

Package ‘wru’

July 22, 2025

Title Who are You? Bayesian Prediction of Racial Category Using Surname, First Name, Middle Name, and Geolocation

Version 3.0.3

Date 2024-05-24

Description Predicts individual race/ethnicity using surname, first name, middle name, geolocation, and other attributes, such as gender and age. The method utilizes Bayes' Rule (with optional measurement error correction) to compute the posterior probability of each racial category for any given individual. The package implements methods described in Imai and Khanna (2016) ``Improving Ecological Inference by Predicting Individual Ethnicity from Voter Registration Records" Political Analysis <DOI:10.1093/pan/mpw001> and Imai, Olivella, and Rosenman (2022) ``Addressing census data problems in race imputation via fully Bayesian Improved Surname Geocoding and name supplements" <DOI:10.1126/sciadv.adc9824>. The package also incorporates the data described in Rosenman, Olivella, and Imai (2023) ``Race and ethnicity data for first, middle, and surnames" <DOI:10.1038/s41597-023-02202-2>.

License GPL (>= 3)

URL <https://github.com/kosukeimai/wru>

BugReports <https://github.com/kosukeimai/wru/issues>

Depends R (>= 4.1.0), utils

Imports cli, dplyr, tidyr, furrr, future, piggyback (>= 0.1.4), PL94171, purrr, Rcpp, rlang

Suggests covr, testthat (>= 3.0.0), tidycensus

LinkingTo Rcpp, RcppArmadillo

Config/testthat/edition 3

Encoding UTF-8

LazyData yes

LazyDataCompression xz

LazyLoad yes

RoxygenNote 7.3.1
NeedsCompilation yes
Author Kabir Khanna [aut],
 Brandon Bertelsen [aut, cre],
 Santiago Olivella [aut],
 Evan Rosenman [aut],
 Alexander Rossell Hayes [aut],
 Kosuke Imai [aut]
Maintainer Brandon Bertelsen <brandon@bertelsen.ca>
Repository CRAN
Date/Publication 2024-05-24 18:00:02 UTC

Contents

format_legacy_data	2
get_census_data	3
predict_race	4
state_fips	7
surnames2000	8
surnames2010	9
voters	9
wru_data_preflight	10
Index	11

format_legacy_data	<i>Legacy data formatting function.</i>
--------------------	---

Description

format_legacy_data formats legacy data from the U.S. census to allow for Bayesian name geocoding.

Usage

```
format_legacy_data(legacyFilePath, state, outFile = NULL)
```

Arguments

- legacyFilePath A character vector giving the location of a legacy census data folder, sourced from https://www2.census.gov/programs-surveys/decennial/2020/data/01-Redistricting_File-PL_94-171/. These file names should end in ".pl".
- state The two letter state postal code.
- outFile Optional character vector determining whether the formatted RData object should be saved. The filepath should end in ".RData".

Details

This function allows users to construct datasets for analysis using the census legacy data format. These data are available for the 2020 census at https://www2.census.gov/programs-surveys/decennial/2020/data/01-Redistricting_File-PL_94-171/. This function returns data structured analogously to data from the Census API, which is not yet available for the 2020 Census as of September 2021.

Examples

```
## Not run:
gaCensusData <- format_legacy_data(PL94171::pl_url('GA', 2020))
predict_race_new(ga.voter.file, namesToUse = 'last, first, mid', census.geo = 'block',
  census.data = gaCensusData)

## End(Not run)
```

get_census_data	<i>Multilevel Census data download function.</i>
-----------------	--

Description

get_census_data returns county-, tract-, and block-level Census data for specified state(s). Using this function to download Census data in advance can save considerable time when running predict_race and census_helper.

Usage

```
get_census_data(
  key = Sys.getenv("CENSUS_API_KEY"),
  states,
  age = FALSE,
  sex = FALSE,
  year = "2020",
  census.geo = c("tract", "block", "block_group", "county", "place", "zcta"),
  retry = 3,
  county.list = NULL
)
```

Arguments

key	A character string containing a valid Census API key, which can be requested from the U.S. Census API key signup page . By default, attempts to find a census key stored in an environment variable named CENSUS_API_KEY.
states	which states to extract Census data for, e.g., c("NJ", "NY").

age	A TRUE/FALSE object indicating whether to condition on age or not. If FALSE (default), function will return $\text{Pr}(\text{Geolocation} \mid \text{Race})$. If TRUE, function will return $\text{Pr}(\text{Geolocation}, \text{Age} \mid \text{Race})$. If sex is also TRUE, function will return $\text{Pr}(\text{Geolocation}, \text{Age}, \text{Sex} \mid \text{Race})$.
sex	A TRUE/FALSE object indicating whether to condition on sex or not. If FALSE (default), function will return $\text{Pr}(\text{Geolocation} \mid \text{Race})$. If TRUE, function will return $\text{Pr}(\text{Geolocation}, \text{Sex} \mid \text{Race})$. If age is also TRUE, function will return $\text{Pr}(\text{Geolocation}, \text{Age}, \text{Sex} \mid \text{Race})$.
year	A character object specifying the year of U.S. Census data to be downloaded. Use "2010", or "2020". Default is "2020". Warning: 2020 U.S. Census data is downloaded only when age and sex are both FALSE.
census.geo	An optional character vector specifying what level of geography to use to merge in U.S. Census 2010 geographic data. Currently "county", "tract", "block", and "place" are supported.
retry	The number of retries at the census website if network interruption occurs.
county.list	A named list of character vectors of counties present in your <i>voter.file</i> , per state.

Value

Output will be an object of class `list` indexed by state. Output will contain a subset of the following elements: `state`, `age`, `sex`, `county`, `tract`, `block_group`, `block`, and `place`.

Examples

```
## Not run: get_census_data(states = c("NJ", "NY"), age = TRUE, sex = FALSE)
## Not run: get_census_data(states = "MN", age = FALSE, sex = FALSE, year = "2020")
```

predict_race	<i>Race prediction function.</i>
--------------	----------------------------------

Description

`predict_race` makes probabilistic estimates of individual-level race/ethnicity.

Usage

```
predict_race(
  voter.file,
  census.surname = TRUE,
  surname.only = FALSE,
  census.geo = c("tract", "block", "block_group", "county", "place", "zcta"),
  census.key = Sys.getenv("CENSUS_API_KEY"),
  census.data = NULL,
  age = FALSE,
  sex = FALSE,
  year = "2020",
```

```

party = NULL,
retry = 3,
impute.missing = TRUE,
skip_bad_geos = FALSE,
use.counties = FALSE,
model = "BISG",
race.init = NULL,
name.dictionaries = NULL,
names.to.use = "surname",
control = NULL
)

```

Arguments

- | | |
|----------------|---|
| voter.file | An object of class data.frame. Must contain a row for each individual being predicted, as well as a field named <i>surname</i> containing each individual's surname. If using geolocation in predictions, <i>voter.file</i> must contain a field named <i>state</i> , which contains the two-character abbreviation for each individual's state of residence (e.g., "nj" for New Jersey). If using Census geographic data in race/ethnicity predictions, <i>voter.file</i> must also contain at least one of the following fields: <i>county</i> , <i>tract</i> , <i>block_group</i> , <i>block</i> , and/or <i>place</i> . These fields should contain character strings matching U.S. Census categories. County is three characters (e.g., "031" not "31"), tract is six characters, block group is usually a single character and block is four characters. Place is five characters. See below for other optional fields. |
| census.surname | A TRUE/FALSE object. If TRUE, function will call <code>merge_surnames</code> to merge in $\text{Pr}(\text{Race} \mid \text{Surname})$ from U.S. Census Surname List (2000, 2010, or 2020) and Spanish Surname List. If FALSE, user must provide a <code>name.dictionary</code> (see below). Default is TRUE. |
| surname.only | A TRUE/FALSE object. If TRUE, race predictions will only use surname data and calculate $\text{Pr}(\text{Race} \mid \text{Surname})$. Default is FALSE. |
| census.geo | An optional character vector specifying what level of geography to use to merge in U.S. Census geographic data. Currently "county", "tract", "block_group", "block", and "place" are supported. Note: sufficient information must be in user-defined <i>voter.file</i> object. If <code>census.geo = "county"</code> , then <i>voter.file</i> must have column named <i>county</i> . If <code>census.geo = "tract"</code> , then <i>voter.file</i> must have columns named <i>county</i> and <i>tract</i> . And if <code>census.geo = "block"</code> , then <i>voter.file</i> must have columns named <i>county</i> , <i>tract</i> , and <i>block</i> . If <code>census.geo = "place"</code> , then <i>voter.file</i> must have column named <i>place</i> . If <code>census.geo = "zcta"</code> , then <i>voter.file</i> must have column named <i>zcta</i> . Specifying <code>census.geo</code> will call <code>census_helper</code> function to merge Census geographic data at specified level of geography. |
| census.key | A character object specifying user's Census API key. Required if <code>census.geo</code> is specified, because a valid Census API key is required to download Census geographic data.

If <code>NULL</code> , the default, attempts to find a census key stored in an environment variable named <code>CENSUS_API_KEY</code> . |

census.data	A list indexed by two-letter state abbreviations, which contains pre-saved Census geographic data. Can be generated using <code>get_census_data</code> function.
age	An optional TRUE/FALSE object specifying whether to condition race predictions on age (in addition to surname and geolocation). Default is FALSE. Must be same as <code>age</code> in <code>census.data</code> object. May only be set to TRUE if <code>census.geo</code> option is specified. If TRUE, <code>voter.file</code> should include a numerical variable <code>age</code> .
sex	optional TRUE/FALSE object specifying whether to condition race predictions on sex (in addition to surname and geolocation). Default is FALSE. Must be same as <code>sex</code> in <code>census.data</code> object. May only be set to TRUE if <code>census.geo</code> option is specified. If TRUE, <code>voter.file</code> should include a numerical variable <code>sex</code> , where <code>sex</code> is coded as 0 for males and 1 for females.
year	An optional character vector specifying the year of U.S. Census geographic data to be downloaded. Use "2010", or "2020". Default is "2020".
party	An optional character object specifying party registration field in <code>voter.file</code> , e.g., <code>party = "PartyReg"</code> . If specified, race/ethnicity predictions will be conditioned on individual's party registration (in addition to geolocation). Whatever the name of the party registration field in <code>voter.file</code> , it should be coded as 1 for Democrat, 2 for Republican, and 0 for Other.
retry	The number of retries at the census website if network interruption occurs.
impute.missing	Logical, defaults to TRUE. Should missing be imputed?
skip_bad_geos	Logical. Option to have the function skip any geolocations that are not present in the census data, returning a partial data set. Default is set to FALSE, in which case it will break and provide error message with a list of offending geolocations.
use.counties	A logical, defaulting to FALSE. Should census data be filtered by counties available in <code>census.data</code> ?
model	Character string, either "BISG" (default) or "fBISG" (for error-correction, fully-Bayesian model).
race.init	Vector of initial race for each observation in <code>voter.file</code> . Must be an integer vector, with 1=white, 2=black, 3=hispanic, 4=asian, and 5=other. Defaults to values obtained using <code>model="BISG_surname"</code> .
name.dictionaries	Optional named list of <code>data.frame</code> 's containing counts of names by race. Any of the following named elements are allowed: "surname", "first", "middle". When present, the objects must follow the same structure as <code>last_c</code> , <code>first_c</code> , <code>mid_c</code> , respectively.
names.to.use	One of 'surname', 'surname, first', or 'surname, first, middle'. Defaults to 'surname'.
control	List of control arguments only used when <code>model="fBISG"</code> , including <ul style="list-style-type: none"> iter Number of MCMC iterations. Defaults to 1000. burnin Number of iterations discarded as burnin. Defaults to half of <code>iter</code>. verbose Print progress information. Defaults to TRUE. me.correct Boolean. Should the model correct measurement error for races geo? Defaults to TRUE. seed RNG seed. If NULL, a seed is generated and returned as an attribute for reproducibility.

Details

This function implements the Bayesian race prediction methods outlined in Imai and Khanna (2015). The function produces probabilistic estimates of individual-level race/ethnicity, based on surname, geolocation, and party.

Value

Output will be an object of class `data.frame`. It will consist of the original user-input `voter.file` with additional columns with predicted probabilities for each of the five major racial categories: *pred.whi* for White, *pred.bla* for Black, *pred.his* for Hispanic/Latino, *pred.asi* for Asian/Pacific Islander, and *pred.oth* for Other/Mixed.

Examples

```
#' data(voters)
try(predict_race(voter.file = voters, surname.only = TRUE))
## Not run:
try(predict_race(voter.file = voters, census.geo = "tract"))

## End(Not run)
## Not run:
try(predict_race(
  voter.file = voters, census.geo = "place", year = "2020"))

## End(Not run)
## Not run:
CensusObj <- try(get_census_data(state = c("NY", "DC", "NJ")))
try(predict_race(
  voter.file = voters, census.geo = "tract", census.data = CensusObj, party = "PID")
)

## End(Not run)
## Not run:
CensusObj2 <- try(get_census_data(state = c("NY", "DC", "NJ"), age = T, sex = T))
try(predict_race(
  voter.file = voters, census.geo = "tract", census.data = CensusObj2, age = T, sex = T))

## End(Not run)
## Not run:
CensusObj3 <- try(get_census_data(state = c("NY", "DC", "NJ"), census.geo = "place"))
try(predict_race(voter.file = voters, census.geo = "place", census.data = CensusObj3))

## End(Not run)
```

Description

Dataset including FIPS codes and postal abbreviations for each U.S. state, district, and territory.

Usage

```
state_fips
```

Format

A tibble with 57 rows and 3 columns:

```
state  Two-letter postal abbreviation
state_code  Two-digit FIPS code
state_name  English name
```

Source

Derived from [tidycensus::fips_codes\(\)](#)

surnames2000	<i>Census Surname List (2000).</i>
--------------	------------------------------------

Description

Census Surname List from 2000 with race/ethnicity probabilities by surname.

Usage

```
surnames2000
```

Format

A data frame with 157,728 rows and 6 variables:

```
surname  Surname
p_whi    Pr(White | Surname)
p_bla    Pr(Black | Surname)
p_his    Pr(Hispanic/Latino | Surname)
p_asia   Pr(Asian/Pacific Islander | Surname)
p_othe   Pr(Other | Surname) #'
```

Examples

```
data(surnames2000)
```

surnames2010	<i>Census Surname List (2010).</i>
--------------	------------------------------------

Description

Census Surname List from 2010 with race/ethnicity probabilities by surname.

Usage

```
surnames2010
```

Format

A data frame with 167,613 rows and 6 variables:

surname Surname

p_whi Pr(White | Surname)

p_bla Pr(Black | Surname)

p_his Pr(Hispanic/Latino | Surname)

p_asl Pr(Asian/Pacific Islander | Surname)

p_oth Pr(Other | Surname) #'

Examples

```
data(surnames)
```

voters	<i>Example voter file.</i>
--------	----------------------------

Description

An example dataset containing voter file information.

Usage

```
voters
```

Format

A data frame with 10 rows and 12 variables:

VoterID Voter identifier (numeric)
surname Surname
state State of residence
CD Congressional district
county Census county (three-digit code)
first First name
last Last name or surname
tract Census tract (six-digit code)
block Census block (four-digit code)
precinct Voting precinct
place Voting place
age Age in years
sex 0=male, 1=female
party Party registration (character)
PID Party registration (numeric) #'

Examples

```
data(voters)
```

wru_data_preflight	<i>Preflight for name data</i>
--------------------	--------------------------------

Description

Checks if namedata is available in the current working directory, if not downloads it from github using piggyback. By default, wru will download the data to a temporary directory that lasts as long as your session does. However, you may wish to set the wru_data_wd option to save the downloaded data to your current working directory for more permanence.

Usage

```
wru_data_preflight()
```

Index

* **datasets**

- state_fips, [7](#)
- surnames2000, [8](#)
- surnames2010, [9](#)
- voters, [9](#)

environment variable, [3](#), [5](#)

format_legacy_data, [2](#)

get_census_data, [3](#)

NULL, [5](#)

predict_race, [4](#)

- state_fips, [7](#)
- surnames2000, [8](#)
- surnames2010, [9](#)

tidycensus::fips_codes(), [8](#)

voters, [9](#)

wru_data_preflight, [10](#)