

2008

# MCM-test: a fuzzy-set-theory-based approach to differential analysis of gene pathways

Lily R. Liang

*University of the District of Columbia, lliang@udc.edu*

Vinay Mandal

*Wayne State University, vinaym@wayne.edu*

Yi Lu

*Prairie View A&M University, yilu@pvamu.edu*

Deepak Kumar

*University of the District of Columbia, dkumar@udc.edu*

---

## Recommended Citation

Liang *et al.* *BMC Bioinformatics* 2008, 9(Suppl 6):S16

doi:10.1186/1471-2105-9-S6-S16

Available at: <http://digitalcommons.wayne.edu/biomedcentral/143>

Research

Open Access

## MCM-test: a fuzzy-set-theory-based approach to differential analysis of gene pathways

Lily R Liang\*<sup>1</sup>, Vinay Mandal<sup>2</sup>, Yi Lu<sup>3</sup> and Deepak Kumar\*<sup>4</sup>

Address: <sup>1</sup>Department of Computer Science and Information Technology, University of the District of Columbia, Washington, D.C., USA, <sup>2</sup>Department of Computer Science, Wayne State University, Michigan, USA, <sup>3</sup>Department of Computer Science, Prairie View A&M University, Texas, USA and <sup>4</sup>Department of Biological and Environmental Sciences, University of the District of Columbia, Washington, D.C., USA

Email: Lily R Liang\* - lliang@udc.edu; Vinay Mandal - vinaym@wayne.edu; Yi Lu - yilu@pvamu.edu; Deepak Kumar\* - dkumar@udc.edu

\* Corresponding authors

from Symposium of Computations in Bioinformatics and Bioscience (SCBB07)  
Iowa City, Iowa, USA. 13–15 August 2007

Published: 28 May 2008

BMC Bioinformatics 2008, 9(Suppl 6):S16 doi:10.1186/1471-2105-9-S6-S16

This article is available from: <http://www.biomedcentral.com/1471-2105/9/S6/S16>

© 2008 Liang et al; licensee BioMed Central Ltd.

This is an open access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

### Abstract

**Background:** Gene pathway can be defined as a group of genes that interact with each other to perform some biological processes. Along with the efforts to identify the individual genes that play vital roles in a particular disease, there is a growing interest in identifying the roles of gene pathways in such diseases.

**Results:** This paper proposes an innovative fuzzy-set-theory-based approach, Multi-dimensional Cluster Misclassification test (MCM-test), to measure the significance of gene pathways in a particular disease. Experiments have been conducted on both synthetic data and real world data. Results on published diabetes gene expression dataset and a list of predefined pathways from KEGG identified OXPHOS pathway involved in oxidative phosphorylation in mitochondria and other mitochondrial related pathways to be deregulated in diabetes patients. Our results support the previously supported notion that mitochondrial dysfunction is an important event in insulin resistance and type-2 diabetes.

**Conclusion:** Our experiments results suggest that MCM-test can be successfully used in pathway level differential analysis of gene expression datasets. This approach also provides a new solution to the general problem of measuring the difference between two groups of data, which is one of the most essential problems in most areas of research.

### Background

Current microarray technologies conduct simultaneous studies of gene expression data of thousands of genes under different conditions. In most of these studies, expression data are analyzed using various standard statistical methods to identify a list of genes responsible for a

particular condition. However, in these approaches, interplay among genes is not taken into account. Since organisms behave as complex systems, functioning through chemical networks and interaction of various molecules (also known as pathways), this interplay of genes may provide invaluable insight to the understanding of various

diseases. Thus, along with the efforts to identify the individual genes that play vital roles in a particular disease, there is a growing interest in identifying the roles of different pathways in such diseases.

Biological pathway is a group of related genes coding for proteins that interact with each other to perform some biological processes. According to the biological processes they are involved with, pathways can be divided into several categories, such as metabolic pathways and regulatory pathways. Metabolic pathways are series of chemical reactions occurring within a cell, catalyzed by enzymes, resulting in either the formation of a metabolic product to be used or stored by the cell, or the initiation of another metabolic pathway. Regulatory pathways represent protein-protein interactions.

During the past few years, many signaling and metabolic pathways have been discovered experimentally and have been integrated into pathway databases, such as KEGG [1] and Biocarta [2]. Various statistical techniques have been developed to analyze microarray expression data for the relevance of predefined pathways to a particular disease. These techniques include gene set enrichment analysis [3,4], pathway level analysis of gene expression using singular value decomposition by Tomfohr et al. [5], and hypothesis testing [6] by Tian et al. These approaches are reviewed in detail in the related works section.

Generally speaking, these approaches can be divided into two categories:

- Conduct statistical differential analysis at the individual gene level, and integrate the result statistics of the genes in the same pathway;
- Obtain activity level indices of each pathway for different sample groups and conduct differential analysis of these indices.

For the first category, when the statistics at individual gene level miss significant genes, the effectiveness of the pathway analysis will be affected. An example is given in the later part of this section. For the second approach, extracting pathway activity level indices from gene expression data may cause loss of information.

Diabetes is a group of diseases characterized by high levels of blood glucose resulting from defects in insulin production, insulin action, or both. It is one of the most common diseases, affecting 18.2 million people in the United States, or 6.3% of the population [7]. Hence, identifying active pathways in diabetes is a critical task for understanding this disease. Several pathway analysis works have been proposed in this area [3,5,6].

In gene set enrichment analysis (GSEA) [3], a differential statistic is calculated first for each gene from their expression data of two different groups of samples. Then the genes are ordered according to the statistic values. A running sum of weights is calculated from the ordered list for a particular pathway. The maximum value of this running sum is called the enrichment score of that pathway. To measure the significance of this score, a null distribution of enrichment scores is generated by permuting the sample labels. This approach falls into the first category stated previously, i.e., statistical analysis at individual gene level is performed followed by an integration of these statistics of genes in the same pathway.

In [5], a hypothesis testing framework for pathway differential analysis is proposed. T-test and Wilcoxon rank test are recommended to measure the difference of expressions of a single gene between two groups of samples. Then this statistic is accumulated over each gene in a particular pathway and standardized by the total number of genes in this pathway. The significance of the result is then interpreted by rejecting two null hypotheses, each with a null population generated by permuting sample labels or gene indices. This approach also belongs to the first category above. Statistical analysis at individual gene level is still required for the pathway analysis in this approach.

In [6], singular value decomposition is used to obtain pathway activity levels from the gene expression matrix. T-test is applied to the pathway activity levels of the two different sample groups to measure the difference. Significance of the measurement is also obtained by permuting the sample labels. In this approach, no differential analysis at individual gene level is required. However, an extraction of pathway activity level prior to the differential analysis is required. During this extraction process, since only the first eigenvector of singular value decomposition is used, some information of expressions is lost. This approach belongs to the second category stated above.

As discussed above, either t-test or rank sum test is used as a core step by [3,6] to identify individual genes which are expressed differently from two different sample groups. Thus these methods inevitably inherit the disadvantage of t-test and rank sum test. While the t-test is very sensitive to extreme values and cannot distinguish two sets with close means even though the two sets are significantly different from each others, the rank sum test is not sensitive to absolute values. In turn, those pathways contain genes which can not be identified by t-test or rank sum test but actually are significantly differently expressed in two different sample groups will be affected. For example, as showed in Table 1, the expressions of Gene 3 are significantly different under two conditions. However this gene was not identified by t-test. Thus, a pathway involving this

**Table 1: An example of five gene pathway**

Gene ID	$S_1$					$S_2$					CM $d$ -value	P-value		
	CM-test	t-test	Rank Sum test											
1	750	559	649	685	636	310	359	135	97	178	1	0.001	0.000	0.008
2	391	379	268	323	380	774	506	416	468	449	1	0.005	0.029	0.008
3	598	424	695	451	141	342	260	266	229	234	0.904	0.018	0.077	0.152
4	233	216	193	394	327	436	980	363	424	416	0.905	0.017	0.071	0.015
5	305	221	241	183	158	201	176	189	177	250	0.812	0.143	0.448	0.693

gene is less likely to be identified by the first category of analysis that uses t-test at the gene level.

In this paper, we propose an innovative fuzzy-set-theory-based approach for differential analysis of gene pathways and apply it on identifying significant pathways for diabetes. In our proposed MCM-test, instead of identifying individual genes first, the differential analysis is done directly at the pathway level without individual gene differential statistic. All expression values of genes which belong to a pathway of a particular patient are treated as a vector. The intuition behind this is based on the fact that genes for each patient interplay with each other. MCM-test does not extract activity level of pathways either. This allows keeping the maximum amount of information for the pathway differential analysis. Moreover, the fuzzy concept makes the approach more tolerant to individual data item noise.

**Results**

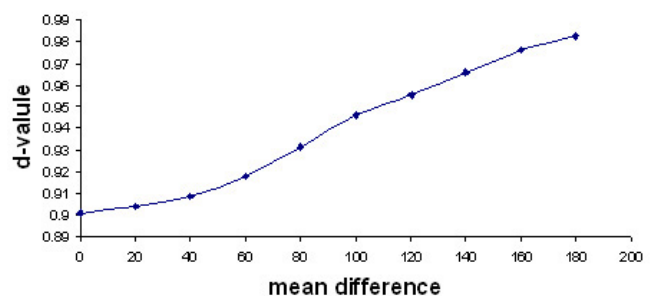
To investigate our approach, we conducted experiments on both synthetic data and real world data. We first conducted a series of experiments on synthetic datasets to find the characteristics of MCM  $d$ -value. We then used the MCM-test on the real world diabetes dataset analyzed by Tomfohr et al. [5] and GSEA [3]. Results on real world diabetes data identified several pathways that were deregulated in diabetes patients. The top three pathways identified were related to mitochondrial functions in accordance with previous diabetes studies. Mitochondrial dysfunction is known to be related to insulin resistance and type-2 diabetes. Our data suggests that the method can be successfully used in pathway level differential analysis of gene expression datasets.

**Relationship between MCM  $d$ -value and mean difference of the distributions**

Suppose two sets  $S_1$  and  $S_2$  are drawn from two different distributions, then a good divergence value will satisfy the following property: the less the overlap, the higher the  $d$ -value. To validate that our MCM-test has this property, we performed the following steps:

1. generated 17 values from Gaussian distribution  $N(\mu, \sigma)$ , where  $\mu$  is the mean and  $\sigma$  is the variance, to use as gene expression data. The number 17 was chosen to mimic the real world diabetes dataset used for the analysis in this paper.
2. repeated Step 1 for 100 times to get expression data of 100 genes
3. generated 17 values from Gaussian distribution  $N(\mu + x, \sigma)$ , with  $x = 0$  at this time.
4. repeated Step 3 for 100 times
5. analyzed these 100 pairs of sets of values with MCM-test and obtained the  $d$ -value.
6. repeated Step 1 to Step 5 for 1000 times and averaged the  $d$ -values over all the iterations.
7. repeated Step 1 to Step 6 for each  $x$ :  $x = 0, 20, 40, 60, 80, 100, 120, 140, 160,$  and  $180$ .

Figure 1 shows the average  $d$ -value verses mean difference. We can see that the MCM-test has the desired property: the



**Figure 1 Relationship between  $d$ -value and mean difference.** Two datasets are generated from two distributions  $N(\mu, \sigma)$  and  $N(\mu + x, \sigma)$ . As the mean difference,  $x$ , increases, the  $d$ -value also increases.

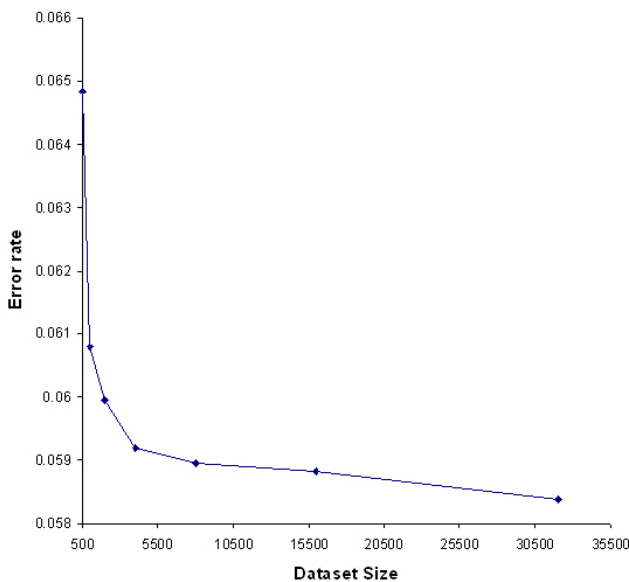
larger the mean difference between two sets, the larger the divergence d-value.

**Impact of population size on standard deviation of MCM d-value**

At this point, a natural question is, what the standard deviation of MCM d-values look like and how population size of d-value influences it. To answer these questions, we generated sample expression datasets and calculate d-values following a process similar to step 1 to 5 in the previous section. Again, we fixed the pathway length to 100. We repeat the process and obtained 500 d-values and calculated the standard deviation of the d-values. This is then repeated for 10 times and the 10 standard deviations are averaged and recorded as error rate of MCM d-value for that population size 500. Similarly we obtained the error rate for various d-value population sizes. As shown in Figure 2, the error rate decreases as the dataset size increases. We also note that the error rate becomes stable after the size of the population becomes greater than 8000.

**Relationship between MCM d-value and empirical p-value**

Suppose two vectors  $S_1$  and  $S_2$  are drawn from same Normal distribution. What is the probability that the MCM d-value of these vectors is greater than a particular  $D$ ? Does the probability increases with the increase of  $D$ ? To answer these questions, we studied the relationship between MCM d-value and empirical p-value as follows:



**Figure 2**  
Impact of number of permutation on the error rate of PDF of MCM d-value. We show the error rate for various numbers of permutations ranging from 500 to 32000. The error rate decreases as the number of permutations increases.

1. We generated 15000 pairs of sets, each set with 15 values from standard normal distribution.

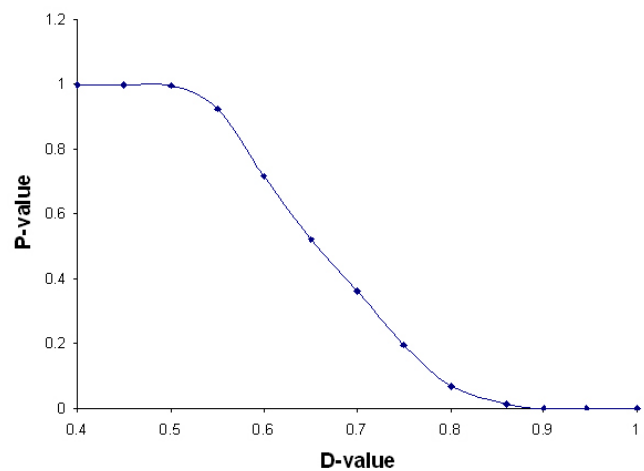
2. From these 15000 pairs of sets, we randomly selected 100 pairs of sets to simulate expression data of a pathway with 100 genes under two conditions. We calculated d-value for this pathway. Since we know that the data size required to obtain stable standard deviation of d-value is 8000 from the previous experiment, this process is repeated 10000 times.

3. For each pathway generated above with d-value  $D$ , we calculated the empirical p-value as  $n+1/10001$ , where  $n$  is the number of d-values generated above that are equal to or greater than  $D$ . The relationship between the d-value and p-value is shown in Figure 3.

In Figure 3 we can see that as the d-value increases, the p-value decreases. In particular, when d-value is greater than 0.809, we have p-value  $\leq 0.05$ .

**Impact of number of samples on error rate of MCM-test d-value**

In order to understand the effect of the number of samples on error rate of MCM d-value, we generated datasets with different sample sizes. For each sample size, we generated 10000 datasets and calculated the corresponding 10000 d-values. The standard deviation of these d-values was calculated. This process is repeated 10 times and the average of the standard deviations is recorded as the error rate. The same is done for the other sample sizes. The relationship between number of samples and the error rate is shown in



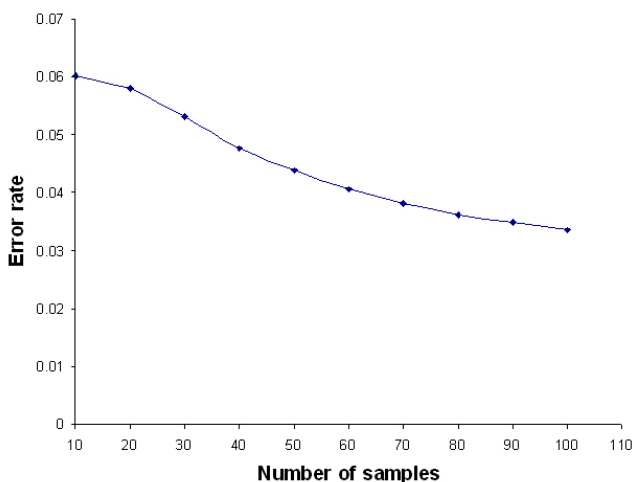
**Figure 3**  
Relationship between MCM d-value and its empirical p-value. As the d-value increases, the corresponding empirical p-value decreases.

Figure 4. As expected, the error rate decreases as the number of samples increases.

#### Analyzing the diabetes dataset with MCM-test

The diabetes dataset contains the transcriptional profiles of smooth muscle biopsies of diabetic and normal individuals. In the expression dataset, for each gene, there are 17 expression values from subjects with type 2 diabetes (DM2), 17 expression values from subjects with normal glucose tolerance (NGT) and 10 expression values from subjects with impaired glucose tolerance (IGT). For our analysis, we only used the 34 expression values from subjects with type 2 diabetes and subjects with normal glucose tolerance. Furthermore, we used about 150 pathways obtained from KEGG (Kyoto Encyclopedia of Genes and Genomes) [1].

The expression values in the dataset which are too small, i.e., less than 100 are considered to be the result of noise. So, to minimize the effect of these low values, we only included the genes which have at least one of the expression values greater than 100. Out of the 22,283 genes in the dataset, 10,983 met the filtering criteria. The d-value for each pathway was calculated as described in the methodology section before. The p-value for the pathway was calculated using permutation test. We permuted the genes 1000 times, each time selecting the same number of genes as that of the pathway under consideration. We then calculated the d-value of each pathway obtained thus and the p-value for the pathway was the fraction of times the d-values of the pathways obtained by 1000 permutation equaled or exceeded the original d-value.



**Figure 4**  
Impact of number of samples on error rate of PDF of MCM-test d-value. As the number of samples in a pathway increases, the error rate of PDF of MCM d-value also decreases.

The pathways are ordered in the ascending order of their p-values. The significant pathways, i.e., the pathways with p-value less than 0.05, are then ordered according to the percentage of the genes in the pathway which were represented in the dataset. Table 2 shows the result after sorting.

Using our method, we identified the deregulation of mitochondrial pathways in the dataset which is in accordance with previous studies. The first cluster of genes involved was from the mitochondrial OXPHOS pathway. The OXPHOS pathway was well represented in the data with 93% of genes (106 out of 114) present in the dataset. Oxidative phosphorylation in mitochondria provides energy in the form of ATP generation. In muscle cells, mitochondrial dysfunction has been linked to insulin resistance and type-2 diabetes [8-10]. The involvement of genes coded by mitochondria in insulin resistance is also well known. The depletion of cellular mitochondrial DNA has been shown to cause insulin resistance in experimental model [11]. Reduced mitochondrial oxidative phosphorylation leads to the accumulation of intracellular triglycerides resulting in insulin resistance. The next 2 clusters, c20\_U133 which is a manually curated cluster of genes coregulated with OXPHOS [3] and the mitochondrial gene cluster human\_mitoDB\_6\_2002 reinforce that muscle mitochondrial dysfunction is linked to type-2 diabetes.

#### Conclusion

In this paper, we propose an innovative fuzzy-set-theory-based approach for differential analysis of gene pathways and apply it on identifying significant pathways for diabetes. Experiments have been conducted on both synthetic datasets and real world dataset. Results on real world diabetes data identified several number of gene pathways. Of note our top significant pathways were related to mitochondrial function which is well known to be involved in insulin resistance and type-2 diabetes. This approach can be used not only for pathway analysis of other diseases but also for other domains. As measuring the difference of two groups of data are essential to most of researches, our approach provides a solution to this general and most critical problem.

#### Methods

In [12-14], we proposed two fuzzy-set-theory based methods, CM-test and FM-test, to identify the individual genes that expressed significant differences under two conditions. In this paper, we extended the cluster misclassification concept to a multi-dimensional space and propose a new approach for pathway analysis, Multi-dimensional Cluster Misclassification test (MCM-test). Comparing with CM-test and FM-test, MCM-test looks for a group of genes significant under two conditions instead of identifying significant individual genes under two conditions. In

**Table 2: The results from MCM-test on diabetes dataset**

Pathway Name	MCM-test p-value	No. of genes hits in dataset	Actual no. of genes in pathway	Percentage of gene hits
OXPHOS	0.04995	106	114	92.98
c20 UI33 probes	0.013	215	270	79.62
human mitoDB	0.029	436	594	73.4
c33 UI33 probes	0.021	245	362	67.67
MAP00252	0.022	23	35	65.71
c34 UI33	0.012	274	452	60.62
c21 UI33	0.026	166	287	57.84
c8 UI33	0.013	164	288	56.94

this approach, the expression values of a group of  $Q$  genes for a particular sample under a particular condition are considered as a  $Q$ -dimension vector. The differential analysis is done at the vector level, without individual gene differential statistic.

In this section, we first introduce the concept of fuzzy membership function of vectors, then the details of MCM-test.

**Fuzzy membership Function of Vectors**

In fuzzy set theory, the degree for one variable to belong to a fuzzy set is defined by a function. For a vector which has two dimensions, the degree that it belongs to a set of vectors can be defined by a three-dimensional function, with the third dimension being the measurement of the membership. Figure 5 shows a sample fuzzy membership function for a vector  $(x, y)$ .

For vectors with  $n$  dimensions, their fuzzy membership function will be  $n+1$ -dimensional, with one dimension measuring the fuzzy membership.

**Our approach**

Consider a pathway that consists of  $Q$  genes, the problem now is to determine how these  $Q$  genes are expressed differently under two conditions. To answer this question, we consider the expression values of the  $Q$  genes for a particular sample under a particular condition as a  $Q$ -dimension vector. Then the expression values of a pathway under one condition  $j$  can be modeled as set  $S_j$  ( $j = 1, 2$ ) of points in a  $Q$ -dimension space. The idea is to consider the two sets of points  $S_1$  and  $S_2$  as samples from two different fuzzy sets. We then examine the membership value of each element with respect to these two fuzzy sets and determine the  $d$ -value between the two sets of samples.

The mean  $\bar{\mu}_j$  of the expression values of set  $S_j$  is:

$$\bar{m}_j = \frac{1}{N_j} \sum_{\bar{x}_n \in S_j} \bar{x}_n \tag{1}$$

where,

$$\bar{m}_j = \begin{bmatrix} m_{j1} \\ m_{j2} \\ \dots \\ m_{jQ} \end{bmatrix} \text{ and } \bar{x}_n = \begin{bmatrix} x_{n1} \\ x_{n2} \\ \dots \\ x_{nQ} \end{bmatrix}$$

$N_j$  is the number of samples in  $S_j$ ,  $\bar{x}_n$  is vector made by the expression values of the  $n$ -th sample under condition  $j$ .

We then characterize set  $S_j$  ( $j = 1, 2$ ) by a fuzzy set  $FS_j$  ( $j = 1, 2$ ) whose fuzzy membership function is defined as:

$$f_{FS_j}(\bar{x}) = \exp\left(-\frac{1}{2}(\bar{x}_n - \bar{m}_j)^T \Sigma_j^{-1}(\bar{x}_n - \bar{m}_j)\right) \tag{2}$$

where,

$$\Sigma_j = \frac{1}{N_j - 1} \sum_{\bar{x}_n \in S_j} (\bar{x}_n - \bar{\mu}_j)(\bar{x}_n - \bar{\mu}_j)^T \tag{3}$$

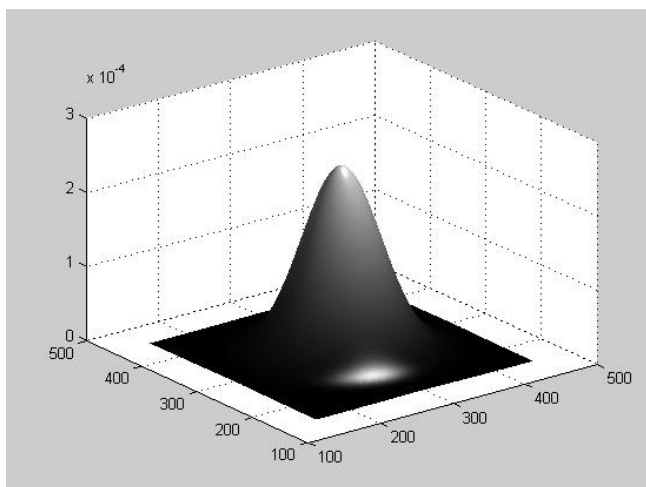
Given an element  $\bar{e}$  in  $S_1$ , we calculate its element misclassification degree with respect to  $FS_2$  as

$$m(\bar{e}, FS_2) = \text{Max}(f_{FS_2}(\bar{e}) - f_{FS_1}(\bar{e}), 0) \tag{4}$$

We denote the misclassified elements in  $S_1$  with respect to  $FS_2$  as  $M_{FS_2}(S_1) = \{\bar{e} \mid \bar{e} \in S_1, m(\bar{e}, FS_2) > 0\}$ . Similarly, we denote the misclassified elements in  $S_2$  with respect to  $FS_1$  as  $M_{FS_1}(S_2) = \{\bar{f} \mid \bar{f} \in S_2, m(\bar{f}, FS_1) > 0\}$ . We denote the number of misclassified elements in  $S_1$  and  $S_2$  with respect to each other as  $\# M(S_1, S_2) = |M_{FS_2}(S_1)| + |M_{FS_1}(S_2)|$ . We then define the convergence degree ( $c$ -value) of  $S_1$  and  $S_2$  as a linear interpolation of the number of misclassified elements and the mutual misclassification degrees as follows.

$$c(S_1, S_2) = \beta * T_1 + (1 - \beta) * T_2 \tag{5}$$





**Figure 5**  
A sample fuzzy membership function of vector (x, y).

where,

$$T_1 = \frac{\# M(S_1, S_2)}{S_1 + S_2} \tag{6}$$

and

$$T_2 = \frac{\sum_{\bar{e} \in S_1} m(\bar{e}, S_2) + \sum_{\bar{f} \in S_2} m(\bar{f}, S_1)}{S_1 + S_2} \tag{7}$$

Then, the divergence between  $S_1$  and  $S_2$  can be calculated using the following:

$$d(S_1, S_2) = 1 - c(S_1, S_2) \tag{8}$$

In our method, to negate the effect of outliers, we used  $\alpha$ -trimmed mean instead of normal mean. The  $\alpha$ -trimmed mean is calculated by ordering the sample under consideration and taking away the smallest and largest  $\alpha$  values from the ordered sample. The mean of the remaining values in the sample is  $\alpha$ -trimmed mean of the sample. For instance, if we have a sample of (3, 17, 25, 29, 23, 53, 22, 31, 45, 81, 90, 1), the 2-trimmed mean is calculated by removing the smallest two values (1, 3), and largest two values (81, 90) from the sample set. The mean of the remaining values (30.625) becomes the 2-trimmed mean of the sample.

For computational simplicity, an *Epanechnikov* function shown as following can be used instead of the Gaussian function of equation (2):

$$f_{FS_j}(\bar{x}_n) = \text{Max}\left\{0, 1 - \frac{||\bar{x}_n - \bar{m}_j||^2}{\sum_{q=1}^Q s_q^2}\right\} \tag{9}$$

where,

$$s_q = \sqrt{\frac{1}{N_j - 1} \sum_{\bar{x}_n \in S_j} (\bar{x}_n - \bar{m}_j)^2} \tag{10}$$

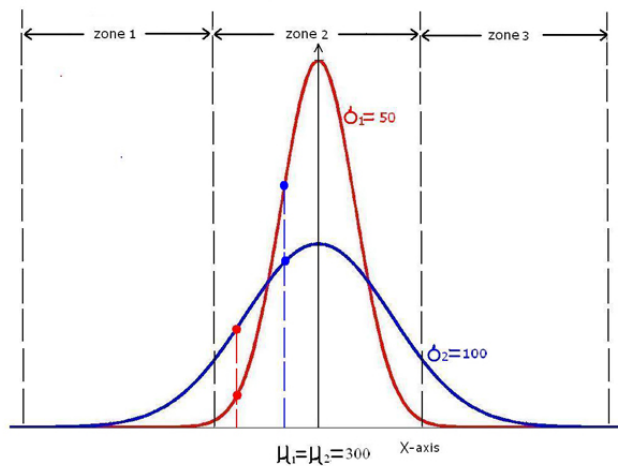
**Analysis of method**

**One dimension: a special case**

In this section we analyze MCM-test for its theoretical justification. For the sake of clarity, we start with one dimension, the simplest and special case of multi-dimension. The one dimensional MCM-test corresponds to differential analysis of a single gene.

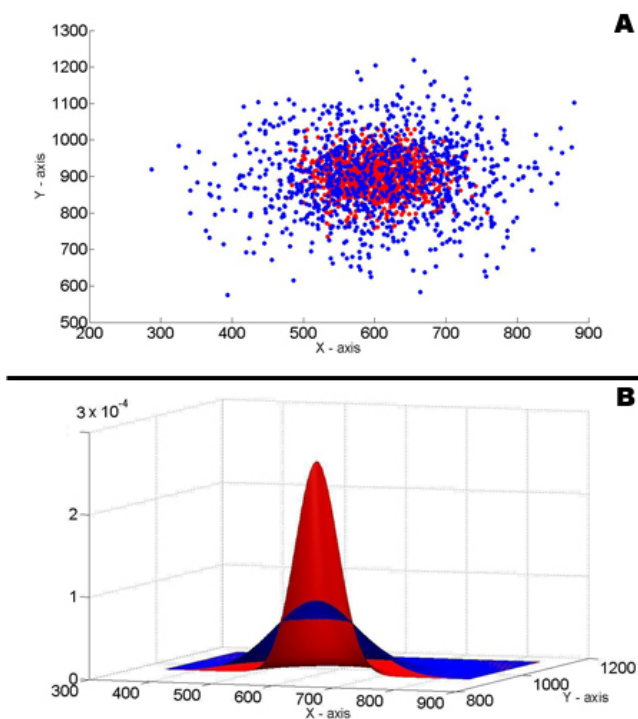
In Figure 6, two distributions,  $D_1$  and  $D_2$  are displayed in blue and red respectively, with mean  $\mu_1 = 600$  and standard deviation  $\sigma_1 = 50$  for  $D_1$  and  $\mu_2 = 700$ ,  $\sigma_2 = 100$  for  $D_2$ . The visualization tells us that they are different as they cover different areas and have different shapes. The CM-test, which can be considered as a special case of the MCM-test differentiates them by measuring the differences on the Y axis, which is a combined result of the location difference together with the difference of the variances.

MCM-test uses the probability distribution functions of these two distributions as their fuzzy membership functions respectively, and takes advantage of the membership differences of "misclassified" samples. As shown in Figure



**Figure 6**  
One dimension Gaussian distributions,  $\mu_1 = 600$ ,  $\mu_2 = 700$ ,  $\sigma_1 = 50$ ,  $\sigma_2 = 100$ .





**Figure 7**  
 (a) 1000 samples of 2-D Gaussian distribution  $\mu_1x = 600$ ,  $\mu_1y = 900$ ,  $\sigma_1x = \sigma_1y = 50$  and 1000 samples of 2-D Gaussian distribution  $\mu_2x = 700$ ,  $\mu_2y = 1000$ ,  $\sigma_2x = \sigma_2y = 100$ . (b) Probability density functions of the two distributions.

6, a sample  $x_1$  of  $D_2$  has a higher degree of belonging to  $D_1$ , thus is "misclassified" in MCM-test. This misclassification degree is aggregated with all the other "misclassified" samples of  $D_2$  that are misclassified. Similarly,  $x_2$  of  $D_1$  has a higher degree for  $D_2$ , thus is also misclassified. This misclassification degree is also aggregated with all the other misclassified samples of  $D_1$ .

MCM-test collects all the misclassified degrees and the number of misclassified samples and form them into an index to measure the divergence of these two distributes. With the mean difference between these two distributions increases, the number of misclassified samples, as well as the aggregated misclassification degree decreases. Thus the MCM  $d$ -value will decrease.

### Two and higher dimensions

Figure 7(a) illustrates samples of two distributions, each of which is a 2-D Gaussian function. In pathway analysis, the X and Y axis can be the expression data of two individual genes respectively. Figure 7(b) shows the probability density functions of these two distributions, which can be used as their fuzzy membership functions after multiplying a constant.

Distributions of higher dimensions are hard to visualize. But the idea of the misclassification test stays the same. In multi-dimension space, each sample is a vector. And their misclassification degrees are used to measure the divergence of their distributions.

### List of abbreviations

MCM-test: multi-dimensional Cluster Misclassification test

CM-test: cluster misclassification test

FM-test: fuzzy membership test

GSEA: gene set enrichment analysis

### Competing interests

The authors declare that they have no competing interests.

### Authors' contributions

LRL designed the algorithm, coordinated the project and wrote part of the manuscript. VM implemented the algorithm, conducted experiments and wrote part of the manuscript. YL designed experiments and wrote part of the manuscript. DK located gene expression and pathway data for experiments, analyzed the results and wrote part of the manuscript.

### Acknowledgements

We would like to thank Ying Wang and Togba Liberty for generating some of the figures used in this paper. This work was supported by the Agriculture Experiment Station at the University of the District of Columbia (Project No.: DC-0LIANG; Accession No.: 0203877).

This article has been published as part of *BMC Bioinformatics* Volume 9 Supplement 6, 2008: Symposium of Computations in Bioinformatics and Bioscience (SCBB07). The full contents of the supplement are available online at <http://www.biomedcentral.com/1471-2105/9?issue=S6>.

### References

1. **Kyoto Encyclopedia of Genes and Genomes** [<http://www.genome.jp/kegg/>]
2. **Biocarta** [<http://www.biocarta.com/>]
3. Mootha VK, Lindgren CM, Eriksson KF, Subramanian A, Sihag S, Lehar J, Puigserver P, Carlsson E, Ridderstrale M, Laurila E, Houstis N, Daly MJ, Patterson N, Mesirov JP, Golub TR, Tamayo P, Spiegelman B, Lander ES, Hirschhorn JN, Altshuler D, Groop LC: **Pgc-1alpha-responsive genes involved in oxidative phosphorylation are coordinately downregulated in human diabetes.** *Nature Genet* 2003, **34**:267-273.
4. Subramanian AS, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR, Lander ES, Mesirov JP: **Gene set enrichment analysis: A knowledge-based approach for interpreting genomewide expression profiles.** *PNAS* 2005, **102**:15545-15550.
5. Tomfohr J, Lu J, Kepler TB: **Pathway level analysis of gene expression using singular value decomposition.** *BMC Bioinformatics* 2005, **6**:225.
6. Tian L, Greenberg SA, Kong SW, Altshuler J, Kohane IS, Park PJ: **Discovering statistically significant pathways in expression profiling studies.** *PNAS* 2005, **102**:13544-13549.
7. **American Diabetes Association Website** [<http://www.diabetes.org/>]

8. Morino K, Petersen KF, Shulman GI: **Molecular mechanisms of insulin resistance in humans and their potential links with mitochondrial dysfunction.** *Diabetes* 2006, **55(Suppl 2):S9-S15.**
9. Takamura T, Honda M, Sakai Y, Ando H, Shimizu A, Ota T, Sakurai M, Misu H, Kurita S, Matsuzawa-Nagata N, Uchikata M, Nakamura S, Matoba R, Tanino M, Matsubara K, Kaneko S: **Gene expression profiles in peripheral blood mononuclear cells reflect the pathophysiology of type 2 diabetes.** *Biochem Biophys Res Commun* 2007, **361(2):379-84.**
10. Skov V, Glinborg D, Knudsen S, Jensen T, Kruse TA, Tan Q, Bruggaard K, Beck-Nielsen H, Højlund K: **Reduced expression of nuclear-encoded genes involved in mitochondrial oxidative metabolism in skeletal muscle of insulin-resistant women with polycystic ovary syndrome.** *Diabetes* 2007, **56(9):2349-55.**
11. Park SY, Lee W: **The depletion of cellular mitochondrial DNA causes insulin resistance through the alteration of insulin receptor substrate-1 in rat myocytes.** *Diabetes Res Clin Pract* 2007, **77(Suppl 1):S165-S171.** 17462778
12. Liang LR, Lu S, Lu Y, Dhawan P, Kumar D: **CM-test: An innovative divergence measurement and its application in diabetes gene expression data analysis.** *Proceedings of the IEEE International Conference on Granular Computing: 262-268 May 2006; Atlanta, Georgia, USA .*
13. Liang LR, Lu S, Wang X, Lu Y, Mandal V, Patacsil D, Kumar D: **FM-test: A fuzzy-set-theory-based approach to differential gene expression data analysis.** *BMC Bioinformatics* 2006, **7(Suppl 4):S7.**
14. Lu Y, Lu S, Liang LR, Kumar D: **FM-test: A fuzzy-set-theory based approach for discovering diabetes genes.** *Proceedings of the IEEE International Symposium of Computations in Bioinformatics and Bioscience: 48-55 June 2006; Hangzhou, Zhejiang, P.R. China .*

Publish with **BioMed Central** and every scientist can read your work free of charge

"BioMed Central will be the most significant development for disseminating the results of biomedical research in our lifetime."

Sir Paul Nurse, Cancer Research UK

Your research papers will be:

- available free of charge to the entire biomedical community
- peer reviewed and published immediately upon acceptance
- cited in PubMed and archived on PubMed Central
- yours — you keep the copyright

Submit your manuscript here:  
[http://www.biomedcentral.com/info/publishing\\_adv.asp](http://www.biomedcentral.com/info/publishing_adv.asp)

