# shinyGEO: a web application for analyzing Gene Expression Omnibus (GEO) datasets using shiny

## Jasmine Dumas[1], Michael A. Gargano[2], Garrett M. Dancik[2]

[1]DePaul University, Chicago, Illinois, USA, 2015 (Johns Hopkins Engineering For Professionals, 2016)
[2]Department of Mathematics and Computer Science, Eastern Connecticut State University, Willimantic, CT, USA

## Introduction

Identifying associations between patient gene expression profiles and clinical data provides insight into the biological processes associated with health and disease. The Gene Expression Omnibus (GEO) is a public repository of gene expression and sequence-based datasets, and currently includes >42,000 datasets with gene expression profiles obtained by microarray. Although GEO has its own analysis tool (GEO2R) for identifying differentially expressed genes, the tool is not designed for survival analysis and does not generate publication-ready graphics. In this work, we describe a web-based, easy-to-use tool for biomarker analysis in GEO datasets, called *shinyGEO*.

## Methods

The tool is developed using *shiny*, a web application framework for R. Specifically, *shinyGEO* allows a user to download the expression and clinical data from a GEO dataset, to modify the dataset correcting for spelling and misaligned data frame columns, to select a gene of interest, and to perform a survival or differential expression analysis using the available data. The tool uses the Bioconductor package *GEOquery* to retrieve the GEO dataset, while survival and differential expression analyses are carried out using the *survival* and *stats* packages, respectively. For both analyses, *shinyGEO* produces publication-ready graphics using *ggplot2* and generates the corresponding R code to ensure that all analyses are reproducible.
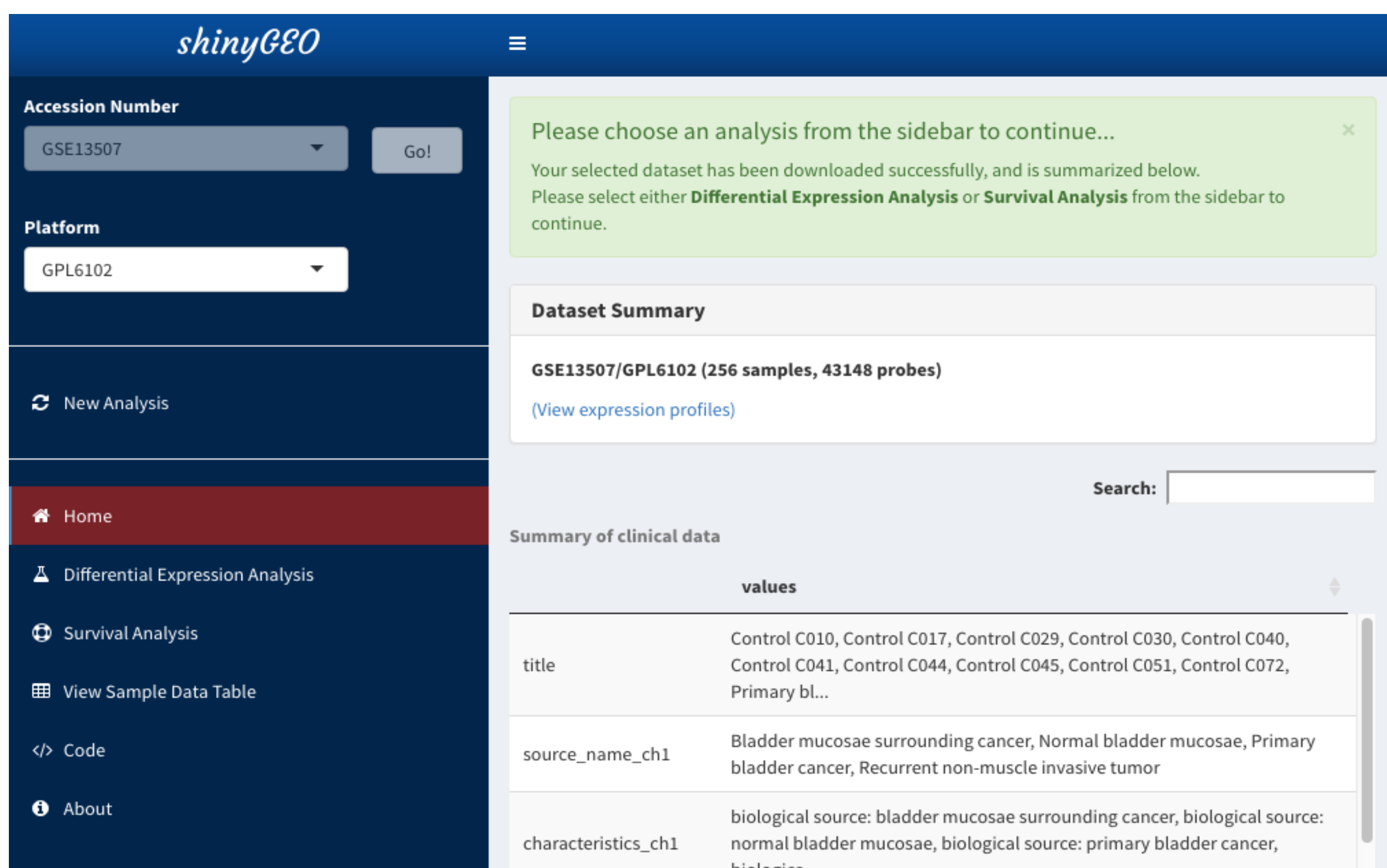
## Results



**Figure 1. Screenshot of *shinyGEO* following dataset selection.** The bladder cancer dataset GSE13507 has been selected and downloaded. From the homepage, the main panel summarizes the dataset. From the sidebar, the user can choose to carry out a differential expression analysis or survival analysis, view the sample data table, view the *R* code used to generate the analyses, or get more information about *shinyGEO*.
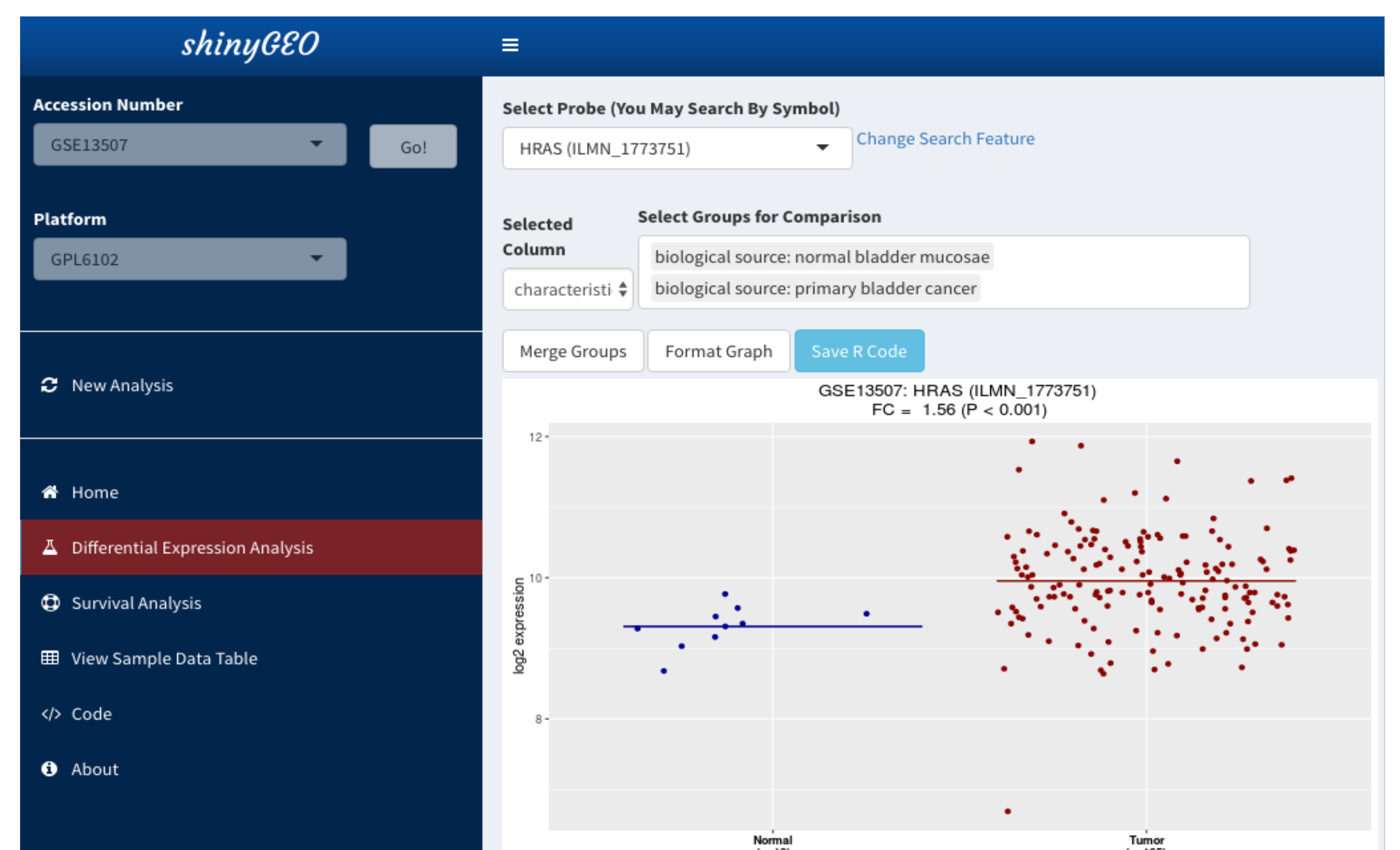


**Figure 2. Screenshot of *shinyGEO* following differential expression analysis.** Results indicate that the oncogene HRAS is significantly up-regulated in bladder tumors (FC = 1.56, P < 0.001) in this cohort. For differential expression analyses, the user selects the probe and the groups to compare. Optionally, the user can format the graph and save the *R* code used to generate the analysis.



**Figure 3. Screenshot of *shinyGEO* following survival analysis.** In this analysis, we select patients with muscle-invasive tumors who did not receive chemotherapy (inset). Results indicate that high TERT expression is significantly associated with poor disease- specific survival (DSS) in these patients (HR = 3.61, P = 0.0118), reproducing the analysis of Borah et al. (Science 347, 1006–1010 (2015)). For survival analyses, the user selects the probe and columns containing time and outcome information.
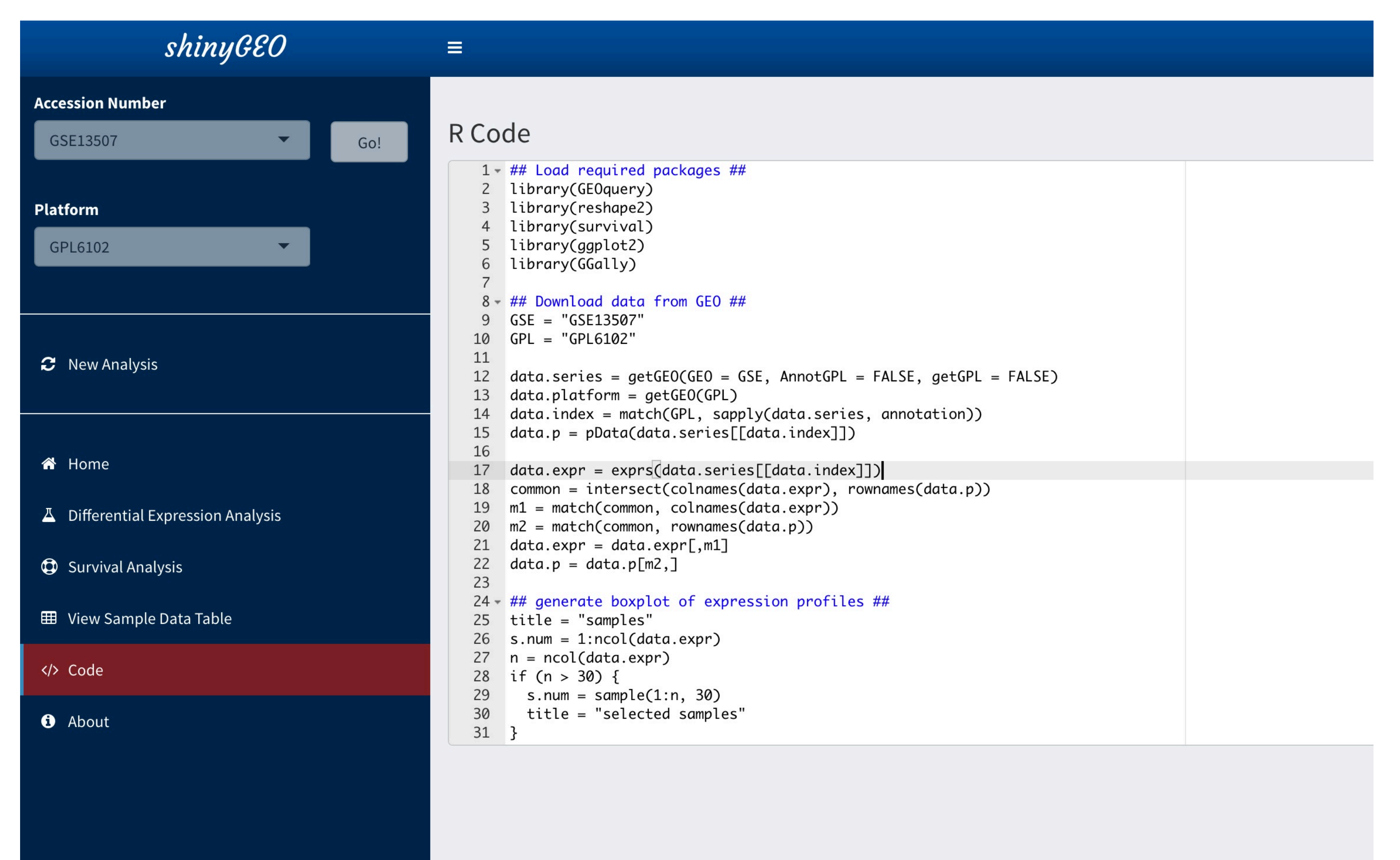


**Figure 4. Screenshot of *shinyGEO* following view of reproducbile *R* code.** R code is automatically generated following dataset download. The user can optionally save *R* code following differential expression or survival analysis.

**Availability:** http://gdancik.github.io/shinyGEO/