

Bioconductor overview

Mikhail Dozmorov
Fall 2016

Bioconductor packages

- Bioconductor software consists of R add-on packages.
- An R package is a structured collection of code (R, C, or other), documentation, and/or data for performing specific types of analyses.
- E.g. **affy**, **limma**, **sva** packages provide implementations of specialized statistical and graphical methods.

Bioconductor Project

- The Bioconductor project started in 2001

Goal: make it easier to conduct reproducible consistent analysis of data from new high-throughput biological technologies

- Core maintainers of the Bioconductor website located at Fred Hutchinson Cancer Research Center
- Updated version released biannually coinciding with the release of R
- Many contributed software packages

Goals of the Bioconductor Project

- Provide access to statistical and graphical tools for analysis of high-dimensional biological data
1. Microarray analysis
 2. High-throughput 'omics' data

Goals of the Bioconductor Project

- Include comprehensive **documentation** describing and providing examples for packages
- Packages have associated **vignettes** that provide examples of how to use functions
- Have additional tools to work with publically available databases and other meta-data

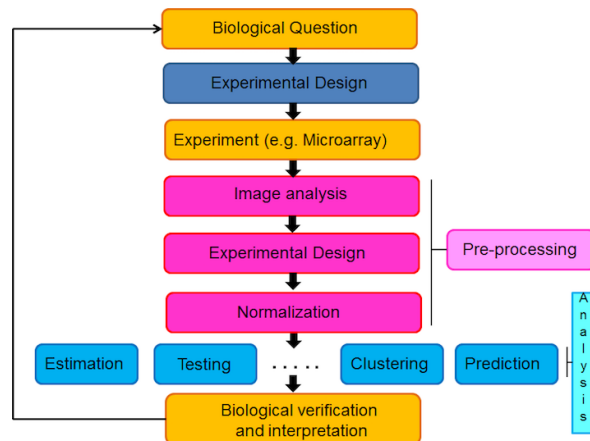
5/15

Vignettes

- Bioconductor has adopted a new documentation paradigm, the vignette.
- A vignette is an executable document consisting of a collection of documentation text and code chunks.
- Vignettes form dynamic, integrated, and reproducible statistical documents that can be automatically updated if either data or analyses are changed.
- Vignettes can be generated using the `sweave` function (or, `roxygen2` package)

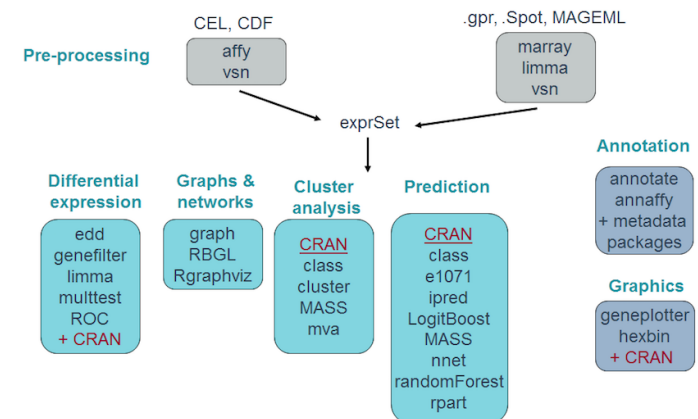
6/15

Microarray data analysis



7/15

Microarray data analysis



8/15

Bioconductor website

Lets take a look at the website...

<http://bioconductor.org/>

marrayRaw class

Pre-normalization intensity data for a batch of arrays

<code>maRf</code>	<code>maGf</code>	Matrix of red and green foreground intensities
<code>maRb</code>	<code>maGb</code>	Matrix of red and green background intensities
<code>maW</code>		Matrix of spot quality weights
<code>maLayout</code>		Array layout parameters - <code>marrayLayout</code>
<code>maGnames</code>		Description of spotted probe sequences - <code>marrayInfo</code>
<code>maTargets</code>		Description of target samples - <code>marrayInfo</code>
<code>maNotes</code>		Any notes

9/15

10/15

AffyBatch class

Probe-level intensity data for a batch of arrays (same CDF)

<code>cdfName</code>	Name of CDF file for arrays in the batch	
<code>nrow</code>	<code>ncol</code>	Dimensions of the array
<code>exprs</code>	<code>se.exprs</code>	Matrices of probe-level intensities and SEs rows → probe cells, columns → arrays.
<code>phenoData</code>	Sample level covariates, instance of class <code>phenoData</code>	
<code>annotation</code>	Name of annotation data	
<code>description</code>	MIAME information	
<code>notes</code>	Any notes	

ExpressionSet class

Processed Affymetrix or spotted array data

<code>exprs</code>	Matrix of expression measures, genes x samples
<code>se.exprs</code>	Matrix of SEs for expression measures, genes x samples
<code>phenoData</code>	Sample level covariates, instance of class <code>phenoData</code>
<code>annotation</code>	Name of annotation data
<code>description</code>	MIAME information
<code>notes</code>	Any notes

11/15

12/15

MIAME

Minimum Information About a Microarray Experiment (MIAME)

<http://fged.org/projects/miame/>

The six most critical elements contributing towards MIAME are:

- 1 - The raw data for each hybridization (e.g., CEL or GPR files)
- 2 - The final processed (normalized) data for the set of hybridizations in the experiment (study) (e.g., the gene expression data matrix used to draw the conclusions from the study)
- 3 - The essential sample annotation including experimental factors and their values (e.g., compound and dose in a dose response experiment)

13/15

MIAME

Minimum Information About a Microarray Experiment (MIAME)

<http://fged.org/projects/miame/>

- 4 - The experimental design including sample data relationships (e.g., which raw data file relates to which sample, which hybridizations are technical, which are biological replicates)
- 5 - Sufficient annotation of the array (e.g., gene identifiers, genomic coordinates, probe oligonucleotide sequences or reference commercial array catalog number)
- 6 - The essential laboratory and data processing protocols (e.g., what normalization method has been used to obtain the final processed data)

14/15

Pre-processing packages

1. `marray`: Spotted DNA microarrays.
 2. `affy`: Affymetrix oligonucleotide chips.
 3. `limma`: all, from spotted arrays to Affy to RNA-seq
- Reading in intensity data, diagnostic plots, normalization, computation of expression measures.
 - The packages start with very different data structures, but produce similar objects of class `ExpressionSet`.
 - One can then use other Bioconductor and CRAN packages for exploratory data analysis and visualization, differential expression detection

15/15