

Multiple Hypothesis Testing

Katie Pollard

BMI 206

docpollard.org/bmi206

October 12, 2016

Testing many hypotheses at once

Large **multiplicity problem**: thousands of hypotheses are tested simultaneously!

Increased chance of **false positives**.

Chance of at least one p-value $< \alpha$ for N independent tests is $1 - (1 - \alpha)^N$

→ converges to one as N increases.

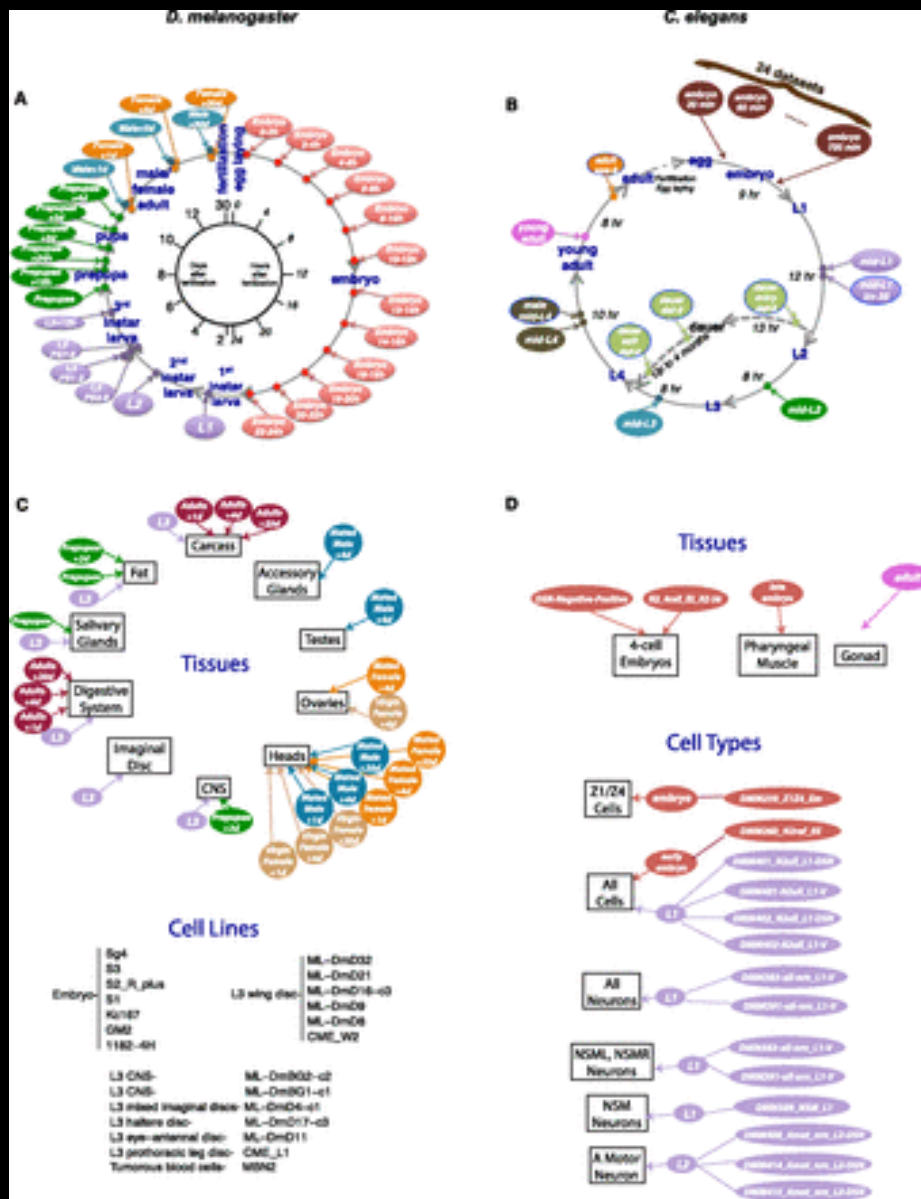
e.g., For N=1,000 and $\alpha = 0.01$, this chance is 0.9999568!

Individual p-values of 0.01 no longer correspond to significant findings.

Need to **adjust for multiple testing** when assessing the statistical significance of the observed associations.

Multiple testing in RNA-seq

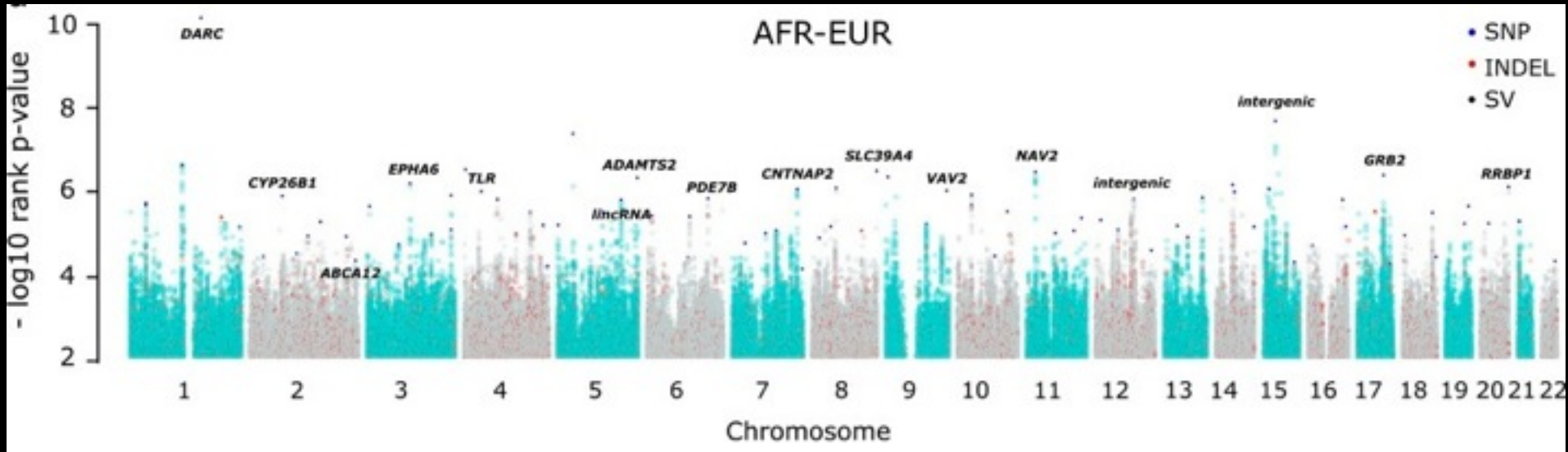
Comparison of fly and worm gene expression across developmental stages



Multiplicity on many levels:

- Two species
- Many stages
- Tissues vs. cell lines

Multiple testing in population genetics



Genomic regions with exceptionally high population differentiation identified in 9 | 1 whole genomes

Multiplicity on many levels:

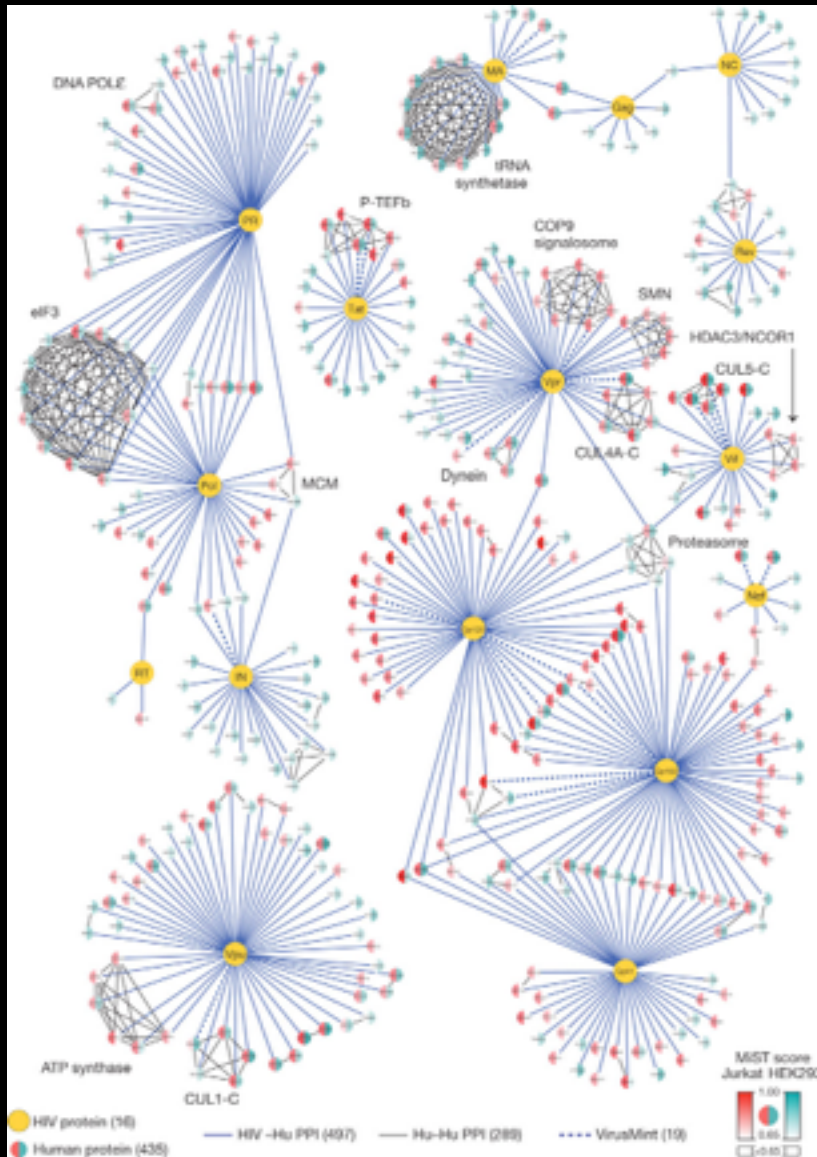
- Genome-wide
- SNPs, indels, SVs
- Several pairs of populations

Multiple testing in mass spec

Identifying human proteins that interact with each protein in the HIV genome

Interactions mean many tests:

- Tens of HIV proteins
- Thousands of human proteins
- Many thousands of potential protein-protein interactions



Components of a Multiple Hypothesis Test

1. **Parameters_u**: quantity of interest
2. **Null and alternative hypotheses**: family of tests; statements about parameter values
3. **Test statistics_u**: quantify evidence
4. **Error rate**: control mistakes
5. **Null distribution**: assess significance (high dim)
6. **Procedure**: decision rule for all tests jointly

Errors in multiple testing

	# non-rejected hypotheses	# rejected hypotheses	
# true null hypotheses	$m_0 - V_n$	V_n Type I error	$m_0 = S_0 $
# false null hypotheses	U_n Type II error	$m_1 - U_n$	$m_1 = S_0^c $
	$m - R_n$	R_n	m

Adapted from Benjamini & Hochberg (1995).

Type I error rates

- **Per family error rate (PFER):** Expected number of false positives.

$$\text{PFER} = E(V_n)$$

- **Per comparison error rate (PCER):** Expected rate of false positives.

$$\text{PCER} = E(V_n)/m$$

Type I error rates

- **Family-wise error rate (FWER):** Probability of at least one false positive.

$$\text{FWER} = P(V_n > 0)$$

- **Generalized FWER (gFWER):** Probability of at least $k+1$ false positives.

$$\text{gFWER}(k) = P(V_n > k)$$

Type I error rates

- **False discovery rate (FDR):** Expected proportion of false positives.

$$\text{FDR} = E(V_n/R_n)$$

- **False discovery proportion (FDP):** Probability that the proportion of false positives is at least q .

$$\text{FDP}(q) = P(V_n/R_n > q)$$

Null distribution for multiple testing

Joint distribution of the vector of test statistics if the null hypotheses were all true.

Used to convert test statistics to p-values.

Multiple testing p-values can be compared across tests, whereas statistics may be in different scales.

Different types:

- same for all tests?

- marginal vs. joint

- parametric vs. non-parametric

Marginal null distributions

- **Parametric** (a.k.a. tabled distributions)

Normal distributions **z-statistics**

Student's t-distribution **t-statistics**

F distribution **F-statistics**

Wilcoxon/Mann-Whitney U **U-statistics**

- **Non-parametric** (i.e., resampling based)

Permutation (2+ groups or continuous)

Bootstrap (various types)

Permutations

- Randomize group labels, positions, locations, ...
 - Estimates a distribution that is the pool of the groups (e.g., same mean, same variance, etc)
 - Usually easy to implement
- Some issues to consider
 - What to permute is not always obvious
 - Permuting into regions that cannot be observed
 - Strict null distribution because all parameters are different from the observed data, potentially including parameters other than in null hypothesis

Implementing a permutation test

- Simulate two vectors of numbers ($n=10$ random normal variables per group).
- Perform a parametric t-test.
- Generate $b=100$ permutations.
- Compute a t-statistic for each permutation.
- Calculate a permutation p-value.
- Compare parametric and permutation results.
- Repeat for different values of n (possibly unbalanced) and b . Also try different means in the two groups.

Bootstrap

- Resampling observed data with replacement estimates the variability in the empirical distribution
- Statistics over bootstrap iterations will have a range of values, providing an empirical test statistics distribution
- If this can be adjusted so the null hypothesis holds, it provides a suitable test statistics null distribution
 - Can be easy, e.g., make means the same in each group by computing sample means and subtracting
 - Need to think explicitly about the null hypothesis to make this adjustment to the bootstrap
 - Does not involve changing the labels, positions, etc.

Implementing a bootstrap test

- Simulate two vectors of numbers ($n=10$ random normal variables per group).
- Generate $b=100$ bootstrap samples. Standardize to have mean zero in each group.
- Compute a t-statistic for each bootstrap.
- Calculate a bootstrap p-value.
- Compare parametric, permutation, and bootstrap results.
- Repeat for different values of n (possibly unbalanced) and b . Also try different means in the two groups.

Joint null distributions

- **Parametric** (a.k.a. tabled distributions)
 - Multivariate Normal distributions
 - Multivariate distribution of F-statistics
- **Non-parametric** (i.e., resampling based)
 - Permutation (2+ groups or continuous)
 - Bootstrap (various types)

multtest package
MTP function

Resampling observations jointly

- Permutations

- Think about the sampling unit
- Permute label, position, location for vector of observed variables for each sampling unit
- Scrambling the variables is a common mistake

- Bootstrap

- Resample vectors of variables with replacement
- Adjust the joint bootstrap distribution so that the null hypothesis holds

Implementing multivariate resampling

- Simulate two vectors of numbers ($n=10$ random normal variables per group) 50 times independently. Store as a 50×20 matrix.
- Generate $b=100$ permutation and bootstrap samples. Standardize the bootstrap data to have mean zero in each group (50 rows).
- Compute a t-statistic for each row.
- Calculate parametric, permutation and bootstrap p-values. Compare results.
- Repeat for different means in the two groups and with correlation between the rows.

Multiple Testing Procedures

Goal: Given test statistics, an error rate, significance level & a high-dimensional null distribution, make a rejection decision for every test.

- Produces a set of **rejected hypotheses**
- Equivalently, compute **adjusted p-values**
 - Related to tail probabilities of the null distribution, but must account for all the other tests so that error rate is controlled
 - Value of multiple testing error rate if reject for all statistics at least this significant

How to get adjusted p-values?

Two different approaches to control multiple testing error rate (e.g., FWER or FDR):

1. Marginal methods that have two steps

- Get usual p-values, i.e., tail probabilities under each test's null distribution (marginal or joint)
- Adjust these probabilities based on the p-values of all other tests

2. Joint methods directly compute adjusted p-values from the joint null distribution

Types of marginal methods

- **Single-step:** Same p-value adjustment for all hypotheses.
- **Step-wise:** Adjustments depend on observed data (test statistics).
 - **Step-down** = start with most significant, reduce adjustment at each step, stop at first null hypothesis not rejected
 - **Step-up** = start with least significant, increase adjustment at each step, stop at first rejected null hypothesis

FWER controlling p-value adjustment

Name	Type	Adjustment
Bonferroni	Single-step	α/m
Sidak (ss)	Single-step	$1-(1-$
Holm	Step-down	α
Sidak (sd)	Step-down	$1-(1-$
Hochberg	Step-up	α

r_j = order statistics (ranks of test statistics)

FDR controlling p-value adjustment

Name	Type	Adjustment
Benjamini & Hochberg	Step-up	r
Benjamini & Yekutieli	Step-up	r
Storey	Step-up	Estimates pFDR and q-value

`qvalue` package

`multtest` package

`mt.rawp2adjp` function

Dependence Assumptions

Independence of test statistics

Bonferroni

Benjamini & Hochberg (or PRD)

Storey

Positive orthant dependent statistics

Sidak (both versions)

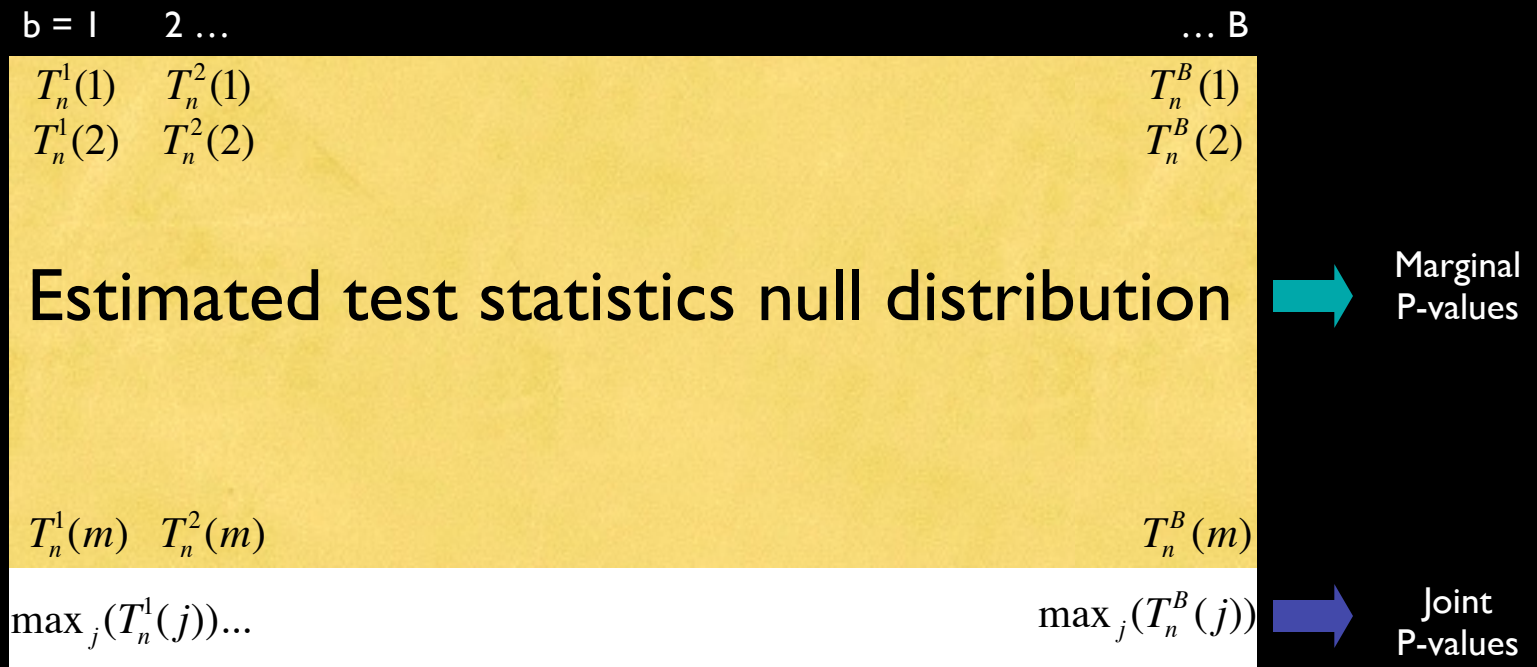
P-values satisfy Simes inequality

$$P(p_{r_j} > \alpha r_j / m) \geq 1 - \alpha$$

Hochberg (also assumes independence)

Joint methods for adjusted p-values

With the joint null distribution of the test statistics, **direct control** of Type I error rates is possible.



Joint methods for adjusted p-values

Name	Error Rate	Type	Details
ss.maxT	FWER	Single-step	Common cut-off: based on quantiles of max statistics
ss.minP	FWER	Single-step	Common quantile: based on quantiles of min p-values
sd.maxT	FWER	Step-down	Gene-specific cut-offs: based on max over subsets of T
sd.minP	FWER	Step-down	Gene-specific qtiles: based on min over subsets of P
ss.T(k+1)	gFWER	Single-step	Common cut-off: based on k + 1st largest T
ss.P(k+1)	gFWER	Single-step	Common qtile: based on k+1st smallest P

Multiple testing summary

- Completely marginal test

Marginal p-values from tabled distribution or resampling one gene at a time

Adjust with a marginal method

- Essentially marginal test

Marginal p-values from joint distribution

Adjust with marginal method

- Completely joint test

Marginal and adjusted p-values from joint distribution (also test statistic cut-offs)



COMPUTATION