

DISNEP simulation

ZHUOHUI Liang

7/15/2021

Simulation Settings

Table 1: Simulation settings

scenario	1.00	2.00	3.00	4.00	5.00
n_signal	20.00	20.00	20.00	20.00	20.00
n_noise	980.00	980.00	980.00	980.00	980.00
strong_signal_mu	10.00	10.00	10.00	11.00	12.00
median_signal_mu	9.00	9.00	9.00	10.00	11.00
noise_mu	8.00	8.00	8.00	8.00	8.00
strong_corr	0.10	0.10	0.10	0.10	0.10
median_corr	0.02	0.02	0.02	0.02	0.02
noise_corr	0.01	0.01	0.01	0.01	0.01
sigma	3.00	2.00	1.00	1.00	1.00
prob_jump	0.50	0.50	0.50	0.50	0.50

For each simulation, there will be two data/matrix simulated, namely gene-gene-interaction network and gene-expression data. Each simulation we will compare t-test, GeneWanderer(random walk with restart on gene-gene-interaction network) and Disnep(random walk with restart on disease-enhanced-gene-gene-interaction network[Ruan,2019])

Gene-gene-interaction network

In each simulation, network is sample from existing database, we have homo sapiens(human) gene-gene(protein-protein) interaction network(PPI) downloaded from V.11 STRING. Since Genewanderer only consider gene-gene-interaction score > 0.4 , we use Cytoscape to pull STRING node information from STRING containing only score > 0.4 and map Gene name to PPI.

And we set the gene from DISGENET V7 renal carcinoma disease(C1378703) as disease/signal gene list and the rest of the gene not in DISGENET list as noise gene.

We sample 20 gene from disease gene list and 980 as noise gene. We keep the edge of the gene from STRING PPI network and set all missing value as 0. And we order the first 20 gene as signal gene just for convenient.

Gene Expression network

In each simulation, we generate 300 case and 50 control sample.

Case

The case set to have 10 strong signal gene, which each gene's mean is generated from $N(\mu_1, 1)$, and use this mean to generate gene expression level. case has another 10 median signal gene, which each gene's mean is generated from $N(\mu_2, 1)$. the rest 980 gene are set as noise gene, which's mean is generate from $N(8, 1)$.

The correlation between strong signal is sampled from $U(0, 0.1)$ and correlation between median signal is sampled from $U(0, 0.02)$ and correlation between noise is sampled from $U(0, 0.01)$.

$$Cov = \sigma^2 * \begin{pmatrix} 1 & \rho_{12} & 0 & 0 \\ \rho_{21} & 1 & \rho_{23} & 0 \\ 0 & \rho_{32} & 1 & \rho_{34} \\ 0 & 0 & \rho_{43} & 1 \end{pmatrix} \quad (1)$$

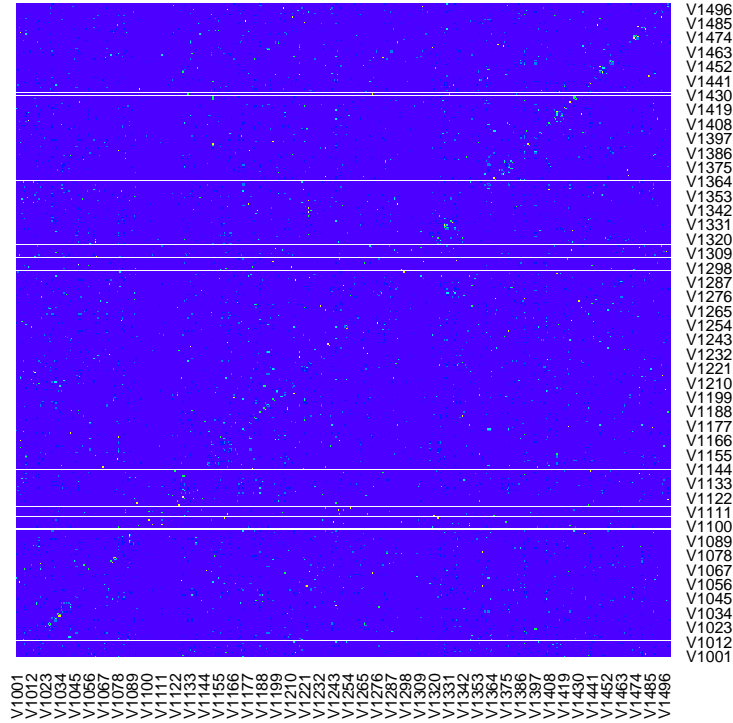
where $\rho_{12} = \rho_{21}^t$ is the triangle block sample from $U(0, 0.1)$ and vise versa for ρ_{23}, ρ_{34}

control

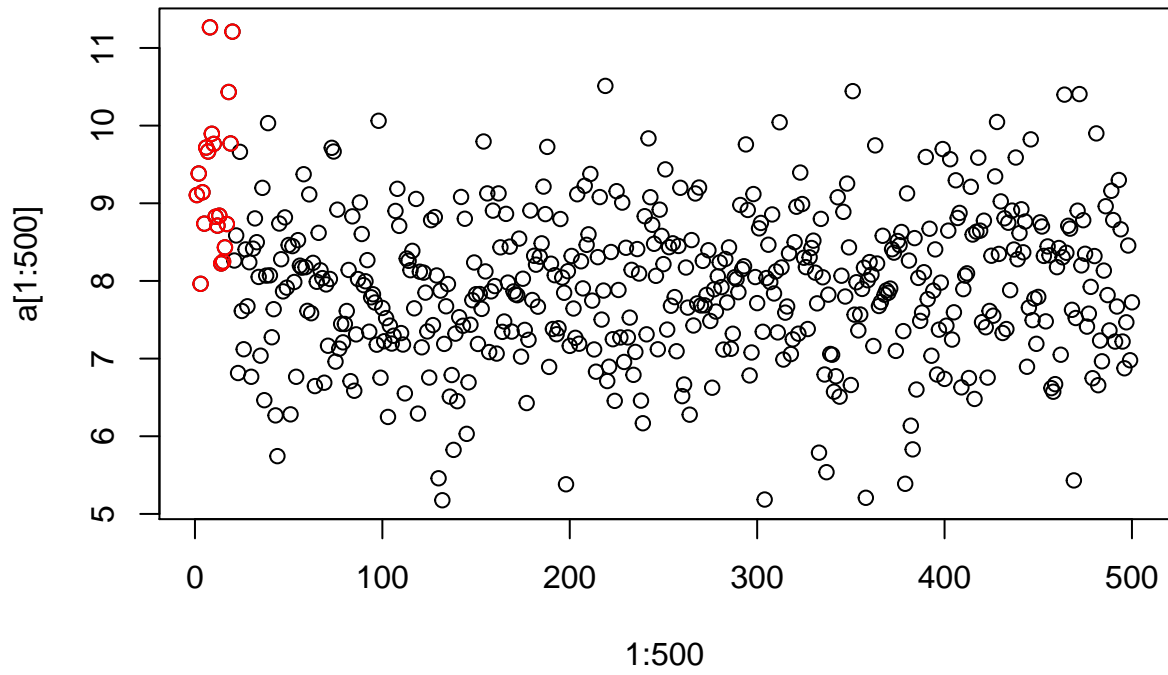
All genes in control is set to be noise and its mean generate from $N(8, 1)$ and correlation sample from $U(0, 0.01)$

all gene expression is generated from multivariate normal, with generated means and covariate as mention above.

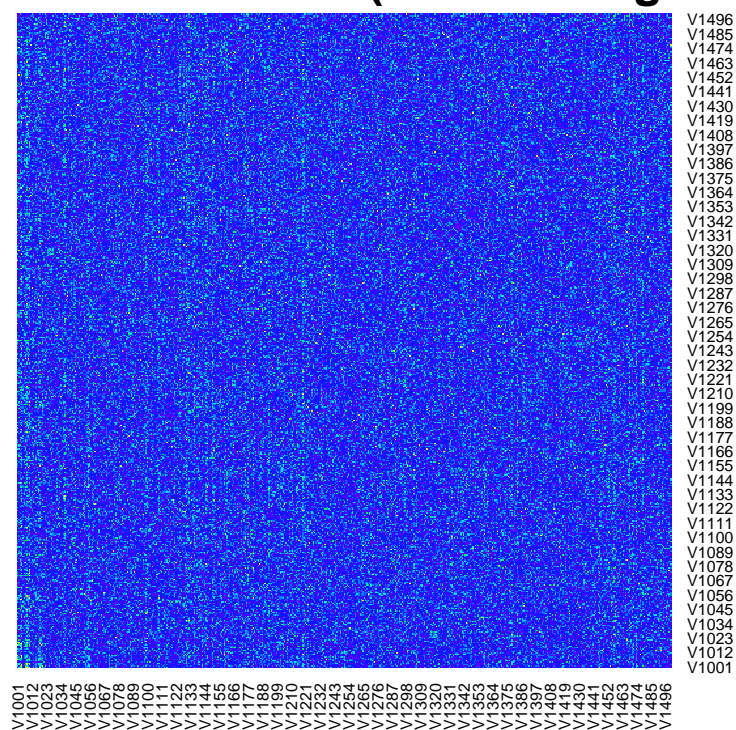
Interaction first 500(included signals)

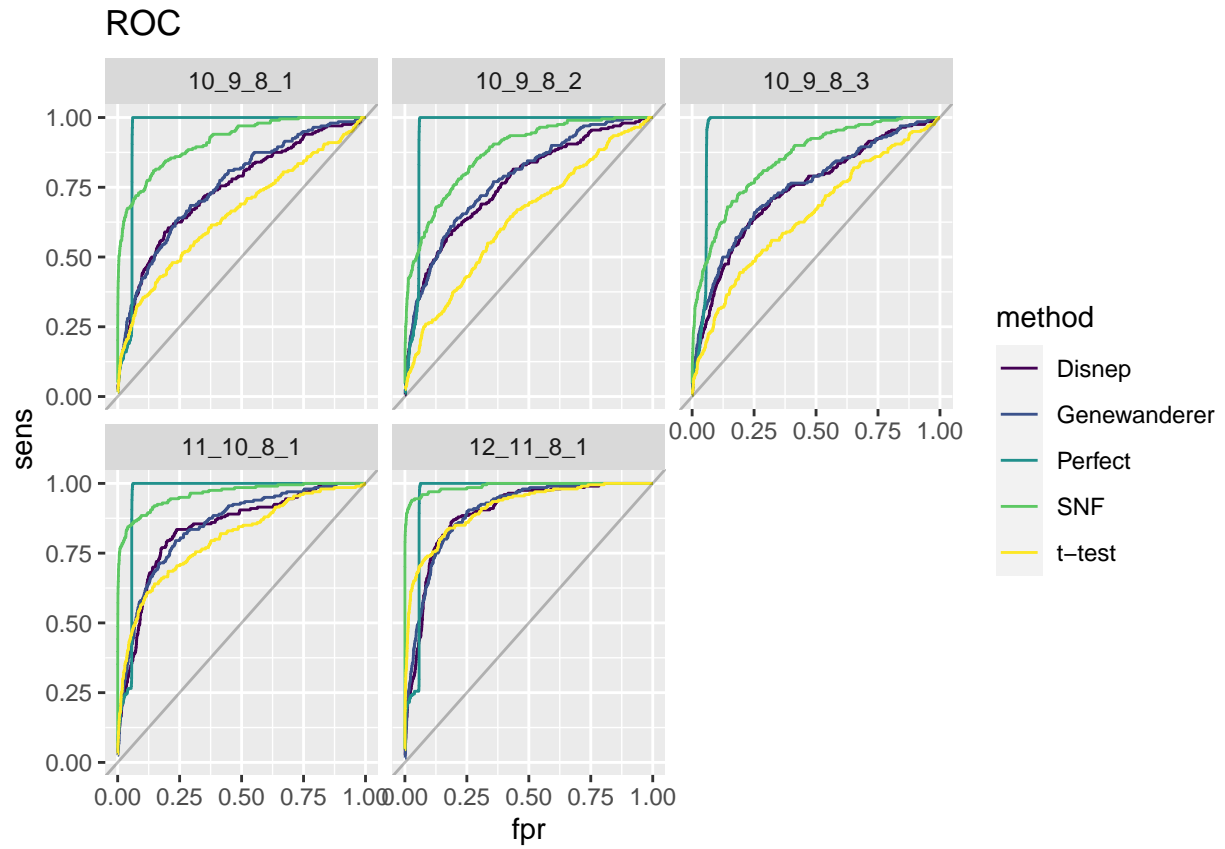


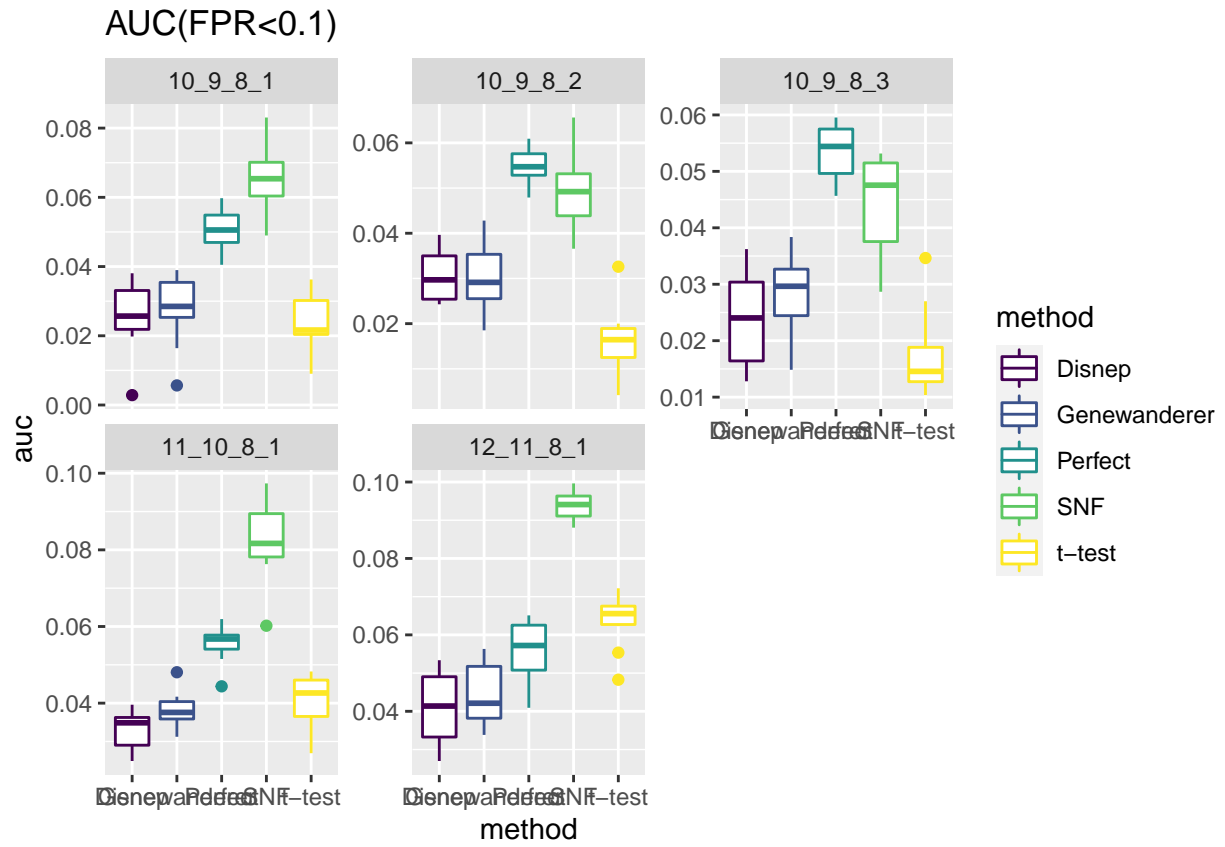
gene expression with mu 10-9-8



Correlation first 500(included signals)







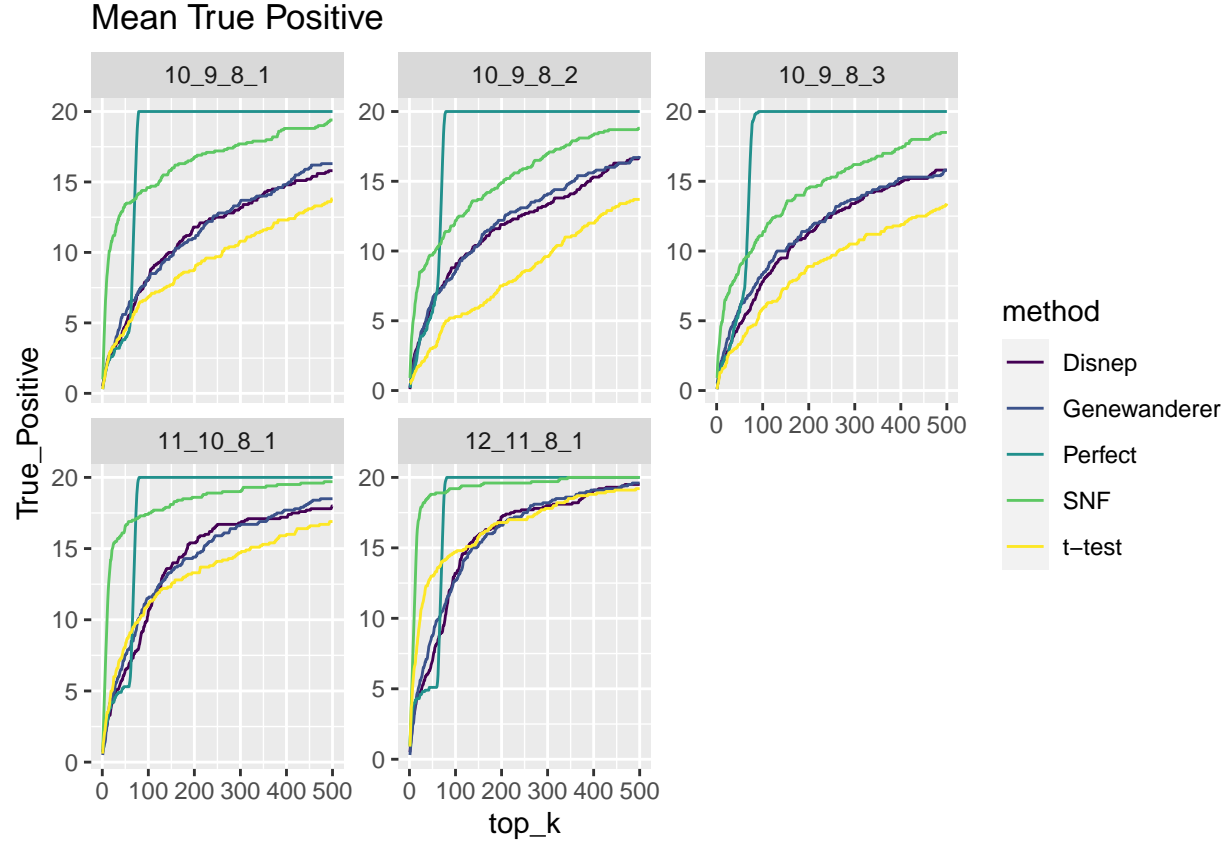


Table 2: mean AUC(FPF<0.1)

mu_sigma	Disnep	Genewanderer	Perfect	SNF	t-test
10_9_8_1	0.026	0.028	0.051	0.066	0.023
10_9_8_2	0.030	0.030	0.055	0.050	0.016
10_9_8_3	0.024	0.028	0.053	0.044	0.018
11_10_8_1	0.033	0.038	0.056	0.083	0.041
12_11_8_1	0.041	0.044	0.056	0.094	0.064