

Using the inSilicoDb package

Jonatan Taminau*

CoMo, Vrije Universiteit Brussel, Pleinlaan 2, 1050 Brussels,
Belgium

1 The inSilico database

With more than 500,000 genomic profiles freely available in the public domain, there is a huge amount of information accessible for computational biologists or bioinformaticians. However, accessibility to these data requires complex computational steps. Manual parsing of annotations and keywords, which is in most cases a necessary evil, tends to be time-consuming and is known to be error-prone. Also the wide variety of normalization and preprocessing methods makes the comparability of different existing studies hard, or even impossible. The inSilico initiative (<http://insilico.ulb.ac.be>) provides an answer to those problems with its freely available web-based database tool: the inSilico database¹. Starting with all public available human Affymetrix studies from Gene Expression Omnibus (GEO) [1] it provides those studies in a consistent and well curated form. With a direct connection to GenePattern [3] and the ability to export the data to different formats, the inSilico database is an efficient mean to re-analyse public datasets and improve reproducibility in genome-wide research.

To further ease the use of this vast amount of genomic data the inSilicoDb R package was developed. This package can be seen as a different front-end to the core inSilico database and, although it provides only limited functionality compared to the web-based tool, it can become very valuable for R programmers or anyone who is interested in large scale analysis using automated scripting.

Similar packages to retrieve gene expression data in R exist [5, 4], but the added value and strength of this package is tightly connected to the innovative concept of the inSilico database and will therefore circumvent common obstacles like incompatibility and missing or malformed annotations.

* jtaminau@vub.ac.be

¹Manuscript in preparation

2 Getting started using inSilicoDb

As this section will show, accessing data from the inSilico database is surprisingly easy and straightforward.

2.1 Simple access

Suppose one is interested in a number of publicly available gene expression studies which he found while browsing Gene Expression Omnibus (GEO) or the inSilico database. Using only the GSE identifier, a completely annotated and formatted dataset can be downloaded in just seconds, without any need for further manual parsing: You need to login to access datasets and datasetinfo. Use your InSilicoDB login and an md5 hash of your password. For this example we're using a restricted test account.

```
> library("inSilicoDb");
> InSilicoLogin("rpackage_tester@insilicodb.com", "5c4d0b231e5cba4a0bc54783b385cc9a");

[1] 5296

> res = getDatasets("GSE4635");
> eset = res[[1]];
```

The result of `getDataset` is a list, containing a Bioconductors ExpressionSet (`eset`) for every platform that exists for this dataset (in the example there is only one platform). An alternative approach to obtain the same data is to specify the platform. In this case no list but the expression set is directly returned:

```
> eset = getDataset("GSE4635", "GPL96");
```

And in case the platform is unknown, the auxiliary function `getPlatforms` is provided:

```
> platforms = getPlatforms("GSE4635");
> print(platforms);

[1] "GPL96"
```

Once an expression set is retrieved, all available Bioconductor packages can be applied for further analysis, as the following code illustrates:

```
> #eset = getDataset("GSE4635", "GPL96", features = "gene");
> #heatmap(exprs(eset)[1:100,]);
> library("limma")
> eset = getDataset("GSE4635", "GPL96",
+                   norm="FRMA", features = "gene");
> # Find 50 most discriminating genes
> f = pData(eset)[ , "Smoker"]
```

Heatmap showing the expression of 48 genes across 32 samples. The samples are grouped into four pairs of 'ever smoker' and 'never smoker'. The genes are listed on the right, including TKT, PDIA4, S100A14, etc. The color scale ranges from yellow (low expression) to red (high expression).

In case only the annotation information is needed and there is no need for the numerical data, the `getAnnotations` function also exists for convenience:

		Smoker	tissue	status	age	race	sex	pkys
GSM15729	never	smoker	Bronchial Epithelium	NILL	NILL	NILL	NILL	NILL
GSM104072	current	smoker	Bronchial Epithelium	NILL	NILL	NILL	NILL	NILL
GSM104074	current	smoker	Bronchial Epithelium	NILL	NILL	NILL	NILL	NILL
GSM104075	current	smoker	Bronchial Epithelium	SNC	45	AFA	Female	14
GSM104076	current	smoker	Bronchial Epithelium	NILL	NILL	NILL	NILL	NILL

GSM104078	never smoker	Bronchial Epithelium	NILL	NILL	NILL	NILL	NILL
GSM104080	never smoker	Bronchial Epithelium	NILL	NILL	NILL	NILL	NILL
GSM104082	never smoker	Bronchial Epithelium	NILL	NILL	NILL	NILL	NILL
		history patient_id	New Column	sdfadsf			
GSM15729		NILL	NILL			NULL	
GSM104072		NILL	NILL			NULL	
GSM104074		NILL	NILL			NULL	
GSM104075	current smoker without cancer		406			NULL	
GSM104076		NILL	NILL			NULL	
GSM104078		NILL	NILL			NULL	
GSM104080		NILL	NILL			NULL	
GSM104082		NILL	NILL			NULL	

2.2 More options

By default all numerical data is retrieved the same way the original authors have submitted the data to GEO and can therefore have been processed by a wide variety of preprocessing methods. However, when combining different studies a consistent preprocessing is required and therefore all studies for which there are CEL files available, are also precomputed by applying the FRMA preprocessing method [2]. The user can retrieve those studies as fast and easy as the original ones, simply by using the optional `norm` parameter.

```
> eset = getDataset("GSE4635", "GPL96", norm="FRMA");
```

All gene expression matrices contain probes as features, although it is also possible to retrieve the genes instead. This probe to gene mapping is precomputed for every dataset and can be selected using the `genes` parameter. By default probes are selected, as this is how the data was submitted to GEO.

```
> eset = getDataset("GSE4635", "GPL96");
> print(nrow(eset));
```

```
Features
22283
```

```
> eset = getDataset("GSE4635", "GPL96", features="gene");
> print(nrow(eset));
```

```
Features
12500
```

2.3 Create your own loop...

One of the advantages of retrieving data through R is the possibility to develop a whole automated workflow in just a few lines of code. The following example illustrates the many opportunities researchers can have using this tool.

In the example code below, we iterate over a list of series GSE identifiers and try to retrieve every dataset from the database. Once retrieved some basic analysis is performed (printing the number of annotations and missing values). Note that the `getDataset` function can throw an error (e.g. no internet connection, dataset is not available, etc.) which is best caught in a try-catch loop, as is shown in the example.

```
> lst = list("GSE4635", "GSExxx", "GSE781");
> gpl = "GPL96";
> for(gse in lst)
+ {
+   catn = function(...) { cat(..., "\n"); }
+   catn("Processing", gse);
+   catn("=====");
+   eset = tryCatch({getDataset(gse, gpl)};
+               error = function(x) { print(as.character(x)); NULL; });
+   if(is.null(eset)) { next; }
+   catn("Number of annotations:");
+   catn(ncol(pData(eset)));
+   catn("Number of missing values:");
+   catn(sum(is.na(exprs(eset))));
+ }
```

```
Processing GSE4635
=====
Number of annotations:
10
Number of missing values:
68
Processing GSExxx
=====
[1] "Error: Stopped because of previous errors\n"
Processing GSE781
=====
Number of annotations:
24
Number of missing values:
0
```

3 Conclusion

This package is built in addition to a very powerful web-based database tool for genomic analysis. Despite its simplicity, it captures many of the benefits of this tool and provides the typical R users efficient means of performing large scale genomic analysis using automated scripting.

4 Session Info

```
> sessionInfo()

R version 3.3.0 alpha (2016-03-28 r70390)
Platform: x86_64-w64-mingw32/x64 (64-bit)
Running under: Windows Server 2008 R2 x64 (build 7601) Service Pack 1

locale:
[1] LC_COLLATE=C
[2] LC_CTYPE=English_United States.1252
[3] LC_MONETARY=English_United States.1252
[4] LC_NUMERIC=C
[5] LC_TIME=English_United States.1252

attached base packages:
[1] parallel stats graphics grDevices utils datasets methods
[8] base

other attached packages:
[1] limma_3.27.14 inSilicoDb_2.7.0 RCurl_1.95-4.8
[4] bitops_1.0-6 Biobase_2.31.3 BiocGenerics_0.17.3
[7] rjson_0.2.15

loaded via a namespace (and not attached):
[1] tools_3.3.0
```

References

- [1] Ron Edgar, Michael Domrachev, and Alex E Lash. Gene expression omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Res*, 30(1):207–10, Jan 2002.
- [2] Matthew N McCall, Benjamin M Bolstad, and Rafael A Irizarry. Frozen robust multiarray analysis (fRMA). *Biostatistics*, 11(2):242–53, Apr 2010.
- [3] Michael Reich, Ted Liefeld, Joshua Gould, Jim Lerner, Pablo Tamayo, and Jill P Mesirov. Genepattern 2.0. *Nat Genet*, 38(5):500–1, May 2006.
- [4] Davis Sean and Paul S Meltzer. GEOquery: a bridge between the gene expression omnibus (GEO) and bioconductor. *Bioinformatics*, 23(14):1846–7, Jul 2007.
- [5] Yuelin Zhu, Sean Davis, Robert Stephens, Paul S Meltzer, and Yidong Chen. GEOmetadb: powerful alternative search engine for the gene expression omnibus. *Bioinformatics*, 24(23):2798–800, Dec 2008.