NJIT
New Jersey Institute
of Technology



# CS636 Data Analytics with R Program

# Project 1

By
Group 6

Duyen Ngyen | Bicheng Xiao | Shubham Gulia

# Scope of the project

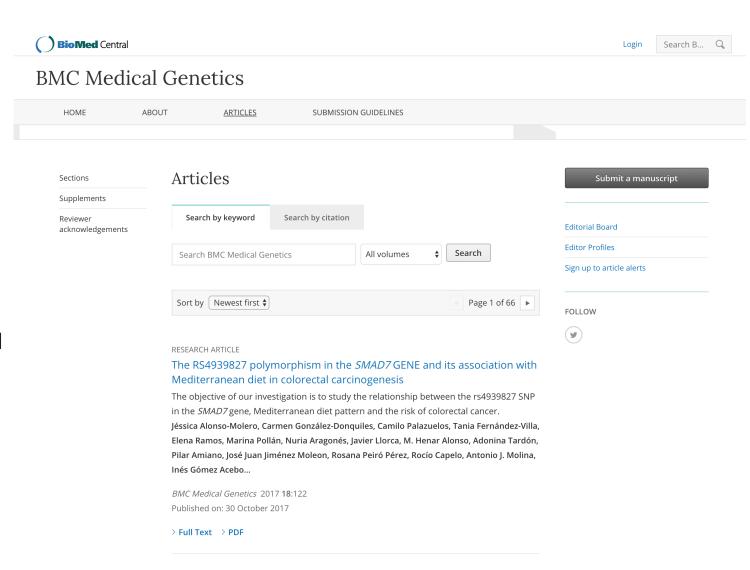**Use R program to extract required information from specified website**

Main page:

https://bmcmedgenet.biomedcentral.com

Total articles: 1642

Required information:

10 fields:

DOI, Title, Authors, Author Affiliations, Corresponding Author, Corresponding Author's Email, Publication Date, Abstract, Keywords, Full Text (Textual format)

---

BioMed Central

Login    Search B...

# BMC Medical Genetics

HOME    ABOUT    ARTICLES    SUBMISSION GUIDELINES

Sections

Supplements

Reviewer acknowledgements

## Articles

Submit a manuscript

Editorial Board

Editor Profiles

Sign up to article alerts

| Search by keyword | Search by citation |

Search BMC Medical Genetics    All volumes    Search

FOLLOW

Sort by  Newest first          Page 1 of 66 ▶

RESEARCH ARTICLE

**The RS4939827 polymorphism in the *SMAD7* GENE and its association with Mediterranean diet in colorectal carcinogenesis**

The objective of our investigation is to study the relationship between the rs4939827 SNP in the *SMAD7* gene, Mediterranean diet pattern and the risk of colorectal cancer.

Jéssica Alonso-Molero, Carmen González-Donquiles, Camilo Palazuelos, Tania Fernández-Villa, Elena Ramos, Marina Pollán, Nuria Aragonés, Javier Llorca, M. Henar Alonso, Adonina Tardón, Pilar Amiano, José Juan Jiménez Moleon, Rosana Peiró Pérez, Rocío Capelo, Antonio J. Molina, Inés Gómez Acebo...

*BMC Medical Genetics* 2017 **18**:122

Published on: 30 October 2017

> Full Text   > PDF

# Contribution of group members

| | | |
|---|---|---|
| Bicheng Xiao | main program: (loadAllArticles.R)<br>Functions(util.R):<br>• analysisArticle<br>• extractAuthors<br>• extractAffliation | Testing and debugging<br>Project report(PDF file) |
| Duyen Ngyen | Function(util.R):<br>• loadArticleList | Testing and debugging<br>Readme.txt |
| Shubham Gulia | Function(util.R):<br>• extract<br>• extracAttribute | Testing and debugging<br>ReadTheResult.R |

## Challenges

| unfamiliar with R packages, like XML | Check the online document and samples |
|---|---|
| extract author, corresponding author, corresponding author's email | • Carefully check the HTML of the article page and find the xpath<br>• Use a 2nd time extraction to get more detailed or related data |

**Thank you**