

Kangar00: Kernel Approaches for Nonlinear Genetic Association Regression

Stefanie Friedrichs

2017-04-26

Introduction

The genetic information collected in genome-wide association studies (GWAS) is represented by the genotypes of various single-nucleotide polymorphisms (SNPs). Testing biological meaningful SNP Sets is a successful strategy for the evaluation of GWAS data, as it may increase power as well as interpretation of results. Via mapping of SNPs to genes forming a network, association between pathways and disease risk can be investigated.

Kernel methods are particularly well suited to cope with the challenges connected to the analysis of large SNP sets from GWAS data. They do not require to model a direct functional relationship between SNPs and effects, while at the same time can deal with high-dimensional data and allow for straightforward incorporation of covariates. The model for a logistic kernel machine regression of a pathway on a binary outcome is given by

$$\text{logit}(P(y_i = 1|x_i, z_i)) = x_i^t \beta + h(z_i)(1)$$

where y_i denotes the case or control status of individual i , x_i is the vector including informative covariates (such as age, sex, etc.) and z_i represents the genotypes of individual i . β is the regression coefficients for the parametric part of the model, while $h(\cdot)$ denotes an unknown function, non-parametrically incorporating the pathway's influence. The intercept is assumed to be included in x_i . For more details see Liu et al (2008).

Different kernels have been proposed that convert the genomic information of two individuals into a quantitative value reflecting their genetic similarity. This package includes the linear kernel as well as two more advanced kernels, adjusting for size bias in the number of SNPs and genes in a pathway or incorporating the network structure of genes within the pathway, respectively. The kernel functions are described in more detail in the instructions below.

A variance component test, constructed around the similarity matrix, can be used to evaluate a pathway's influence on disease risk. In **kangar00** p-values can be calculated with the Satterthwaite approximation or Davies method as described in Schaid (2010) and Davies (1980), respectively.

Data extraction and preparation

Pathways

The **kangar00** package offers several functions for data extraction from internet databases. In the following they will be explained using the Circadian rhythm pathway as an example.

- In the KEGG database (Kanehisa et al 2014) this pathway is identified with the id *hsa04710*.
- The function `pathway_info()` can use this id to create a table listing all genes included in Circadian rhythm. For each gene the startpoint, endpoint and the chromosome are listed.
- Gene membership is obtained directly from KEGG, while startpoints, endpoints and chromosome information is extracted from Ensembl (Cunningham et al 2015). The database is accessed via the

function `getBM()` in the `biomaRt` package. This means that the gene boundaries given will equal the current build used in Ensembl. An internet connection is required for this step.

```
pathway_info('hsa04710')
```

will return a `pathway_info` object containing a data frame of the form

| pathway | gene_start | gene_end | chr | gene |
|----------|------------|----------|-----|---------|
| hsa04710 | 13276652 | 13387266 | 11 | ARNTL |
| hsa04710 | 4979116 | 4985323 | 3 | BHLHE40 |
| hsa04710 | 26120026 | 26125127 | 12 | BHLHE41 |
| hsa04710 | ... | ... | ... | ... |

listing information on all genes KEGG assigned to Circadian rhythm.

Pathway object

In `kangaroo` all information on a specific pathway is combined in a `pathway` object. It includes

- The pathway's ID as used in KEGG.
- The adjacency **matrix**, which equals the network matrix without signs.
- A **vector** giving the signs for the interactions.

The following example creates a new pathway object, to which gene-interaction information has yet to be added

```
pathw <- pathway(id='hsa04710', adj=matrix(0), sign=as.vector(matrix(0)[matrix(0)!=0]))
```

Networkmatrix

The gene-gene interactions within pathways are represented by a network matrix. This quadratic matrix is of dimension equal to the number of genes in the corresponding pathway. It includes entries equal to 1 (representing an activation interaction), -1 (denoting an inhibiting interaction) or 0 (no interaction).

A network matrix can be created using the function `get_network_matrix()`. Gene interaction information for a specific pathway is extracted from the KEGG database. It is accessed via the function `retrieveKGML()` from the `KEGGgraph` package. An internet connection is required for this step.

```
pathw_complete <- get_network_matrix(pathw, directed=FALSE)
```

will download the KEGG XML file for the pathway with ID 'hsa04710' and save it in the working directory. The function will convert the data into a network matrix and add it to the given pathway object. The expanded pathway object will be returned. The user can specify whether the gene-interaction matrix should be given directed (`directed=TRUE`) or undirected (`directed=FALSE`).

SNP positions

Kangaroo offers a function to download positions of the SNPs available in your GWAS dataset from the Ensembl database.

- `snp_info()` will take a vector of rs-numbers and give the corresponding base pair positions.

- Positions are extracted from the Ensembl database and thus equal the current build used on the website. The database is accessed via the function `getBM()` from the package `biomaRt`. This requires an internet connection.

```
snp_info("rs234")
```

will return a `snp_info` object containing the data frame

| chr | position | rsnumber |
|-----|-----------|----------|
| 7 | 105920689 | rs234 |

Pathway Annotation

To define SNP sets representing a pathway, the function `get_anno()` can be used.

- Input arguments are a `pathway_info` as well as a `snp_info` object.
- If you do not want to change positions in your SNP file using the `snp_info()` function, you will have to transform it into a `snp_info` object including a data frame listing all SNPs to be annotated. This data frame must include the columns 'chr', 'position' and 'rsnumber', giving for each SNP the chromosome it lies on, its base pair position on the chromosome and the rs-numbers identifier, respectively. See also the output description of `snp_info()`.
- For annotation the package `sqldf` is used.

```
get_anno(snp_info, pathway_info)
```

will return a data frame listing all SNPs that lie inside the boundaries of one or more genes in the pathway. That means that genes can appear several times, depending on the number of SNPs mapped to them. A SNP can and will be mapped to multiple genes if they overlap. The data frame will have the following format

| pathway | gene | chr | snp | position |
|----------|--------|-----|------------|----------|
| hsa04710 | CSNK1E | 22 | rs11089885 | 38413480 |
| hsa04710 | CSNK1E | 22 | rs13054361 | 38336819 |
| hsa04710 | CSNK1E | 22 | rs135757 | 38307648 |
| hsa04710 | ... | ... | ... | ... |

GWAS data

Data from a case control study is needed to test a pathways influence on disease risk with the logistic kernel machine test in `kangaroo`. Here, GWAS data is represented by the `GWASdata` object. It includes

- Genotype data for each individual.
 - Genotype data needs to be a matrix with one line per individual and one column for each SNP.
 - Rownames give ID numbers for the individuals while columnnames give the rs-numbers corresponding to the SNPs genotyped in the study.
 - Note that missing values are not allowed and SNPs with missing genotypes have to be imputed or excluded from the sample prior to creation of the `GWASdata` object.
- Phenotype data for each individual.
 - Phenotypes need to be given in a **data frame** with the first column including the individual IDs as in the genotype sample.

- Further columns can contain informative covariates (such as age, sex, ...) to be used in the logistic regression model.
- Annotation of study SNPs to pathways created by `get_anno`.
 - This **data frame** defines the SNP set representing a specific pathway. It can be created using the function `get_anno()`.
- A **character** describing the data can be added to the **GWASdata** object. This could for instance be the name of the study.

A **GWASdata** object can be constructed as

```
my_gwas <- GWASdata(pheno=pheno, geno=geno, anno=anno, desc="study xy")
```

Calculation of Kernel Matrices

Once a **GWASdata** object is created, we can start to calculate kernel matrices to test a pathways influence on disease risk. **kangaroo** offers three different kernel functions to compute a similarity matrix for the individuals in analysis. They will be explained in the following.

Linear Kernel (Lin)

The linear kernel assumes additive SNP effects. It is calculated as

$$ZZ^t(2)$$

where Z denotes the genotype matrix (See also Liu et al, 2010). In **kangaroo** a linear kernel can be created using the function `kernel_lin()`. It requires as arguments

- A **GWASdata** object containing the genotype information.
- A **pathway** object specifying the pathway to be tested.
- A value for argument **calculation** to decide how the kernel should be calculated. Options are **cpu** for calculation on cpu and **gpu** for gpu calculation.

```
K_lin <- lin_kernel(gwas, p, calculation='cpu')
```

will return a quadratic matrix of dimension equal to the number of individuals in the **GWASdata** object.

Size-adjusted Kernel (Sia)

The size-adjusted kernel takes into consideration the numbers of SNPs and genes in a pathway to correct for size bias. It is calculated as

$$K_{i,j} = \exp\left(-\sqrt{\frac{1}{r_p}} \sum_g \left(\frac{\|z_i^g - z_j^g\|}{\mu_g k_g^{eff}}\right)^{\delta_g}\right)(3)$$

Here z_i^g is the vector of individual i 's genotypes in gene g and r_p the number of genes in pathway p . Scaling parameters k_g^{eff} , μ_g and δ_g adjust for the number of genes in the pathway and the number of SNPs within these genes (for more details refer to Freytag et al. 2012).

A kernel of this type can be calculated using the function `kernel_sia()` with the following arguments

- A **GWASdata** object containing the genotype information.
- A **pathway** object specifying the pathway to be tested.
- A value for argument **calculation** to decide how the kernel should be calculated. Currently only **cpu** for cpu calculation is available.

```
K_sia <- sia_kernel(gwas, p, calculation='cpu')
```

will return a quadratic **matrix** of dimension equal to the number of individuals in the **GWASdata** object.

Network Kernel (Net)

The network kernel incorporates information about gene-gene interactions into the model. It is defined as

$$K = ZANA^tZ^t(4)$$

where matrix A maps SNPs to genes, N represents the underlying network structure, and Z is the genotype matrix. The network based kernel matrix for a pathway can be calculated with the function **kernel_net()**. Following arguments are needed

- A **GWASdata** object containing the genotype information.
- A **pathway** object specifying the pathway to be tested.
- A value for argument 'calculation' to decide how the kernel should be calculated.

```
K_net <- net_kernel(gwas, p, calculation='cpu')
```

will return a quadratic **matrix** of dimension equal to the number of individuals in the **GWASdata** object.

Alternatively, kernel matrices can be calculated using the function **calc_kernel()**. Here the kernel type is specified via an additional argument **type**. It can be set to **lin**, **sia** or **net**.

```
K <- calc_kernel(gwas, p, type='lin', parallel='none')
```

This function will simply call the suitable kernel function as described above and therefore has the same output.

Variance Component Test

A pathways influence on the probability of being a case is evaluated in a variance component test. The test statistic is

$$Q = \frac{1}{2}(y - \mu)^t K(y - \mu)(5)$$

with μ the vector of null model estimators given by $\mu_i = \text{logit}^{-1}(x_i^t \beta)$ for an individual i and K a kernel matrix of the pathway to be tested. Q follows a mixture of X^2 distributions which can be approximated using the Satterthwaite procedure (Schaid 2012) or Davies method as implemented in the R package **QuadCompForm** (Davies 1980). More details on the test can be found in Wu et al (2010).

In **kangaroo00** the logistic kernel machine test can be applied to a SNP set defining a pathway with the function **lkmt**. It needs the following arguments

- A formula specifying the null model to be used in the test. The dependent variable is the case control status of the individual (in the example denoted as 'pheno') and is explained by an intercept and optional covariates.
- A linear, size-adjusted or network kernel matrix calculated by one of the kernel functions **kernel_lin()**, **kernel_sia()** or **kernel_net()**.

- A **GWASdata** object including the genotype based on which the test should be performed.
- A **character** specifying which method should be used to calculate the p-value. Available are ‘satt’ for the Satterthwaite approximation (Schaid

2010) or ‘davies’ for Davies method (Davies 1980).

```
pval_net <- lkmt(pheno ~ 1+sex+age, K_mat, my_gwas, method='satt')
```

will return an object of type **lkmt** giving the test result for the pathway on which the kernel matrix ‘K_mat’ was calculated. The **GWASdata** object ‘my_gwas’ has to be the same as used to calculate the kernel matrix. The formula above would for example fit for a phenotype file of the following format (IDs in first column are always required in phenotype file)

| ID | pheno | sex | age | smoker |
|------|-------|-----|-----|--------|
| ind1 | 1 | 1 | 41 | 1 |
| ind2 | 0 | 0 | 38 | 0 |
| ind3 | 1 | 1 | 56 | 1 |
| ... | ... | ... | ... | ... |

note, that the columns to be used in the model are specified in the formula given to the **lkmt()** function and not all covariates have to be used.

References

- Liu D, Ghosh D, Lin X. Estimation and testing for the effect of a genetic pathway on a disease outcome using logistic kernel machine regression via logistic mixed models. *BMC Bioinformatics* 2008 9:292.
- Schaid DJ: Genomic similarity and kernel methods I: advancements by building on mathematical and statistical foundations. *Hum Hered* 2010, 70:109-131.
- Davies R: Algorithm as 155: the distribution of a linear combination of chi-2 random variables. *J R Stat Soc Ser C* 1980, 29:323-333.
- Wu MC, Kraft P, Epstein MP, Taylor DM, Chanock SJ, Hunter DJ, Lin X: Powerful SNP-Set Analysis for Case-Control Genome-Wide Association Studies. *Am J Hum Genet* 2010, 86:929-42
- Cunningham F, Amode MR, Barrell D et al. Ensembl 2015. *Nucleic Acids Research* 2015 43 Database issue:D662-D669
- Kanehisa, M., Goto, S., Sato, Y., Kawashima, M., Furumichi, M., and Tanabe, M.; Data, information, knowledge and principle: back to metabolism in KEGG. *Nucleic Acids Res.* 42, D199-D205 (2014).
- Freytag S, Bickeboeller H, Amos CI, Kneib T, Schlather M: A Novel Kernel for Correcting Size Bias in the Logistic Kernel Machine Test with an Application to Rheumatoid Arthritis. *Hum Hered.* 2012, 74(2):97-108.
- Freytag S, Manitz J, Schlather M, Kneib T, Amos CI, Risch A, Chang-Claude J, Heinrich J, Bickeboeller H: A network-based kernel machine test for the identification of risk pathways in genome-wide association studies. *Hum Hered.* 2013, 76(2):64-75.