

## Data and text mining

# hypeR: an R package for geneset enrichment workflows

Anthony Federico<sup>1,2,\*</sup> and Stefano Monti<sup>1,2,\*</sup>

<sup>1</sup>Bioinformatics Department, Boston University, Boston, MA 02118, USA and <sup>2</sup>Section of Computational Biomedicine, Boston University School of Medicine, Boston, MA 02118, USA

\*To whom correspondence should be addressed.

Associate Editor: Jonathan Wren

Received on June 7, 2019; revised on August 16, 2019; editorial decision on September 2, 2019; accepted on September 4, 2019

## Abstract

**Summary:** Geneset enrichment is a popular method for annotating high-throughput sequencing data. Existing tools fall short in providing the flexibility to tackle the varied challenges researchers face in such analyses, particularly when analyzing many signatures across multiple experiments. We present a comprehensive R package for geneset enrichment workflows that offers multiple enrichment, visualization, and sharing methods in addition to novel features such as hierarchical geneset analysis and built-in markdown reporting. hypeR is a one-stop solution to performing geneset enrichment for a wide audience and range of use cases.

**Availability and implementation:** The most recent version of the package is available at <https://github.com/montilab/hypeR>.

**Contact:** [anfed@bu.edu](mailto:anfed@bu.edu) or [smonti@bu.edu](mailto:smonti@bu.edu)

## 1 Introduction

Geneset enrichment is an important step in biological data analysis workflows, particularly in bioinformatics and computational biology. At a basic level, one is performing a hypergeometric or Kolmogorov-Smirnov test to determine if a group of genes is over-represented or enriched, respectively, in pre-defined sets of genes, which suggests some biological relevance. The R package hypeR brings a fresh take to geneset enrichment, focusing on the analysis, visualization and reporting of enriched genesets. Although similar tools exist—such as Enrichr (Kuleshov *et al.*, 2016), fgsea (Sergushichev, 2016) and clusterProfiler (Yu *et al.*, 2012), among others—hypeR excels in the downstream analysis of geneset enrichment workflows—in addition to sometimes overlooked upstream analysis methods such as allowing for a flexible background population size or reducing genesets to a background distribution of genes. Finding relevant biological meaning from a large number of often obscurely labeled genesets may be challenging for researchers. hypeR overcomes this barrier by incorporating hierarchical ontologies—also referred to as relational genesets—into its workflows, allowing researchers to visualize and summarize their data at varying levels of biological resolution. All analysis methods are compatible with hypeR's markdown features, enabling concise and reproducible reports easily shareable with collaborators. Additionally, users can import custom genesets that are easily defined, extending the analysis of genes to other areas of interest such as proteins, microbes, metabolites etc. The hypeR package goes beyond performing basic enrichment, by providing a suite of

methods designed to make routine geneset enrichment seamless for scientists working in R.

## 2 Implementation

hypeR is a Bioconductor package written completely in R. The core function hypeR() accepts one or more signatures and a list of genesets to test for either over-representation or enrichment, depending on whether the signature is a vector of unranked genes, or a ranked vector of genes with or without weights. The former case is applicable to clusters of genes, such as those identified through co-expression analysis, while the latter is useful when signatures of genes can be ranked, such as through differential expression analysis. Despite its flexibility, hypeR() always returns one or more hyp objects that are defined using R6 (Mailund, 2017), which is an implementation of encapsulated object-oriented programming for R. A hyp object contains all information relevant to the enrichment analysis, including a data frame of results, enrichment plots for each geneset tested, as well as the arguments used to perform the analysis. All downstream functions used for analysis, visualization and reporting recognize hyp objects and utilize their data. Adopting an object-oriented framework brings modularity to hypeR, enabling flexible workflows. Additionally, most of hypeR's functionalities are applicable after enrichment results have been calculated. Therefore, users can perform enrichment with other popular tools, and use hypeR to analyze the results, by formatting the output into a hyp object. As an example, the documentation includes a tutorial

