

Bayesian multiple testing and q-values for genetics applications

Oleg Makhnin
New Mexico Tech
Dept. of Mathematics

September 28, 2012

Bioinformatics applications require simultaneous testing of multiple quantitative traits. For example, the studies of gene expression are common, for which two or more groups are compared (say, Normal and High Fat diets) and we wish to identify which genes are expressed differently in these groups.

A Bayesian framework is developed using which we can borrow information between different traits tested, and perform multiple testing more efficiently than using a simple Bonferroni approach. Additionally, we discuss the concept of q-values, or "False Discovery Rates" that can be obtained naturally using our approach.

Multiple testing

Bayesian approach and q-values

Application to microbial genome testing

Bioinformatics applications require simultaneous testing of multiple quantitative traits. For example, the studies of gene expression are common, for which two or more groups are compared (say, Normal and High Fat diets) and we wish to identify which genes are expressed differently in these groups.

A Bayesian framework is developed using which we can borrow information between different traits tested, and perform multiple testing more efficiently than using a simple Bonferroni approach. Additionally, we discuss the concept of q-values, or "False Discovery Rates" that can be obtained naturally using our approach.

Multiple testing

Bayesian approach and q-values

Application to microbial genome testing

Bioinformatics applications require simultaneous testing of multiple quantitative traits. For example, the studies of gene expression are common, for which two or more groups are compared (say, Normal and High Fat diets) and we wish to identify which genes are expressed differently in these groups.

A Bayesian framework is developed using which we can borrow information between different traits tested, and perform multiple testing more efficiently than using a simple Bonferroni approach. Additionally, we discuss the concept of q-values, or "False Discovery Rates" that can be obtained naturally using our approach.

Multiple testing

Bayesian approach and q-values

Application to microbial genome testing

Bioinformatics applications require simultaneous testing of multiple quantitative traits. For example, the studies of gene expression are common, for which two or more groups are compared (say, Normal and High Fat diets) and we wish to identify which genes are expressed differently in these groups.

A Bayesian framework is developed using which we can borrow information between different traits tested, and perform multiple testing more efficiently than using a simple Bonferroni approach. Additionally, we discuss the concept of q-values, or "False Discovery Rates" that can be obtained naturally using our approach.

Multiple testing

Bayesian approach and q-values

Application to microbial genome testing

Hypothesis testing

- ▶ Comparing Sample 1 and Sample 2. For example, Sample 1 = rats on High-fat diet (HFD, "treatment"), Sample 2 = rats on Regular diet ("control").

Some genes might be involved in Metabolic Syndrome that is caused by high-fat diet. Which genes?

Hypothesis testing

- ▶ Comparing Sample 1 and Sample 2. For example, Sample 1 = rats on High-fat diet (HFD, "treatment"), Sample 2 = rats on Regular diet ("control").

Some genes might be involved in Metabolic Syndrome that is caused by high-fat diet. Which genes?

- ▶ **Hypothesis test** tries to separate the random fluctuations in data from the real changes (presumably) caused by HFD.

Testing if some condition (in our case, **gene expression level**) has changed from Sample 1 to Sample 2.

P-value

P-value is a measure of how extreme the change in each gene is.

$$P\text{-value}_i = P(\textit{observed difference}_i | \textit{no change})$$

Declare the change (in gene i) **statistically significant** when **$p\text{-value}_i < \alpha$** , for a given threshold α .

α = false positive rate (FPR). For N experiments with no change at all, in approximately $N\alpha$ we will declare that a change has occurred!

Multiple tests: corrections needed

- ▶ ● **Why a correction is needed:** too many false positives?
If, say, $N = 1000$ tests are run, each with a false positive rate (FPR) of 0.05, then we'll observe $\approx 1000(0.05) = 50$ false positives for the entire experiment

Multiple tests: corrections needed

- ▶ ● **Why a correction is needed:** too many false positives?
If, say, $N = 1000$ tests are run, each with a false positive rate (FPR) of 0.05, then we'll observe $\approx 1000(0.05) = 50$ false positives for the entire experiment
- ▶ ● **Bonferroni correction**
Thus, run each test (gene) with a smaller positive rate. If the final FPR of α is needed, run each test with a FPR of α/N .

Q-value

Q-value (Storey and Tibshirani, 2003) is the opposite of p-value:

P-value \approx P(Test positive | no change) \mapsto False Positive Rate

Q-value = P(no change | Test positive)
 \mapsto False Discovery Rate (FDR)

Like with p-value, declare genes with q-value $< \alpha$ “statistically significant”.

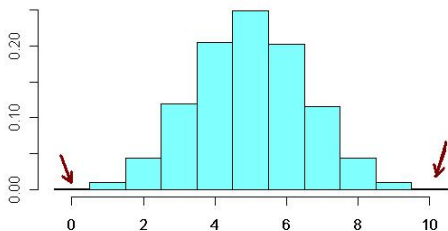
Q-value may be more desirable to practitioners: “Which portion of genes I have ‘discovered’ are bogus?”

P-value example

- ▶ Suppose we toss a coin 10 times, and observe 0 Heads. Are we inclined to believe that $P(\text{Heads}) = 1/2$?

P-value example

- ▶ Suppose we toss a coin 10 times, and observe 0 Heads. Are we inclined to believe that $P(\text{Heads}) = 1/2$?



In our case, p-value

$$= P(0 \text{ Heads}) + P(10 \text{ Heads}) = 2 * 0.5^{10} \approx 0.002$$

Since this is unlikely, we claim that the coin must be unfair.

In 1000 tosses, we'll have 2 false positives.

Q-value example

On the other hand, q-value will depend on what we know about the person tossing. If we trust the person, we estimate $P(\text{fair coin}) = 0.9$

A little calculation:

$$\begin{aligned} P(0 \text{ or } 10 \text{ Heads}) &= P(0 \text{ or } 10 \text{ Heads} | \text{fair}) * P(\text{fair}) + \\ &+ P(0 \text{ or } 10 \text{ Heads} | \text{unfair}) * P(\text{unfair}) = 0.002 * 0.9 + 0.1 * 0.1 \\ &\approx 0.012 \end{aligned}$$

$$P(\text{unfair} | 0 \text{ or } 10 \text{ Heads}) \approx 0.002 / 0.012 = 1/6$$

Hidden variable method

Idea 1: introduce *hidden variables* h_i which indicate whether the change occurred.

$$\begin{aligned}h_i &= 1, && \text{if change in Gene } i \\h_i &= 0, && \text{if no change in Gene } i\end{aligned}$$

Bayesian **Markov Chain Monte Carlo** computation through **Gibbs sampler** produces samples of all unknown variables and parameters. Our method allows for easy estimation of q-values:

$$q_i = \#\{h_{i,m} = 0, m = 1, \dots, M\} / M$$

where $h_{i,m}$, is the m th sample from the hidden variable h_i .

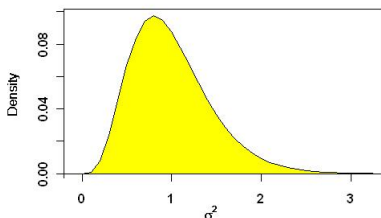
M is the Monte Carlo sample size (do not confuse with the sample size n for our data and with number of genes N). Higher M leads to more precise estimates of q_i .

Dealing with σ

Each gene has (potentially) different variability σ_i . For small sample sizes, it is hard to estimate σ_i well. Leads to highly inefficient Bonferroni method.

Earlier work (Dapo) assumed that σ_i 's are the same for all genes $i, i = 1, \dots, N$, but it's an oversimplification.

Idea 2: pool the estimates of σ_i from different genes i . This means that σ_i have some common distribution to be estimated.



Microbial study

PCE (perchloroethylene) used in dry cleaning

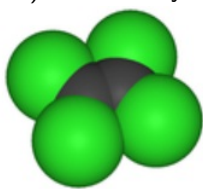
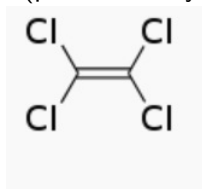


Image: Wikipedia

A well contaminated with PCE was injected with oil to stimulate microbial growth. The microbes break down PCE. We measure the species' abundance before and after remediation.

In this example, we are concerned not with individual genes, but rather whole microbial species. The data are counts of RNA segments (extracted from well water) that are attributed to different species.

Model

Let $L_{i,k}^{\text{before}}$ and $L_{i,k}^{\text{after}}$ be the expression levels (counts of RNA segments) for species i , sample k , from before and after the remediation.

The data here come from *matched pairs* design, so we'll take as model input $X_{i,k} := \log(L_{i,k}^{\text{after}} / L_{i,k}^{\text{before}})$
Logging is done to make data more Normally distributed.

Model

Let $L_{i,k}^{\text{before}}$ and $L_{i,k}^{\text{after}}$ be the expression levels (counts of RNA segments) for species i , sample k , from before and after the remediation.

The data here come from *matched pairs* design, so we'll take as model input $X_{i,k} := \log(L_{i,k}^{\text{after}} / L_{i,k}^{\text{before}})$
Logging is done to make data more Normally distributed.

- $X_{i,k}$ are Normal with mean μ_i and variance σ_i^2 ;
shortly $X_{i,k} \sim \mathcal{N}(\mu_i, \sigma_i^2)$

Model

Let $L_{i,k}^{\text{before}}$ and $L_{i,k}^{\text{after}}$ be the expression levels (counts of RNA segments) for species i , sample k , from before and after the remediation.

The data here come from *matched pairs* design, so we'll take as model input $X_{i,k} := \log(L_{i,k}^{\text{after}} / L_{i,k}^{\text{before}})$
Logging is done to make data more Normally distributed.

- $X_{i,k}$ are Normal with mean μ_i and variance σ_i^2 ;
shortly $X_{i,k} \sim \mathcal{N}(\mu_i, \sigma_i^2)$
- $\mu_i = 0$ if no change for species i (i.e. $h_i = 0$)
and $\mu_i \sim \mathcal{N}(0, \tau^2)$ otherwise

Model

Let $L_{i,k}^{\text{before}}$ and $L_{i,k}^{\text{after}}$ be the expression levels (counts of RNA segments) for species i , sample k , from before and after the remediation.

The data here come from *matched pairs* design, so we'll take as model input $X_{i,k} := \log(L_{i,k}^{\text{after}} / L_{i,k}^{\text{before}})$
Logging is done to make data more Normally distributed.

- $X_{i,k}$ are Normal with mean μ_i and variance σ_i^2 ;
shortly $X_{i,k} \sim \mathcal{N}(\mu_i, \sigma_i^2)$
- $\mu_i = 0$ if no change for species i (i.e. $h_i = 0$)
and $\mu_i \sim \mathcal{N}(0, \tau^2)$ otherwise
- σ_i^2 have Inverse Chi-square distribution, shortly
 $\sigma_i^2 \sim \text{Inv}\chi^2(\nu_0, \sigma_0^2)$

Microbial study: results

1594 out of 1900 most abundant species had q-value below 0.05,
1415 had q-value below 0.01.

Other model parameters: (estimates \pm standard errors)

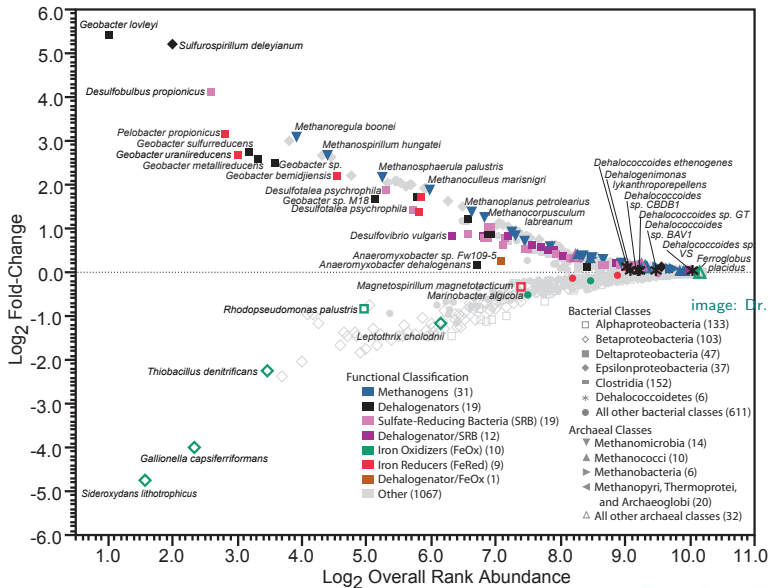
Center of σ -distribution $s_0 = 0.109 \pm 0.004$

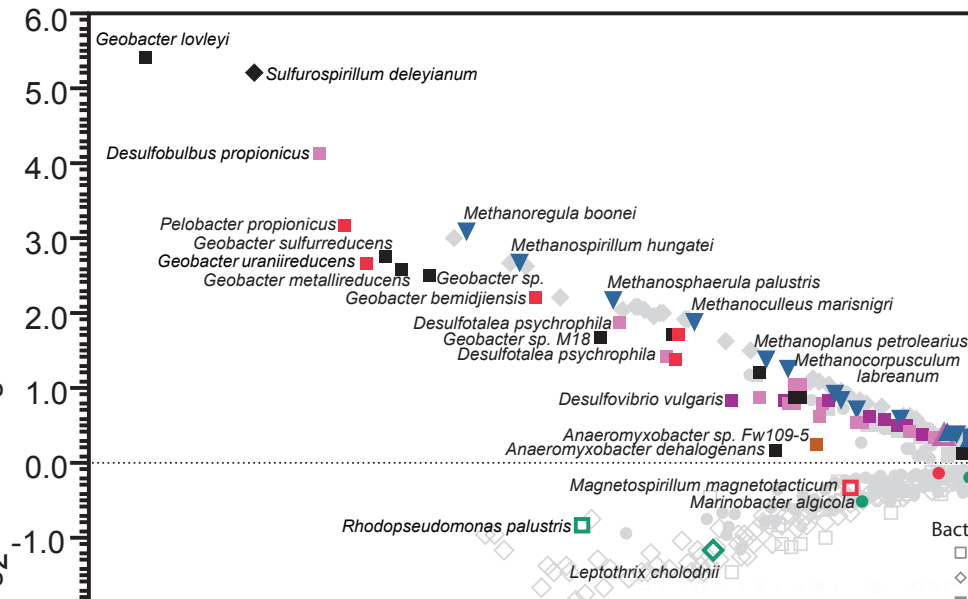
Degrees of freedom (tightness) of σ -distribution

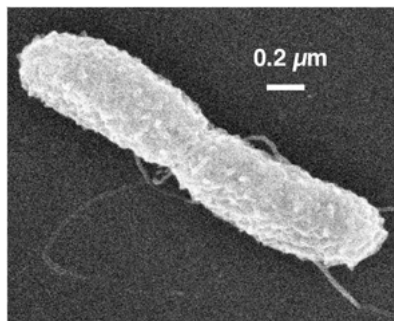
$\nu_0 = 2.67 \pm 0.22$

Probability of change $p_{DE} = 0.949 \pm 0.008$

Average change size (\log scale) $\tau = 1.34 \pm 0.02$







Geobacter lovleyi

QUESTIONS?

THANK YOU!

see `www.nmt.edu/~olegm/talks/BMT`