

# isomiRs

Lorena Pantano<sup>1\*</sup>, Georgia Escaramis<sup>2\*</sup>, Eulalia Martin<sup>2\*</sup>

Harvard H Chan School of Public Health, Boston, US;

<sup>2</sup> Center of Genomic Regulation, Barcelona, Spain;

\*lpantano (at) iscb.org

Modified: 3 Feb, 2015. Compiled: July 20, 2015

## Contents

---

<b>1</b>	<b>Citing isomiRs</b>	<b>2</b>
<b>2</b>	<b>Input format</b>	<b>2</b>
<b>3</b>	<b>IsomirDataSeq class</b>	<b>2</b>
3.1	Access data . . . . .	2
3.2	isomiRs annotation . . . . .	3
<b>4</b>	<b>Quick start</b>	<b>3</b>
4.1	Reading input . . . . .	3
4.2	Descriptive analysis . . . . .	3
4.3	Differential expression analysis . . . . .	4
4.4	Count data . . . . .	6
4.5	Supervised classification . . . . .	7

## Introduction

miRNA are small RNA fragments (18-23 nt long) that influence gene expression during development and cell stability. Morin et al [1], discovered isomiRs first time after sequencing human stem cells.

IsomiRs are miRNAs that vary slightly in sequence, which result from variations in the cleavage site during miRNA biogenesis (5'-trimming and 3'-trimming variants), nucleotide additions to the 3'-end of the mature miRNA (3'-addition variants) and nucleotide modifications (substitution variants)[2].

There are many tools designed for isomiR detection, however the majority are web application where user can not control the analysis. The two main command tools for isomiRs mapping are SeqBuster and sRNAbench[3]. *isomiRs* package is designed to analyze the output of SeqBuster tool or any other tool after converting to the desire format.

## 1 Citing isomiRs

If you use the package, please cite this paper [4].

## 2 Input format

The input should be the output of SeqBuster-miraligner tool (\*.mirna files) for each sample in the following format:

seq	name	freq	mir	start	end	mism	add	t5	t3	s5	s3	DB	am
TGTAAACATCCTACACTCAGCTGT				seq_100014_x23	23		hsa-miR-30b-5p	17	40	0	0	0	0
TGTAAACATCCCTGACTGGAA			seq_100019_x4	4		hsa-miR-30d-5p	6	26	13TC	0	0	0	u-
TGTAAACATCCCTGACTGGAA			seq_100019_x4	4		hsa-miR-30e-5p	17	37	12CT	0	0	0	u-
CAAATTCGTATCTAGGGGATT			seq_100049_x1	1		hsa-miR-10a-3p	63	81	0	u-TT	0	0	u-
TGACCTAGGAATTGACAGCCAGT			seq_100060_x1	1		hsa-miR-192-5p	25	47	8GT	0	d-C	d-	d-

This is the standard output of SeqBuster-miraligner tool, but can be converted from any other tool having the mapping information on the precursors. Read more on [miraligner manual](#)

## 3 IsomirDataSeq class

This object will store all raw data from the input files and some processed information used for visualization and statistical analysis. It is a subclass of SummarizedExperiment with colData and counts methods. Beside that, the object contains raw and normalized counts from miraligner allowing to update the summarization of miRNA expression.

### 3.1 Access data

The user can access the normalized count matrix with `counts(object, norm=TRUE)`.

You can browse for the same miRNA or isomiRs in all samples with `isoSelect` method.

```
library(isomiRs)
data(isomiRexp)
head(isoSelect(isomiRexp, mirna="hsa-let-7a-5p"))
```

##		nb1	nb2	nb3	o1	o2	o3
##	hsa-let-7a-5p.iso.t5:0.seed:0.t3:0.ad:0.mm:12TG	21785	27357	25998	31263	16492	40180
##	hsa-let-7a-5p.iso.t5:0.seed:0.t3:0.ad:0.mm:18CT	7412	6454	11449	14130	5088	17159

```
## hsa-let-7a-5p.iso.t5:0.seed:0.t3:0.ad:0.mm:19TA      8668  5195      0      0      0      0
## hsa-let-7a-5p.iso.t5:0.seed:0.t3:0.ad:u-A.mm:0      36208 34981 42310 35404 23139 37970
## hsa-let-7a-5p.iso.t5:0.seed:0.t3:d-T.ad:0.mm:0      35572 25121 36082 22035 13302 22766
## hsa-let-7a-5p.iso.t5:0.seed:0.t3:d-T.ad:0.mm:9GT     4430  4992 12003      0      0      0
```

## 3.2 isomiRs annotation

IsomiR names follows this structure:

- miRNA name
- type: ref if the sequence is the same than the miRNA reference. 'iso' if the sequence has variations.
- t5 tag: indicates variations at 5' position. The naming contains two words: 'direction - nucleotides', where direction can be 'u' (changes upstream of the 5' reference position) or 'd' (changes downstream of the 5' reference position). '0' indicates no variation, meaning the 5' position is the same than the reference. After 'direction', it follows the nucleotide/s that are added (for upstream changes) or deleted (for downstream changes).
- t3 tag: indicates variations at 3' position. The naming contains two words: 'direction - nucleotides', where direction can be 'u' (upstream of the 3' reference position) or 'd' (downstream of the 3' reference position). '0' indicates no variation, meaning the 3' position is the same than the reference. After 'direction', it follows the nucleotide/s that are added (for downstream changes) or deleted (for upstream changes).
- ad tag: indicates nucleotides additions at 3' position. The naming contains two words: 'direction - nucleotides', where direction is 'u' (upstream of the 5' reference position). '0' indicates no variation, meaning the 3' position has no additions. After 'direction', it follows the nucleotide/s that are added.
- mm tag: indicates nucleotides substitutions along the sequences. The naming contains three words: 'position-nucleotideATsequence-nucleotideATreference'.
- seed tag: same than 'mm' tag, but only if the change happens between nucleotide 2 and 8.

## 4 Quick start

We are going to use a small RNAseq data from human frontal cortex samples [5] to give some basic examples of isomiRs analyses.

In this data set we will find two groups:

- b: 3 individuals with less than a year
- o: 3 individuals in the elderly.

```
library(isomiRs)
data(isomiRexp)
```

### 4.1 Reading input

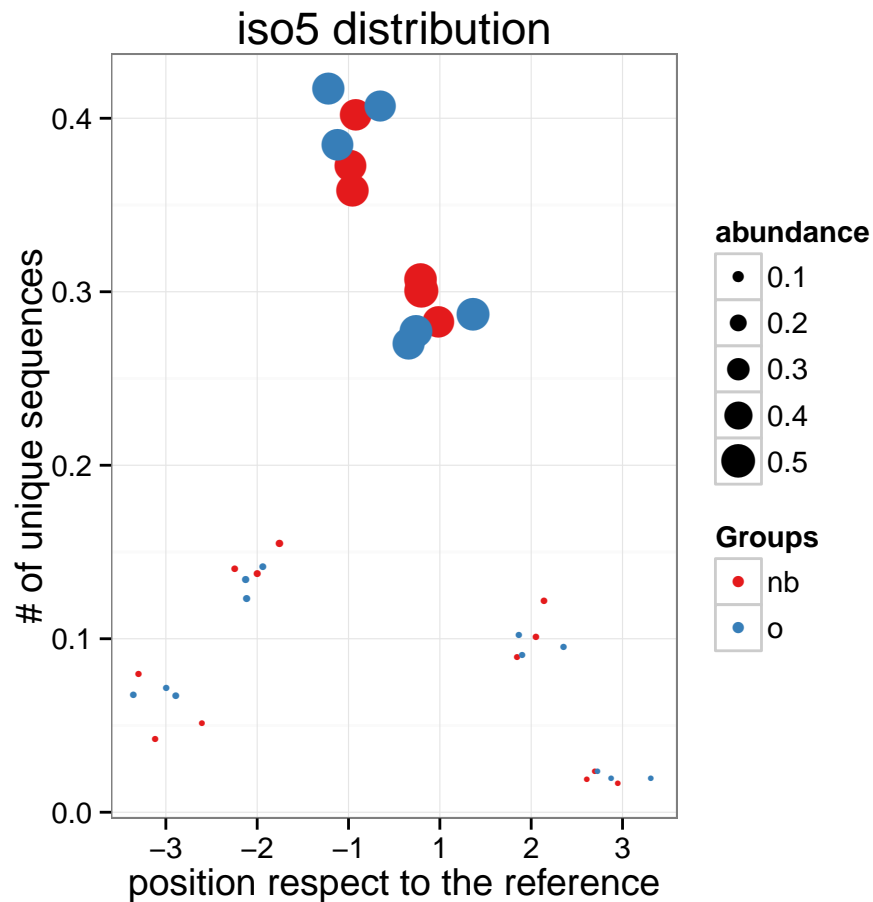
The function `loadIso` needs a vector with the paths for each file and a data frame with the design experiment similar to the one used for a mRNA differential expression analysis.

```
obj <- IsomirDataSeqFromFiles(fn_list, design=de)
```

### 4.2 Descriptive analysis

You can plot isomiRs expression with `isoPlot`. In this figure you will see how abundant is each type of isomiRs at different positions considering the total abundance and the total number of sequences. The type parameter controls what type of isomiRs to show. It can be trimming (iso5 and iso3), addition (add) or substitution (subs) changes.

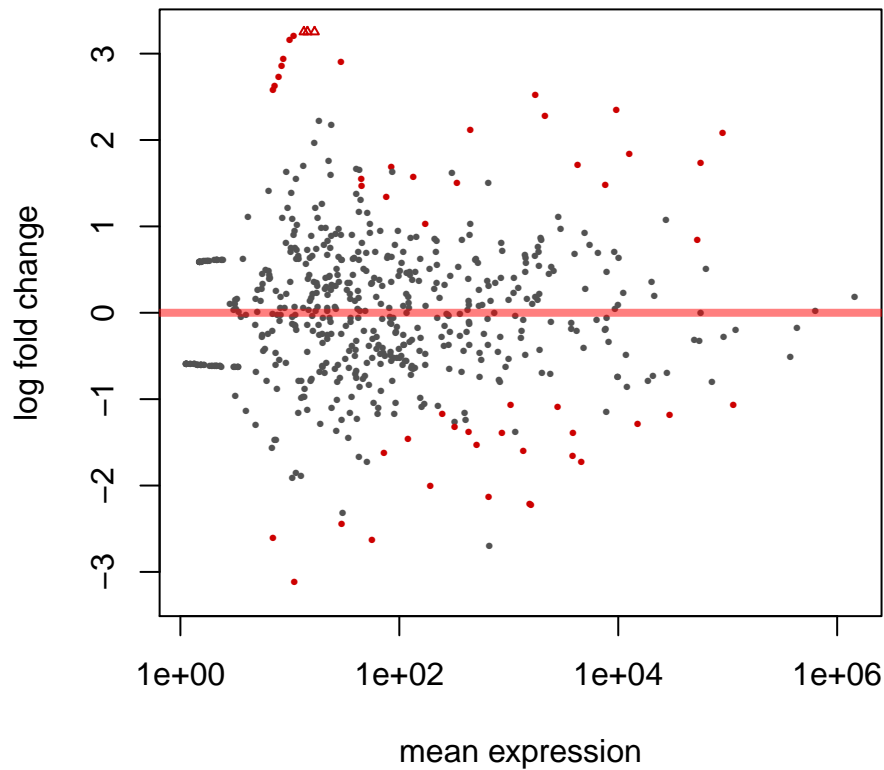
```
obj <- isoPlot(isomiRexp, type="iso5")
```



### 4.3 Differential expression analysis

The `isoDe` uses functions from [DESeq2](#) package. This function has parameters to create a matrix using only the reference miRNAs, all isomiRs, or some of them. This matrix and the design matrix are the inputs for `DESeq2`. The output will be a `DESeqDataSet` object, allowing to generate any plot or table explained in `DESeq2` package vignette.

```
dds <- isoDE(obj, formula=~condition)
library(DESeq2)
plotMA(dds)
```



```
head(results(dds, tidy=TRUE))
```

##	row	baseMean	log2FoldChange	lfcSE	stat	pvalue	padj
## 1	hsa-let-7a-3p	39.28208	0.1325786	0.5975562	0.2218680	0.8244166	0.9541789
## 2	hsa-let-7a-5p	429526.75771	-0.1697043	0.3917262	-0.4332217	0.6648537	0.8804819
## 3	hsa-let-7b-3p	23.82048	0.2636501	0.6501903	0.4054968	0.6851124	0.8917654
## 4	hsa-let-7b-5p	118492.78690	-0.1999231	0.3232163	-0.6185428	0.5362176	0.8141445
## 5	hsa-let-7c-5p	373488.14001	-0.5151406	0.4461972	-1.1545132	0.2482898	0.6998768
## 6	hsa-let-7d-3p	11.95405	0.6544099	0.8588832	0.7619312	0.4461010	0.7398849

You can differentiate between reference sequences and isomiRs at 5' end with this command:

```
dds = isoDE(obj, formula=~condition, ref=TRUE, iso5=TRUE)
```

```
head(results(dds, tidy=TRUE))
```

##	row	baseMean	log2FoldChange	lfcSE	stat	pvalue
## 1	hsa-let-7a-3p.iso.t5:0	2.254768e+01	-0.11279202	0.7201078	-0.15663212	0.8755348
## 2	hsa-let-7a-3p.iso.t5:d-C	4.819743e+00	0.36981376	0.8924843	0.41436445	0.6786072
## 3	hsa-let-7a-3p.iso.t5:d-CT	1.794593e+00	0.69047324	0.7594203	0.90921089	NA
## 4	hsa-let-7a-5p.iso.t5:0	1.481710e+05	-0.36374581	0.4192563	-0.86759762	0.3856146
## 5	hsa-let-7a-5p.ref.t5:0	2.796471e+05	-0.03897592	0.3932878	-0.09910281	0.9210566
## 6	hsa-let-7b-3p.iso.t5:0	2.343947e+01	0.28720908	0.6815838	0.42138485	0.6734741
##	padj					
## 1		0.9806026				
## 2		0.8870062				

```
## 3      NA
## 4 0.7102917
## 5 0.9933185
## 6 0.8840507
```

## 4.4 Count data

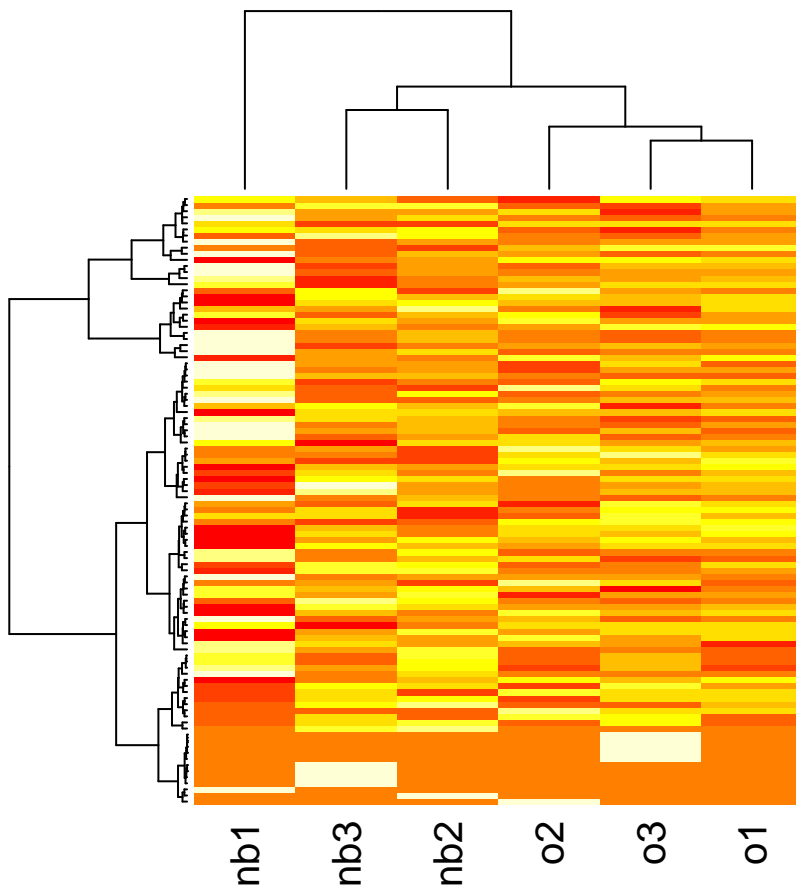
`isoCounts` gets the count matrix that can be used for many different downstream analyses changing the way isomiRs are collapsed. The following command will merge all isomiRs into one feature: the reference miRNA.

```
obj = isoCounts(obj)
head(counts(obj))
```

##		nb1	nb2	nb3	o1	o2	o3
##	hsa-let-7a-3p	24	23	70	47	26	65
##	hsa-let-7a-5p	364764	373809	485644	495146	284300	552624
##	hsa-let-7b-3p	12	17	38	24	27	33
##	hsa-let-7b-5p	90429	107361	162601	131970	89807	136599
##	hsa-let-7c-5p	406504	314874	404699	375977	235729	379519
##	hsa-let-7d-3p	0	13	15	20	15	17

The normalization uses `rlog` from [DESeq2](#) package and allows quick integration to another analyses like heatmap, clustering or PCA.

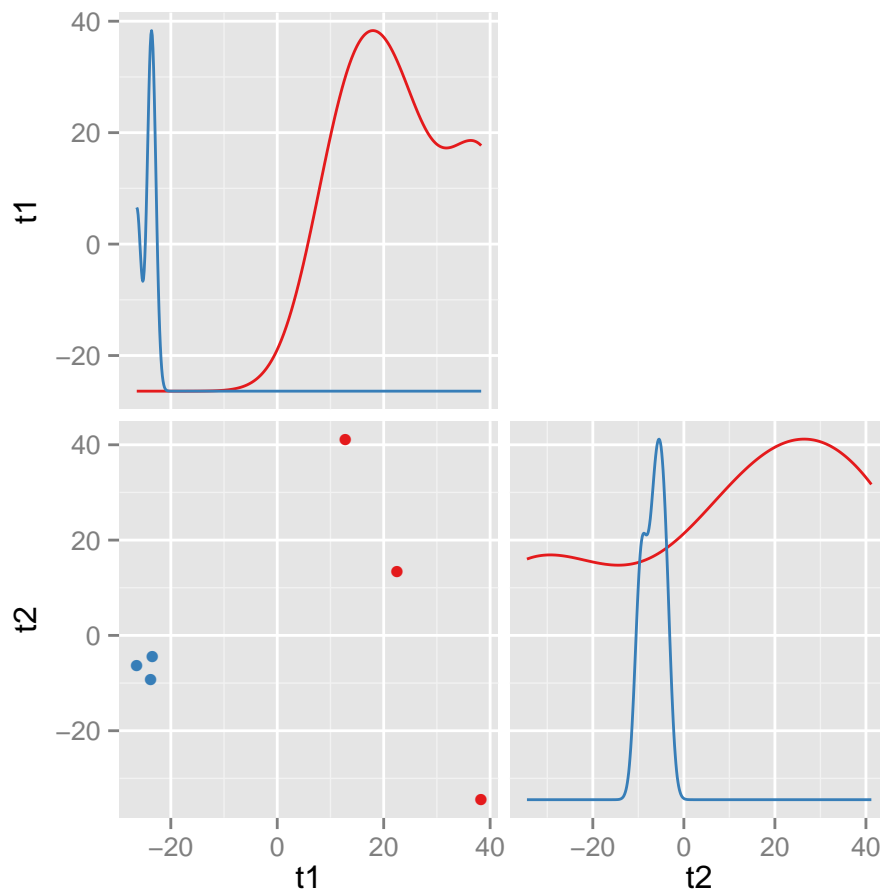
```
obj = isoNorm(obj)
heatmap(counts(obj, norm=TRUE)[1:100,], labRow = "")
```



## 4.5 Supervised classification

Partial Least Squares Discriminant Analysis (PLS-DA) is a technique specifically appropriate for analysis of high dimensionality data sets and multicollinearity [6]. PLS-DA is a supervised method (i.e. makes use of class labels) with the aim to provide a dimension reduction strategy in a situation where we want to relate a binary response variable (in our case young or old status) to a set of predictor variables. Dimensionality reduction procedure is based on orthogonal transformations of the original variables (isomiRs) into a set of linearly uncorrelated latent variables (usually termed as components) such that maximizes the separation between the different classes in the first few components [7]. We used sum of squares captured by the model ( $R^2$ ) as a goodness of fit measure. We implemented this method using the *Discriminer* into *isoPLSDA* function. The output p-value of this function will tell about the statistical significant of the group separation using miRNA expression data. Moreover, the function *isoPLSDAplot* helps to visualize the results. It will plot the samples using the significant components (t1, t2, t3 ...) from the PLS-DA analysis and the samples distribution along the components.

```
obj = isoCounts(obj, iso5=TRUE, iso3=TRUE, add=TRUE, ref=TRUE)
obj = isoNorm(obj)
pls.obj = isoPLSDA(obj, "condition", nperm = 10)
## [1] "pval:0.1"
isoPLSDAplot(pls.obj$component, colData(obj)[,"condition"])
```



The analysis can be done again using only the most important discriminant isomiRS from the PLS-DA models based on the analysis. We used Variable Importance for the Projection (VIP) criterion to select the most important features, since takes into account the contribution of a specific predictor for both the explained variability on the response and the explained variability on the predictors.

```
pls.obj = isoPLSDA(obj, "condition", refinement = FALSE, vip = 0.8)
```



## Session info

---

Here is the output of `sessionInfo` on the system on which this document was compiled:

- R version 3.2.1 (2015-06-18), x86\_64-unknown-linux-gnu
- Locale: LC\_CTYPE=en\_US.UTF-8, LC\_NUMERIC=C, LC\_TIME=en\_US.UTF-8, LC\_COLLATE=en\_US.UTF-8, LC\_MONETARY=en\_US.UTF-8, LC\_MESSAGES=en\_US.UTF-8, LC\_PAPER=en\_US.UTF-8, LC\_NAME=C, LC\_ADDRESS=C, LC\_TELEPHONE=C, LC\_MEASUREMENT=en\_US.UTF-8, LC\_IDENTIFICATION=C
- Base packages: base, datasets, graphics, grDevices, methods, parallel, stats, stats4, utils
- Other packages: Biobase 2.29.1, BiocGenerics 0.15.3, DESeq2 1.9.21, DiscrMiner 0.1-29, GenomeInfoDb 1.5.8, GenomicRanges 1.21.16, IRanges 2.3.14, isomiRs 0.99.6, knitr 1.10.5, Rcpp 0.11.6, RcppArmadillo 0.5.200.1.0, S4Vectors 0.7.10, SummarizedExperiment 0.3.2
- Loaded via a namespace (and not attached): acepack 1.3-3.3, annotate 1.47.0, AnnotationDbi 1.31.17, assertthat 0.1, BiocParallel 1.3.31, BiocStyle 1.7.4, bitops 1.0-6, caTools 1.17.1, cluster 2.0.2, colorspace 1.2-6, DBI 0.3.1, digest 0.6.8, dplyr 0.4.2, evaluate 0.7, foreign 0.8-65, formatR 1.2, Formula 1.2-1, futile.logger 1.4.1, futile.options 1.0.0, gdata 2.17.0, genefilter 1.51.0, geneplotter 1.47.0, GGally 0.5.0, ggplot2 1.0.1, gplots 2.17.0, grid 3.2.1, gridExtra 0.9.1, gtable 0.1.2, gtools 3.5.0, highr 0.5, Hmisc 3.16-0, KernSmooth 2.23-15, labeling 0.3, lambda.r 1.1.7, lattice 0.20-31, latticeExtra 0.6-26, lazyeval 0.1.10, locfit 1.5-9.1, magrittr 1.5, MASS 7.3-42, munsell 0.4.2, nnet 7.3-10, plyr 1.8.3, proto 0.3-10, R6 2.1.0, RColorBrewer 1.1-2, reshape 0.8.5, reshape2 1.4.1, rpart 4.1-10, RSQLite 1.0.0, scales 0.2.5, splines 3.2.1, stringi 0.5-5, stringr 1.0.0, survival 2.38-3, tools 3.2.1, XML 3.98-1.3, xtable 1.7-4, XVector 0.9.1

## References

---

- [1] R. D. Morin, M. D. O'Connor, M. Griffith, F. Kuchenbauer, A. Delaney, A.-L. Prabhu, Y. Zhao, H. McDonald, T. Zeng, M. Hirst, C. J. Eaves, and M. A. Marra. Application of massively parallel sequencing to microRNA profiling and discovery in human embryonic stem cells. *Genome Res.*, 18:610–621, 2008. doi:10.1101/gr.7179508, PMID:18285502.
- [2] Eulàlia Martí, Lorena Pantano, Mónica Bañez Coronel, Franc Llorens, Elena Miñones Moyano, Sílvia Porta, Lauro Sumoy, Isidre Ferrer, and Xavier Estivill. A myriad of miRNA variants in control and Huntington's disease brain regions detected by massively parallel sequencing. *Nucleic Acids Res.*, 38:7219–35, 2010. doi:10.1093/nar/gkq575, PMID:20591823.
- [3] Barturen Guillermo, Rueda Antonio, Hamberg Maarten, Alganza Angel, Lebron Ricardo, Kotsyfakis Michalis, Shi BuJun, KoppersLalic Danijela, and Hackenberg Michael. sRNAbench: profiling of small RNAs and its sequence variants in single or multi-species high-throughput experiments. *Methods in Next Generation Sequencing*, 1(1):2084–7173, 2014. doi:10.2478/mngs-2014-0001.
- [4] Estivil X Pantano L, Marti E. SeqBuster. *Nucleic Acids Res.*, 38:e34, 2010. doi:10.1093/nar/gkp1127, PMID:20008100.
- [5] Mehmet Somel, Song Guo, Ning Fu, Zheng Yan, Hai Yang Hu, Ying Xu, Yuan Yuan, Zhibin Ning, Yuhui Hu, Corinna Menzel, Hao Hu, Michael Lachmann, Rong Zeng, Wei Chen, and Philipp Khaitovich. MicroRNA, mRNA, and protein expression link development and aging in human and macaque brain. *Genome Research*, 20(9):1207–1218, 2010. doi:10.1101/gr.106849.110.
- [6] Miguel Pérez-Enciso and Michel Tenenhaus. Prediction of clinical outcome with microarray data: a partial least squares discriminant analysis (PLS-DA) approach. *Human genetics*, 112:581–592, 2003. doi:10.1007/s00439-003-0921-9, PMID:12607117.
- [7] Jianguo Xia and David S Wishart. Web-based inference of biological patterns, functions and pathways from metabolomic data using MetaboAnalyst. *Nature protocols*, 6:743–760, 2011. doi:10.1038/nprot.2011.319, PMID:21637195.