

Public Opinion Question Database Documentation

December 20, 2019

Table of Contents

Project Description.....	1
Summary.....	2
Data	2
Universe	2
Citation	2
Funding	2
Basic Concepts	2
User Guide.....	3
01 Get All File Info.....	3
02 Load and Parse.....	3
03 Import Datasets	4
04 Copy Data to Dropbox And Merge	4
05 Question Identifications	4
06 Wording Flags.....	5
07 Abort Classes.....	5
07 Abort Classes B	6
08 Demog Classes.....	6
09 Demog Imputations	6
10 Combine Datasets	6
11 Standardize Metrics.....	7
12 ClassifyTyper	7

Project Description

The purpose of this Github repository is to allow researchers to merge numerous public opinion datasets from a variety of sources into a single dataset. The repository houses all of the R code needed to load, parse, import, and merge files, as well as to identify questions and analyze the data. The repository serves as an example using abortion, but could be applied to other projects as well.

Summary

The files in the Github repository originated as part of a larger study on abortion attitudes and opinions conducted from 2017-2020. A portion of the project, conducted at the University of Michigan, involved creating a single file that comprehensively archived survey questions asked on abortion attitudes over time. For this part of the project, data were collected from representative surveys or polls that asked a question about abortion. The research team combed data repositories for all publicly available data, relied on restricted use data, and included proprietary data where it could be made available. Datasets were in a variety of formats, including fwf, dat, data, sav, por, portable, rdata, tab, rda, sps, and dta. This Github repository archives the code associated with combining and organizing the datasets.

Data

The data associated with this project are restricted access; data were gathered from collections and sources that require membership and can only be made available directly from the organizations responsible for archiving the information. Data come from a variety of archives including Dataverse, the Roper Center, ICPSR, RAND, Quinnipiac, and the Pew Center, and sources such as GSS, ANES, ABC News-Washington Post Polls, etc.

Universe

1,837 studies with relevant abortion attitude questions consisting of data from 9.3 million responses to 5,559 abortion attitude questions.

Citation

Pasek, J. and Lippman, J. R. 2019 Public Opinion on Abortion Question Archive. Ann Arbor, MI, GitHub repository.

—NOTE: FOR DOI <https://help.github.com/en/github/creating-cloning-and-archiving-repositories/referencing-and-citing-content>

Funding

Funding for the project was provided by a foundation grant.

Basic Concepts

Data Directory

Top level data directory for all surveys. Includes folders from each survey organization which contain every survey that the organization has run on the topic of interest. Within each individual survey folder, there are typically several files—often the data files in several formats and related survey documentation.

Keyword Files

Metadata file in csv format providing study information at the dataset level. Provided by Roper.

Parser Files

Files manually created to read dat format from Roper. Question-level CSV file that lists question by topic, question wording and response categories, and contains other metadata, such as survey begin and end dates, N, sample type, survey organization, etc.

Topic Files

Metadata file in csv format providing information at the question level. Provided by Roper. Catalogues each question by topic, and contains other metadata, such as survey begin and end dates, N, survey organization, survey sponsor, sample type, etc.

Topline File

For each survey question, percentage of respondents selecting each answer category. Usually csv, pdf, txt, rtf, doc, docx. Roper topline files are csv files, usually called “report.csv”. Roper is the only source of topline files in the abortion study. Roper csv topline files are set up so that each case is a response category with variables related to the question and the percentage of respondents selecting the response.

User Guide

This section provides an overview of all of the R files in the repository. To use the R files, the user begins with a top level data directory that includes folders for each survey containing several files, often with data files in several formats and related survey documentation.

01 Get All File Info

This R file (01_GetAllFileInfo.r) identifies all the files located within the data directory to be included in the user’s project and gives the user information about what is being imported. Specifically, this R routine produces a csv and rdata file containing information about all of the files located within the data directory, looking within each survey folder and identifying files by extension type. The routine searches the data directory for data except for “donotread.” It deals with duplicate names by assuming that they are from different filing schemes. This information is then saved to a folder called “Summary Information” as a csv file (StudyFiles.csv) and also produces an rdata file (located in EachFileExports/01_FileInformation.rdata). It can, for example, tell the user how many csv, doc, dta, or other file types are located in the data directory.

02 Load and Parse

The purpose of the 02_LoadAndParse.r file is to catalogue the data by producing a single csv and rdata file that contain the metadata from each survey or topline. This includes study ID, start and end date for the study, country where data collection took place, study

sponsor, survey organization, type of sample, N, mode, notes, study title, and other information. For the abortion study, topic files and topline files were used that were created by Roper. These types of files are typically available for any study, and can be exported. For the abortion study, some metadata for surveys from other sources were manually recorded, and these are also imported at this step. The routine also produces error files for topline files and parsers, indicating duplicate question numbers, missing metadata, parsers or topline files that cannot be used or have parsing failures. Error files list the studies or topline files with problems, and these can be investigated for correction or follow-up.

03 Import Datasets

The 03_ImportDatasets.r file imports all of the data from data files including rdata, sps, spss and portable (por) files, stata, and tab and saves them collectively based on file type. It also parses and then imports all data from files originally from dat format (i.e., the Roper files), linking the parser files with the dat files and exporting the numeric data. For all of the file formats, errors are flagged where necessary, and troubleshooting files are produced in csv format that list all of the files that cannot be imported, categorized by problem type. Note that the parser files will work on any kind of fixed-width data.

04 Copy Data to Dropbox And Merge

This file (04_CopyDataToDropboxAndMerge.r) builds an rdata file that provides the labels for the question wording and response options for the datasets from the previous step. It does this by first extracting and adding the wording and response options using the parser files and adding them to the previously parsed datasets, and provides error files that list datasets where the number of response options in the parser file and the number of values in the data do not match. It also lists errors from the process of building the parsed datasets in a csv file. Next, this R file then combines all the datasets that were imported using the previous file (i.e., 03_ImportDatasets.r), listing errors in csv files that deal with duplicate level names, corrupt data imports, and empty data frames. Finally, it identifies question wording and response options for all of the variables based on the labels provided in the original sources of data.

05 Question Identifications

The Question Identification R file (05_QuestionIdentifications.r) adds the question wording and response options from the topline files to the rdata file in the previous step so that there is a single, comprehensive rdata file of question wording and response options. It does this primarily by extracting the question wording and response options from all of the topline files. The topline files, from Roper, are arranged such that each response category is a separate case. Only topline files in csv format are imported. It also produces a csv troubleshooting file that lists all of the topline files where there are duplicate question names so that the user can go back and fix these in the original topline files.

06 Wording Flags

The purpose of this routine is to find the variables and recode and relabel them. The R file goes through all of the question and response options and converts as many of the labels as possible to a single standard. To do this, it first calls 99_TranslateResponses.r which orders demographic variables and cleans label abbreviations and misspellings. It looks at the variables sex, race, hispanic origin, state, region, education, religion (denomination, fundamentalism, attendance frequency), weight, partisanship, party strength, presidential approval, age, marital status, and—for the abortion study—abortion questions.

The routine writes three types of csv files for each variable that is being incorporated into the final dataset: CheckTranslate files, CheckWords files, and CheckDrops files. The check translate files show how well the routine did at relabeling the answer categories. The user can examine any of the csv files in CheckTranslate to compare the original categories with the relabeled answer categories. The column in this csv file, labeled “unmatched” shows any unusual labels that can’t be relabeled by the routine. The “n” column shows how many times the labeling scheme from the “original” column appears in the datasets. The csv files in the CheckDrops folder show variables that have been dropped from the final dataset because they are unmatched after translation; since they can’t be appropriately recoded or relabeled they are dropped.

The routine produces an rdata file with cleaned question and response options located at EachFileExports/06_SelectedQuestionInfo.rdata. Please note that the 06_WordingFlags.r file takes a long time to run.

—writes “CheckWords/RaceWords.csv” question wording - how often each type appears?

07 Abort Classes

07_AbortClasses.r pulls in the raw question-level data from the previous step (06 Wording Flags) and categorizes abortion-relevant questions by question and response type. It excludes questions on abortion that are not relevant to abortion attitudes (e.g., recognizing the position of candidate’s positions on abortion). For the remaining questions, it merges them with the metadata from the second r file (i.e., 02_MetaInformation.rdata) so that more information can be appended to the question categorizations.

The routine then writes three types of item banks as csv files: Item Bank by Wordings, Item Banks by Question Subtype, and Item Bank By Question Type. The first of these lists every question asked on abortion with the full question wording, the responses, the number of times asked, when it was last asked, the total N, the survey organization, and sponsor, the data type, and provides the classification of the question by type and subtype (e.g., Attitude Strength and Importance). The latter two files collapse the Item Bank by Wording file into question subtype and type. Overall, there are 9 subtypes (Importance, Access, Maternal Circumstances, Morality, Pregnancy Circumstances, Pregnancy Term, Regulation, Pro-Choice/Pro-Life, and NA). The Item Bank by Question Type tells us, for example, that Importance questions were asked 1,342 times of 1,786,826 respondents between 1955 and

2018 (when the datasets were incorporated into the project), and provides us with further metadata related to the survey organizations and sponsors that asked that type of question.

The routine writes 07-01_UncategorizedAbortionQs.csv, a troubleshooting file, to list the questions (and number of times they were asked) that appear in the NA category of the Item Bank files. Finally, the routine writes the rdata file for the categorizations, EachFileExports/07-01_AbortionVariablesCategorized.rdata.

07 Abort Classes B

07_AbortClassesB.r replicates this process using an alternate coding scheme.

08 Demog Classes

08_DemogClasses.r pulls in the cleaned questions and response options from the demographic variables in step 06 Wording Flags, and categorizes variables from that step by type. It also further cleans the demographic variables.

09 Demog Imputations

The file 09_DemogImputations.r imports the abortion variables as individual level data. It puts the responses from factor variables in order based on their ordinality, using the first factor from principal components analysis to order the responses. It also builds several study level datasets. Finally, the routine uses Amelia to multiply impute missing data.

10 Combine Datasets

The purpose of 10_CombineDatasets.r is to build a dataset where each case is an answer to an abortion attitude question. Each case contains “respID” (the respondent identifier), the study ID (“StudyID”), and a “UniqueID” for the abortion attitude question. In this construction of the dataset, each person may be represented more than once (for example, if they answered more than one question about abortion attitudes).

The goal of this file is to identify sets of response options that can be coded in an ordinal manner. To do so, we consider all response options that were used by at least 5% of (weighted) respondents. If more than 20% of respondents chose low-probability responses to a question, we did not code ordinal responses to that question. The sets of response options chosen at least this frequently are exported into two files: one for measures with three or more such response options and one for measures with only two response options. Unique response option sets for questions with 3+ response options were then exported and manually reordered to identify ordinality. They were then manually divided into those where the direction of the ordinality was known (e.g. always illegal to always legal) that could serve as a training set and those where the direction of the ordinality was unknown (e.g. strongly favor to strongly oppose) where directionality would need to be determined by the data. Items that did not have a clear ordering of response options were excluded as part of this step. Direction was not denoted for items with two categories as these could be assigned using the training set.

It creates csv files listing the uncategorized response options and the untrained items, and writes an rdata file “10_IndividualMasterPreReversals.rdata.”

11 Standardize Metrics

The 11_StandardizeMetrics.r file takes the ordinal directional questions and strength questions (from the training dataset) recodes them from 0 to 1 and standardizes them. Next, it runs a prediction on the standardized data and imputes demographics for cases. Then, it predicts probabilities with the imputed demographics on the standardized versions in the training set (separately for measures classified as either directional attitudes or those classified as importance-based). After that, it takes the predicted probabilities for the trained and untrained data, assesses which variables are likely to be directional vs. importance measures, and determines which ordinal variables need to be reverse-coded. To do this, it correlates the predicted probability from the training data (based on a regression including time interacted with imputed demographics) with the observed values within each variable to see whether they are positive or negative with the predicted probability. Negative correlations are then reversed. This process is repeated, reversing items with slightly smaller (in magnitude) correlations for a few iterations. The routine also produces a figure: “CorrelationsBetweenVariableTypes.pdf” that shows these correlations. Finally, it writes a file, “10-02_FullDatasetWithReversals.rdata” in which all ordinal variables are coded in the correct ordinal direction and each case is an individual’s abortion measure answer.

12 ClassifyTyper

The 12_ClassifyTyper.r uses the microdata to produce estimates of the means and standardized means for each question within each demographic group of interest. This export into the file “12-01_DemographicsByDataset.rdata,” which has one row per question asked and columns for each outcome-demographic combination. Outcomes include: “Mean01_” files which are the average value when coded from 0 to 1; “AdjMn01_” files, which shift the Mean01 values based on how the global mean for a question compares to what might be expected on that date for a measure of that type; “MeanStdz_” files, which show the group’s mean on a standardized scale; and “Ns01_” files, which show the number of respondents in each group.