

COVID-19 Specimen Data  
Processing:  
manifest\_code.R

IMPORTANT: All manifest file names should follow the same formatting:  
-SourceLocation-,<YYYYMMDD>\_<n>.csv  
SourceLocation = Where the file is from, should correspond to the Manifest: submitters (CSR, Main, CSTP)  
YYYYMMDD = When the manifest and samples on the list were received  
n = 1, 2, 3, ..., n. If more than one manifest & batch of samples are received on the same day, this should differentiate the lists

Libraries Needed:  
library(dplyr)  
library(tidyverse)  
library(jsonlite)  
library(readr)  
library(ggplot2)

Overview: This code file puts in all manifest files received from sample sources. Data checks are completed, and additional columns are added.

Fill in starting\_path. This is the path from your own machine to Box/DropBox  
Ex: "C:\Users\jake@Box Sync"

Fill in Manifest folder paths (there should be a separate folder for each location that sends sample manifests)

chr\_manifest\_fp = "SampleMetadataOrganization\Manifest\CSR"  
main\_manifest\_fp = "SampleMetadataOrganization\Manifest\Main"  
cstp\_manifest\_fp = "SampleMetadataOrganization\Manifest\CSTP"

List of all the manifest folder paths (manifest\_folder\_list) is created. This is iterated over later to ensure all files from all locations are completed.

Fill in output\_LOC, the output location of the manifest file completion

"SampleMetadataOrganization\Manifests\ManifestComplete"

Created as empty dataframe first, then gets filled in later in our process.

manifest\_storage

Iterate through every folder in manifest\_folder\_list

Store the name of every .csv file in the folder in file\_list

IMPORTANT: This will need to be altered if manifests begin coming in a different file format (.xlsx, .xls, etc.)

Iterate through every file in file\_list

Read in the file as a dataframe (file\_all) with all columns as character data type

We're trying to avoid losing leading zeros in our subject\_ids (MIDs and/or UMIDs)

Any cells that contain only " " or "" are converted to NA, so that remove\_empty() works to take away any completely blank rows

This sometimes occurs with .csv files that have been opened/manipulated in Excel

Store the names of the columns of file, in matrix column\_name\_check

Acting under the assumption that the manifest files have been pre-screened for column names, true\_columns = "position", "sample\_id", "subject\_id", "coll\_date", "flag"

If true\_columns and column\_name\_check are identical

If they aren't identical, but there are still 5 columns

If they aren't identical and there are not 5 columns

The difference between true\_columns and column\_name\_check will be printed to the console, along with a warning and a file location

Code will STOP executing if this occurs

If any sample\_id's are NA

Code will STOP executing and print a warning to the console

If any subject\_id's are NA

Code will STOP executing and print a warning to the console

If any coll\_date's are NA

Code will STOP executing and print a warning to the console

The first 4 characters of the coll in row 1, column n of file\_in are stored in test\_date\_format

If the manifest file structure changes order, this will need to be altered. This should correspond with the first row entry of coll\_date

If test\_date\_format sums into NA when attempting to correct it into a number

This may need additions if more alternate date structures are found in the manifests

Assume the coll\_date column is in "YYYYMMDD" format and correct it into a date format (NA=NA/NA)

The file name is split on "." and the second item is kept as rec\_date, a character string with all leading and trailing whitespace removed

rec\_date is then reformatted into first four characters (characters 5 and 6) (characters 7 and 8) to form a date

rec\_date then gets assigned as the value for every row in a new column of file\_in, called received\_date

The file name is again split on "." and the first item is kept as rec\_source, a character string with all leading and trailing whitespace removed

rec\_source then gets assigned as the value for every row in a new column of file\_in, called received\_source

file\_in is now bound to manifest\_storage and kept as manifest\_storage to continue to collect new rows

Only unique, distinct rows are kept in manifest\_storage

At this point in the code, we have all our unique specimen receipt rows from all our sources, compiled into manifest\_storage. Following this, we're going to flag some rows that are potentially of concern, write out our final compiled manifest receipt file, and write out a report of what we've got for record-keeping

Store the number of unique combinations of sample\_id, subject\_id, and coll\_date from manifest\_storage in unique\_ids

If unique\_ids is not equal to the number of rows in manifest\_storage

If we assume that every row is a unique instance of sample\_id, subject\_id, and coll\_date, then we expect the row count of manifest\_storage to match unique\_ids. If they don't match, we need to identify the duplicates

Get the count of each combination of the three variables, and use those to identify the duplicate rows as duplicate\_ssc

Alter flag for these rows in manifest\_storage to mark them with "Duplicate Sample - Subject - Collection"

Create subject\_id length on manifest\_storage, to look at the number of characters in each subject\_id

Current assumptions are that samples from CSR should be 9 characters long (as MIDs) and samples from CSTP should be 8 characters long (as UMIDs). If more information is learned about the samples, additional checks should be added for these next steps.

Add "MID < 9 digits + leading 0s restored" to flag where received\_source is "CSR" and subject\_id length is less than 9

Leading zeros are added to subject\_id for those marked rows

Add "UMID < 8 digits + leading 0s restored" to flag where received\_source is "CSTP" and subject\_id length is less than 8

Leading zeros are added to subject\_id for those marked rows

manifest\_storage is written out as a .csv at the output\_LOC as sample\_id\_manifest\_list.csv

manifest\_output\_report\_template.xlsx is read in from output\_LOC

The date, number of columns of manifest\_storage, and number of rows of manifest\_storage are added to the SUMMARY tab. If there are duplicate rows, then duplicate\_ssc is written out to the SUMMARY tab. Any instances of added zeros to MIDs or UMIDs are written out to the RESTORE\_ZEROS tab.

The populated report is then written out to output\_LOC as manifest\_output\_report\_YYYYMMDD.xlsx