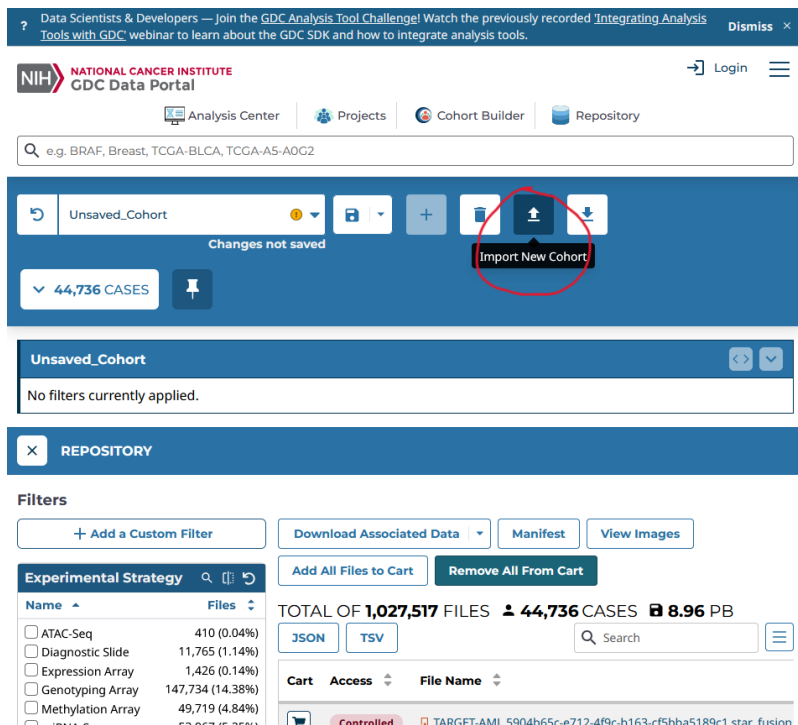# README

Jose Angel Sanchez Gomez
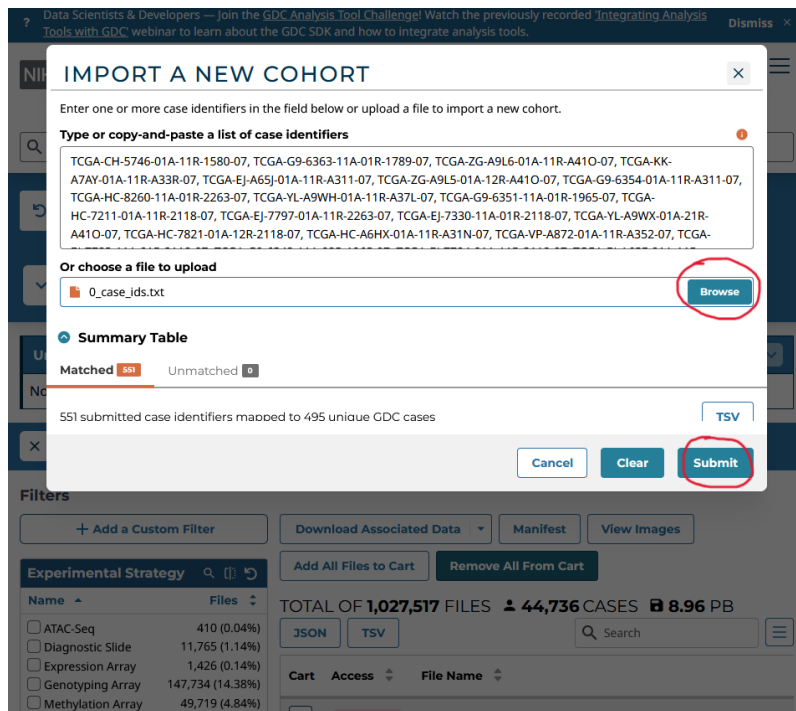
2024-10-01

INSTRUCTIONS FOR RECOVERING CLINICAL DATA:
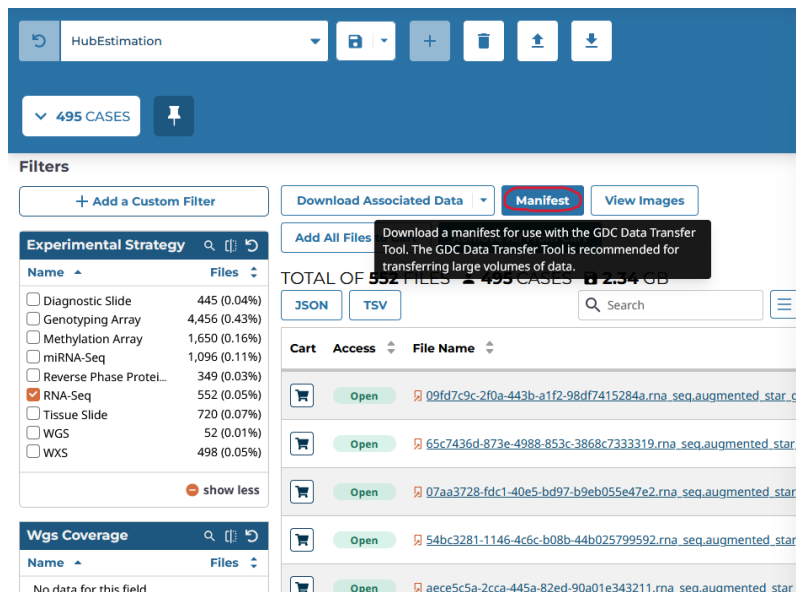
0) Download the github repository "HubEstimationCodeSubmission" into your local environment.

1) Access the TCGA online data repository: https://portal.gdc.cancer.gov/analysis_page?app=Downloads

2) Select the option "Import New Cohort."



3) Select the option "Choose a file to upload" and browse through your directories to upload the file "./realdata_4/200_ProcessingClinicalData/200_case_ids.txt"

4) Click on "cases" which provides "additional cases details and features." Here, you will be provided the option to download clinical data files in ".tsv" format.



5) After downloading, you will have a zip file containing "clinical.tsv", "exposure.tsv", "family_history.tsv", "follow_up.tsv", and "pathology_detail.tsv". Extract them so these files are located in the directory "./realdata_4/200_ProcessingClinicalData/".

6) In order to use the data in the file "clinical.tsv" we first require some processing. To do this, load the file "clinical.tsv" into R using the RMarkdown "./realdata_4/200_ProcessingClinicalData/211_CleaningClinicalData.Rmd". This file provides explanations on how the data is being cleaned and pre-processed for further analysis.

7) The resulting processed clinical data obtained through the file "211_ImportingClinicalData.Rmd" is saved in the file "./realdata_4/200_ProcessingClinicalData/212_clinical_reduced.csv".