



# Statistical Evaluation for Methods of Gene-set Analysis with Multivariate Non-normal Scenarios

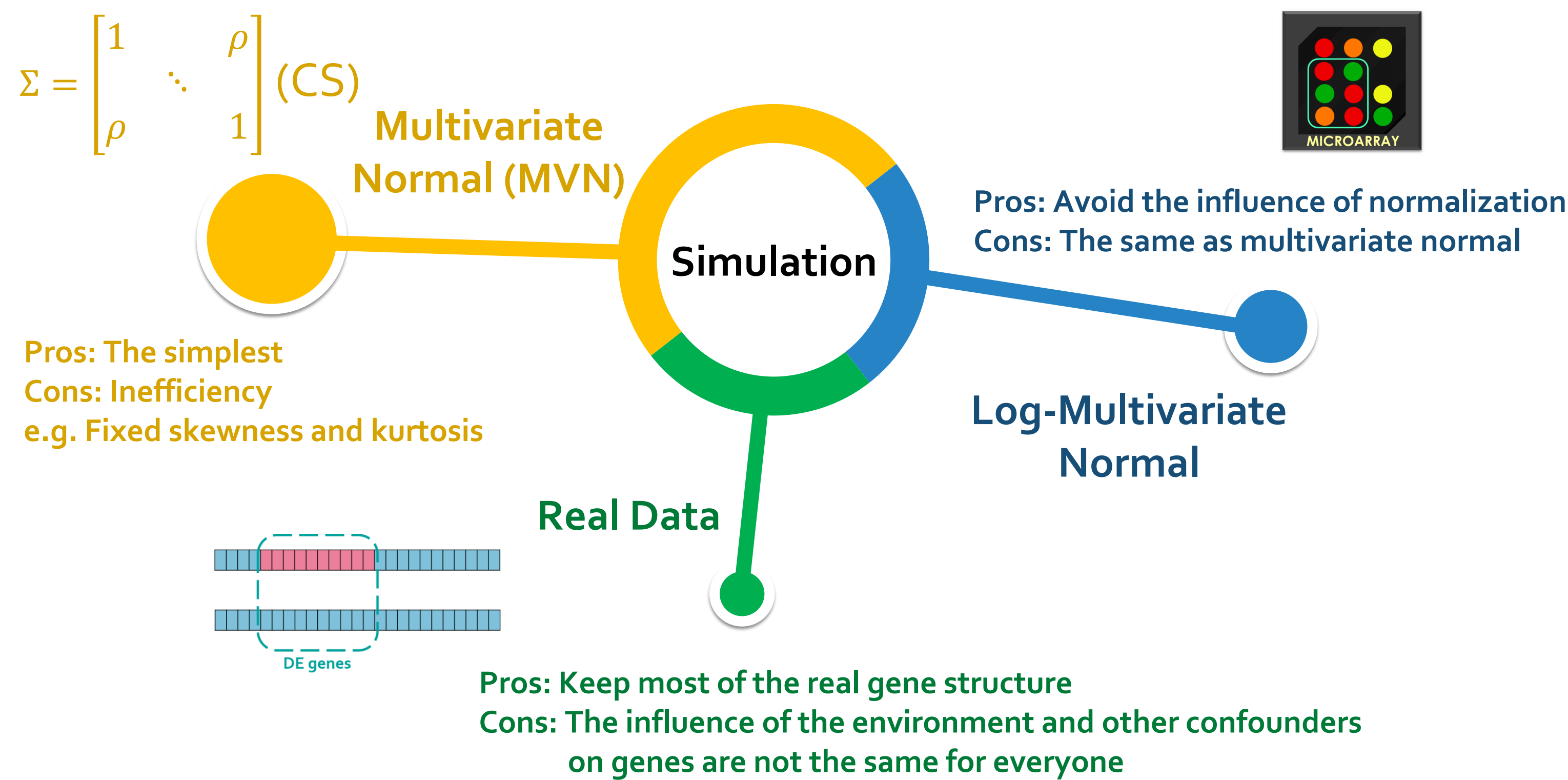
## 以非常態情境評估基因集合分析方法在真實基因資料下之表現研究

Presenter: Chi-Hsuan Ho Advisor: Prof. Chuhsing Kate Hsiao

### Introduction

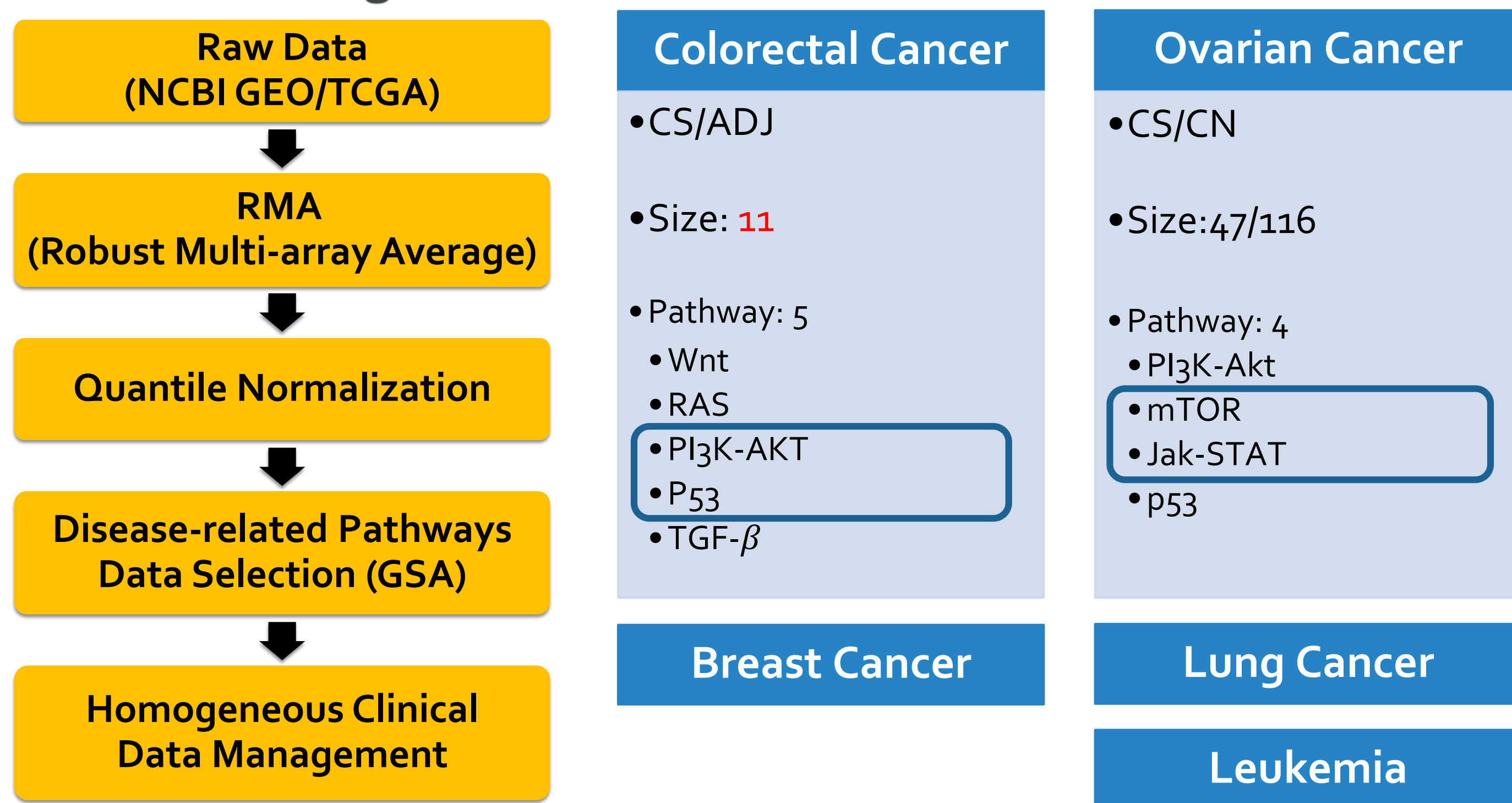
As the technology improves, more and more statistical methods for gene-set analysis (GSA) are developed to find pathogenic pathways and genes. In recent years, many studies use the multivariate normal (or lognormal) distribution in simulation studies to evaluate the performance of these GSA methods. However, the normality assumption for the gene expression data has been questionable. Therefore, here we first reveal the non-normality of the real data, and then propose a statistical evaluation method to compare five well-known GSA methods via radar plots.

#### Simulation methods used in current GSA

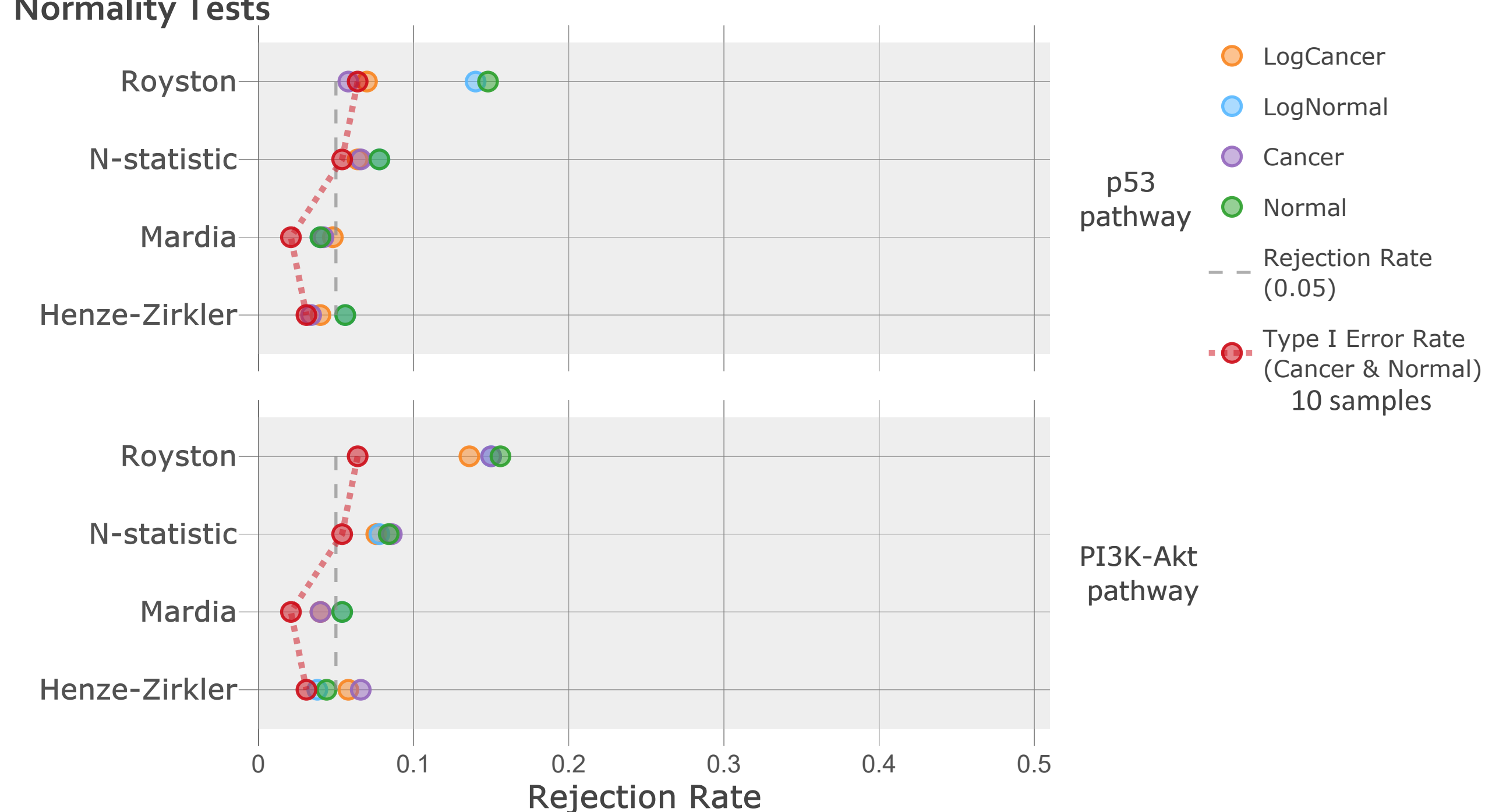


### Non-normality of Real Data

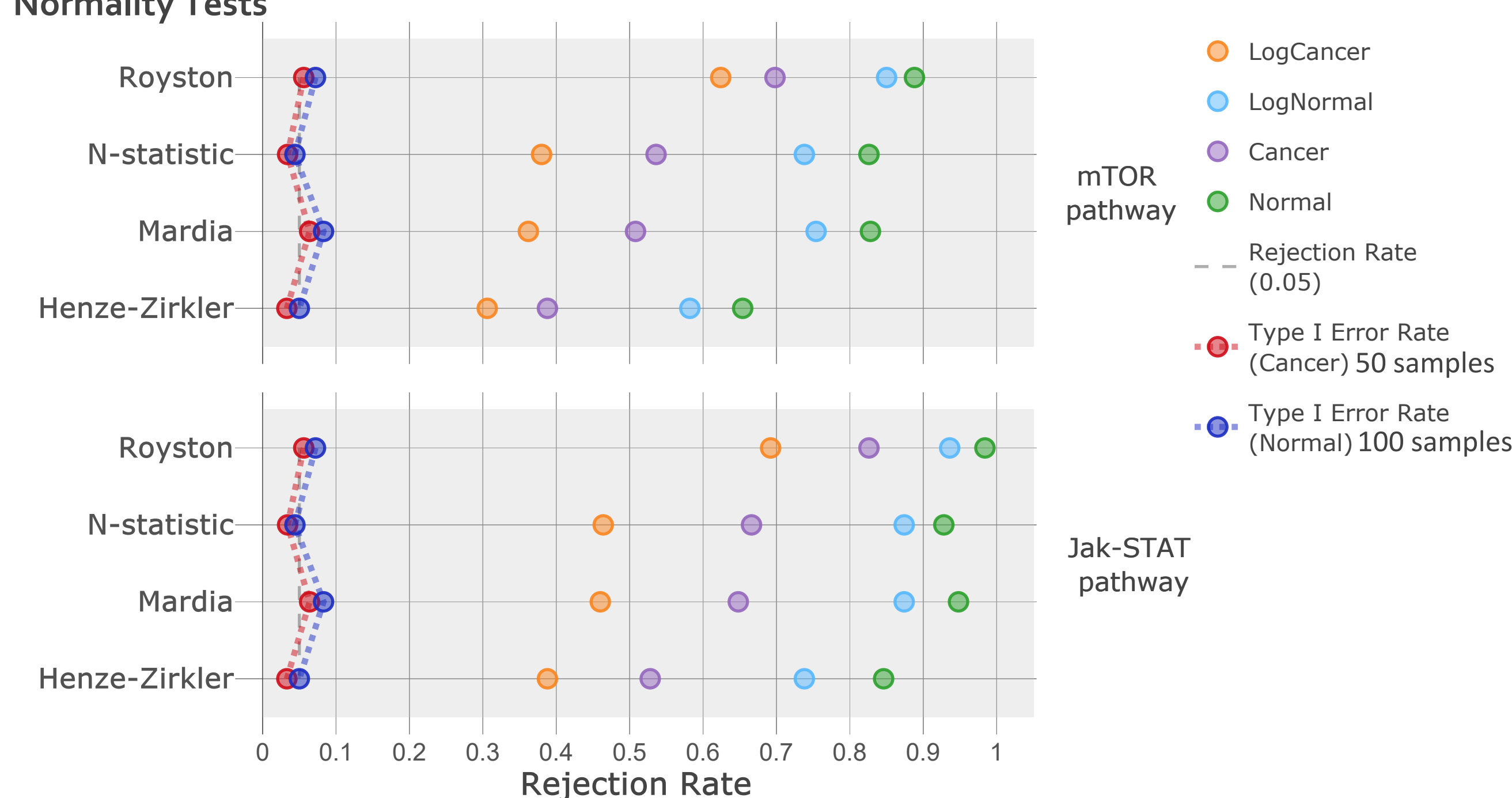
#### Data Processing



#### Colorectal Cancer Data Result (Dot Plot)



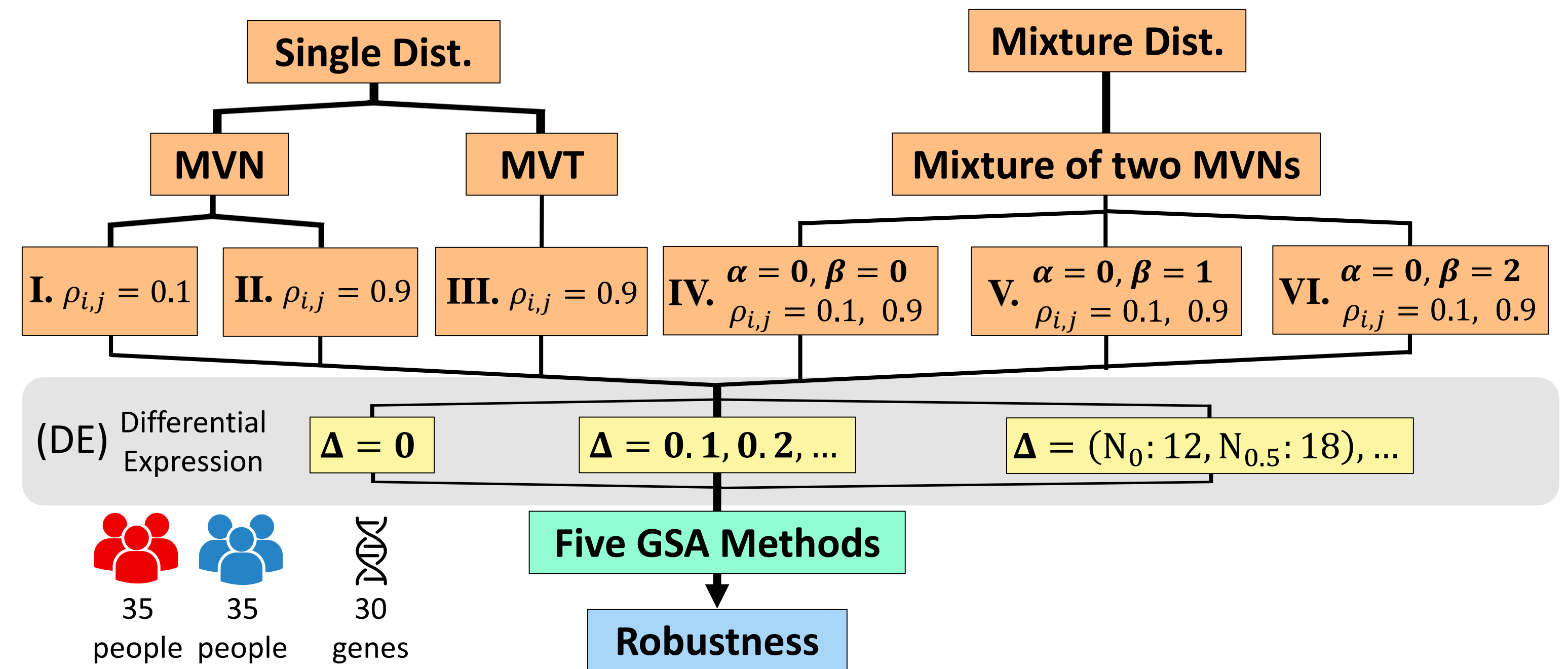
#### Ovarian Cancer Data Result (Dot Plot)



- We sampled 4 genes from each pathway to test normality (repeat 500 times).
- Gene expression data are unlikely normally (log-normally) distributed.
- The test result of colorectal cancer data set is influenced by sample size.

### Statistical Evaluation for GSA Methods

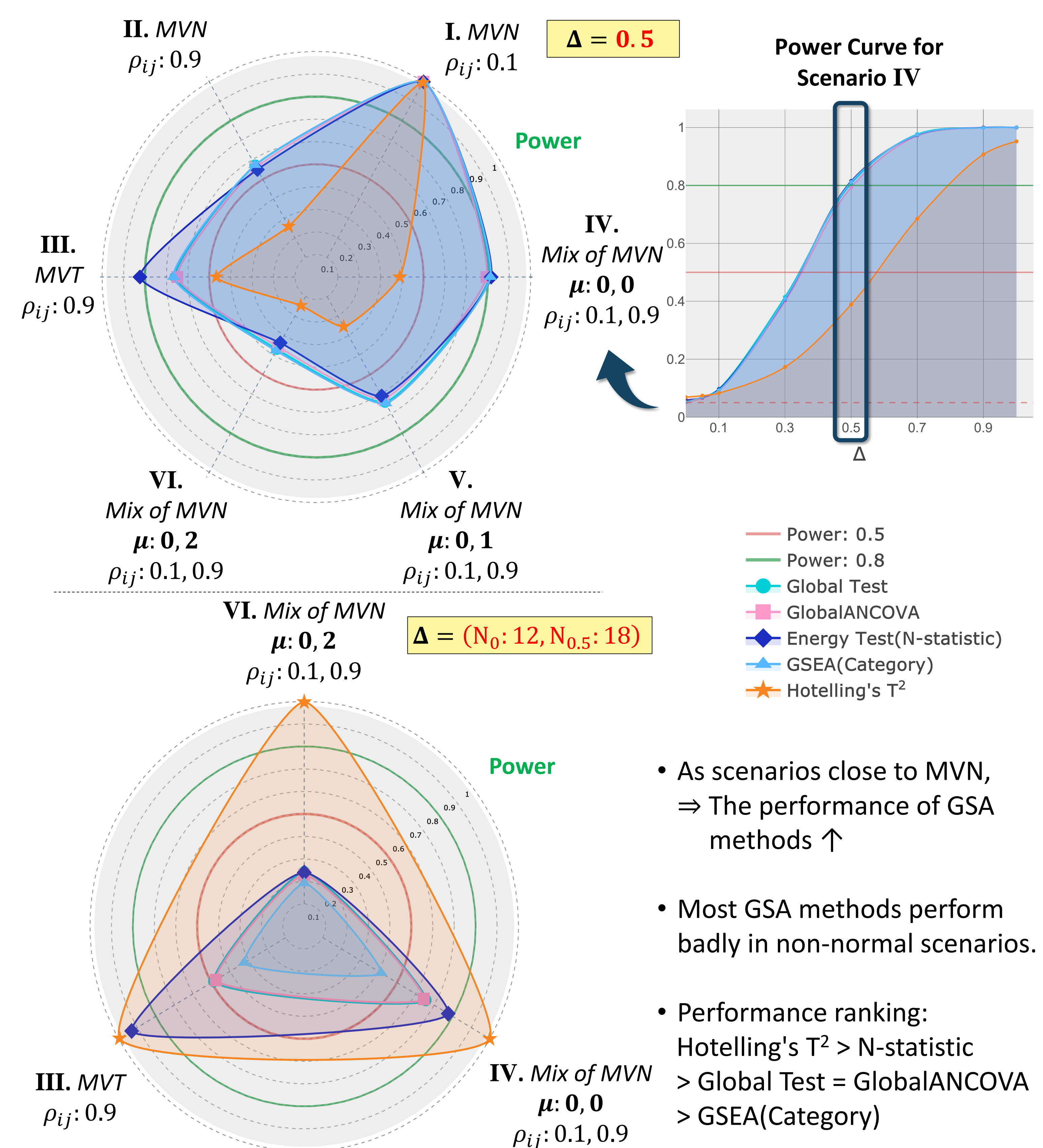
#### Evaluation Procedure:



#### Methods of Gene-set Analysis (GSA Methods)

Distance-based	Regression-based	Other
$D(dist_{CS}, dist_{CN})$ • <b>Hotelling's T<sup>2</sup></b> → Multidimensional version of two-sample t-test (Mahalanobis distance)	• <b>The Global test</b> → Score test for detecting the random effect of each gene (Y: Phenotype, X: Genes) • <b>The GlobalANCOVA</b> → Linear regression model for each gene (RSS) (Y: Genes, X: Confounder)	• <b>GSEA (Category)</b> → Self-contained version of GSEA (Running-sum statistic: Per-gene t-statistic)

#### Evaluation Results (Radar Plots)



- As scenarios close to MVN,  $\Rightarrow$  The performance of GSA methods  $\uparrow$
- Most GSA methods perform badly in non-normal scenarios.
- Performance ranking:  
Hotelling's T<sup>2</sup> > N-statistic  
> Global Test = GlobalANCOVA  
> GSEA(Category)

### Conclusions & Discussion

- Three elements in our evaluation method are important:
  - Correlations between genes ( $\rho_{ij}$ )
  - The choice of scenarios
  - The structure of difference vector ( $\Delta$ )
- Some advice for researchers to develop GSA methods:
  - Distribution-free methods are preferred.
  - Binary phenotypes: Multivariate distance-based methods perform better than regression-based methods.
- Some settings can be considered in the future:
  - Real data scenarios
  - Copula-based scenarios
  - DE  $\Rightarrow$  Correlation matrix or distribution