

PCOT2: Principal Coordinates and Hotelling's T^2 for the analysis of microarray data

Sarah Song and Mik Black

October 3, 2006

1 Overview

`pcot2` is an R-package for the analysis of groups of genes in microarray experiments. It utilizes inter-gene correlation information to detect significant alterations in the activities of gene sets. Incorporating additional (usually functional) information into the data analysis process allows gene interactions to be investigated in a statistical framework. One of the reasons that gene set analysis is becoming important is that it is suitable for detecting small coordinated changes in expression of groups of genes which are functionally related, which may not be considered significant in a single gene analysis. This vignette gives a tutorial-style introduction to the functions in the `pcot2` package. These functions are used for testing and visualizing changes in expression activity for groups of genes.

2 Example: ALL/AML data

In this example the ALL/AML leukemia data set of Golub *et al.*(1999) is used to illustrate the functionality of the `pcot2` package. This data set contains 38 bone marrow samples obtained from adult leukemia patients, 11 relating to acute myeloid leukemia (AML, class 1) and 27 relating to acute lymphoblastic leukemia (ALL, class 0). Gene expression levels were measured using Affymetrix high density oligonucleotide arrays containing 6817 human genes, of which 3051 genes were considered suitable for analysis by Golub *et al.*(1999) after pre-processing. This data set is available as part of the `multtest` package and gene sets are defined as KEGG pathways using the `hu6800` annotation package. Both packages can be downloaded from www.bioconductor.org.

```
> library(pcot2)
> library(multtest)
> library(hu6800)
> set.seed(1234567)
```

3 The `pcot2` function

The `pcot2` function implements the PCOT2 testing method, which is a two-stage permutation-based approach for testing changes in activity in pre-specified


```
> results$res.sig
```

```
[1] Num          T2          P.nor          P.adj          P.permu        P.permu.adj  
<0 rows> (or 0-length row.names)
```

```
> results$res.all
```

	Num	T2	P.nor	P.adj	P.permu	P.permu.adj
KEGG01032	11	16.083606	1.559644e-03	1.229864e-02	0.1	0.5632532
KEGG04010	107	37.799126	3.502656e-06	4.852223e-05	0.1	0.5632532
KEGG04060	86	50.940814	1.990119e-07	5.666977e-06	0.1	0.5632532
KEGG04350	26	23.712919	1.425931e-04	1.379011e-03	0.1	0.5632532
KEGG04520	37	23.992055	1.314175e-04	1.320773e-03	0.1	0.5632532
KEGG00190	43	14.212036	2.959080e-03	2.166725e-02	0.1	0.5632532
KEGG00193	18	47.397496	4.122077e-07	9.186144e-06	0.1	0.5632532
KEGG04810	89	54.942820	9.054124e-08	3.569835e-06	0.1	0.5632532
KEGG04514	72	25.591713	8.291894e-05	8.673666e-04	0.1	0.5632532
KEGG04670	64	102.797497	5.541345e-11	1.420137e-08	0.1	0.5632532
KEGG00564	13	52.401665	1.486743e-07	4.762786e-06	0.1	0.5632532
KEGG00590	19	47.576677	3.970128e-07	9.186144e-06	0.1	0.5632532
KEGG04370	35	18.656730	6.707025e-04	5.729592e-03	0.1	0.5632532
KEGG04664	38	61.379812	2.735820e-08	1.558081e-06	0.1	0.5632532
KEGG04730	33	47.690379	3.876787e-07	9.186144e-06	0.1	0.5632532
KEGG04912	38	15.919109	1.648417e-03	1.261065e-02	0.1	0.5632532
KEGG04510	88	45.140440	6.662542e-07	1.219627e-05	0.1	0.5632532
KEGG00280	22	45.846527	5.725182e-07	1.086852e-05	0.1	0.5632532
KEGG00240	32	58.978644	4.234863e-08	2.170623e-06	0.1	0.5632532
KEGG04360	32	39.327865	2.446638e-06	3.800151e-05	0.1	0.5632532
KEGG03022	13	23.820339	1.381779e-04	1.362010e-03	0.1	0.5632532
KEGG00220	12	7.580525	3.529195e-02	2.032501e-01	0.1	0.5632532
KEGG00260	12	9.014209	2.002974e-02	1.222197e-01	0.1	0.5632532
KEGG00330	22	67.699684	9.103087e-09	5.832352e-07	0.1	0.5632532
KEGG03320	20	53.201493	1.269935e-07	4.339457e-06	0.1	0.5632532
KEGG04310	42	37.040273	4.197126e-06	5.661264e-05	0.1	0.5632532
KEGG04330	15	14.413820	2.758517e-03	2.049140e-02	0.1	0.5632532
KEGG05120	39	72.098995	4.399616e-09	3.758448e-07	0.1	0.5632532
KEGG00230	54	18.234794	7.681197e-04	6.454225e-03	0.1	0.5632532
KEGG04612	56	41.303014	1.555451e-06	2.559465e-05	0.1	0.5632532
KEGG00561	17	58.313748	4.788850e-08	2.231432e-06	0.1	0.5632532
KEGG04512	33	49.878215	2.467977e-07	6.657828e-06	0.1	0.5632532
KEGG04340	11	6.073128	6.534459e-02	3.640549e-01	0.1	0.5632532
KEGG04640	70	115.400542	1.210776e-11	6.205958e-09	0.1	0.5632532
KEGG04650	65	46.542476	4.936711e-07	9.732163e-06	0.1	0.5632532
KEGG04662	39	46.757474	4.717016e-07	9.732163e-06	0.1	0.5632532
KEGG04610	15	73.363867	3.589230e-09	3.758448e-07	0.1	0.5632532
KEGG04070	31	25.776064	7.869364e-05	8.403176e-04	0.1	0.5632532
KEGG00980	13	69.188212	7.093654e-09	5.194180e-07	0.1	0.5632532
KEGG00350	15	4.806800	1.115512e-01	6.082631e-01	0.1	0.5632532
KEGG04660	43	33.338131	1.042993e-05	1.243249e-04	0.1	0.5632532
KEGG00380	22	74.739625	2.883968e-09	3.758448e-07	0.1	0.5632532
KEGG04110	49	51.495005	1.780683e-07	5.368869e-06	0.1	0.5632532

KEGG01510	25	9.919144	1.413848e-02	9.058529e-02	0.1	0.5632532
KEGG00360	11	38.952479	2.670169e-06	4.025361e-05	0.1	0.5632532
KEGG04940	34	7.779280	3.259065e-02	1.920078e-01	0.1	0.5632532
KEGG04020	62	42.299762	1.243043e-06	2.197015e-05	0.1	0.5632532
KEGG04540	43	10.907225	9.740883e-03	6.401014e-02	0.1	0.5632532
KEGG03050	16	41.184124	1.597917e-06	2.559465e-05	0.1	0.5632532
KEGG00120	10	8.026618	2.953082e-02	1.760038e-01	0.1	0.5632532
KEGG04080	68	35.854072	5.589777e-06	7.162746e-05	0.1	0.5632532
KEGG05040	21	13.809328	3.406880e-03	2.459481e-02	0.1	0.5632532
KEGG04210	44	27.299352	5.138148e-05	5.603428e-04	0.1	0.5632532
KEGG04620	48	41.219369	1.585201e-06	2.559465e-05	0.1	0.5632532
KEGG04920	30	57.385632	5.693675e-08	2.431960e-06	0.1	0.5632532
KEGG05110	15	13.112465	4.359517e-03	3.060980e-02	0.1	0.5632532
KEGG00500	17	17.798230	8.848180e-04	7.314881e-03	0.1	0.5632532
KEGG00010	36	8.520856	2.429027e-02	1.464733e-01	0.1	0.5632532
KEGG00030	15	13.506746	3.790243e-03	2.698234e-02	0.1	0.5632532
KEGG00051	18	18.678736	6.659944e-04	5.729592e-03	0.1	0.5632532
KEGG00710	12	6.022369	6.673974e-02	3.678295e-01	0.1	0.5632532
KEGG04910	56	27.699624	4.601491e-05	5.241205e-04	0.1	0.5632532
KEGG04630	56	38.326364	3.092383e-06	4.402869e-05	0.1	0.5632532
KEGG00860	12	38.451145	3.002926e-06	4.397660e-05	0.1	0.5632532
KEGG05210	38	27.327734	5.097998e-05	5.603428e-04	0.1	0.5632532
KEGG00071	18	36.921968	4.317896e-06	5.674826e-05	0.1	0.5632532
KEGG04720	35	16.143909	1.528379e-03	1.224041e-02	0.1	0.5632532
KEGG04930	18	17.466958	9.858206e-04	8.020517e-03	0.1	0.5632532
KEGG00310	15	33.903265	9.048742e-06	1.131226e-04	0.1	0.5632532
KEGG04150	18	11.009560	9.376387e-03	6.241513e-02	0.1	0.5632532
KEGG04742	10	9.165107	1.889037e-02	1.166561e-01	0.1	0.5632532
KEGG04530	41	33.404722	1.025619e-05	1.243249e-04	0.1	0.5632532
KEGG05130	27	9.739467	1.514257e-02	9.465221e-02	0.1	0.5632532
KEGG05131	27	9.739467	1.514257e-02	9.465221e-02	0.1	0.5632532
KEGG00150	10	11.673533	7.335679e-03	4.947340e-02	0.1	0.5632532
KEGG05050	11	7.680916	3.389911e-02	1.974471e-01	0.1	0.5632532
KEGG00252	15	20.668382	3.563113e-04	3.261269e-03	0.1	0.5632532
KEGG00760	11	53.597908	1.175117e-07	4.302274e-06	0.1	0.5632532
KEGG00562	14	19.021299	5.970409e-04	5.276198e-03	0.1	0.5632532
KEGG00480	11	72.739759	3.967321e-09	3.758448e-07	0.1	0.5632532
KEGG04740	10	14.887889	2.341753e-03	1.765132e-02	0.1	0.5632532
KEGG00052	15	19.849740	4.596460e-04	4.133269e-03	0.1	0.5632532
KEGG00650	16	15.980992	1.614410e-03	1.253762e-02	0.1	0.5632532
KEGG00410	13	46.661226	4.814060e-07	9.732163e-06	0.1	0.5632532
KEGG00340	10	29.388643	2.910738e-05	3.390748e-04	0.1	0.5632532
KEGG00620	16	21.669123	2.622921e-04	2.444374e-03	0.1	0.5632532
KEGG00640	17	48.892360	3.020600e-07	7.741198e-06	0.1	0.5632532
KEGG00510	14	10.796492	1.015222e-02	6.586868e-02	0.1	0.5632532
KEGG01030	18	12.723491	5.010394e-03	3.470445e-02	0.1	0.5632532
KEGG00970	16	23.403392	1.561698e-04	1.482342e-03	0.1	0.5632532
KEGG05030	13	25.111697	9.508626e-05	9.747490e-04	0.1	0.5632532
KEGG00020	12	12.207513	6.036512e-03	4.125436e-02	0.2	1.0000000
KEGG05010	13	3.332612	2.123831e-01	1.000000e+00	0.2	1.0000000

KEGG04120	12	6.182388	6.244511e-02	3.517240e-01	0.2	1.0000000
KEGG00251	13	6.428522	5.640018e-02	3.212055e-01	0.2	1.0000000
KEGG01031	10	1.390378	5.152243e-01	1.000000e+00	0.4	1.0000000
KEGG00530	11	1.535339	4.814905e-01	1.000000e+00	0.4	1.0000000
KEGG01430	35	2.677797	2.849138e-01	1.000000e+00	0.5	1.0000000
KEGG04320	10	1.201290	5.630263e-01	1.000000e+00	0.7	1.0000000

In the `pcot2` function, the T^2 statistic can be calculated in two ways, using either a pooled estimate of correlation for the two classes (default) or an un-pooled estimate. And users can set `var.equal=F` if the correlation structure is assumed to differ across the two classes.

In the first step of the PCOT2 analysis, the dimensionality of the gene expression data is reduced via principle coordinates. The default dimensionality in the `pcot2` function is set as `ncomp=2`. In the second step of the PCOT2 analysis, the distances between the transformed groups are calculated via euclidean distances by default. Other distances (e.g., correlation or Spearman distances) can also be used by defining `dist.method` in the function. A permutation p -value for each category is calculated by re-arranging the sample labels. The permutations can also be performed by permuting rows (genes), using `permu='ByRow'`.

Table 1 lists computation times (in minutes) required to run 1000 permutations of the `pcot2` function on the AML/ALL data under various parameter configurations. The two machines used were a 3.2GHz Pentium 4 with 1Gb RAM running Microsoft Windows XP and R 2.1.0 (PC), and a 1.70GHz Pentium M with 256Mb of RAM running Fedora Core 3 and R 2.2.0 (Unix).

Table 1: *Computation times (minutes, 1000 permutations)*

Changes	PC machine	UNIX machine
default setting	5.6	6.8
var.equal=F	5.5	6.8
comp=8	6	7.6
dist.method="euclidean"	4.8	6
permu="ByRow"	5.6	6.8

4 The `corplot` and `corplot2` functions

The `corplot` and `corplot2` functions enable visualization of both correlation and gene expression information for a particular gene category, in particular the groups identified as being differentially expressed. The plot produced by the `corplot` function displays the pooled correlation calculated from the two classes, while the `corplot2` function produces a plot based on un-pooled correlation. Gene names can be added to the plot using `add.name=T` (default). The font size can be changed by setting the `font.size` argument. The `main` option specifies the title of the plot.

```
> sel <- c("04620", "04120")
> pvalue <- c(0.001, 0.72)
> library(KEGG)
> pname <- unlist(mget(sel, env = KEGGPATHID2NAME))
```

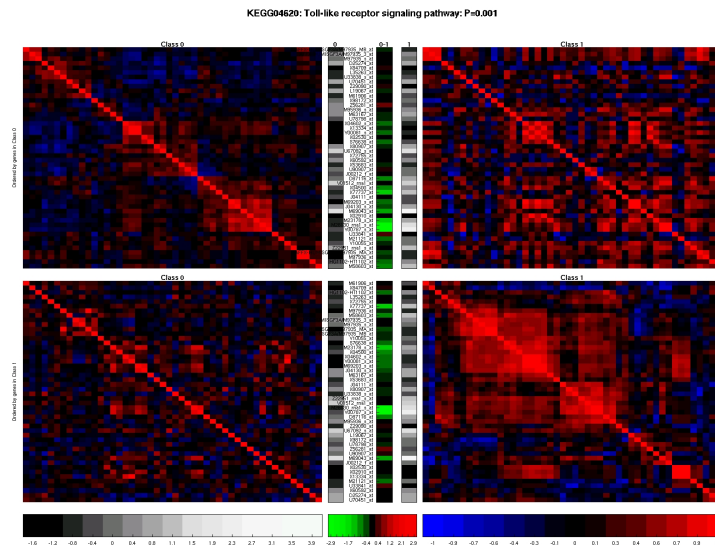


Figure 1: KEGG04620

```

> main <- paste("KEGG", sel, ": ", pname, ": ", "P=", pvalue, sep = "")
> for (i in 1:length(sel)) {
+   fname <- paste("corplot2-KEGG", sel[i], ".jpg", sep = "")
+   jpeg(fname, width = 1600, height = 1200, quality = 100)
+   selgene <- rownames(imat)[imat[, match(paste("KEGG", sel,
+     sep = "")[i], colnames(imat))] == 1]
+   corplot2(golub, selgene, golub.cl, main = main[i])
+   dev.off()
+ }

```

The argument *inputP* allows users to input the *p*-values of individual genes calculated using other approaches, such as the *limma* package (Smyth *et al.*, 2004), allowing the results from both per-gene and per-pathway analysis to be printed on a single plot. To allow users to identify genes from in correlation image plots, the argument *gene.locator=T* allows the selection of interesting (e.g., highly correlated and differentially expressed between two classes) genes by clicking beginning and end points on the main diagonal of the image plots. This prints the identifiers for the selected genes. Further details of this functionality are provided in the *HowToUseGeneLocator.pdf* document. The usage of *corplot2* is similar to that for the *corplot* function.

5 The aveProbes function

In Affymetrix gene expression data, a unique gene can often link to multiple probe sets, with such genes then having a greater influence on the pathway analysis (particularly if the gene is differentially expressed). In order to solve this problem, the *aveProbe* function is provided to change the multiple probe data to the unique gene data by taking the median of the probe values. This

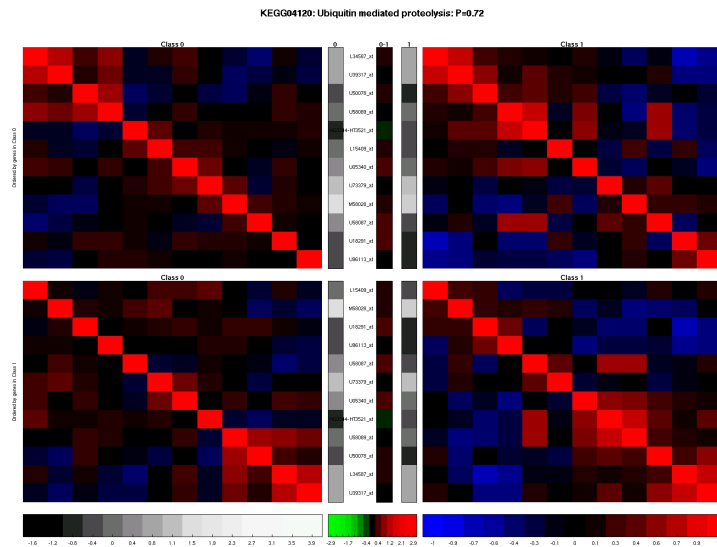


Figure 2: KEGG04120

function can be used to transform both expression data and the indicator matrix by providing a vector of unique gene identifiers.

```
> pathlist <- as.list(hu6800PATH)
> pathlist <- pathlist[match(rownames(golub), names(pathlist))]
> ids <- unlist(mget(names(pathlist), env = hu6800SYMBOL))
> newdata <- aveProbe(x = golub, ids = ids)$newx
> output <- aveProbe(x = golub, imat = imat, ids = ids)
> newdata <- output$newx
> newimat <- output$newimat
> newimat <- newimat[, apply(newimat, 2, sum) >= 10]
> dim(newdata)

[1] 2764 38

> dim(newimat)

[1] 2764 90
```

After the multiple probe data set has been changed to the unique gene symbol data, further analysis such as testing and visualizing pathways can be done on the new data set.

References

- [1] Benjamini, B.Y. and Yekutieli, D. (2001) The control of the false discovery rate in multiple testing under dependency. *The Annals of Statistics*, **29**, 1165-1188.

- [2] Gentleman,R.C., Carey,V.J., Bates,D.M., Bolstad,B., Dettling,M., Dudoit,S., Ellis,B., Gautier,L., Ge,Y., Gentry,J. *et al.* (2004) Bioconductor: open software development for computational biology and bioinformatics. *Genome Biology*, **5**, R80.
- [3] Golub,T.R., Slonim,D.K., Tamayo,P., Huard,C., Gaasenbeek,M., Mesirov,J.P., Coller,H., Loh,M.L., Downing,J.R., Caligiuri,M.A. *et al.* (1999) Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring, *Science*, **286**, 531-537.
- [4] Smyth,G.K. (2004) Linear models and empirical Bayes methods for assessing differential expression in microarray experiments. *Statistical Applications in Genetics and Molecular Biology*, **3**, No.1, Article 3.