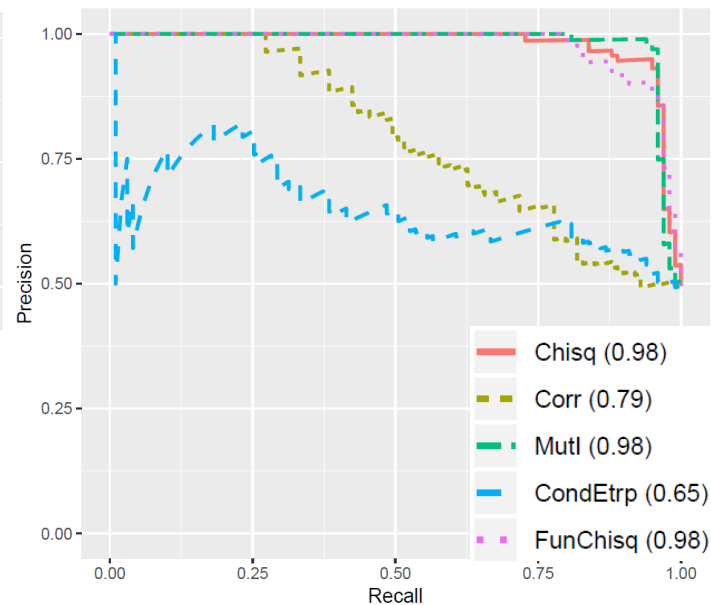
ČESKÉ VYSOKÉ
UČENÍ TECHNICKÉ
V PRAZE

Results 1st Configuration

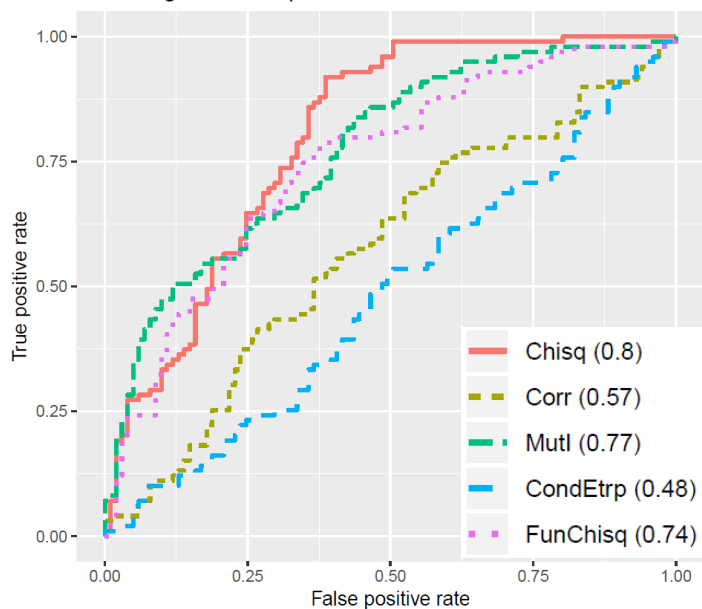
1st configuration dropout: 0.8



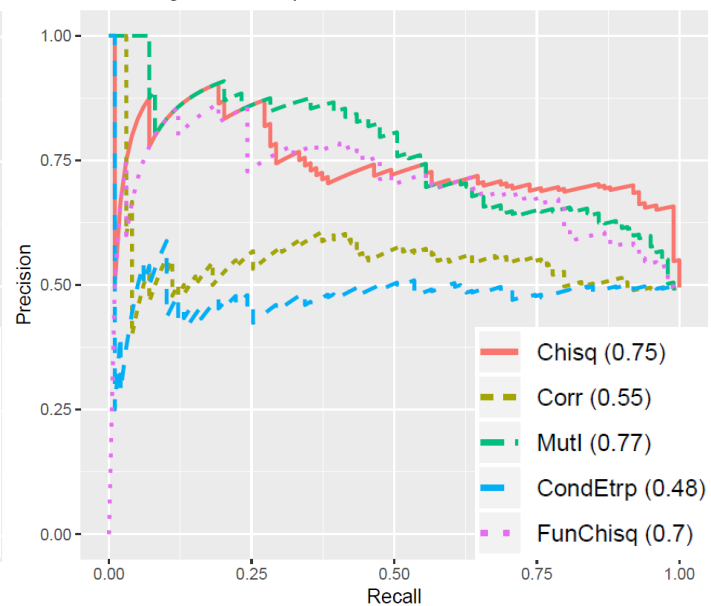
1st configuration dropout: 0.8



1st configuration dropout: 0.9



1st configuration dropout: 0.9



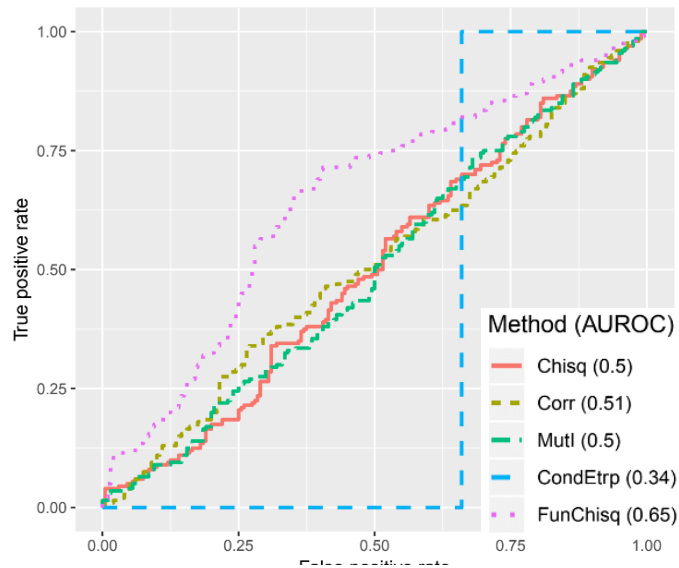


ČVUT

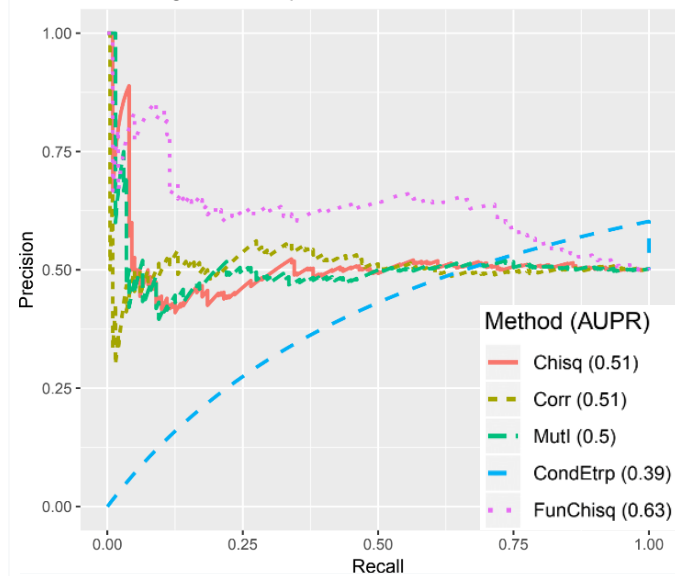
ČESKÉ VYSOKÉ
UČENÍ TECHNICKÉ
V PRAZE

Results 2nd Configuration

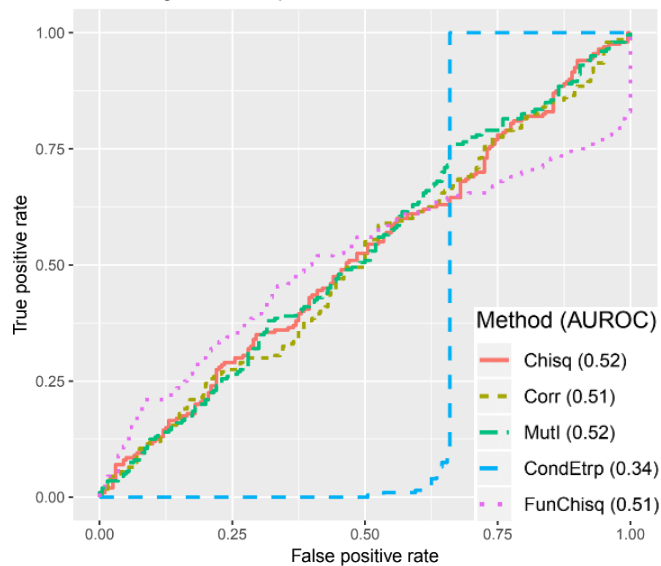
2nd configuration dropout: 0.8



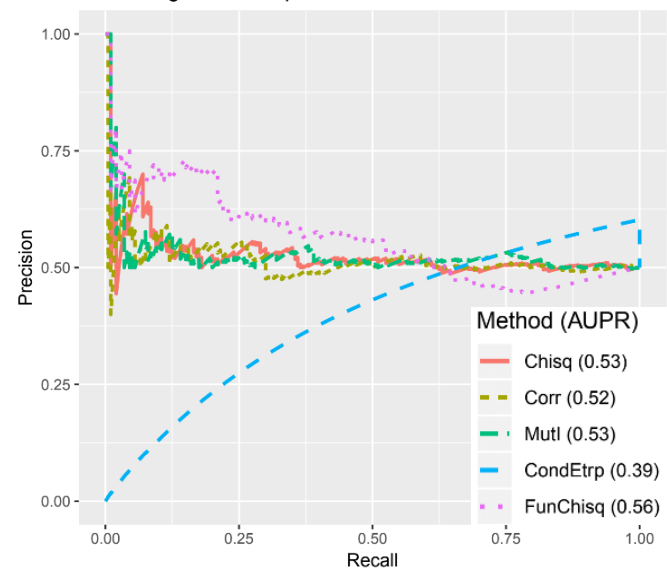
2nd configuration dropout: 0.8



2nd configuration dropout: 0.9



2nd configuration dropout: 0.9



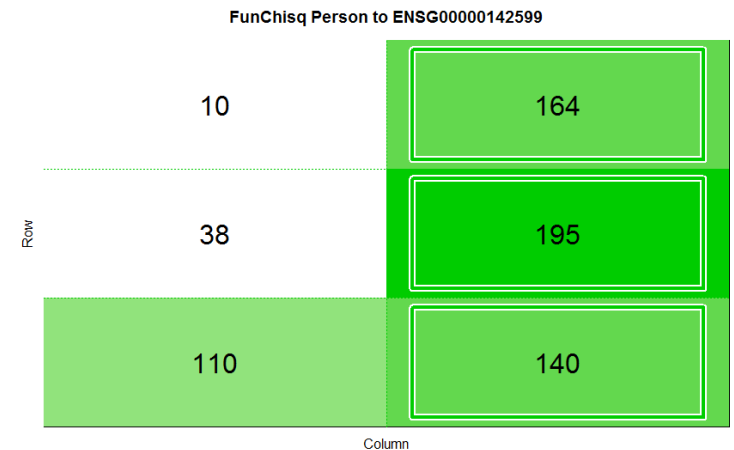


ČVUT

ČESKÉ VYSOKÉ
UČENÍ TECHNICKÉ
V PRAZE

Real dataset biological background

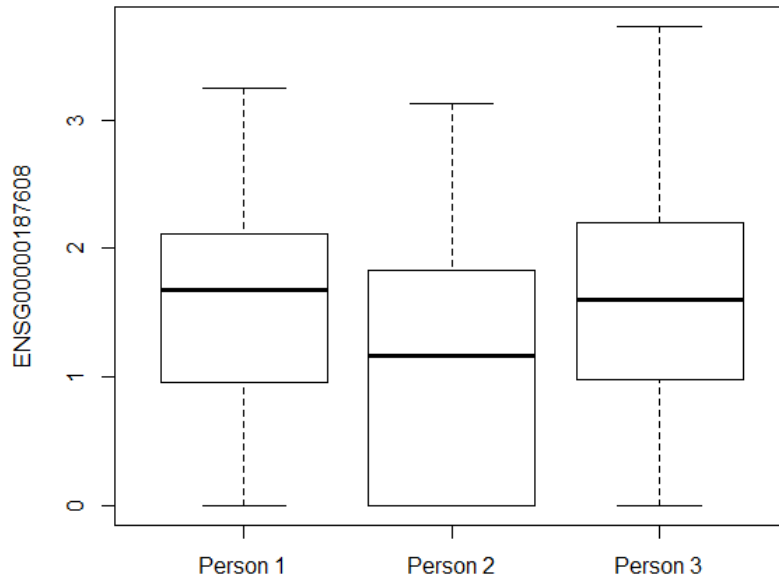
- **Dataset gained from Tung et al. [2]**
- **single-cell Fluidigm C1 platform**
- **three C1 replicates from three human induced pluripotent stem cell (iPSC) lines**
- **unique molecular identifiers (UMI) to all samples**
- **The paper is focused on finding source of variation in gene expression data.**
 - **Genotype**
 - **UMI counts are biased**



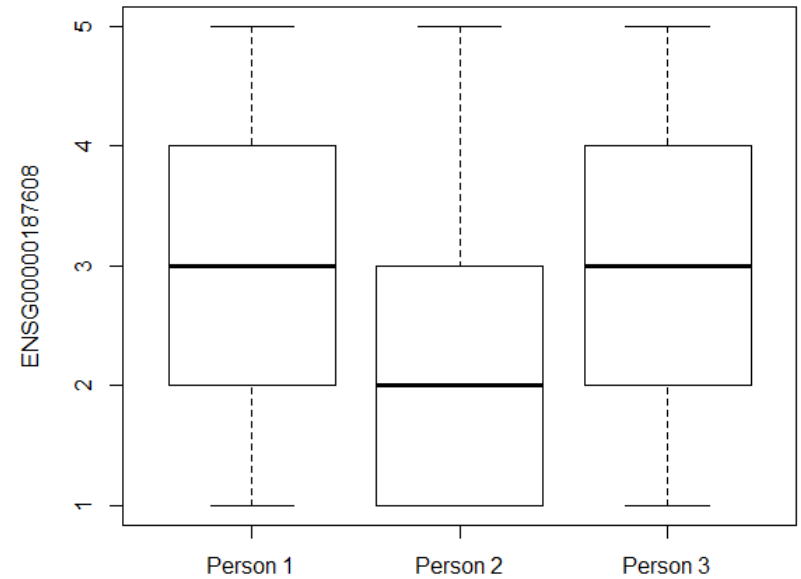


Data discretization

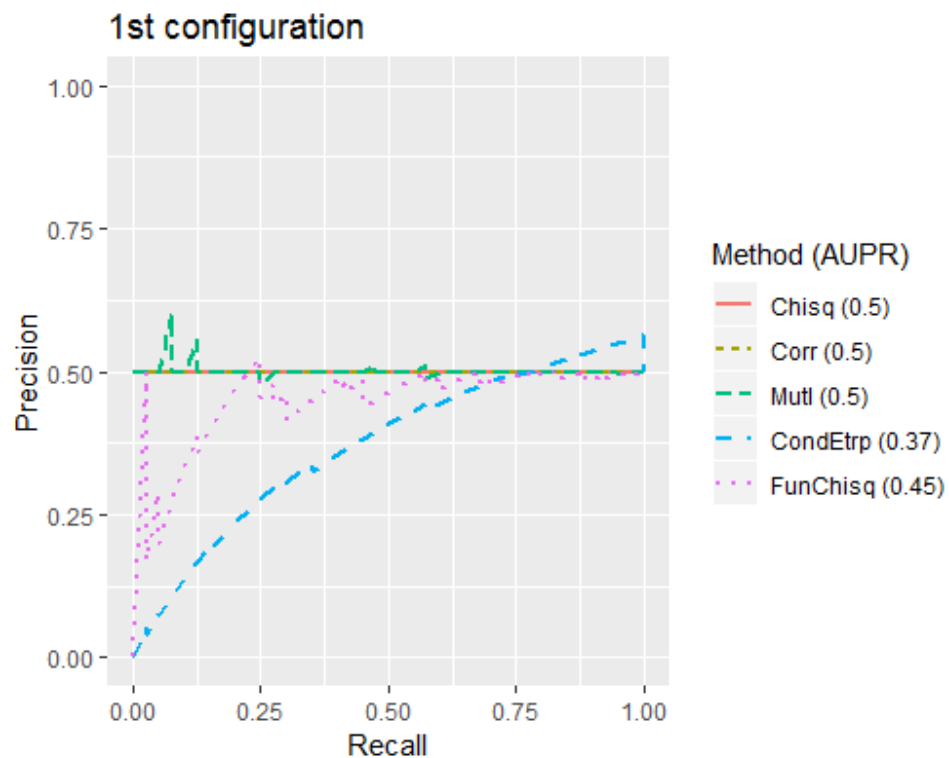
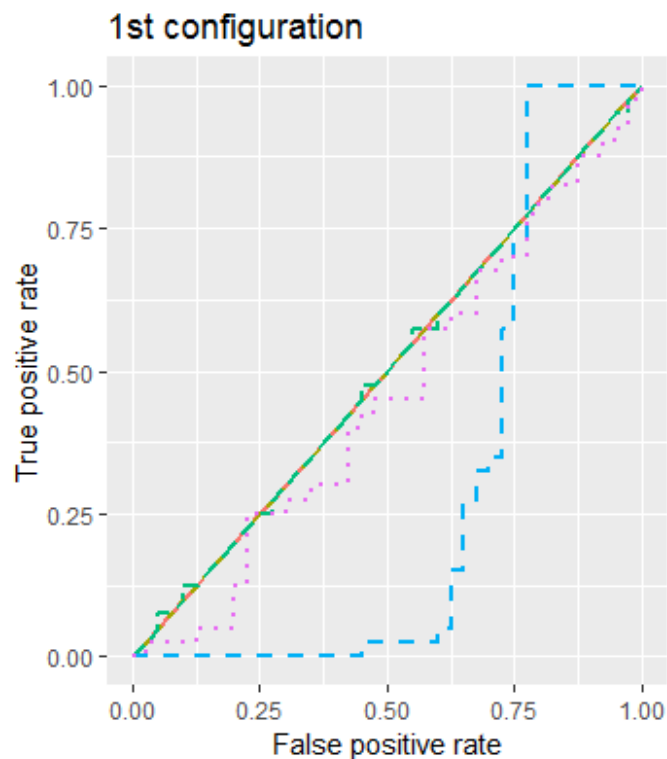
Before discretization



Discretized



Results 1st configuration



Tables with strongest inference for each method

Corr Person to ENSG00000142599

	10	164
38		195
110		140

Row

Column

Chisq Person to ENSG00000142871

	15	84	69	6
3	17	93	120	
	20	92	111	27

Row

Column

Mutl Person to ENSG00000142871

	15	84	69	6
Row	3	17	93	120
	20	92	111	27

Column



ČVUT

ČESKÉ VYSOKÉ
UČENÍ TECHNICKÉ
V PRAZE

Tables with strongest inference for each method

CondEtrp ENSG00000215910 to Person

Row	159	214	237
	15	19	13
Column			

FunChisq ENSG00000142871 to Person

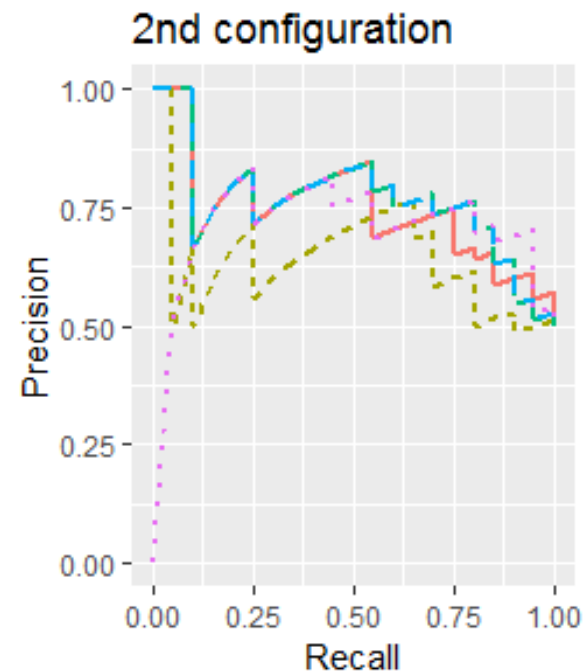
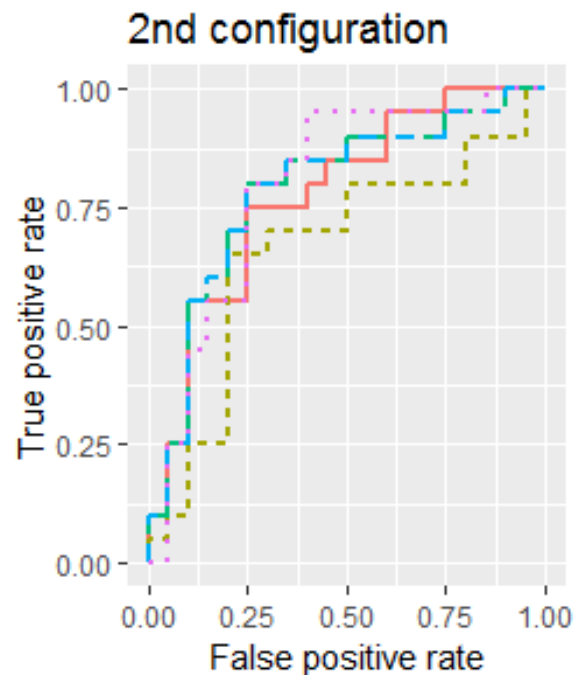
Row	15	3	20
	84	17	92
69	93	111	
	6	120	27
Column			



ČVUT

ČESKÉ VYSOKÉ
UČENÍ TECHNICKÉ
V PRAZE

Results 2nd configuration



Conclusion

- **Simulated dataset**
- **Real dataset**
- **Future work**
 - **Answer How many samples are needed to gain at least 70% accuracy?**
 - **Add test for multiple noise**
 - **Process whole real dataset**

References

- [1] Zhang, Y., & Song, M. (2013). Deciphering interactions in causal networks without parametric assumptions. *arXiv preprint arXiv:1311.2707*.**
- [2] Tung, P.Y., Blischak, J. D., Hsiao, C. J., Knowles, D. A., Burnett, J. E., Pritchard, J. K., & Gilad, Y. (2017). Batch effects and the effective design of single cell gene expression studies. *Sci Rep*, 7, 39921.**