

# Computational Genomics

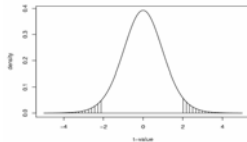
10-810/02-710, Spring 2009

## Differential Analysis of Microarray Gene Expression Data

Eric Xing

Lecture 16, March 16, 2009

Reading: class assignment



© Eric Xing @ CMU, 2005-2009

1

## Outline

- Motivation & examples
- Univariate hypothesis testing
- Multiple hypothesis testing
- Results for the two examples
- Discussion

© Eric Xing @ CMU, 2005-2009

2

# Introduction



- Many microarray experiments are carried out to find genes which are differentially expressed between two (or more) samples of cells:
  - cells (from the liver, say), in a mouse with a gene knocked out, compared with liver cells in a normal mouse of the same strain
  - cells in one region of the brain (say cerebellum), compared with cells in a different region (say the anterior cingulate region)
  - tumor cells in some organ (say the liver), compared with normal cells from the same organ
  - cells from an organism (say yeast) after a treatment (say by heat, or cold, or a drug) compared with cells of the same kind in the untreated state
  - cells from some part of a developing organ or organism at one time, compare with cells of the same kind at a later time, and so on
  - ...

© Eric Xing @ CMU, 2005-2009

3

# Motivation



- SCIENTIFIC: To determine which genes are differentially expressed between two sources of mRNA (trt, ctrl).
- STATISTICAL: To assign appropriately adjusted  $p$ -values to thousands of genes, and/or make statements about false discovery rates.
- We will discuss the issues in the context of two experiments, one which fits the aims above, and one which doesn't, but helps make a number of points.

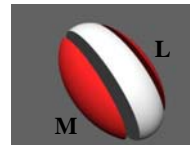
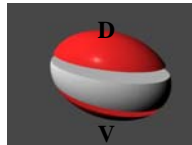
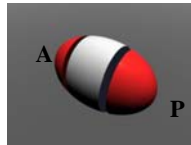
© Eric Xing @ CMU, 2005-2009

4

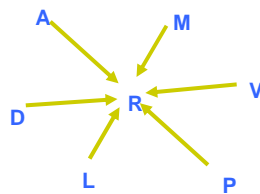
## Preliminary: experimental design



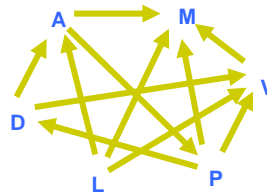
- Comparisons



- Two Ways to Do the Comparisons



Compare all samples to a common reference sample



Multiple direct comparisons between different samples

© Eric Xing @ CMU, 2005-2009

5

## Preliminary: differential analysis with one slide



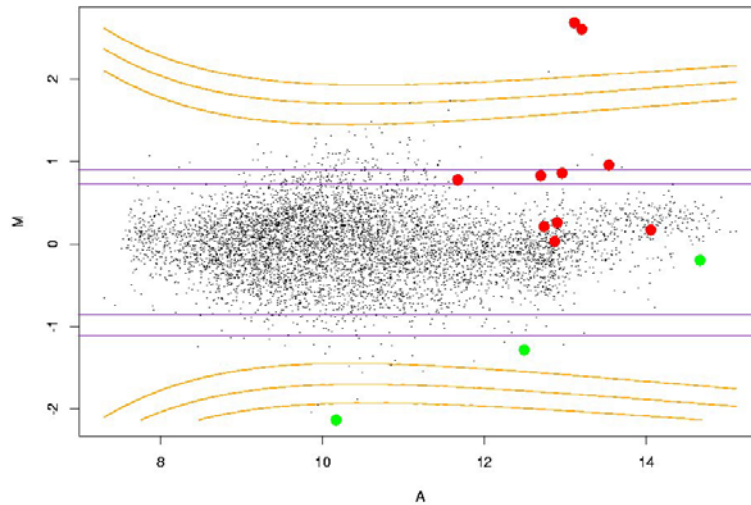
- The simplest cDNA microarray data analysis problem is identifying differentially expressed genes using **one slide**
  - This is a common enough hope
  - Efforts are frequently successful
  - It is not hard to do by eye
  - The problem is probably beyond formal statistical inference (valid p-values, etc) for the foreseeable future....why?
- In the next two panels, genes found to be **up-** or **down-**regulated in an 8 treatment (Srb1 over-expression) versus 8 control comparison are indicated in **red** and **green**, respectively, on plots of the data from **single hybridizations**.
- Also depicted are “confidence lines” determined by different methods and/or different “confidence” levels, which claim to be able to delineate differentially expressed genes using just one hybridization (slide).

© Eric Xing @ CMU, 2005-2009

6

## Matt Callow's Srb1 dataset (#5).

Newton's and Chen's single slide method

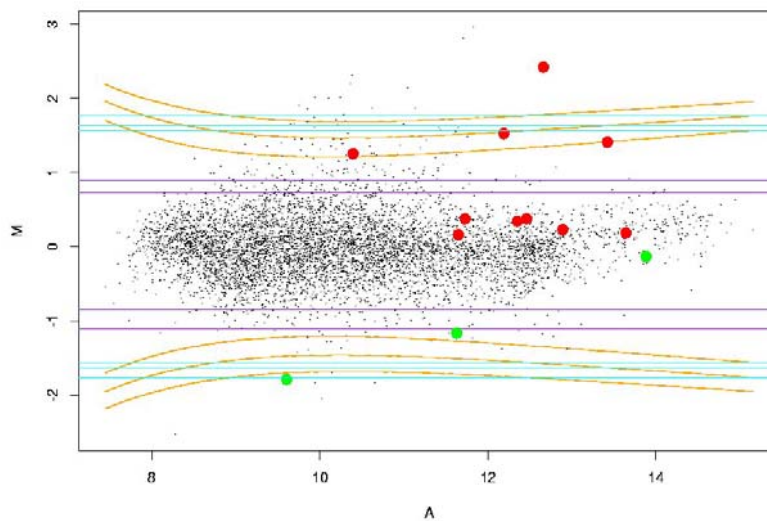


© Eric Xing @ CMU, 2005-2009

7

## Matt Callow's Srb1 dataset (#8).

Newton's, Sapir & Churchill's and Chen's single slide method



© Eric Xing @ CMU, 2005-2009

8

## Differential analysis with replicated hybridizations



- The second simplest cDNA microarray data analysis problem is identifying differentially expressed genes using replicated hybridizations
  - There are a number of different aspects:
    - First, between-slide normalization;
    - Then, what should we look at: averages, SDs, t-statistics, other summaries?
    - How should we look at them?
    - Can we make valid probability statements?
- We will discuss the issues in the context of two experiments, one which fits the aims above, and one which doesn't, but helps make a number of points.

© Eric Xing @ CMU, 2005-2009

9

## Apo AI experiment: (Matt Callow)



- **Goal:** To identify genes with altered expression in the livers of Apo AI knock-out mice (T) compared to inbred C57Bl/6 control mice (C).
  - 8 treatment mice and 8 control mice
  - 16 hybridizations: liver mRNA from each of the 16 mice ( $T_i$ ,  $C_i$ ) is labelled with Cy5, while pooled liver mRNA from the control mice ( $C^*$ ) is labelled with Cy3.
  - Probes: ~ 6,000 cDNAs (genes), including 200 related to lipid metabolism.



© Eric Xing @ CMU, 2005-2009

10

## Golub *et al* (1999) experiments



- **Goal.** To identify genes which are differentially expressed in acute lymphoblastic leukemia (ALL) tumors in comparison with acute myeloid leukemia (AML) tumors.
  - 38 tumor samples: 27 ALL, 11 AML.
  - Data from Affymetrix chips, some pre-processing.
  - Originally 6,817 genes; 3,051 after reduction.
  - Data therefore a  $3,051 \times 38$  array of expression values.
- Comment: this wasn't really the goal of Golub *et al*.

© Eric Xing @ CMU, 2005-2009

11

## Data



- The gene expression data can be summarized as follows

$$X = \begin{array}{c} \begin{array}{cc} \text{treatment} & \text{control} \end{array} \\ \left[ \begin{array}{ccc|ccc} x_{1,1} & \cdots & x_{1,n_1} & x_{1,n_1+1} & \cdots & x_{1,n} \\ x_{2,1} & \cdots & x_{2,n_1} & x_{2,n_1+1} & \cdots & x_{2,n} \\ \vdots & & \vdots & \vdots & & \vdots \\ x_{i,1} & \cdots & x_{i,n_1} & x_{i,n_1+1} & \cdots & x_{i,n} \\ \vdots & & \vdots & \vdots & & \vdots \\ x_{m,1} & \cdots & x_{m,n_1} & x_{m,n_1+1} & \cdots & x_{m,n} \end{array} \right] \end{array}$$

- Here  $x_{i,j}$  is the (relative) expression value of gene  $i$  in sample  $j$ . The first  $n_1$  columns are from the treatment (T); the remaining  $n_2 = n - n_1$  columns are from the control (C).

© Eric Xing @ CMU, 2005-2009

12

## Test strategy



- Which genes have changed? When permutation testing possible.
- 1. For each gene and each hybridization (8 ko + 8 ctl), use  $M = \log_2(R/G)$ .
- 2. For each gene form the  $t$ -statistic:
 
$$\frac{\text{average of 8 ko Ms} - \text{average of 8 ctl Ms}}{\sqrt{(1/8 (\text{SD of 8 ko Ms})^2 + 1/8 (\text{SD of 8 ctl Ms})^2)}}$$
- 3. Form a histogram of 6,000  $t$ -values.
- 4. Do a normal qq-plot; look for values “off the line”.
- 5. Compute the raw p-values for each gene by permutation procedures or from distribution models.
- 6. Adjust for multiple testing.

© Eric Xing @ CMU, 2005-2009

13

## Univariate hypothesis testing



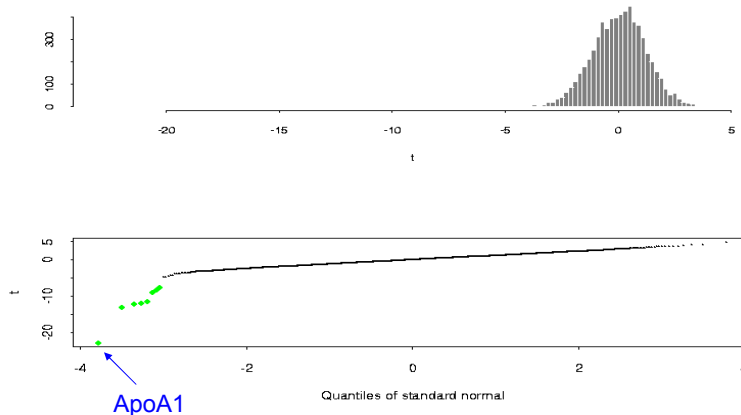
- Initially, focus on one gene only.
- We wish to test the null hypothesis  $H$  that the gene is not differentially expressed.
- In order to do so, we use a two sample  $t$ -statistic:

$$t = \frac{\text{average of } n_1 \text{ trtx} - \text{average of } n_2 \text{ ctlx}}{\sqrt{\frac{1}{n_1} (\text{SD of } n_1 \text{ trtx})^2 + \frac{1}{n_2} (\text{SD of } n_2 \text{ ctlx})^2}}$$

© Eric Xing @ CMU, 2005-2009

14

## Histogram & normal qq-plot of $t$ -statistics



© Eric Xing @ CMU, 2005-2009

15

## What is a normal qq-plot?



We have a random sample, say  $t_i, i=1, \dots, n$ , which we believe might come from a normal distribution. If it did, then for suitable  $\mu$  and  $\sigma$ ,  $\Phi((t_i - \mu)/\sigma)$ ,  $i=1, \dots, n$  would be **uniformly distributed** on  $[0, 1]$  (why?), where  $\Phi$  is the standard normal c.d.f.. Denoting the order statistics of the  $t$ -sample by  $t_{(1)}, t_{(2)}, \dots, t_{(n)}$  we can then see that  $\Phi((t_{(i)} - \mu)/\sigma)$  should be approximately  $i/n$  (why?). With this in mind, we'd expect  $t_{(i)}$  to be about  $\sigma\Phi^{-1}(i/n) + \mu$  (why?).

Thus if we plot  $t_{(i)}$  against  $\Phi^{-1}((i+1/2)/(n+1))$ , we might expect to see a straight line of slope about  $\sigma$  with intercept about  $\mu$ . (The  $1/2$  and  $1$  in numerator and denominator of the  $i/n$  are to avoid problems at the extremes.)

This is our normal quantile-quantile plot, the  $i/n$  being a quantile of the uniform, and the  $\Phi^{-1}$  being that of the normal.

© Eric Xing @ CMU, 2005-2009

16



## Why do a normal q-q plot?



One of the things we want to do with our  $t$ -statistics is roughly speaking, to identify the *extreme* ones.

It is natural to rank them, but how extreme is extreme? Since the sample sizes here are not too small (two samples of 8 each gives 16 terms in the difference of the means), approximate normality is not an unreasonable expectation for the null marginal distribution.

Converting ranked  $t$ 's into a normal qq-plot is a great way to see the extremes: they are the ones that are "off the line", at one end or another. This technique is particularly helpful when we have thousands of values. Of course we can't expect all differentially expressed genes to stand out as extremes: many will be masked by more extreme random variation, which is a big problem in this context. See later in the class for a discussion of these issues.

© Eric Xing @ CMU, 2005-2009

17

gene index	t statistic	
2139	-22	
4117	-15	
5330	-12	
1731	-11	
538	-11	
1489	-9.1	
2526	-8.3	
4916	-7.7	
941	-4.7	
2000	+3.1	
5867	-4.2	
4608	+4.8	
948	-4.7	
5577	-4.5	

**Gene annotation**

Apo AI

EST, weakly sim. to STEROL DESATURASE

CATECHOL O-METHYLTRANSFERASE

Apo CIII

EST, highly sim. to Apo AI

EST

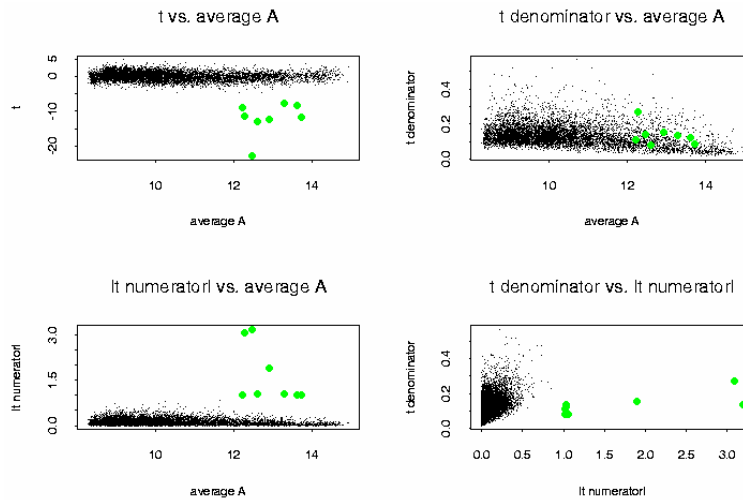
Highly sim. to Apo CIII precursor

similar to yeast sterol desaturase

© Eric Xing @ CMU, 2005-2009

18

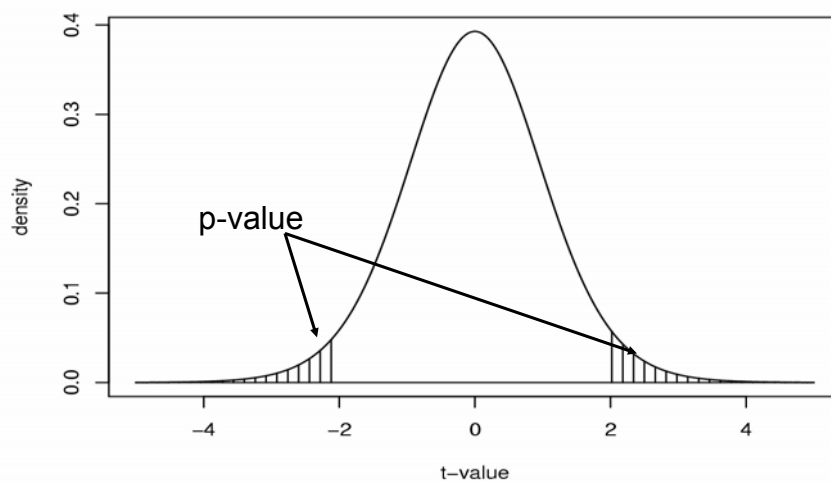
## Useful plots of $t$ -statistics



© Eric Xing @ CMU, 2005-2009

19

## The $p$ -values for two sample $t$ -stat



© Eric Xing @ CMU, 2005-2009

20

## $p$ -values



- The  **$p$ -value** or **observed significance level**  $p$  is the chance of getting a test statistic as or more extreme than the observed one, under the null hypothesis  $H$  of no differential expression.
- Although the previous test statistic is denoted by  $t$ , it would be unwise to assume that its null distribution is that of *Student's  $t$* . We have another way to assign  $p$ -values which is more or less valid: using permutations.

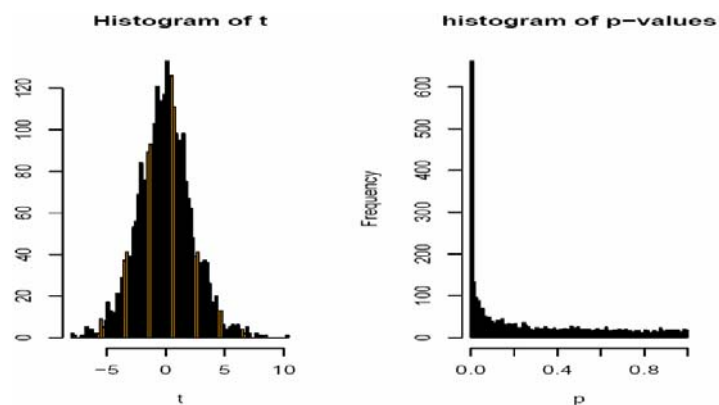
© Eric Xing @ CMU, 2005-2009

21

## Example



- Golub data, 27 ALL vs. 11 AML samples, 3,051 genes.



t-test: 1045 genes with  $p < 0.05$ .

© Eric Xing @ CMU, 2005-2009

22

## Computing $p$ -values by permutations



- We focus on one gene only. For the  $b$ th iteration,  $b = 1, \dots, B$ ;
- Permute the  $n$  data points for the gene ( $x$ ). The first  $n_1$  are referred to as “treatments”, the second  $n_2$  as “controls”.
- For each gene, calculate the corresponding two sample t-statistic,  $t_b$ .
- After all the  $B$  permutations are done;
- Put  $p = \#\{b: |t_b| \geq |t|\}/B$  (plover if we use  $>$ ).

With **all** permutations in the Apo A1 data,  $B = n!/n_1!n_2! = 12,870$ ;  
for the leukemia data,  $B = 1.2 \times 10^9$ .

© Eric Xing @ CMU, 2005-2009

23

## Many tests: a simulation study



- Simulation of this process for 6,000 genes with 8 treatments and 8 controls.
- **All** the gene expression values were simulated *i.i.d* from a  $N(0,1)$  distribution, i.e. **NOTHING** is differentially expressed in our simulation.
- We now present the 10 smallest raw (unadjusted) permutation  $p$ -values.

© Eric Xing @ CMU, 2005-2009

24

## Unadjusted $p$ -values



gene index	t value	$p$ -value (unadj.)
2271	4.93	$2 \times 10^{-4}$
5709	4.82	$3 \times 10^{-4}$
5622	-4.62	$4 \times 10^{-4}$
4521	4.34	$7 \times 10^{-4}$
3156	-4.31	$7 \times 10^{-4}$
5898	-4.29	$7 \times 10^{-4}$
2164	-3.98	$1.4 \times 10^{-3}$
5930	3.91	$1.6 \times 10^{-3}$
2427	-3.90	$1.6 \times 10^{-3}$
5694	-3.88	$1.7 \times 10^{-3}$

Clearly we can't just use standard  $p$ -value thresholds of .05 or .01.

© Eric Xing @ CMU, 2005-2009

25

## Discussion



- What assumptions on the null distributions of the gene expression values  $x_i = (x_{i,1}, x_{i,2}, \dots, x_{i,n})$  are necessary or sufficient for the permutation-based  $p$ -values just described to be *valid*? And, are they applicable in our examples?
  - First,  $p$ -values are *valid* if their distribution is *uniform*(0,1) under the null hypothesis.
  - Secondly, if the null distribution of  $x_i$  is exchangeable, i.e. invariant under permutations of  $1, \dots, n$ , then, we could reasonably hope (and actually prove) that the distribution of the permutation-based  $p$ -values is indeed uniform on  $1, \dots, n$ .
  - We also noted that having the joint distribution i.i.d. would be sufficient, as this implied exchangeability.
- Considered the ApoAI experiment.
  - Because the 16 log-ratios for each gene involved a term from the pooled control mRNA, called  $C^*$  above, it seems clear that an i.i.d. assumption is unreasonable.
  - Had the experiment been carried out by using pooled control mRNA from mice other than the controls in the experiment, an exchangeability assumption under the null hypothesis would have been quite reasonable.
  - Unfortunately,  $C^*$  did come from the same mice as the  $C_i$ , so exchangeability is violated, and the assumption is at best an approximation.

© Eric Xing @ CMU, 2005-2009

26

## Multiple testing: the problem



**Multiplicity problem: thousands of hypotheses are tested simultaneously.**

- Increased chance of false positives.
- E.g. suppose you have 10,000 genes on a chip and not a single one is differentially expressed. You would expect  $10000 \times 0.01 = 100$  of them to have a p-value  $< 0.01$ .
- Individual p-values of e.g. 0.01 no longer correspond to significant findings.

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

© Eric Xing @ CMU, 2005-2009

27

## Multiple testing: Counting errors



Assume we are testing  $H^1, H^2, \dots, H^m$ .

$m_0$  = # of true hypotheses     $R$  = # of rejected hypotheses

	# true null hypo.	# false null hypo.	
# accepted	$U$	$T$	$m - R$
# rejected	$V$	$S$	$R$
	$m_0$	$m - m_0$	

$V$  = # Type I errors [false positives]

$T$  = # Type II errors [false negatives]

© Eric Xing @ CMU, 2005-2009

28

## Type I error rates



- **Per comparison error rate** (PCER): the expected value of the number of Type I errors over the number of hypotheses,  
$$\text{PCER} = E(V)/m.$$
- **Per-family error rate** (PFER): the expected number of Type I errors,  
$$\text{PFER} = E(V).$$
- **Family-wise error rate**: the probability of at least one type I error  
$$\text{FEWR} = \text{pr}(V \geq 1)$$
- **False discovery rate** (FDR) is the expected proportion of Type I errors among the rejected hypotheses  
$$\text{FDR} = E(V/R; R > 0) = E(V/R \mid R > 0)\text{pr}(R > 0).$$
- **Positive false discovery rate** (pFDR): the rate that discoveries are false  
$$\text{pFDR} = E(V/R \mid R > 0).$$

© Eric Xing @ CMU, 2005-2009

29

## Multiple testing: Controlling a type I error rate



- Aim:  
For a given type I error rate, use a procedure to select a set of “significant” genes that guarantees a type I error rate  $\leq \alpha$ .

© Eric Xing @ CMU, 2005-2009

30

# Multiple testing

## Family-wise error rates



- Definition:

$$\begin{aligned}\text{FWER} &= \Pr(\# \text{ of false discoveries} > 0) \\ &= \Pr(V > 0)\end{aligned}$$

Bonferroni (1936)

Tukey (1949)

Westfall and Young (1993) discussed resampling

.....

- FWER and microarrays

- maxT step-down procedure

- Dudoit et al (2002)

- Westfall et al (2001)

- minP step-down procedure

- Ge et al (2003), a novel fast algorithm

© Eric Xing @ CMU, 2005-2009

31

# Multiple testing

## False discovery rates



- Definition:

$$Q = \frac{\# \text{ of false discoveries}}{\# \text{ of discoveries}} = \frac{V}{R}$$

- Q is set to be 0 when R=0

- FDR = expectation of Q = E(V/R; R>0)

- Seeger (1968)

- Benjamini and Hochberg (1995)

- Caution with FDR

- Cheating:

- Adding known diff. expressed genes reduces FDR

- Interpreting:

- FDR applies to a set of genes in a global sense, not to individual gene

© Eric Xing @ CMU, 2005-2009

32



## Types of control of Type I error



- **strong control:** control of the Type I error whatever the true and false null hypotheses. For FWER, strong control means controlling

$$\max_{M_0 \subset H_0^C} pr(V \geq 1 | M_0)$$

where  $M_0$  = the set of true hypotheses (note  $|M_0| = m_0$ );

- **exact control:** under  $M_0$ , even though this is usually unknown.
- **weak control:** control of the Type I error only under the **complete null hypothesis**  $H_0^C = \cap_i H_i$ . For FWER, this is control of  $pr(V \geq 1 | H_0^C)$ .

© Eric Xing @ CMU, 2005-2009

33

## Adjustments to $p$ -values



- For strong control of the FWER at some level  $\alpha$ , there are procedures which will take  $m$  unadjusted  $p$ -values and modify them separately, so-called *single step* procedures, the **Bonferroni** adjustment or correction being the simplest and most well known. Another is due to Sidák.
- Other, more powerful procedures, adjust sequentially, from the smallest to the largest, or vice versa. These are the *step-up* and *step-down* methods, and we'll meet a number of these, usually variations on single-step procedures.
- In all cases, we'll denote adjusted  $p$ -values by  $\pi$ , usually with subscripts, and let the context define what type of adjustment has been made. Unadjusted  $p$ -values are denoted by  $p$

© Eric Xing @ CMU, 2005-2009

34

## p-value adjustments: single-step



- Suppose we conduct a hypothesis test for each gene  $g = 1, \dots, m$ , producing
  - an observed test statistic:  $T_i$
  - an unadjusted p-value:  $p_i$ .
- Define adjusted p-values  $\pi_i$ , such that the FWER is controlled at level  $\alpha$  where  $H_i$  is rejected when  $\pi_i \leq \alpha$ .

$$\text{Bonferroni: } \pi_i = \min(m p_i, 1)$$

- Bonferroni always gives strong control.

© Eric Xing @ CMU, 2005-2009

35

## Proof for Bonferroni (single-step adjustment)



$$\begin{aligned} & \text{pr (reject at least one } H_i \text{ at level } \alpha \mid H_0^C) \\ &= \text{pr (at least one } \pi_i \leq \alpha \mid H_0^C) \\ &\leq \sum_{i=1}^m \text{pr}(\pi_i \leq \alpha \mid H_0^C), && \text{by Boole's inequality} \\ &= \sum_{i=1}^m \text{pr}(P_i \leq \alpha/m \mid H_0^C), && \text{by definition of } \pi_i \\ &= m \times \alpha/m, && \text{assuming } P_i \sim U[0,1]) \\ &= \alpha. \end{aligned}$$

### Notes:

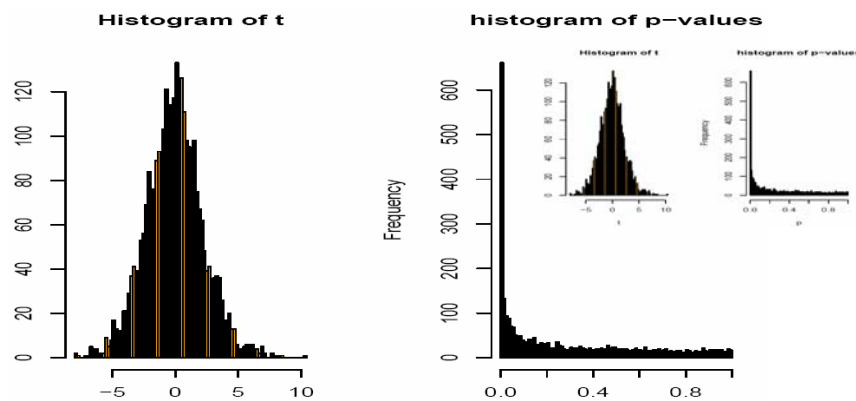
1. We are testing  $m$  genes,  $H_0^C$  is the complete null hypothesis, **that no gene is differentially expressed**.
2.  $P_i$  is the unadjusted p-value for gene  $i$ , while  $\pi_i$  here is the Bonferroni adjusted p-value.
3. We use lower case letters for observed p-values, and upper case for the corresponding random variables.

© Eric Xing @ CMU, 2005-2009

36

## Example

- Golub data, 27 ALL vs. 11 AML samples, 3,051 genes.



98 genes with Bonferroni-adjusted  $\pi_i < 0.05 \Leftrightarrow p_i < 0.000016$  (t-test)

© Eric Xing @ CMU, 2005-2009

37

## More is not always better

- Suppose you produce a small array with 500 genes you are particularly interested in.
- If a gene on this array has an unadjusted  $p$ -value of 0.0001, the Bonferroni-adjusted  $p$ -value is still 0.05.
- If instead you use a genome-wide array with, say, 50,000 genes, this gene would be much harder to detect, because roughly 5 genes can be expected to have such a low  $p$ -value by chance.

© Eric Xing @ CMU, 2005-2009

38

## p-value adjustments: single-step



- Suppose we conduct a hypothesis test for each gene  $g = 1, \dots, m$ , producing
  - an observed test statistic:  $T_i$
  - an unadjusted p-value:  $p_i$ .
- Define adjusted p-values  $\pi_i$ , such that the FWER is controlled at level  $\alpha$  where  $H_i$  is rejected when  $\pi_i \leq \alpha$ .

$$\text{Sidák: } \pi_i = 1 - (1 - p_i)^m$$

- Sidák is less conservative than Bonferroni. When the genes are independent, it gives strong control exactly (FWER =  $\alpha$ ), proof later. It controls FWER in many other cases, but is still conservative.

© Eric Xing @ CMU, 2005-2009

39

## Proof for Sidák's method (single-step adjustment)



$$\begin{aligned} & \text{pr}(\text{reject at least one } H_i \mid H_0^C) \\ &= \text{pr}(\text{at least one } \pi_i \leq \alpha \mid H_0^C) \\ &= 1 - \text{pr}(\text{all } \pi_i > \alpha \mid H_0^C) \\ &= 1 - \prod_{i=1}^m \text{pr}(\pi_i > \alpha \mid H_0^C) \quad \text{assuming independence} \end{aligned}$$

Here  $\pi_i$  is the Sidák adjusted p-value, and so  $\pi_i > \alpha$  if and only if  $P_i > 1 - (1 - \alpha)^{1/m}$  (check), giving

$$\begin{aligned} & 1 - \prod_{i=1}^m \text{pr}(\pi_i > \alpha \mid H_0^C) \\ &= 1 - \prod_{i=1}^m \text{pr}(P_i > 1 - (1 - \alpha)^{1/m} \mid H_0^C) \\ &= 1 - \{ (1 - \alpha)^{1/m} \}^m \text{ since all } P_i \sim U[0, 1], \\ &= \alpha \end{aligned}$$

© Eric Xing @ CMU, 2005-2009

40

## Single-step adjustments (ctd)



### FWER: Improvements to Bonferroni

- The *minP* method of Westfall and Young:

$$\Pi_i = \text{pr}(\min_{1 \leq l \leq m} P_l \leq p_i | H)$$

- Based on the joint distribution of the  $p$ -values  $\{P_l\}$ . This is the most powerful of the three single-step adjustments.
- If  $P_i \sim U[0,1]$ , it gives a FWER exactly  $= \alpha$  (see next page).
- It always confers weak control, and gives strong control under subset pivotality (definition next but one slide).

## Proof for (single-step) minP adjustment



Given level  $\alpha$ , let  $c_\alpha$  be such that

$$\text{pr}(\min_{1 \leq i \leq m} P_i \leq c_\alpha | H_0^C) = \alpha.$$

Note that  $\{\Pi_i \leq \alpha\} \equiv \{P_i \leq c_\alpha\}$  for any  $i$ .

$$\begin{aligned} & \text{pr}(\text{reject at least one } H_i \text{ at level } \alpha | H_0^C) \\ &= \text{pr}(\text{at least one } \Pi_i \leq \alpha | H_0^C) \\ &= \text{pr}(\min_{1 \leq i \leq m} \Pi_i \leq \alpha | H_0^C) \\ &= \text{pr}(\min_{1 \leq i \leq m} P_i \leq c_\alpha | H_0^C) \\ &= \alpha. \end{aligned}$$

## Strong control and subset pivotality



- The above proofs are under  $H_0^C$ , which is what we term weak control
- In order to get strong control, we need the condition of *subset pivotality*.
- The distribution of the unadjusted  $p$ -values  $(P_1, P_2, \dots, P_m)$  is said to have the *subset pivotality* property if for all subsets  $L \subseteq \{1, \dots, m\}$  the distribution of the subvector  $\{P_i; i \in L\}$  is identical under the restrictions  $\cap \{H_i; i \in L\}$  and  $H_0^C$ .
- Using the property, we can prove that for each adjustment under their conditions, we have
 
$$\begin{aligned} & \text{pr (reject at least one } H_i \text{ at level } \alpha, i \in M_0 \mid H_{M_0})} \\ &= \text{pr (reject at least one } H_i \text{ at level } \alpha, i \in M_0 \mid H_0^C) \\ &\leq \text{pr (reject at least one } H_i \text{ at level } \alpha, \text{ for all } i \mid H_0^C) \\ &\leq \alpha. \end{aligned}$$
- Therefore, we have proved strong control for the previous three adjustments, assuming subset pivotality.

© Eric Xing @ CMU, 2005-2009

43

## Permutation-based single-step minP adjustment of $p$ -values



- For the  $b$ th iteration,  $b = 1, \dots, B$ ;
- Permute the  $n$  columns of the data matrix  $X$ , obtaining a matrix  $X_b$ . The first  $n_1$  columns are referred to as “treatments”, the second  $n_2$  columns as “controls”.
- For each gene, calculate the corresponding unadjusted  $p$ -values,  $p_{i,b}$ ,  $i = 1, 2, \dots, m$ , (e.g. by further permutations) based on the permuted matrix  $X_b$ .
- After all the  $B$  permutations are done.
- Compute the adjusted  $p$ -values  $\pi_i = \#\{b: \min_i p_{i,b} \leq p_i\} / B$ .

© Eric Xing @ CMU, 2005-2009

44

## Example



- Suppose  $p_{\min} = 0.0003$  (the minimal unadjusted  $p$ -value).
- Among the randomized data sets (permuted sample labels), count how often the minimal  $p$ -value is smaller than 0.0003. If this appears e.g. in 4% of all cases,  $\pi_{\min} = 0.04$ .

## The computing challenge: iterated permutations



- The procedure is quite computationally intensive if  $B$  is very large (typically at least 10,000) **and** we estimate all unadjusted  $p$ -values by further permutations.

- Typical numbers:

To compute one unadjusted  $p$ -value  $B = 10,000$

# unadjusted  $p$ -values needed  $B = 10,000$

# of genes  $m = 6,000$ . In general run time is  $O(mB^2)$ .

- How to avoid computational difficulty?

## Single-step minP adjustment



- *maxT* method: (Chapter 4 of Westfall and Young)

$$\pi_i = \Pr( \max_{1 \leq l \leq m} |T_l| \geq |t_i| \mid H_0^C )$$

- needs  $B = 10,000$  permutations only.
- However, if the distributions of the test statistics are not identical, it will give more weight to genes with heavy tailed distributions (which tend to have larger  $t$ -values)
- There is a fast algorithm which does the minP adjustment in  $O(m \log B + m \log m)$  time.

© Eric Xing @ CMU, 2005-2009

47

## Proof for the single-step maxT adjustment



Given level  $\alpha$ , let  $c_\alpha$  such that  $\Pr(\max_{1 \leq i \leq m} |T_i| \leq c_\alpha \mid H_0^C) = \alpha$ .

Note the  $\{P_i \leq \alpha\} \equiv \{|T_i| \leq c_\alpha\}$  for any  $i$ . Then we have (cf. min P)

$$\begin{aligned} & \Pr(\text{reject at least one } H_i \text{ at level } \alpha \mid H_0^C) \\ &= \Pr(\text{at least one } P_i \leq \alpha \mid H_0^C) \\ &= \Pr(\min_{1 \leq i \leq m} P_i \leq \alpha \mid H_0^C) \\ &= \Pr(\max_{1 \leq i \leq m} |T_i| \leq c_\alpha \mid H_0^C) \\ &= \alpha. \end{aligned}$$

To simplify the notation we assumed a two sided test by using the statistic  $T_i$ . We also assume  $P_i \sim U[0, 1]$ .

© Eric Xing @ CMU, 2005-2009

48



## More powerful methods: step-down adjustments



The idea: S Holm's modification of Bonferroni.

Also applies to Sidák, maxT, and minP.

## S Holm's modification of Bonferroni



- Order the unadjusted  $p$ -values such that  $p_{r_1} \leq p_{r_2} \leq \dots \leq p_{r_m}$ .

The indices  $r_1, r_2, r_3, \dots$  are fixed for given data.

- For control of the FWER at level, the step-down Holm adjusted  $p$ -values are

$$\pi_{r_j} = \max_{k \in \{1, \dots, j\}} \{\min((m-k+1)p_{r_k}, 1)\}.$$

- The point here is that we don't multiply every  $p_{r_k}$  by the same factor  $m$ , but only the smallest. The others are multiplied by successively smaller factors:  $m-1, m-2, \dots$ , down to multiplying  $p_{r_m}$  by 1.
- By taking successive maxima of the first terms in the brackets, we can get monotonicity of these adjusted  $p$ -values.
- Holm's adjusted  $p$ -values deliver strong control.

## Step-down adjustment of minP



- Order the unadjusted  $p$ -values such that  $p_{r1} \leq p_{r2} \leq \dots \leq p_{rm}$ .
- Step-down adjustment: it has a complicated formula, see below, but in effect is

1. Compare  $\min\{P_{r1}, \dots, P_{rm}\}$  with  $p_{r1}$  ;
2. Compare  $\min\{P_{r2}, \dots, P_{rm}\}$  with  $p_{r2}$  ;
3. Compare  $\min\{P_{r3}, \dots, P_{rm}\}$  with  $p_{r3}$  ;
- ...
- m. Compare  $P_{rm}$  with  $p_{rm}$  .

- Enforce **monotonicity** on the adjusted  $p_{ri}$  . The formula is

$$\pi_{rj} = \max_{k \in \{1, \dots, j\}} \{pr(\min_{l \in \{rk, \dots, rm\}} P_l \leq p_{rk} | H_0^c)\}.$$

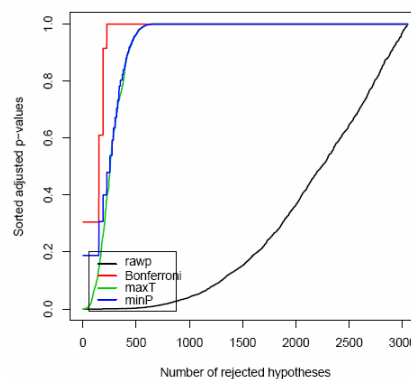
© Eric Xing @ CMU, 2005-2009

51

## FWER: Comparison of different methods



- Golub data, 27 ALL vs. 11 AML samples, 3,051 genes.



The FWER is a conservative criterion: many interesting genes may be missed.

© Eric Xing @ CMU, 2005-2009

52

## False Discovery Rate



- FDR:  $= E(V/R)$   
 $= E(F/S | S > 0) P(S > 0)$

In case  $R=0$ , define  $F/R=0$  if  $R=0$ .

- Alternatively, define  $pFDR = E(V/R | R > 0)$ .
  - When  $m$  is large,  $P(S > 0)$  is approx. 1 and FDR is approx. equal to pFDR.
- FDR is a measure of the overall accuracy of a set of significant features.

## False discovery rate (Benjamini and Hochberg 1995)



### Steps:

- Select desired limit  $\alpha$  on  $E(FDR)$
- Rank the p-values  $p_{r1} \leq p_{r2} \leq \dots \leq p_{rm}$ .
- The adjusted p-values are to control FDR when  $P_i$  are **independently** distributed are given by the step-up formula:

$$\pi_{ri} = \min_{k \in \{i, \dots, m\}} \{ \min (mp_{rk}/k, 1) \}.$$

- We use this as follows: reject  $H_{r1}, H_{r2}, \dots, H_{rk^*}$  where  $k^*$  is the largest  $k$  such that  $p_{rk} \leq (k/m)\alpha$ . This keeps the  $FDR \leq \alpha$  under independence
  - Thus the FDR Adjusted p-value = lowest level of FDR for which the hypothesis is first included in the set of rejected hypothesis
  - Compare the above with Holm's adjustment to control FWE, the step-down version of Bonferroni, which is  $\pi_i = \max_{k \in \{1, \dots, i\}} \{ \min (kp_{rk}, 1) \}$ .

## Positive false discovery rate (Storey, 2001, independent case)



- A new definition of FDR, called positive false discovery rate (pFDR)

$$\text{pFDR} = E(V/R \mid R > 0)$$

- The logic behind this is that in practice, at least one gene should be expected to be differentially expressed.
- The adjusted  $p$ -value (called  $q$ -value in Storey's paper) are to control pFDR.

$$\Pi_i = \min_{k \in \{1, \dots, i\}} \{ (m p_k / k) \pi_0 \}$$

- Note  $\pi_0 = m_0 / m$  can be estimated by the following formula for suitable  $\beta$

$$\pi_0 = \# \{ p_i > \beta \} / \{ (1 - \beta) m \}.$$

© Eric Xing @ CMU, 2005-2009

55

## Estimation of the FDR



- Idea: Depending on the chosen cutoff-value(s) for the test statistic  $T_i$ , estimate the expected proportion of false positives in the resulting gene list through a permutation scheme.
  - Estimate the number  $m_0$  of non-diff. genes:  $m_0 = \# \{ p_i > \beta \} / (1 - \beta)$ .
  - Compute the average number of significant genes under permutations of the sample labels.
    - For  $b = 1, \dots, B$ , (randomly) permute the sample labels – this corresponds to the complete null hypothesis. Compute test statistics  $T_{ib}$  for each gene.
    - For any threshold  $t_0$  of the test statistic, compute the numbers  $V_b$  of genes with  $T_{ib} > t_0$  (numbers of false positives).
    - compute the mean of the  $V_b$ .
  - Estimate the FDR

$$E(V/R) = \hat{V} \frac{m_0}{m} / R$$

© Eric Xing @ CMU, 2005-2009

56

## FWER or FDR?

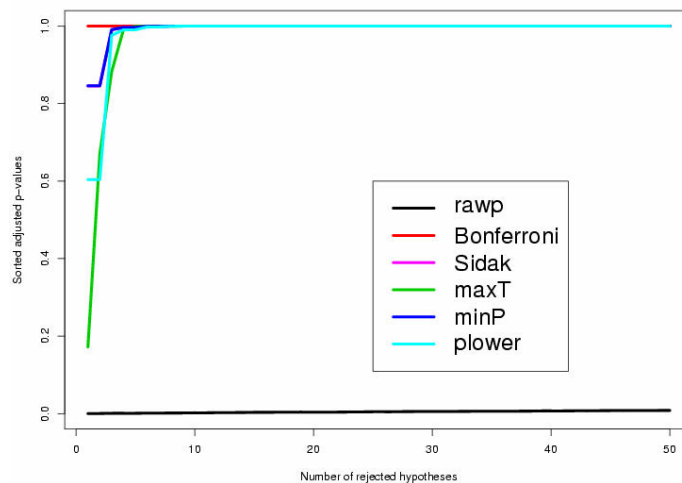
- Chose control of the FWER if high confidence in all selected genes is desired. Loss of power due to large number of tests: many differentially expressed genes may not appear as significant.
- If a certain proportion of false positives is tolerable: Procedures based on FDR are more flexible; the researcher can decide how many genes to select, based on practical considerations.

© Eric Xing @ CMU, 2005-2009

57

## Results: random data

Random data---complete permutations

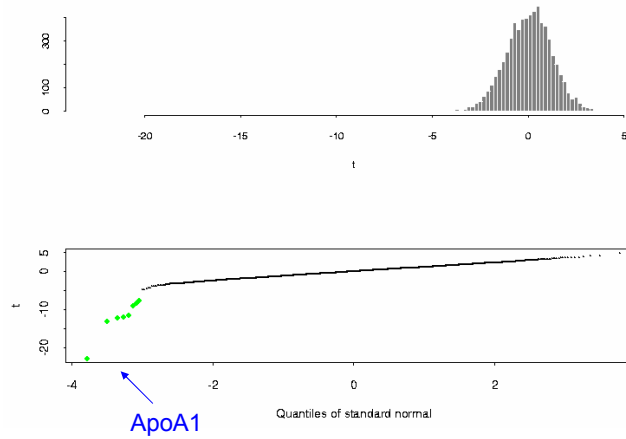


© Eric Xing @ CMU, 2005-2009

58

## Results: Apo AI data

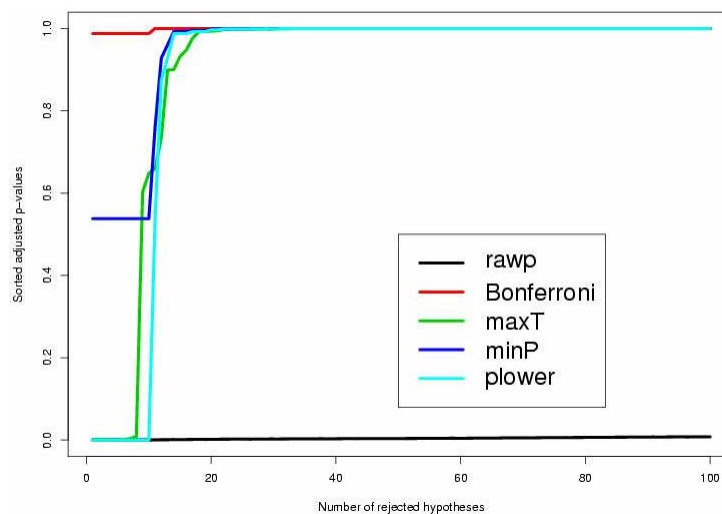
- Histogram & normal q-q plot of t-statistics



© Eric Xing @ CMU, 2005-2009

59

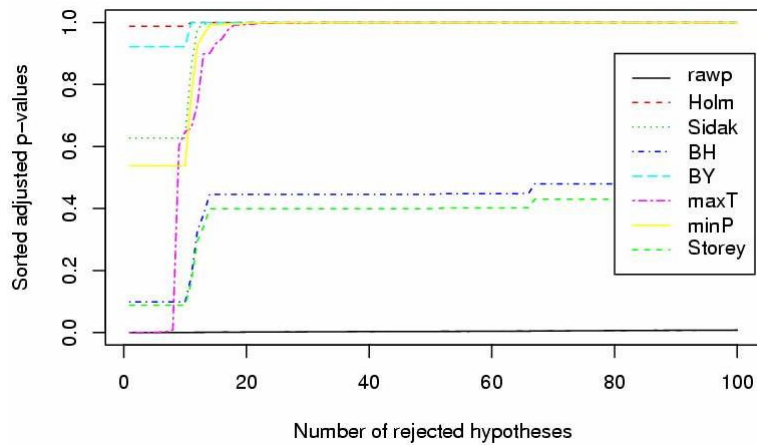
## Callow's AI ko data – complete permutation



© Eric Xing @ CMU, 2005-2009

60

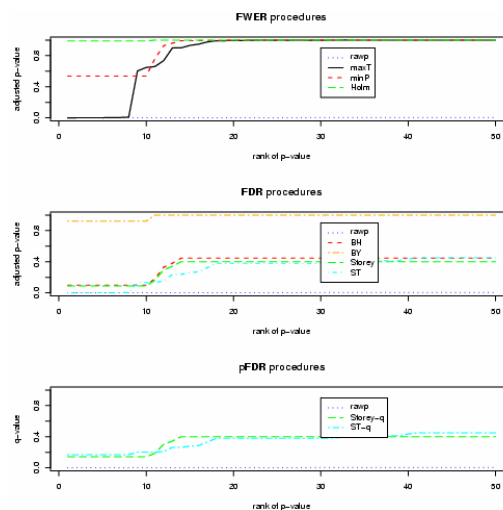
## Callow data with some FDR values included



© Eric Xing @ CMU, 2005-2009

61

## Comparison



© Eric Xing @ CMU, 2005-2009

62

## Comparison



gene index	t statistic	unadj. p ( $\times 10^4$ )	minP adjust.	plower	maxT adjust.
2139	-22	1.5	.53	$8 \times 10^{-5}$	$2 \times 10^{-4}$
4117	-13	1.5	.53	$8 \times 10^{-5}$	$5 \times 10^{-4}$
5330	-12	1.5	.53	$8 \times 10^{-5}$	$5 \times 10^{-4}$
1731	-11	1.5	.53	$8 \times 10^{-5}$	$5 \times 10^{-4}$
538	-11	1.5	.53	$8 \times 10^{-5}$	$5 \times 10^{-4}$
1489	-9.1	1.5	.53	$8 \times 10^{-5}$	$1 \times 10^{-3}$
2526	-8.3	1.5	.53	$8 \times 10^{-5}$	$3 \times 10^{-3}$
4916	-7.7	1.5	.53	$8 \times 10^{-5}$	$8 \times 10^{-3}$
941	-4.7	1.5	.53	$8 \times 10^{-5}$	0.65
2000	+3.1	1.5	.53	$8 \times 10^{-5}$	1.00
5867	-4.2	3.1	.76	0.54	0.90
4608	+4.8	6.2	.93	0.87	0.61
948	-4.7	7.8	.96	0.93	0.66
5577	-4.5	12	.99	0.93	0.74

© Eric Xing @ CMU, 2005-2009

63

## The gene names



**Index      Name**

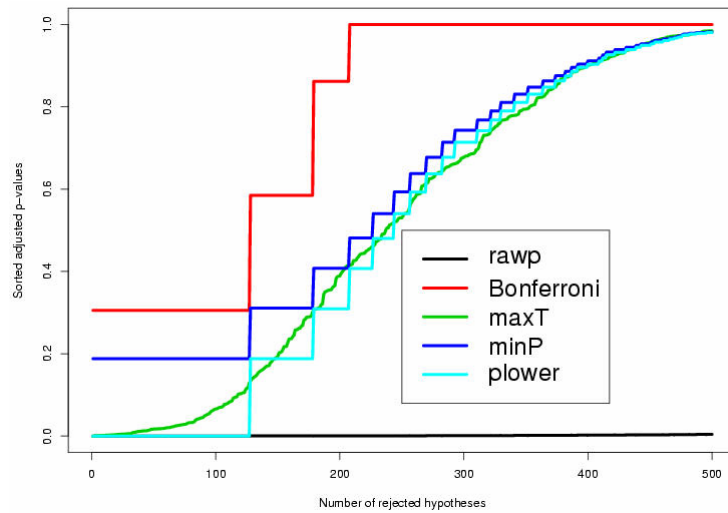
- 2139      Apo AI
- 4117      EST, weakly sim. to STEROL DESATURASE
- 5330      CATECHOL O-METHYLTRANSFERASE
- 1731      Apo CIII
- 538 EST, highly sim. to Apo AI
- 1489      EST
- 2526      Highly sim. to Apo CIII precursor
- 4916      similar to yeast sterol desaturase

© Eric Xing @ CMU, 2005-2009

64



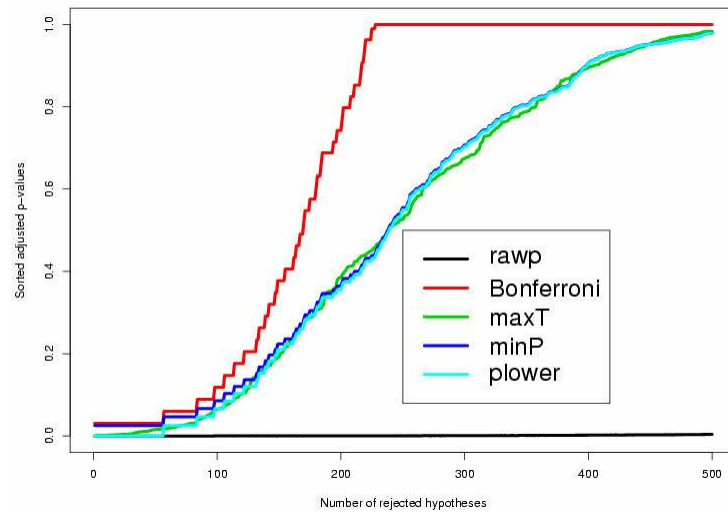
## Golub's data --- 10K simulations



© Eric Xing @ CMU, 2005-2009

65

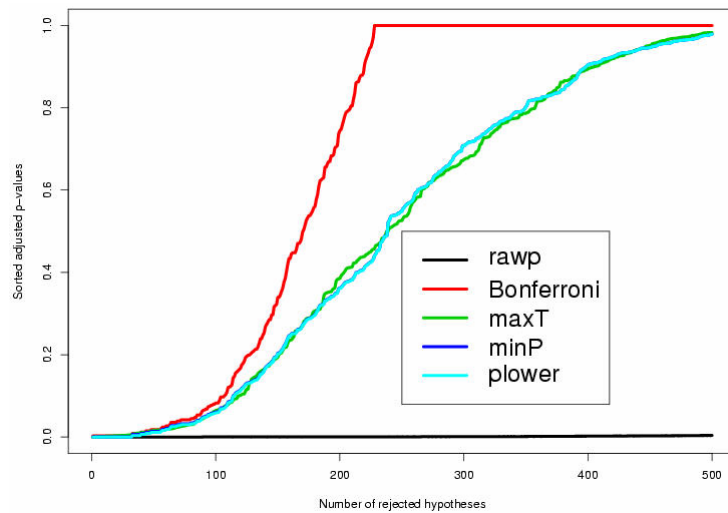
## Golub's data --- 100K simulations



© Eric Xing @ CMU, 2005-2009

66

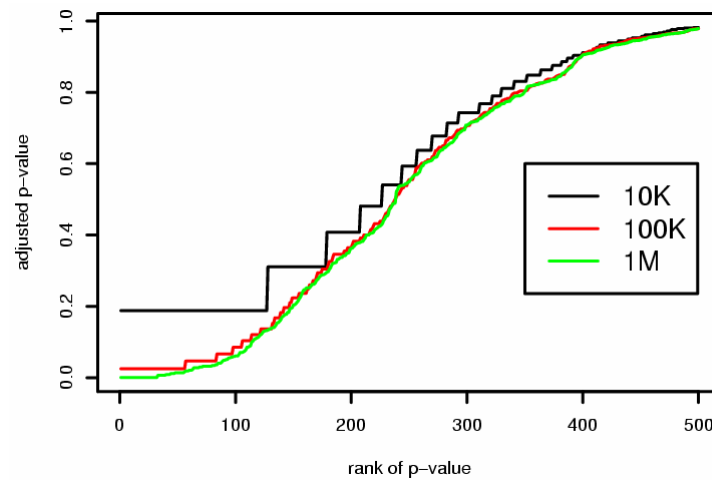
## Golub's data --- 1M simulations



© Eric Xing @ CMU, 2005-2009

67

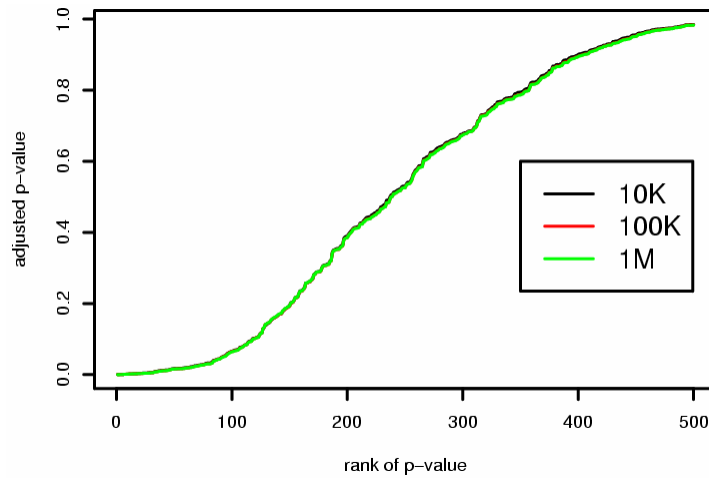
## Golub data with minP



© Eric Xing @ CMU, 2005-2009

68

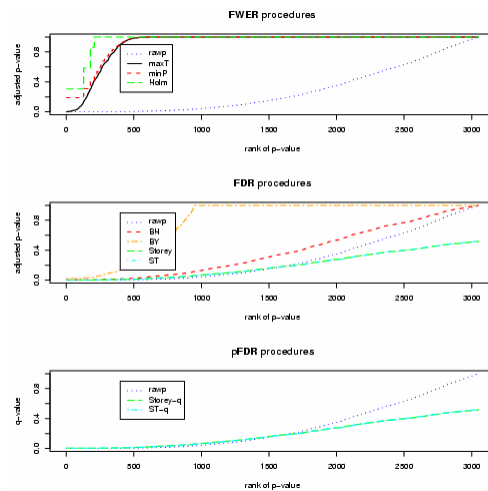
## Golub data with maxT



© Eric Xing @ CMU, 2005-2009

69

## Comparisons



© Eric Xing @ CMU, 2005-2009

70

## What should one look for in a multiple testing procedure?



- There is a bewildering variety of multiple testing procedures. How can we choose which to use? There is no simple answer here, but each can be judged according to a number of criteria:
  - **Interpretation:** does the procedure answer a relevant question for you?
  - **Type of control:** strong, exact or weak?
  - **Validity:** are the assumptions under which the procedure applies clear and definitely or plausibly true, or are they unclear and most probably not true?
  - **Computability:** are the procedure's calculations straightforward to calculate accurately, or is there possibly numerical or simulation uncertainty, or discreteness?

© Eric Xing @ CMU, 2005-2009

71

## Discussion



- The minP adjustment seems more conservative than the maxT adjustment, but is essentially model-free.
- With the Callow data, we see that the adjusted minP values are very discrete; it seems that 12,870 permutations are not enough for 6,000 tests.
- With the Golub data, we see that the number of permutations matters. Discreteness is a real issue here to, but we do have enough permutations.
- The same ideas extend to other statistics: Wilcoxon, paired t, F, blocked F.
- Same speed-up works with the bootstrap.

© Eric Xing @ CMU, 2005-2009

72

## Selected references



- Westfall, PH and SS Young (1993) *Resampling-based multiple testing: Examples and methods for p-value adjustment*, John Wiley & Sons, Inc
- Benjamini, Y & Y Hochberg (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing *JRSS B* 57: 289-300
- J Storey (2001): 3 papers (some with other authors), [www-stat.stanford.edu/~jstorey/](http://www-stat.stanford.edu/~jstorey/)
  - The positive false discovery rate: a Bayesian interpretation and the q-value.
  - A direct approach to false discovery rates
  - Estimating false discovery rates under dependence, with applications to microarrays
- Y Ge et al (2001) Resampling-based multiple testing for microarray data analysis, *Test* (to appear), see #633 in <http://www.stat.Berkeley.EDU/tech-reports/index.html>
- Software
  - C and R code available for different tests: multtest in <http://www.bioconductor.org>

## Acknowledgements



Slides are adapted from Lecture  
notes of  
Terry Speed