

## Gene expression

The `tspair` package for finding top scoring pair classifiers in R

Jeffrey T. Leek

Department of Oncology, Johns Hopkins School of Medicine, Baltimore, MD 21287, USA

Received on January 11, 2009; revised on February 16, 2009; accepted on March 1, 2009

Advance Access publication March 10, 2009

Associate Editor: Joaquin Dopazo

## ABSTRACT

**Summary:** Top scoring pairs (TSPs) are pairs of genes whose relative rankings can be used to accurately classify individuals into one of two classes. TSPs have two main advantages over many standard classifiers used in gene expression studies: (i) a TSP is based on only two genes, which leads to easily interpretable and inexpensive diagnostic tests and (ii) TSP classifiers are based on gene rankings, so they are more robust to variation in technical factors or normalization than classifiers based on expression levels of individual genes. Here I describe the R package, `tspair`, which can be used to quickly identify and assess TSP classifiers for gene expression data.

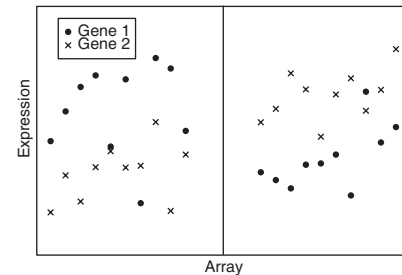
**Availability:** The R package `tspair` is freely available from Bioconductor: <http://www.bioconductor.org>

**Contact:** [jtleek@jhu.edu](mailto:jtleek@jhu.edu)

## 1 INTRODUCTION

Classification of patients into disease groups or subtypes is the most direct way to translate microarray technology into a clinically useful tool (Quackenbush, 2006). A small number of tests based on microarrays have even been approved for clinical use, for example, for diagnosing breast cancer subtypes (Ma *et al.*, 2004; Marchionni *et al.*, 2008; Paik *et al.*, 2004; van't Veer *et al.*, 2002). But standard microarray classifiers are based on complicated functions of many gene expression measurements. This type of classifier is both hard to interpret and depends critically on the platform, pre-processing and normalization steps to be effective (Quackenbush, 2006). Identifying biologically interpretable, robust and cheap classifiers based on small subsets of genes would greatly speed progress in the development of clinical tests from microarray experiments.

Top scoring pairs (TSPs) are pairs of genes that accurately classify patients into clinically relevant groups based on their ranks (Geman *et al.*, 2004; Tan *et al.*, 2005; Xu *et al.*, 2005). The basic idea is to search among all pairs of genes, and look for genes whose ranking most consistently switches between two groups. To understand how the classification scheme works, consider the simulated gene expression data in Figure 1. In this figure there are two groups of arrays, separated by the black line. These groups could represent healthy patients versus cancer patients, or two distinct subtypes of cancer. For all but one array in Group 1, Gene 1 has higher expression than Gene 2, and the reverse is true in Group 2. In this case, Genes 1 and 2 form a classifier based on their relative levels of expression. A new sample where the gene expression for Gene 1 was higher than the gene expression for Gene 2 would be classified as Group 1.



**Fig. 1.** An Example of a TSP. In this simulated example, the expression for Gene 1 is higher than the expression for Gene 2 for almost all of the arrays in the group on the left and this relationship reverses for the group on the right.

The TSP approach has been successfully applied to identify subtypes of sarcoma, resulting in a RT-PCR-based test that correctly classified 20 independent tumors with perfect accuracy (Price *et al.*, 2007). This early success suggests that it may be possible to identify TSP classifiers for other important diseases and quickly develop new inexpensive diagnostic tests.

## 2 THE TSPAIR PACKAGE

Calculating the TSP for a gene expression dataset is relatively straightforward, but computationally intensive. I have developed an R package `tspair` that can rapidly calculate the TSP for typical gene expression datasets, with tens of thousands of genes. The TSP can be calculated both in R or with an external C function, which allows both for rapid calculation and flexible development of the `tspair` package. The `tspair` package includes functions for calculating the statistical significance of a TSP by permutation test, and is fully compatible with Bioconductor expression sets. The R package is freely available from the Bioconductor web site ([www.bioconductor.org](http://www.bioconductor.org)).

## 3 AN EXAMPLE SESSION

Here I present an example session on a simple simulated dataset included in the `tspair` package. I calculate the TSP, assess the strength of evidence for the classifier with a permutation test, plot the output and show how to predict outcomes for a new dataset. The main function in the `tspair` package is `tspcalc()`. This function accepts either (i) a gene expression matrix or an expression set and a group indicator vector, or (ii) an expression set object and a column number, indicating which column of the annotation data to use as the group indicator. The result is a `tsp` object which gives

