# **Pathway and Gene Set Analyses**

Self contained and competitive tests

Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017

## What is a pathway?

No clear definition!

Wikipedia: "In <u>biochemistry</u>, **metabolic pathways** are series of <u>chemical</u> reactions occurring within a cell. In each pathway, a principal chemical is modified by <u>chemical reactions</u>." These pathways describe enzymes and metabolites.

But often the word "pathway" is also used to describe gene regulatory networks or protein interaction networks.

In all cases a pathway describes a biological function very specifically.

### What is a gene set?

Just what it says: a set of genes!

All genes involved in a pathway are an example of a gene set. All genes corresponding to a Gene Ontology term are a gene set. All genes mentioned in a paper of Smith et al might form a gene set.

A gene set is a much more general and less specific concept than a pathway.

Still: we will sometimes use two words interchangeably, as the analysis methods are mainly the same.

Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017

[LM]

## Gene set / pathway analysis

The aim is to give one number (e.g. a p-value, a score) to a gene set / pathway:

Are many genes in the pathway differentially expressed (upregulated / down-regulated)?

Can we give a number (a p-value) as the probability of observing changes of this magnitude (or larger) just by chance?



Conversion can lead to errors!

There are many more resources out there. BioCarta, BioPax, NetPath, ...

Commercial packages often use their own pathway / gene set definitions. Ingenuity, Metacore, Genomatix, ...

### Self contained versus competitive tests

The distinction between "self-contained" and "competitive" methods goes back to Goeman and Buehlman (2007). PMID 17303618.

A **self-contained** method only uses the values for the genes of a gene set.

A **competitive method** compares the genes within the gene set with the other genes interrogated (array or sequencing).

Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017



> gsids[131:135] \$chryq11 [1] 2260 3798 3963 3849 4862 4629 4826 4827 4629 4826 4827 4629 4826 4827 4629 4826 4827 4856 4862 4855 2611 4842 3963 4856 5968 3742 4842 [28] 4842 4842 4842 4842 4842 4842 4470 \$chr15q13 [1] 5667 7986 5667 1355 6917 6464 7436 1952 6623 881 \$chr3q22 [1] 5145 7328 7840 1042 5016 6193 1235 7016 5465 3678 3752 3160 1368 5574 7077 \$chr6q15 [1] 8416 1607 8212 8082 \$chr4q31

[1] 2635 3295 4107 7827 6266 4581 5543 6509 5949 4313 2138 5751 2295 3258 118 481 2930 4522 812 118 481 1908 2863 7970 7084 5880 2785 [28] 2799 2907

## Enrichment

### Are genes on chromosome yq11 enriched?

	Not in gene set	In gene set
Not differentially expressed	8796	9
Differentially expressed	11	4

**Self contained**: Under the null, we would expect 13\*0.05 < 1 gene to be significant. Using a Binomial(13,0.05) distribution, the probability that 4 or more genes are differentially expressed is 0.003.

**Competitive**: Under the null, we would expect that virtually all differentially expressed genes are not in the gene set. Fisher's exact test gives a p-value of less than 1e-8.

Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017





### Wilcoxon test for differential expression



#### [ RI ]

### A simple test for a mean shift

The simplest statistic to test for a mean shift is the average difference in mean. One way to summarize this difference is to average the t-statistics.

$$ar{t} = rac{1}{N}\sum_{i\in G} t_i$$
 ( with N being the size of gene set G ).

Under the null, the t statistics have mean 0 and SD 1.

If they are independent, then:

$$\sqrt{N} \ \bar{t} \sim N(0,1)$$



## Variance of average of correlated variables

$$\operatorname{var}(\bar{t}) = \frac{1}{N^2} \operatorname{var}\{(1\dots 1)(t_1\dots t_N)'\}$$
$$= \frac{1}{N^2}(1\dots 1) \begin{pmatrix} 1 & \rho & \dots & \rho \\ \rho & 1 & \dots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \dots & \rho & 1 \end{pmatrix} (1\dots 1)'$$
$$= \frac{1}{N}\{1 + (N-1)\rho\}$$

Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017

**Correction factor** 

$$\frac{\sqrt{N}}{\sqrt{1+(N-1)\bar{r}}}\,\bar{t}$$

Here,  $\bar{r}$  is the average pairwise correlation within gene set.

Note that there are also plenty of other approaches for dealing with correlation!

[ RI ]









Ingo Ruczinski | Asian Institute in Statistical Genetics and Genomics | July 21-22, 2017

PMID 27070863

