# 약 력

## 1. 인적사항

| | | |
|---|---|---|
| 성  명 | 김 선 영 | |
| 소속기관 | 생명공학연구원 | |
| 직   위 | 선임연구원 | |
| 전자메일 | kimsy@kribb.re.kr | |

## 2. 학력/경력

| 연    도 | 학교 / 기관 | 전공 / 직위 | 학위 / 비고 |
|---|---|---|---|
| 1988. 3 ~ 1992. 2 | 서울대학교 / 미생물학과 | 미생물학 | 학사 |
| 1992. 3 ~ 1994. 2 | 서울대학교 / 미생물학과 | 미생물학 | 석사 |
| 1994. 3 ~ 1998. 8 | 서울대학교 / 미생물학과 | 미생물학 | 박사 |
| 1998. 9 ~ 2001. 6 | 생명공학연구소/ 세포생물학 실험실 | 박사 후 연구원 | |
| 2001. 7 ~ 2004.12 | Dept. Pathology, Columbia University | 박사 후 연구원 | |
| 2005. 1 ~ 현재 | 생명공학연구원/의약유전체연구센터 | 선임 연구원 | |

## 3. 주요연구실적(개조식, 간단하게)

• *Kim S.Y. and *Volsky DJ (2005) PAGE: Parametric Analysis of Gene Set Enrichment, *BMC Bioinformatics*, 6:144, *(equal corresponding authors)

• Kim S.Y., Li J, Bentsman G, Brooks AI, Volsky DJ. (2004) Microarray analysis of changes in cellular gene expression induced by productive infection of primary human astrocytes: implications for HAD, *J Neuroimmunol*, 157(1-2):17-26.

• Kim, S. Y., Choi, S. Y., Wei, C., and D. J. Volsky, (2003) Transcriptional regulation of human excitatory amino acid transporter 1 (EAAT1): Cloning of the EAAT1 promoter and characterization of its basal and inducible activity in human astrocyte. *J. Neurochemistry*, 87: 1485-1498

• Kim, S. Y., Wei, C., Choi, S.Y., and D. J. Volsky, (2003) Cloning and Characterization of the 3'-untranslated region of the human excitatory amino acid transporter 2 transcript, J. *Neurochemistry*, 86(6): 1458-67

# Gene Set and Pathway Analysis of Microarray Data

생명공학연구원 선임연구원 **김 선 영**

최근의 microarray 기술의 발달로 인해 점점 더 많은 양의 mRNA 발현 데이터가 쌓여 가고 있다. 이제는 데이터를 만드는 단계보다는 데이터로부터 중요한 생물학적 의미를 끌어내는 것이 더욱 중요한 일이 되었다. micorarray 기술이 처음 도입된 이후로, 많은 앨고리즘과 소프트웨어가 개발되어, 실험자들이 microarray 데이터로부터 생물학적 의미를 끌어내는 작업을 도와주어 왔다. 그런데, 이전의 데이터 마이닝 방법들은 거의 예외 없이 전체 데이터로부터 선택된 몇 십, 몇 백 개의 유전자 리스트로부터 출발한다. 그런데, 이러한 방법 (over-representation analysis, ORA로 줄임)은 몇 가지 한계를 가지고 있어서, 최근에는 전체 데이터로부터 의미 있는 유전자 세트 (gene set)를 찾아내는 방법들이 도입되었다. 본 세미나는 이런 방법들, 줄여서 gene set analysis라 함, 에 사용되는 앨고리즘들과 소프트웨어들을 비교, 검토하고자 한다.

# Gene Set and Pathway Analysis of Microarray Data

---

**Outline**

- Existing GO and Pathway Analysis Tools

- Disadvantages of Over-Representation Analysis (ORA)

- Gene Set Analysis or Functional Class Scoring (FCS)

- PAGE (Parametric Analysis of Gene Set Enrichment)

- Comparison of Current Gene Set Analysis Tools

---

## Interpretation of Microarray Data



Raw Data → preprocessing

clustering

gene selection:
t-test, ANOVA, ....

List of a few hundred genes

Functional annotation, pathway analysis, ...

Identify coregulated genes or samples

---

## Evolution history of GO-based functional analysis software



*Khatri & Draghici, 2005, Bioinformatics*

---

## Review of GO analysis tools: statistical model



*Khatri S Draghici, 2005, Bioinformatics*

---

## Review of GO analysis tools: user interface



*Khatri S Draghici, 2005, Bioinformatics*

## Existing pathway analysis software



Curtis et al., 2005. Trends Biotech.

---

## ORA (Over-Representation Analysis)

### Contingency Table



count genes with GO term in set

count genes without GO term in set

count in set (e.g. differentially expressed genes)

Count in reference set (e.g. all genes on array)

P-value

$8\times10^{-52}$

Fisher's exact test or $\chi^2$-square test

| 51 | 416 | 467 |
| 125 | 8588 | 8713 |
| 173 | 9004 | 9177 |

---

## Problems of Over-Representation Analysis (ORA)

•Arbitrary Cut-off (or threshold)

•Actual values are ignored after gene selection

•Many small changed, but important genes are ignored.

---

## Problem 1. Arbitrary Cut-off

Hypothetical example          Cut-off: |log2(FC)| > 1, p-value < 0.05

| Gene | FC (log2) | P-value (t-test) | Selected? |
|------|-----------|------------------|-----------|
| Gene1 | 3 | 0.001 | Yes |
| Gene2 | 0.02 | 0.98 | No |
| Gene3 | 1.01 | 0.049 | Yes |
| Gene4 | 0.99 | 0.002 | No |
| Gene5 | 1.01 | 0.051 | No |
| Gene6 | 4.2 | 0.051 | No |

---

## Problem 2: Actual values are ignored after selection

| Gene | FC (log2) | P-value (t-test) |
|------|-----------|------------------|
| Gene1 | 3 | 0.001 |
| Gene2 | 1.1 | 0.02 |
| Gene3 | 1.01 | 0.049 |
| Gene4 | 3.99 | 0.00002 |
| Gene5 | 1.01 | 0.01 |
| Gene6 | 4.2 | 0.00001 |

Gene1, Gene2, Gene3, Gene4, Gene5, Gene6... AS USER INPUT

---

## Problem 3: If No or Few genes pass selection?

No genes to input into statistical test.

Data mining becomes impossible.

## Slide 1

### Gene Set Enrichment Analysis

Compare two experimental groups at the level of pre-defined *gene sets* rather than *individual genes*

Biological change often occurs by *co-regulation* of related genes

By examining groups of genes, GSEA often *detects biological changes that were hidden at individual gene level.*

GSEA is especially useful when *individual* gene expression changes are *minimal* or *moderate.*

First introduced by Mootha et al. (2003), *Nat. Genet.*

## Slide 2

GSEA: procedures



*Mootha et*

## Slide 3

### GSEA with 149 gene sets



*Mootha et al., 2003*

## Slide 4



*Mootha et al. 2003*

## Slide 5

PAGE: Parametric Analysis of Gene set Enrichment

*Kim and Volsky, 2005, BMC Bioinformatics*

## Slide 6

### Enrichment Score (ES) vs. Z score

**ES**



N: total genes

G: gene set

**Z score**

$$Z = \frac{X - \mu}{\sigma}$$

X: Mean of test statistics of *a gene set G*

$\mu$: population N (total genes in a microarray)'s mean of test statistics

$\delta$: population's standard deviation divided by square root of the size of a gene set

Test statistic: fold change, t-value, S2N value...

## Slide 1

"PAGE is a *parametric* version of GSEA"

### Comparison of PAGE with GSEA

|  | PAGE | GSEA |
|---|---|---|
| Nature of statistical test | Parametric Statistical Test | Nonparametric Statistical Test |
| test statistic | fold change, correlation coefficient etc. | Rank of genes |
| Result | Z score | Enrichment Score |
| Background Distribution | *Standard Normal Distribution* | Calculated from permutation of data set |
| Computation Time | Fast | Slow |
| Sensitivity | Sensitive | Less sensitive |

## Slide 2

PAGE                    GSEA

1. Collect gene sets

2. Calculate test statistic by expression difference

3. Calculate Z score for each gene sets

4. Evaluate significance of Z scores against *Standard Normal* distribution

*Mootha et al. 2003, Nat. Gen.*

## Slide 3

*"Is PAGE statistically well-based?"*          *"Yes"*

### Statistical Background of PAGE

**Central Limit Theorem:**

**"The distribution of an average tends to be Normal, even when the distribution from which the average is computed is decidedly non-Normal."**

N = 1          N = 3          N = 8          N= 32

## Slide 4

### Distribution Pattern of Microarray Data

**N=1**

Kolmogorov-Smirnov Normality Test
$H_0$: Distribution is Normal
D = 0.08, p-value < 2.2e-16
=> *Non-Normal*

**N=10**

Kolmogorov-Smirnov Normality Test
$H_0$: Distribution is Normal
D = 0.0239, p-value = 0.1783
=> *Normal*

## Slide 5

### Comparison of PAGE with GSEA

| Gene Set | PAGE Z score | p-value | GSEA Gene Set | ES | p-value |
|---|---|---|---|---|---|
| OXPHOS_HG-U133A | -10.5835 | 1.0E-11 | OXPHOS_HG-U133A | 348.8827 | 0.003 |
| human_mitoDB_6_2002_HG-U133A | -6.7213 | 1.81E-11 | human_mitoDB_6_2002_HG-U133A | 215.0424 | 0.011 |
| mitochond_HG-U133A | -6.4761 | 9.44E-11 | mitochond_HG-U133A | 207.9381 | 0.017 |
| MAP00150_Oxidative_phosphorylation | -4.5745 | 4.78E-06 | c20_U133 | 181.1569 | 0.042 |
| c20_U133 | 3.7461 | 0.0002 | MAP00190_Oxidative_phosphorylation | 148.9051 | 0.084 |
| c25_U133 | -2.7817 | 0.0058 | c22_U133 | 142.9006 | 0.026 |
| c21_U133 | -2.1116 | 0.0347 | c29_U133 | 131.4732 | 0.026 |

***PAGE* is statistically more *sensitive* than GSEA**

## Slide 6

### Simulation

Hypothetical data set of 20 samples (10 controls vs. 10 treatment) with 2,000 genes generated.

**10 Treatment Samples**

20
1,980

20 from rnorm, then + α
1,980 from rnorm

**10 Controls**

20
1,980

2,000 from rnorm

## Simulation study comparing PAGE with GSEA

| | PAGE | | GSEA | |
|---|---|---|---|---|
| Mean Difference | Z score | p-value | ES | p-value |
| 1 | 8.227 | 2.22E-16 | 128.54 | 0.001 |
| 0.95 | 8.055 | 8.88E-16 | 152.8663 | 0.001 |
| 0.9 | 8.522 | 0 | 112.3632 | 0.001 |
| 0.85 | 7.338 | 2.18E-13 | 133.268 | 0.003 |
| 0.8 | 6.706 | 2.01E-11 | 138.0922 | 0.001 |
| 0.75 | 5.232 | 1.68E-07 | 127.3383 | 0.001 |
| 0.7 | 6.932 | 4.18E-12 | 91.96096 | 0.001 |
| 0.65 | 7.067 | 1.59E-12 | 86.73476 | 0.001 |
| 0.6 | 5.078 | 3.82E-07 | 138.2932 | 0.001 |
| 0.55 | 6.071 | 1.28E-09 | 41.20655 | 0.014 |
| 0.5 | 4.481 | 7.44E-06 | 79.599 | 0.002 |
| 0.45 | 5.203 | 1.96E-07 | 20.00025 | 0.136 |
| 0.4 | 3.967 | 7.59E-05 | 60.70428 | 0.003 |
| 0.35 | 2.603 | 0.009241 | 23.71889 | 0.0555 |
| 0.3 | 2.132 | 0.033007 | 3.316825 | 0.7625 |
| 0.25 | 3.216 | 0.0013 | 7.939739 | 0.415 |
| 0.2 | 1.408 | 0.159131 | 25.9288 | 0.0675 |
| 0.15 | 1.711 | 0.087081 | 0 | 0.994 |
| 0.1 | 0.753 | 0.45146 | 9.04534 | 0.34 |
| 0.05 | -0.269 | 0.772582 | 9.54786 | 0.3785 |

---

## Application of PAGE to different Affymetrix probe level analysis methods

| MAS5 | | MBEI | | RMA | |
|---|---|---|---|---|---|
| Gene Set | Z score p-value | Gene Set | Z score p-value | Gene Set | Z score p-value |

---

## Comparison of PAGE results from data sets produced using different microarray platforms

| U95A | | U133A | | Agilent | |
|---|---|---|---|---|---|
| Gene Set | Z score p-value | Gene Set | Z score p-value | Gene Set | Z score p-value |

---

## Comparison of two microarray data sets at gene set level

| | gds287 | | gds472 | |
|---|---|---|---|---|
| Gene Set | Z score | p value | Z score | p value |

---

## Comparison of multiple data sets: Gene level vs. Gene Set level

Exp1, Exp2: Two identical experiments except cells (different batches of primary cells)

**Gene Level**
*|FC| > 2 & p < 0.05*

153 | 9 | 192   (9+9)/(162+201) = 5%

**Gene Set Level**
*|Z score| > 2*

186 | 98 | 166   (98+98)/(284+264)=36%
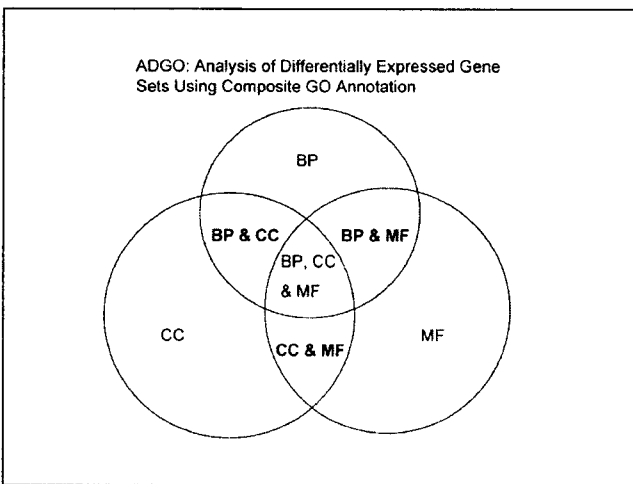
---

## Conclusion I

1. PAGE is statistically *more sensitive* than GSEA

2. PAGE is *easier* and *faster* to compute than GSEA

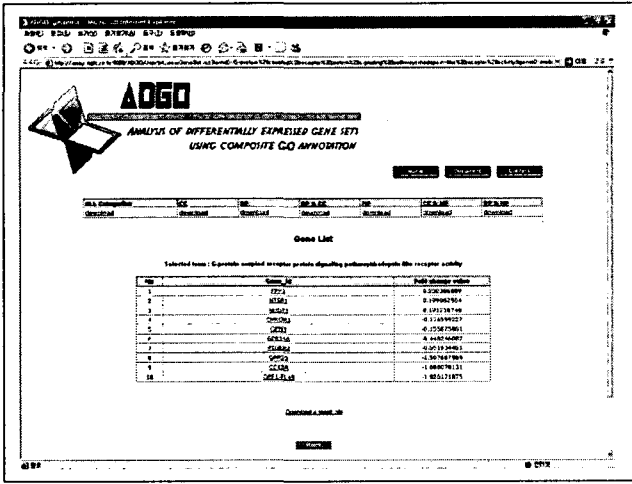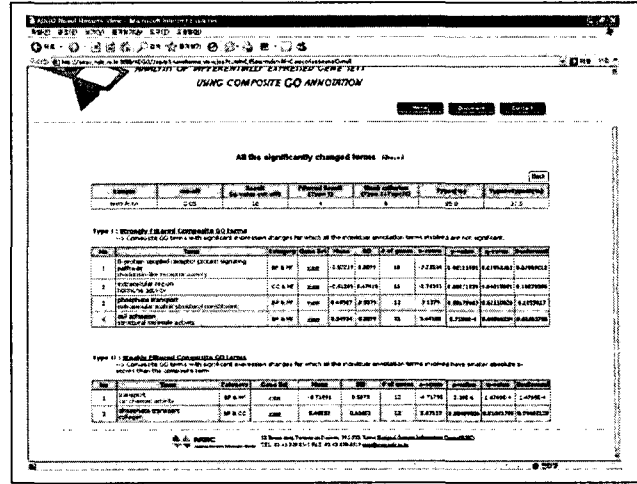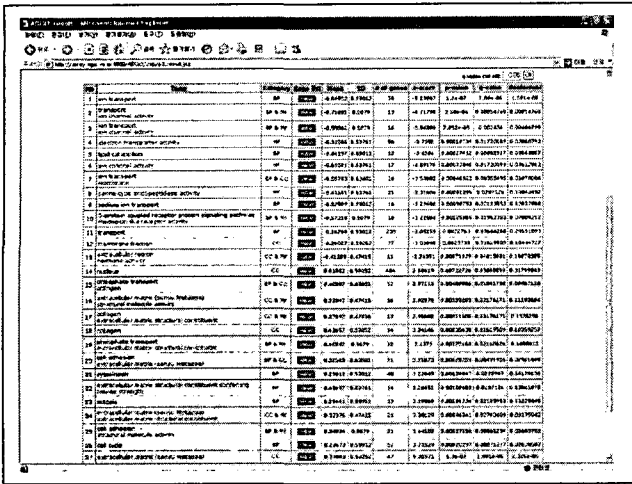3. PAGE produces consistent result over multiple platforms and data preprocessing methods

ADGO: Analysis of Differentially Expressed Gene Sets Using Composite GO Annotation

## Recent Works in FCS (Functional Class Scoring) Methods

Pavlidis, P. et al., 2002, *Pac. Symp. Biocomput.* 474-485

Mootha et al., 2003, *Nat. Gen.* 34: 267-273

Breslin et al., 2004, *BMC Bioinformatics*, 5: 193⁻          => Catmap

Al-Shahrour et al., 2005, *Bioinformatics*, 21: 2988-93

Tu et al., 2005, *Biotechniques*, 2: 277-83          => MEGO

Kim & Volsky, 2005, *BMC Bioinformatics*, 6: 144          => PAGE, GAzer

Boorsman et al., 2005, *Nuc. Acids Res.*, 33          => T-Profiler

Tian et al., 2005, *PNAS*, 38: 13544-9

Subramanian et al., 2005, *PNAS*, 38: 15545-50          => GSEA

Lee et al., 2005, BMC *Bioinformatics*, 6: 269          => ErmineJ

## Comparison of Gene Set Analysis Softwares

|  | PAGE | GSEA | ErmineJ | MEGO | Catmap | T-profiler |
|---|---|---|---|---|---|---|
| Used Statistics | Fold Change | Rank | P-value | Fold change | P-value | T-statistic |
| Statistical test | Z-test | permutation | permutation |  | permutation | t-test |
| Speed | Fast | Slow | Moderate | Fast | Slow | Fast |
| Standalone software | YES | YES | YES | YES | YES | NO |
| GUI | NO | YES | YES | YES | NO | NO |
| Web server | YES | NO | NO | NO | NO | YES |
| Organism | H, M, R, Y | H, M | H, M, R | H, | | Y |
| Gene Sets | GO, Pathways, Chromosomes | GO, Pathways, Chromosoma s | GO, Pathways, Chromosome s | GO, Pathways, Chromosome s | GO, Pathways, Chromosome s | GO, Pathways, Chromosome s |

## ACKNOWLEDGEMENTS

NGIC (National Genome Information Center)          Columbia University

In-Sun Chu          David J Volsky
Dougu Nam
Sang-Bae Kim
Sang-Cheol Kim
Seon-Kyu Kim
Seong-Jin Yang
Hyun-Goo Woo