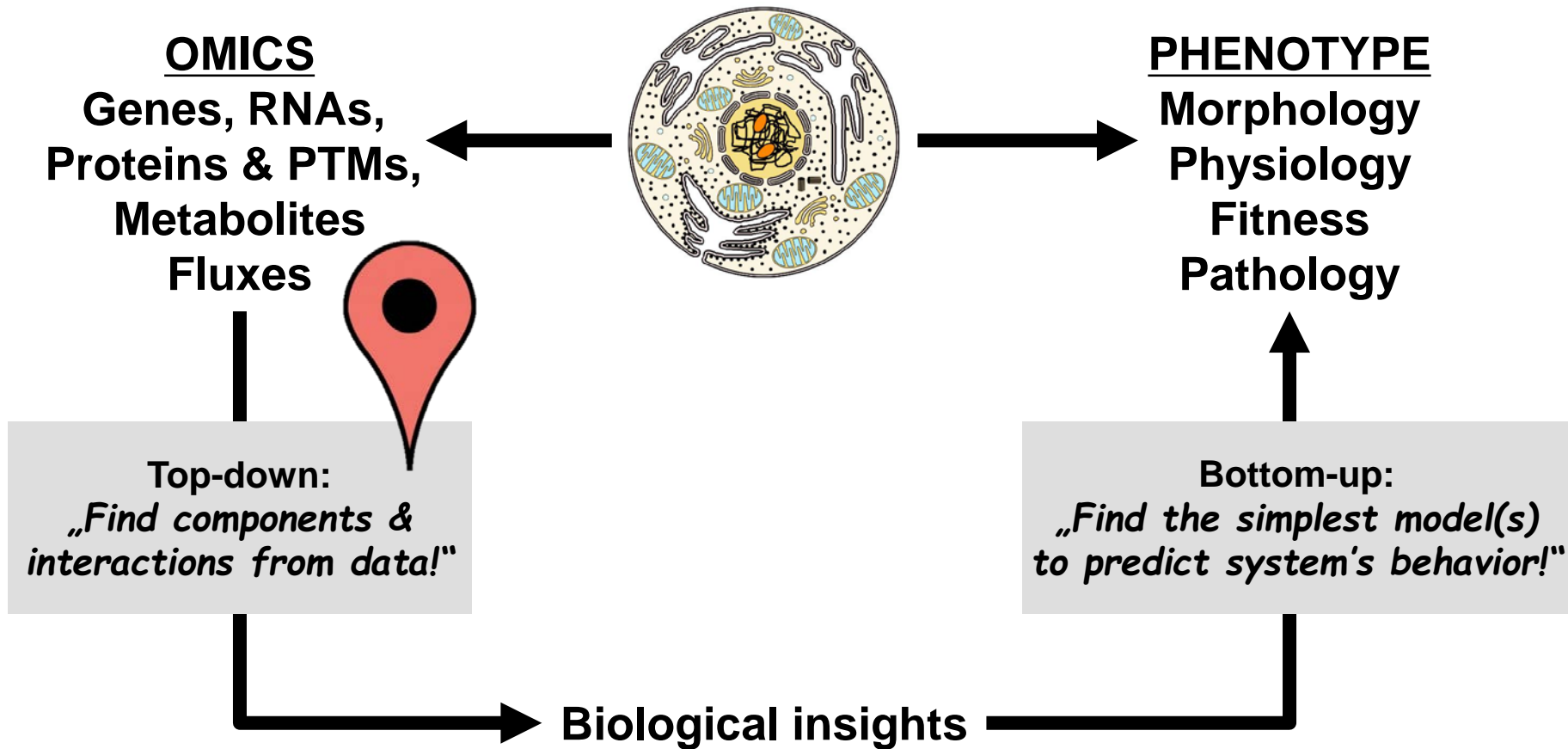# Learning Goals Lecture 8
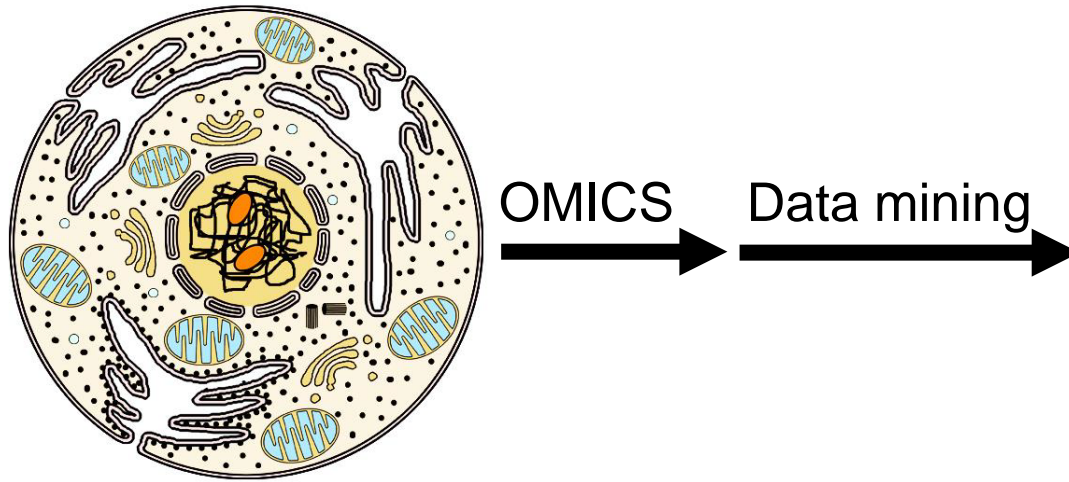
- **Intro on data mining of omics data**

  - Know what biological insights are to be obtained by data mining of omics data
  - Describe the typical problems associated with omics data?
  - Describe the role of statistics in biological context

- **Differential, univariate analysis between two groups**

  - Understand the basic form of univariate analysis
  - Understand pathway Enrichment analysis
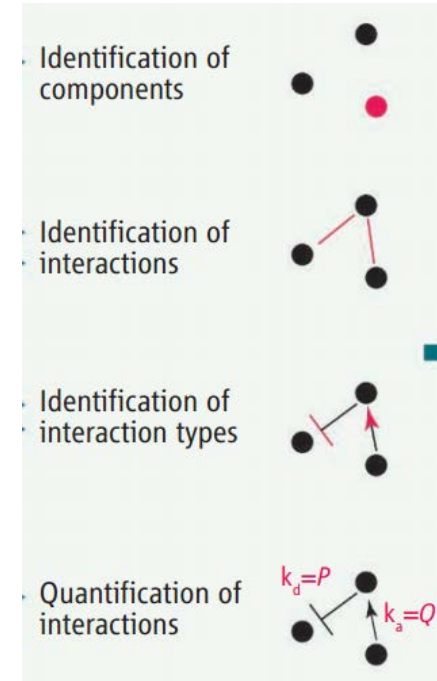  - Correctly assess statistical significance

# The big picture



OMICS
Genes, RNAs,
Proteins & PTMs,
Metabolites
Fluxes

PHENOTYPE
Morphology
Physiology
Fitness
Pathology

Top-down:
„Find components &
interactions from data!"

Bottom-up:
„Find the simplest model(s)
to predict system's behavior!"

Biological insights

# What are *biological insights*?

## 1) Molecular level



Identification of components

Identification of interactions

Identification of interaction types

Quantification of interactions $k_d=P$ $k_a=Q$

OMICS → Data mining →

## 2) Phenotype

e.g. growth, disease, …

# Exemplary BIO questions

COMPONENTS ID

- What is the composition of a cell?

- What is the function of a gene / gene product?

- What are the variants of a protein?

NETWORK RECONSTRUCTION (LINK PREDICTION)

- What cellular regulators are active?

- Infer a transcriptional network

CLASSIFICATION, FEATURE SELECTION

- Group cells/organs/patients based on molecular features (markers)

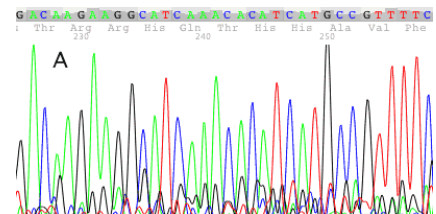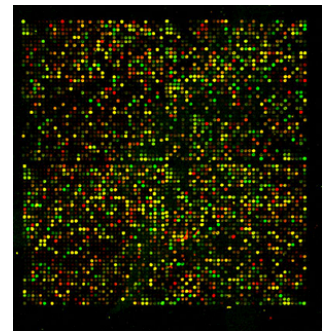- Describe molecular response of cells to treatments
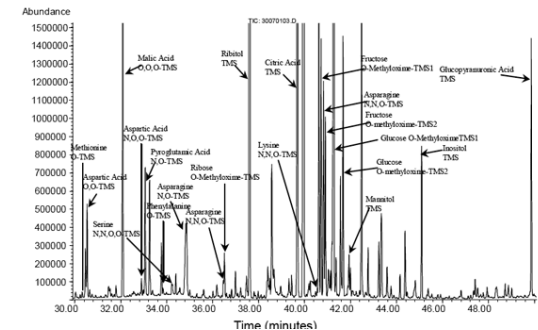
- Predict behavior of cells to treatments

…

# OMICS data are large

It doesn't have to be omics data all the time. However, with the current state of the art in technology, it's fairly <u>common</u> to obtain data on

- **100s** of metabolites/lipids
  - Mass spectrometry
- **1'000s** of proteins
  - Mass spectrometry or antibodies
- **1'000s – 10'000s** genes and RNAs
  - Microarrays and NextGen Sequencing

… **for <u>EACH</u> tested condition.**

# Big data, low sample number

Most frequent problems in analyzing OMICS data:

1) **# of detected features >>> # of samples**

2) Data are noisy

3) Technical reproducibility of experiments

There is a serious risk of **overfitting** the data, i.e. to identify patterns that capture NOISE of the data and are hardly reproducible. Such models are excessively complicated and lead to poor predictions on new data sets.

# Good statistics is a must but not sufficient

Opportunity and quality of statistical tests are under debate…



**SOCIAL SELECTION** — *Popular articles on social media*

## Psychology journal bans *P* values

A controversial statistical test has met its end, at least in one journal. Earlier this month, the editors of *Basic and Applied Social Psychology* (*BASP*) announced that the journal would no longer publish papers containing *P* values, because the values were too often used to support lower-quality research.

Authors are still free to submit papers to *BASP* with *P* values and other statistical measures that form part of 'null hypothesis significance testing' (NHST), but the numbers will be removed before publication. "Basic and Applied Social Psychology just went science rogue and banned NHST from their journal. Awesome," tweeted Nerisa Dozo, a PhD student in psychology at the University of Queensland in Brisbane, Australia. But Jan de Ruiter, a cognitive scientist at Bielefeld University in Germany, tweeted: "NHST is really problematic", adding that banning all inferential statistics is "throwing away the baby with the p-value".

*Basic Appl. Soc. Psych.* 37, **1–2** (2015)

Based on data from altmetric.com. Altmetric is supported by Macmillan Science and Education, which owns Nature Publishing Group.

⟳ **NATURE.COM** For more on popular papers: go.nature.com/ynfi49

---

*Open access, freely available online*

**Essay**

## Why Most Published Research Findings Are False

John P. A. Ioannidis

DOI: 10.1371/journal.pmed.0020124

**Summary**

There is increasing concern that most current published research findings are false. The probability that a research claim is true may depend on study power and

factors that influence this problem and some corollaries thereof.

**Modeling the Framework for False Positive Findings**

Several methodologists have

is characteristic of the field and can vary a lot depending on whether the field targets highly likely relationships or searches for only one or a few true relationships among thousands and millions of hypotheses that may

---

# STATISTICAL ERRORS

*P* values, the 'gold standard' of statistical validity, are not as reliable as many scientists assume.

**BY REGINA NUZZO**

150 | NATURE | VOL 506 | 13 FEBRUARY 2014

For a brief moment in 2010, Matt Motyl was on the brink of scientific glory: he had discovered that extremists quite literally see the world in black and white.

The results were "plain as day", recalls Motyl, a psychology PhD student at the University of

It turned out that the problem was not in the data or in Motyl's analyses. It lay in the surprisingly slippery nature of the *P* value, which is neither as reliable nor as objective as most scientists assume. "*P* values are not doing their job, because they can't," says Stephen Ziliak, an

Goodman, a physician and statistician at Stanford. "Then 'laws' handed down from God are no longer handed down from God. They're actually handed down to us by ourselves, through the methodology we adopt."

DALE EDWIN MURRAY

# Recommended resource

http://www.nature.com/collections/qghhqm/

- Practical guides
- Primers for biologists

# The biologist's view on statistics

- **Research won't proceed if solely based on absolutely proven facts.** It's ok to allow for uncertain hypotheses as long as we are honest about that (to us and others).
    - Depending on the study, it might be more important to minimize false positives before maximizing true positives. This is particularly true if the follow up experiments are very tedious/complex.

- **Biologists use data mining to generate hypotheses.** Some of these are for sure wrong.
    - Hypotheses are a start for research, not an endpoint.
    - Hypotheses are tested/proven with ad-hoc follow-up experiments.

- Good stats are an indicator of good lab work (reproducibility).

# Best practice in (bio) data mining

- **Start from a biological question!**
  - Explorative analyses are difficult and, typically, lead nowhere.

    Nice reading: 10.1126/science.aaa6146

- Don't expect an answer at all costs!
  Because of lack of data, noise, lack of response, wrong question

- Always think about **positive and negative controls**

- It's a good habit to (try to) keep data mining **simple**!

- Use different techniques and parameters

# Self check

- What do we expect to obtain from omics data?

- What are the typical problems associated with omics data?

- Why do we need computers?

- What is the role of statistics?

# Two groups problem

**Samples are divided <u>a priori</u> in two groups/classes.**
Typical examples:

- – mutant vs. wild-type
- – disease vs. health
- – diet A vs. diet B
- – treated vs. untreated (by a drug)
- – resistant vs. sensitive (to a drug)
- – old vs. young
- – …

BIO questions:

- – **Identify significantly changed features**
- – **Identify significantly changed cellular processes**



Phenotype Classes

A    B

Features (eg mRNA, proteins, metabolites)

How would you identify significantly changed features between A and B?



Phenotype Classes
A   B

Features (eg mRNA, proteins, metabolites)

# Univariate analysis

The simplest approach is to use univariate analysis as we know it from low-dimensional data.

That is: consider 1 feature at the time and compare the levels in the two groups. The test is repeated independently for all features.

Typical approach:

- Calculate <u>magnitude</u> of change > fold-change between group means or medians
- Verify statistical <u>significance</u> > p-value for null hypothesis

# Calculate <u>magnitude</u> of change

Fold-change:

$$FC = \frac{mean(GroupA)}{mean(GroupB)}$$

$[0\ldots1\ldots+\infty]$

For better (symmetric) visualization:

$$log2FC = log2\left(\frac{mean(GroupA)}{mean(GroupB)}\right)$$

$[-\infty\ldots0\ldots+\infty]$

# Verify statistical <u>significance</u>

Statistical hypothesis testing is often used to the likelihood to obtain the measured data given the hypothesis that the two groups originate from the same distribution = they are NOT different (= null hypothesis).

Common variants:

- **T-test** (usually using the two-tailed, unequal variance form)
  assumes normal distribution
- **Mann-Whitney/Wilcoxon/Ranksum test**
  uses ranking (no normality/distribution assumed)
- **Permutation test**
  For studies with many samples, it counts how frequently the result is better to that obtained with randomized sample labels.

Such tests produce a <u>p-value</u>. It indicates that under the assumption that the null hypothesis is true (= no difference), there are [p-value] chances to obtain the same results. The p-value is NOT the probability that the null hypothesis is true.

Traditionally in biology:

p-value < 0.05 : acceptable!

p-value < 0.01 : nice to have!

# Volcano plot

It's a common visualization method to evaluate the FC and statistical significance of a two-group comparison



Each point is one feature (gene, RNA, metabolite, …)

# How to choose thresholds?

# Fold Change threshold

**Pick an arbitrary value:** given e.g. overall reproducibility, object, sample size (example: |log2FC| > 1 in transcriptomics)

**Empirical procedure**: choose value such that things that are supposed to be identical are <u>not</u> called significantly different.

Example of negative controls:

- all wild-types
- all healthy samples
- all non-treated samples

**Thresholds can be determined by estimating the distribution of fold-changes between random (all) subsets of negative controls.**



log2(FC) between CONTROLS

# p-value threshold

**Multiple testing problem**: Choosing a significance cutoff α = 0.05 means that in a single test there is a 5% chance that the result is a false positive. This is ok if we do only 1 test, but…

Probability of
making at least 1 error

$$1 - (1 - \alpha)^{m}$$



$\alpha = 0.05$

$\alpha = 0.01$

Number of tests

# p-value threshold

Example: 821 features, significance cutoff < 0.05



**Case 1: Group A vs. B**

(t-test)
**p-value**
distribution

**56% significant**

**Case 2: Group C vs. C**

**14% significant**

# Error types (more in future…)

*False negatives*

## "TRUTH"

"DECISION"

| | $H_0$ true | $H_0$ false | Total |
|---|---|---|---|
| Do not reject $H_0$ | Correct<br>U<br>$1 - \alpha$ | **Type II Error**<br>T<br>$\beta$ | m-R |
| **Reject $H_0$** | **Type I Error**<br>V<br>$\alpha$ | Correct<br>S<br>$1 - \beta$ | R |
| | $m_0$ | m-$m_0$ | m |

*False positives*

*False Discovery Rate (FDR) = V/R*
*False Positive Rate (FPR) = V/$m_0$*

# Correcting for multiple testing

**Control Family-Wise Error Rate**

Guard against ANY false positives FWER = P(V ≥ 1).

**Bonferroni**: adjust p-value by number of tests **$m$**: reject $H_i$ if $p_i \leq \frac{\alpha}{m}$   <span style="color:red">very strict!!</span>

*Less stringent, sequential procedures exist (e.g. Holm), but don't change the issue that in biology we don't really want to control for FWER. In discovery, we can live with some false positives. Hence, all FWER tend to be <u>overly restrictive</u>.*

Better approaches:

**Control False Discovery Rate** (**Benjamini and Hochberg** 1995)
Controls the proportion of false positives among the set of rejected hypotheses

p-values > **FDR-adjusted p-values**, reject $H_i$ if $adj. p_i \leq \alpha$

<span style="color:red">OK if differences are rare ($pi_0 \approx 1$)</span>

**Control positive False Discovery Rate** (**Storey and Tibshirani** 2003)

Controls the rate that discoveries are false (pFDR). Uses all p-values to derive **$q$-values** (minimum FDR that can be attained when calling that feature significant)

<span style="color:red">OK if differences are frequent ($pi_0 < 1$)</span>

*(example: if gene X has q = 0.02 it means that 2% of genes that show p-values at least as small as gene X are false positives).*

[More info: Noble, *How does multiple testing correction work?* Nature Biotech 2009]
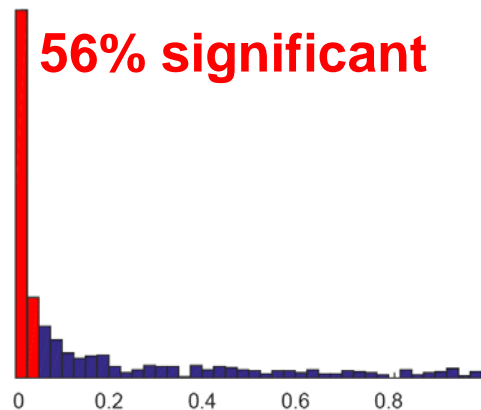
# The effect of MT correction

Example: 821 features, significance < 0.05
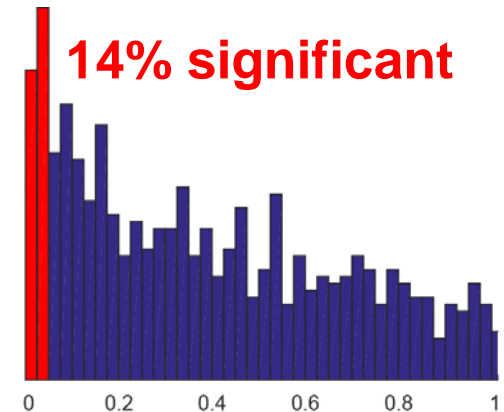


**Group A vs. B**

**56% significant**

(t-test)
**p-value** distribution

**Group C vs. C**

**14% significant**

**Benjamini-Hochberg correction**

**adjusted p-value** distribution
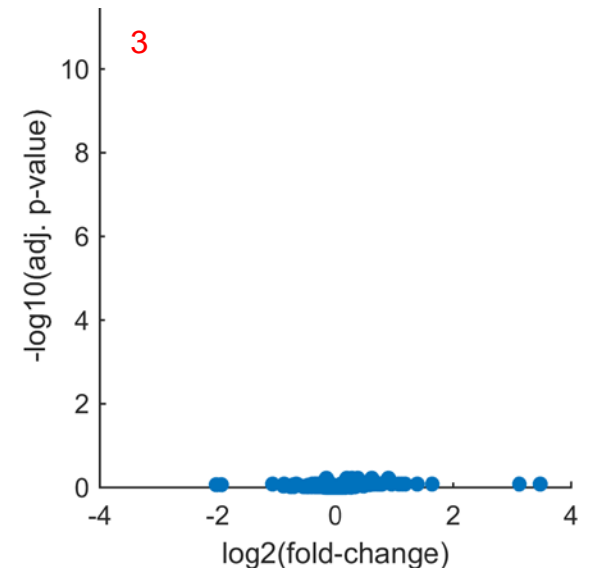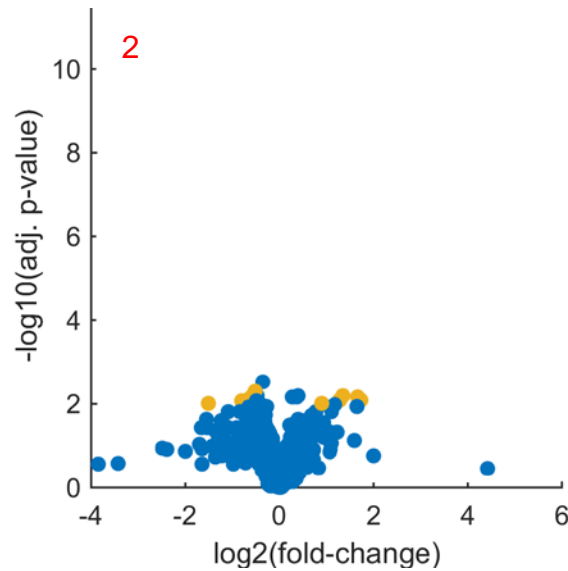
**51% significant**

**0.2% significant**

# Volcano plot

It's a common visualization method to evaluate the FC and statistical significance of a two-group comparison

# What do we conclude from these volcano plots?

# What do we learn from that?

| | |
|---|---|
| **Zero** significant changes<br><br>3 | Check data quality<br>Low p-values, high FC > increase number of replicates |
| Very **few** significant changes or very few "strong" changes<br><br>2 | Simplest case<br>> What is the identity of the markers? |
| **A lot** of significant changes<br>No clear ranking<br><br>1 | Quite common, but hard to generate hypotheses<br>> Too many hits<br>> Combination of primary and/or secondary effects<br><br>> **Enrichment analysis!** |

# Enrichment analysis

*Provided that __many__ significantly changed features were found (between two groups),*

how would you identify significantly changed cellular processes?

# Cellular "Processes"

**In a (systems-) biological context, we frequently identify processes with the molecular features that we know to be involved.**
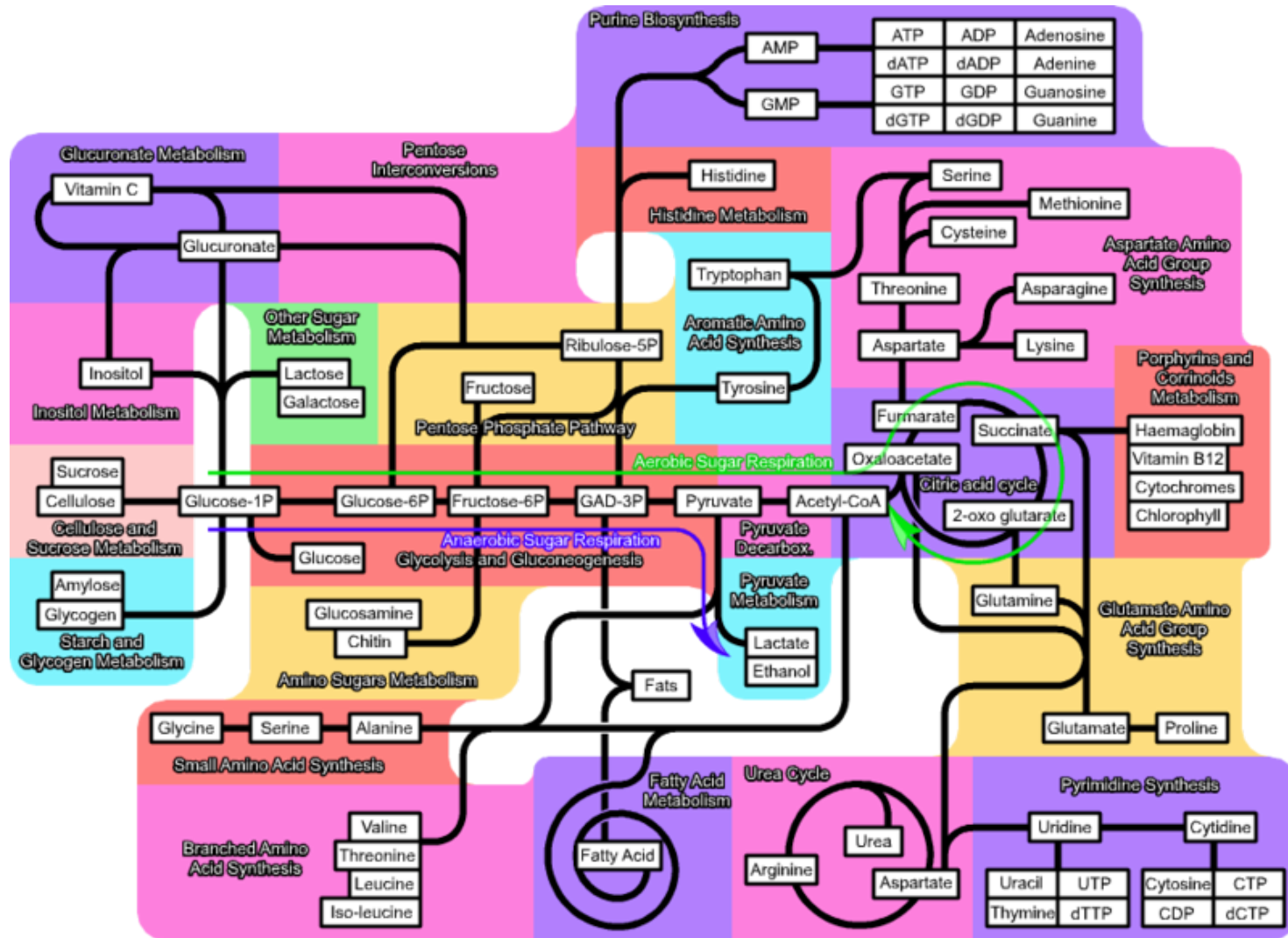
Examples:

- Enzymes or metabolite belonging to a pathway

- Proteins belonging to a complex (eg Ribosome)

- Genes controlled by a transcription factor

- GO terms

- Molecular Signatures Database (GSEA)

Important: A feature can belong to several groups/processes

# Biochemical pathways

# Transcriptional Network

Example: RegulonDB *(E. coli)*

*E. coli* TF-Gene Network

# Gene Ontology & Collections

Reference of controlled vocabularies and annotations representing gene product functions (10'000s entries, hierarchical).

# Enrichment analysis (EA)

*We want to test whether in the set of significantly changed features (i.e. genes, proteins, metabolites), the members of a specific group/process/pathway are over-represented.*

**Volcano Plot** > significance threshold > count > **Contingency table:**

*From definitions (GO term)*



| | In group (G+) 🟡 | Not in group (G-) 🔵 |
|---|---|---|
| Significant (S+) | A | B |
| Not significant (S-) | C | D |

# [Fisher's Exact Test]

|  | In group (G+) | Not in group (G-) |
|---|---|---|
| Significant (S+) | A | B |
| Not significant (S-) | C | D |

From the contingency table, we can directly calculate the likelihood that the same results could be obtained assuming the null hypothesis according to a Fisher's exact test (or an hypergeometric distribution)

In MATLAB:     *(all formula give the same result)*

```
>> [~,p] = fishertest([A B;C D],'Tail','right')
or      >> p = sum(hygepdf(A:A+B, A+B+C+D, A+C, A+B))
or      >> p = 1-hygecdf(A-1, A+B+C+D, A+C, A+B)
```

# Problem: How to choose cutoffs?

The contingency table depends on pre-defined cutoffs. These



Arguments against using the "normal" cutoffs:
[Subramanian PNAS 2005]

*After correcting for multiple hypotheses testing, no individual gene may meet the threshold for statistical significance, because the __relevant biological differences are modest relative to the noise.__*

*Single-gene analysis may miss important effects on pathways. Cellular processes often affect sets of genes acting in concert. __An increase of 20% in all genes encoding members of a metabolic pathway may dramatically alter the flux through the pathway and may be more important than a 20-fold increase in a single gene__.*

# Gene Set Enrichment Analysis

Subramanian et al, PNAS 2005

*(A) Genes are ranked in a list L based on the correlation between their expression and the class distinction by using any suitable metric, e.g. log2(fold-change) or t-test p-value.*

*(B) Given an a priori defined set of genes S, the goal of GSEA is to determine whether the members of S are randomly distributed or primarily found at the top or bottom.*

*We expect that sets related to the phenotypic distinction will tend to show the latter distribution.*

# The GSEA procedure

1. **Decide on what signs do we want to include:**



POS or NEG        POS only        NEG only

2. **Procedure for 1 pathway/group/process:**
    1. Choose very permissive thresholds (e.g. log2FC > 0.2 and p-value < 0.1)
    2. Rank all hits on either p-value (preferred) or log2FC
    3. Build contingency tables with either 2, 3, 4, …., ALL top hits
    4. For each table, calculate p-value using Fisher's exact test
    5. Keep lowest p-value = "best" enrichment

3. **Repeat for all pathways**
4. **FDR-correction (Benjamini-Hochberg or Storey)**

# Example

Out of 650 features detected, we know that 20 belong to a pathway/GO term.
Out of 650, 36 features passed our criteria abs(log2FC) > 0.2 and adj. p-value < 0.1.



*Features ranked by Fold-Change (high to low)*

#1 #2 #3 #4

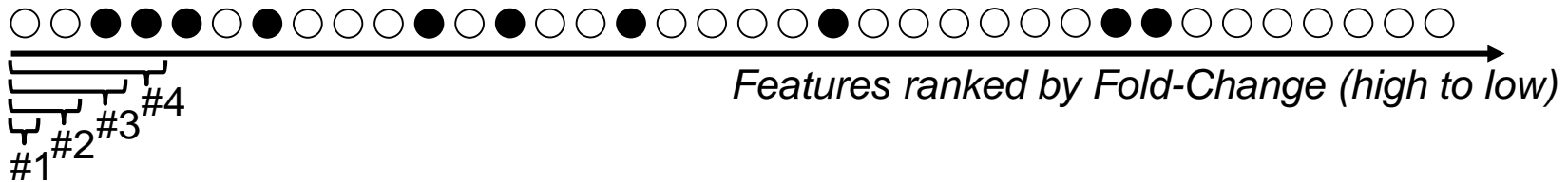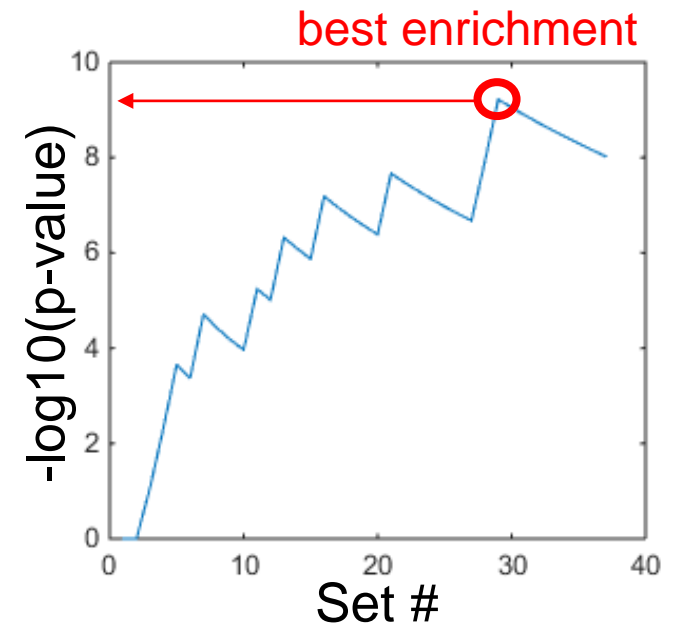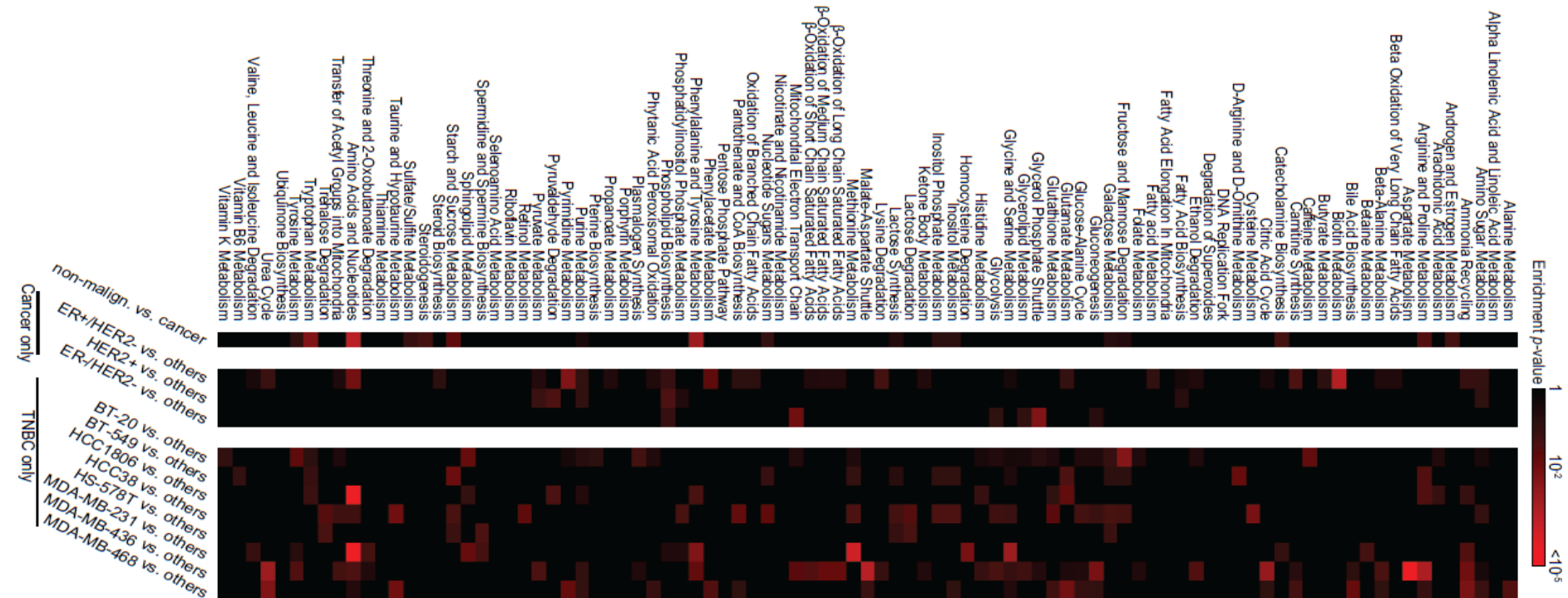| set | A+B+C+D | A+C | A | B | C | D | Enrichment p-value (Fisher) |
|-----|---------|-----|---|---|---|---|------------------------------|
| #1 | 650 | 20 | 0 | 1 | 20 | 629 | 1.00 |
| #2 | 650 | 20 | 0 | 2 | 20 | 628 | 1.00 |
| #3 | 650 | 20 | 1 | 2 | 19 | 628 | 0.087 |
| #4 | 650 | 20 | 2 | 2 | 18 | 628 | 0.049 |
| #5 | 650 | 20 | 3 | 2 | 17 | 628 | 0.00022 |
| #6 | 650 | 20 | 3 | 3 | 17 | 627 | 0.00043 |
| #7 | 650 | 20 | 4 | 3 | 16 | 627 | 0.000019 |
| #8 | 650 | 20 | 4 | 4 | 16 | 626 | 0.00038 |
| … | | | | | | | |

(A ● B ○)

best enrichment

-log10(p-value) vs Set #

`[~,p] = fishertest([A B;C D],'Tail','right')`

# EA across pathways

A full analysis is typically repeated on tens/hundreds of pathways/GO terms.

1. For each pathway, we keep the lowest p-value
2. FDR correction of p-values (e.g. Benjamini-Hochberg or Storey)
3. Enrichments are significant if adj. p-value or q-value < 0.01 (0.05)

# Self check

- What if no enrichment is found?



- What happens if the enrichment analysis is done <u>without</u> any log2FC or p-value cutoff?

# Summary

- **Define the question in biological and data mining terms**.

- Data mining of omics data is about hypothesis generation.

- Seeking for absolute certainty is likely to prevent discovery of testable hypotheses.

- In the simplest approach, significant features are found by univariate analysis of all features (independently)
> t-test > FDR correction

- Significantly changed "processes" are identified by enrichment analysis
> rank by p > test subsets > Fisher's exact test > FDR correction

# Exercises

- This week:
  - two groups problem, p-values, multiple hypothesis testing

  *What are the significantly changed features?*

- Next week:
  - gene set enrichment analysis

  *What are the significantly changed processes?*