

Moment based gene set enrichment testing – the npGSEA package

Jessica L. Larson^{1*} and Art Owen²

¹ Department of Bioinformatics and Computational Biology, Genentech, Inc.

² Department of Statistics, Stanford University

*larson.jessica (at) gene.com

May 2, 2014

Contents

1	Introduction	2
2	Example workflow for GSEA	2
2.1	Preparing our gene sets and our dataset for analysis	2
2.2	Running npGSEA	4
2.3	Running npGSEA with the beta and chi-sq approximations	5
2.4	Adding weights to the model	6
2.5	Adding covariates to model	7
2.6	Running npGSEA with multiple gene sets	8
3	Methods in brief	9
3.1	Disadvantages to a permutation approach	9
3.2	Test statistics	9
3.3	Moment based reference distributions	10
4	Session Info	11
5	References	11

1 Introduction

Gene set methods are critical to the analysis of gene expression data. The npGSEA package provides methods to run permutation-based gene set enrichment analyses without the typically computationally expensive permutation cost. These methods allow users to adjust for covariates and approximate corresponding permutation distributions. We are currently evaluating the applicability and accuracy of our method for RNA-seq expression data.

Our methods find the exact relevant moments of a weighted sum of (squared) test statistics under permutation, taking into account correlations among the test statistics. We find moment-based gene set enrichment p -values that closely approximate the permutation method p -values.

This vignette describes a typical analysis workflow and includes some information about the statistical theory behind npGSEA. For more technical details, please see Larson and Owen, 2014 .

2 Example workflow for GSEA

2.1 Preparing our gene sets and our dataset for analysis

For our example, we will use the ALL dataset. We begin by loading relevant libraries, subsetting the data, and running `featureFilter` on this data set. For details on these methods, please see the `limma` manual.

```
> library(ALL)
> library(hgu95av2.db)
> library(genefilter)
> library(limma)
> library(GSEABase)
> library(npGSEA)
> data(ALL)
> ALL <- ALL[, ALL$mol.biol %in% c('NEG','BCR/ABL') &
+   !is.na(ALL$sex)]
> ALL$mol.biol <- factor(ALL$mol.biol,
+   levels = c('NEG', 'BCR/ABL'))
> ALL <- featureFilter(ALL)
```

We adjust the feature names of the ALL dataset so that they match the names of our gene sets below. We convert them to entrez ids.

```
> featureNames(ALL) <- select(hgu95av2.db, featureNames(ALL),
+   "ENTREZID", "PROBEID")$ENTREZID
```

We now make four arbitrary gene sets by randomly selecting from the genes in our universe.

```
> xData <- exprs(ALL)
> geneEids <- rownames(xData)
```

```

> set.seed(12345)
> set1 <- GeneSet(geneIds=sample(geneEids,15, replace=FALSE),
+               setName="set1",
+               shortDescription="This is set1")
> set2 <- GeneSet(geneIds=sample(geneEids,50, replace=FALSE),
+               setName="set2",
+               shortDescription="This is set2")
> set3 <- GeneSet(geneIds=sample(geneEids,100, replace=FALSE),
+               setName="set3",
+               shortDescription="This is set3")
> set4 <- GeneSet(geneIds=sample(geneEids,500, replace=FALSE),
+               setName="set4",
+               shortDescription="This is set4")

```

As a positive control, we also make three gene sets that include our top differentially expressed genes.

```

> model <- model.matrix(~mol.biol, ALL)
> fit <- eBayes(lmFit(ALL, model))
> tt <- topTable(fit, coef=2, n=200)
> ttUp <- tt[which(tt$logFC >0), ]
> ttDown <- tt[which(tt$logFC <0), ]
> set5 <- GeneSet(geneIds=rownames(ttUp)[1:20],
+               setName="set5",
+               shortDescription="This is a true set of the top 20 DE
+               genes with a positive fold change")
> set6 <- GeneSet(geneIds=rownames(ttDown)[1:20],
+               setName="set6",
+               shortDescription="This is a true set of the top 20 DE genes
+               with a negative fold change")
> set7 <- GeneSet(geneIds=c(rownames(ttUp)[1:10], rownames(ttDown)[1:10]),
+               setName="set7",
+               shortDescription="This is a true set of the top 10 DE genes
+               with a positive and a negative fold change")

```

We then collapse all of our gene sets into a GeneSetCollection. For more information on GeneSets and GeneSetCollections, see the GSEABase manual.

```

> gsc <- GeneSetCollection( c(set1, set2, set3, set4, set5, set6, set7) )
> gsc

```

```

GeneSetCollection
  names: set1, set2, ..., set7 (7 total)
  unique identifiers: 6530, 8608, ..., 2625 (693 total)
  types in collection:
    geneIdType: NullIdentifier (1 total)
    collectionType: NullCollection (1 total)

```

2.2 Running npGSEA

Now that we have both our gene sets and experiment, we are ready to run npGSEA and determine the level of enrichment in our experiment. We can use npGSEA with our eset or expression data (xData) directly. We call npGSEASummary to get a summary of the results. T_Gw is explained in more detail in Section 3.2.

```
> yFactor <- ALL$mol.biol
> res1 <- npGSEA(x = ALL, y = yFactor, set = set1) ##with the eset
> res1

Normal Approximation for set1
T_Gw = 0.206
var(T_Gw) = 0.00427
pLeft = 0.999, pRight = 8e-04, pTwoSided = 0.0016

> res2_exprs <- npGSEA(xData, ALL$mol.biol, gsc[[2]]) ##with the expression data
> res2_exprs

Normal Approximation for set2
T_Gw = 0.0986
var(T_Gw) = 0.0168
pLeft = 0.776, pRight = 0.224, pTwoSided = 0.447
```

npGSEA has several built in accessor functions to gather more information about the analysis of your set of interest in your experiment.

```
> res3 <- npGSEA(ALL, yFactor, set3)
> res3

Normal Approximation for set3
T_Gw = -0.326
var(T_Gw) = 0.0993
pLeft = 0.15, pRight = 0.85, pTwoSided = 0.301

> geneSetName(res3)

"set3"

> stat(res3)

[1] -0.3260773

> sigmaSq(res3)

[1] 0.09929326

> zStat(res3)

[1] -1.03481

> pTwoSided(res3)
```

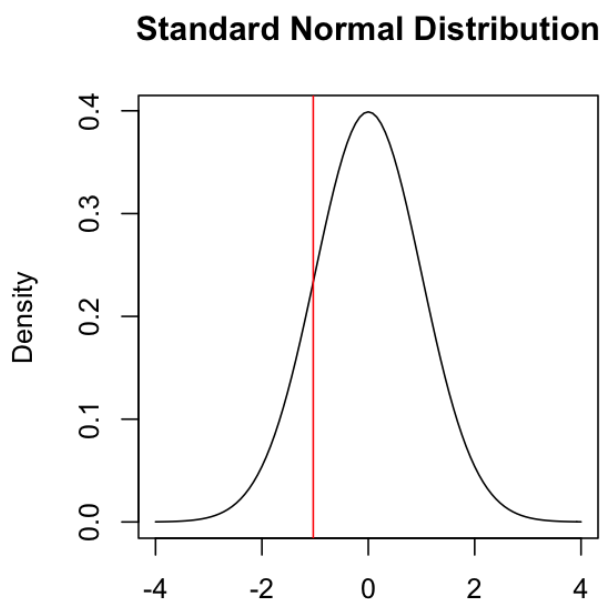


Figure 1: **Set3 normal approximation results** This plot displays the standard normal curve and our observed zStat for set3 in this analysis.

```
| [1] 0.3007577
> pLeft(res3)
| [1] 0.1503788
> dim(xSet(res3))
| [1] 100 109
```

There is also a npGSEA specific plot function (`npGSEAPlot`) to visualize the results of your analysis. Highlighted in red on the plot is the corresponding zStat of our analysis.

```
> npGSEAPlot(res3)
```

2.3 Running npGSEA with the beta and chi-sq approximations

There are three types of approximation methods in npGSEA: "norm", "beta", and "chiSq". Each method is discussed in brief in Section 3. The "norm" approximation method is the default.

```
> res5_norm <- npGSEA(ALL, yFactor, set5, approx= "norm")
> res5_norm
```

```

| Normal Approximation for set5
| T_Gw = 5.81
| var(T_Gw) = 0.533
| pLeft = 1, pRight = 8.72e-16, pTwoSided = 1.74e-15
> npGSEAPlot(res5_norm)

```

The beta approximation yields results quite similar to the normal approximation.

```

> res5_beta <- npGSEA(ALL, yFactor, set5, approx= "beta")
> res5_beta

| Beta Approximation for set5
| T_Gw = 5.81
| var(T_Gw) = 0.533
| pLeft = 1, pRight = 5.24e-29, pTwoSided = 1.05e-28
> npGSEAPlot(res5_beta)

```

The chi-sq approximation method is only available for the two-sided test. Here we call npGSEA and then show how the chiSqStat is related to C_Gw. C_Gw is explained in more detail in Section 3.2.

```

> res5_chiSq <- npGSEA(ALL, yFactor, set5, approx= "chiSq")
> res5_chiSq

| Chi-sq Approximation for set5
| C_Gw = 2.06
| df = 2.42, sigmaSq = 0.0232
| pTwoSided = 0
> chiSqStat(res5_chiSq)
| [1] 88.81123
> stat(res5_chiSq)
| [1] 2.05971
> stat(res5_chiSq)/sigmaSq(res5_chiSq)
| [1] 88.81123
> npGSEAPlot(res5_chiSq)

```

Note that, as we expected, set5 is a significantly enriched in all three methods. In each of the three corresponding plots, the observed statistic is a very rare event.

2.4 Adding weights to the model

Sometimes we do not want to weigh each gene in our set equally. We want to assign a larger weight to genes that are of a particular interest, and a lower weight to genes that we know may behave poorly. In this example, we up-weight the genes in set7 that had a positive fold change, and down-weight those

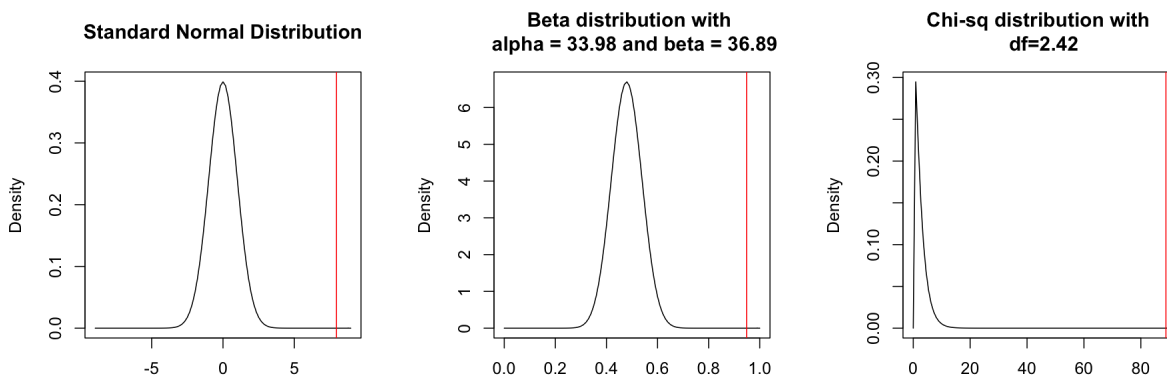


Figure 2: **Set5 normal, beta, and chi-sq approximation results** These plots displays the reference normal, beta, and chi-sq curves, and our observed zStat, betaStat, and chiSqStat for set5 in this analysis.

with a negative fold change. We assign the first 10 genes in our set a weight of 2, and the second 10 a weight of 0.5.

```
> w <- c( rep(2, 10), rep (0.5, 10) )
> res7_nowts <- npGSEA(x = ALL, y= yFactor, set = set7)
> res7_nowts

Normal Approximation for set7
T_Gw = 1.69
var(T_Gw) = 0.0772
pLeft = 1, pRight = 6.39e-10, pTwoSided = 1.28e-09

> res7_wts <- npGSEA(x = ALL, y = yFactor, set = set7, w = w)
> res7_wts

Normal Approximation for set7
T_Gw = 1.87
var(T_Gw) = 0.104
pLeft = 1, pRight = 3.28e-09, pTwoSided = 6.55e-09
```

By adding these weights, we get a slightly more significant result. We can add weights for the beta and chi-sq approximations, too. By default, npGSEA assigns a weight of 1 for all genes.

2.5 Adding covariates to model

Often we want to correct for confounders in our model. In this example, we correct for the age and sex of the subjects in our experiment. For more details on model selection and its relation to inference, please see the `limma` manual.

```

> res3_age <- npGSEA(x = ALL, y = yFactor, set = set3, z = ALL$age)
> res3_age

Normal Approximation for set3
T_Gw = -0.386
var(T_Gw) = 0.0806
pLeft = 0.0872, pRight = 0.913, pTwoSided = 0.174

> res3_agesex <- npGSEA(x = ALL, y = yFactor, set = set3, z = cbind(ALL$age, ALL$sex))
> res3_agesex

Normal Approximation for set3
T_Gw = -0.374
var(T_Gw) = 0.0791
pLeft = 0.0915, pRight = 0.908, pTwoSided = 0.183

```

By adjusting for these variables, we get a slight different result than above. Note that we can adjust for covariates in the beta and chi-sq approximation methods, too.

2.6 Running npGSEA with multiple gene sets

To explore multiple gene sets, we let `set` be a `GeneSetCollection`. This returns a list of `npGSEAResultNorm` objects, called a `npGSEAResultNormCollection`. We can access statistics for each `GeneSet` in our analysis through accessors of `npGSEAResultNormCollection`.

```

> resgsc_norm <- npGSEA(x = ALL, y = yFactor, set = gsc)
> unlist( pLeft(resgsc_norm) )

      set1      set2      set3      set4      set5      set6
9.992001e-01 7.764258e-01 1.503788e-01 9.999824e-01 1.000000e+00 4.247294e-11
      set7
1.000000e+00

> unlist( stat (resgsc_norm) )

      set1      set2      set3      set4      set5      set6      set7
0.20616765 0.09856167 -0.32607725 3.01145204 5.80773112 -3.85135976 1.68604828

> unlist( zStat (resgsc_norm) )

      set1      set2      set3      set4      set5      set6      set7
3.1559313 0.7601778 -1.0348100 4.1369937 7.9583579 -6.4915710 6.0700222

```

Note how quick our method is. We get results as accurate as permutation methods in a fraction of the time, even for multiple gene sets.

Using the `ReportingTools` package, we can publish these results to a HTML page for exploration. We first adjust for multiple testing.


```

> pvals <- p.adjust( unlist(pTwoSided(resgsc_norm)), method= "BH" )
> library(ReportingTools)
> npgseaReport <- HTMLReport (shortName = "npGSEA",
+                             title = "npGSEA Results", reportDirectory = "./reports")
> publish(gsc, npgseaReport, annotation.db = "org.Hs.eg",
+         setStats = unlist(zStat (resgsc_norm)), setPValues = pvals)
> finish(npgseaReport)

```

3 Methods in brief

3.1 Disadvantages to a permutation approach

There are three main disadvantages to permutation-based analyses: cost, randomness, and granularity.

Testing many sets of genes becomes computationally expensive for two reasons. First, there are many test statistics to calculate in each permuted version of the data. Second, to allow for multiplicity adjustment, we require small nominal p -values to draw inference about our sets, which in turn requires a large number of permutations.

Permutations are also subject random inference. Because permutations are based on a random shuffling of the data, there is a chance that we will obtain a different p -value for our set of interest each time we run our permutation analysis.

Permutations also have a granularity problem. If we do M permutations, then the smallest possible p -value we can attain is $1/(M + 1)$. When it is necessary to adjust for multiplicity, the permutation approach becomes very computationally expensive. Another aspect of the granularity problem is that permutations give us no basis to distinguish between two gene sets that both have the same p -value $1/(M + 1)$. There may be many such gene sets, and they have meaningfully different effect sizes.

Because of each of these limitations of permutation testing, there is a need to move beyond permutation-based GSEA methods. The methods we present in npGSEA and discuss in brief below are not as computationally expensive, random, or granular than their permutation counterparts. More details on our method can be found in Larson and Owen (2014).

3.2 Test statistics

We present our notation using the language of gene expression experiments.

Let g and h denote individual genes and G be a set of genes. Our experiment has n subjects. The subjects may represent patients, cell cultures, or tissue samples. The expression level for gene g in subject i is X_{gi} , and Y_i is the target variable on subject i . Y_i is often a treatment, disease, or genotype. We center the variables so that $\sum_{i=1}^n Y_i = \sum_{i=1}^n X_{gi} = 0, \forall g$.

Our measure of association for gene g on our treatment of interest is

$$\hat{\beta}_g = \frac{1}{n} \sum_{i=1}^n X_{gi} Y_i.$$

We consider the linear statistic

$$T_{G,w} = \sum_{g \in G} w_g \hat{\beta}_g$$

and the quadratic statistic

$$C_{G,w} = \sum_{g \in G} w_g \hat{\beta}_g^2,$$

where w_g corresponds to the weight given to gene g in set G .

3.3 Moment based reference distributions

To avoid the issues discussed above, we approximate the distribution of the permuted test statistics $T_{G,w}$ by Gaussian or by rescaled beta distributions. For the quadratic statistic $C_{G,w}$ we use a distribution of the form $\sigma^2 \chi_{(\nu)}^2$.

For the Gaussian treatment of $T_{G,w}$ we calculate $\sigma^2 = \text{Var}(T_{G,w})$ under permutation, and then report the p -value

$$p = \Pr(N(0, \sigma^2) \leq T_{G,w}).$$

The above is a left tail p -value. Two-sided and right tailed p -values are analogous.

When we want something sharper than the normal distribution, we can use a scaled Beta distribution, of the form $A + (B - A)\text{Beta}(\alpha, \beta)$. The $\text{Beta}(\alpha, \beta)$ distribution has a continuous density function on $0 < x < 1$ for $\alpha, \beta > 0$. We choose A , B , α and β by matching the upper and lower limits of $T_{G,w}$ under permutation, as well as its mean and variance. The observed left tailed p -value is

$$p = \Pr\left(\text{Beta}(\alpha, \beta) \leq \frac{T_{G,w} - A}{B - A}\right).$$

For the quadratic test statistic $C_{G,w}$ we use a $\sigma^2 \chi_{(\nu)}^2$ reference distribution reporting the p -value

$$\Pr(\sigma^2 \chi_{(\nu)}^2 \geq C_{G,w}),$$

after matching the first and second moments of $\sigma^2 \chi_{(\nu)}^2$ to $E(C_{G,w})$ and $E(C_{G,w}^2)$ under permutation, respectively.

Additional details on how σ^2 , A , B , α , β , $E(C_{G,w})$, $E(C_{G,w}^2)$, and ν are derived can be found in Larson and Owen (2014).

4 Session Info

- R version 3.1.0 (2014-04-10), x86_64-apple-darwin13.1.0
- Locale: C/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
- Base packages: base, datasets, grDevices, graphics, methods, parallel, stats, utils
- Other packages: ALL 1.6.0, AnnotationDbi 1.26.0, Biobase 2.24.0, BiocGenerics 0.10.0, DBI 0.2-7, GSEABase 1.26.0, GenomInfoDb 1.0.2, RSQLite 0.11.4, annotate 1.42.0, genefilter 1.46.0, graph 1.42.0, hgu95av2.db 2.14.0, limma 3.20.1, npGSEA 1.0.0, org.Hs.eg.db 2.14.0
- Loaded via a namespace (and not attached): BiocStyle 1.2.0, IRanges 1.22.5, XML 3.98-1.1, splines 3.1.0, stats4 3.1.0, survival 2.37-7, tools 3.1.0, xtable 1.7-3

5 References

Larson and Owen. (2014). Moment based gene set tests. Submitted.