covidcast_practice



Data Preparation

by: Joanne Kim, Michael La, Joshua Harasaki, Tyler Chia

```
!pip install covidcast
Collecting covidcast
   Downloading covidcast-0.1.5-py3-none-any.whl (12.3 MB)
                                           12.3 MB 8.3 MB/s
Requirement already satisfied: matplotlib in /shared-libs/python3.7/py/lib/python3.7/site-packages (from covidcast) (3.4.3)
Collecting geopandas
   Downloading geopandas-0.10.2-py2.py3-none-any.whl (1.0 MB)
                                                  1.0 MB 19.2 MB/s
Requirement already satisfied: numpy in /shared-libs/python3.7/py/lib/python3.7/site-packages (from covidcast) (1.19.5)
Collecting descartes
   Downloading descartes-1.1.0-py3-none-any.whl (5.8 kB)
Requirement already satisfied: requests in /shared-libs/python3.7/py/lib/python3.7/site-packages (from covidcast) (2.26.0)
Collecting imageio-ffmpeg
   Downloading imageio_ffmpeg-0.4.5-py3-none-manylinux2010_x86_64.whl (26.9 MB)
                                                    26.9 MB 29.3 MB/s
Requirement already satisfied: pandas in /shared-libs/python3.7/py/lib/python3.7/site-packages (from covidcast) (1.2.5)
Collecting imageio
   Downloading imageio-2.13.0-py3-none-any.whl (3.3 MB)
        3.3 MB 9.3 MB/s
Collecting delphi-epidata>=0.0.11
   Downloading delphi_epidata-0.3.1-py3-none-any.whl (6.8 kB)
Requirement already satisfied: tqdm in /shared-libs/python3.7/py/lib/python3.7/site-packages (from covidcast) (4.62.3)
Collecting epiweeks
   Downloading epiweeks-2.1.3-pv3-none-anv.whl (5.9 kB)
Requirement already satisfied: pyparsing>=2.2.1 in /shared-libs/python3.7/py-core/lib/python3.7/site-packages (from matplotlib->covidcast) (2.4.7)
Requirement already satisfied: cycler>=0.10 in /shared-libs/python3.7/py/lib/python3.7/site-packages (from matplotlib->covidcast) (0.11.0)
Requirement already satisfied: pillow>=6.2.0 in /shared-libs/python3.7/py/lib/python3.7/site-packages (from matplotlib->covidcast) (8.4.0)
Requirement already satisfied: kiwisolver>= 1.0.1 in /shared-libs/python 3.7/py/lib/python 3.7/site-packages (from matplotlib->covid cast) (1.3.2) in /shared-libs/python 3.7/py/lib/python 3.7/site-packages (from matplotlib->covid cast) (1.3.2) in /shared-libs/python 3.7/py/lib/python 3.7/site-packages (from matplotlib->covid cast) (1.3.2) in /shared-libs/python 3.7/site-packages (fr
Requirement already satisfied: python-dateutil>=2.7 in /shared-libs/python3.7/py-core/lib/python3.7/site-packages (from matplotlib->covidcast) (2.8.2)
Collecting pyproj>=2.2.0
   Downloading pyproj-3.2.1-cp37-cp37m-manylinux2010_x86_64.whl (6.3 MB)
                                                            ■1 6.3 MB 24.0 MB/s
from datetime import date
import covidcast
import pandas as pd
import numpy as np
{\tt ca\_counties = covid cast.fips\_to\_name("^06.*", ties\_method="all")}
ca\_counties = list(ca\_counties[0].values())
counties string = []
for i in ca_counties:
       string = "
        for element in i:
              string += element
       counties_string.append(string)
counties_string = counties_string[1:59] # removing 'california'
ca_counties_fips = covidcast.name_to_fips(counties_string)
ca_counties_fips
/root/venv/lib/python3.7/site-packages/covidcast/geography.py:314: UserWarning: Some inputs were not uniquely matched; returning only the first match in each case. To return all match
   warnings.warn("Some inputs were not uniquely matched; returning only the first match '
```

```
['06001',
  '06003',
  '06005',
  '06007',
  '06009',
  '06011',
  '06013',
 '06015'.
  '06017',
  '06019'
  '06021'.
  '06023'.
  '06025',
data = covidcast.signal("indicator-combination", "confirmed_incidence_num",
geo_values= google_sum_fips)
data.head()
  geo_value object | signal object
                                              time_value datetime64[ns]
                                                                         issue datetime64[ns]
                                                                                                     lag int64 missing_value int64 missing_stderr int64 missing_sample_size int64 value
                                              2020-02-20T00:00:00.00000
                                                                          2020-07-10T00:00:00.00000
                    confirmed_incidence_num
                                                                                                                                                           5
@ 06001
                                              2020-02-20T00:00:00.00000
                                                                          2020-07-10T00:00:00.00000
11 06013
                    confirmed_incidence_num
                                                                                                      141
                                                                                                                 0
                                                                                                                                     5
                                                                                                                                                           5
                                                                                                                                                                                     0
                                              2020-02-20T00:00:00.00000
                                                                         2020-07-10T00:00:00.00000
                    confirmed_incidence_num
                                              2020-02-20T00:00:00.00000
                                                                         2020-07-10T00:00:00.00000
38 06029
                    confirmed incidence num
                                                                                                                                                           5
                                              2020-02-20T00:00:00.00000
                                                                         2020-07-10T00:00:00.00000
41 06037
                    confirmed_incidence_num
                                                                                                                0
                                                                                                                                                          5
                                                                                                                                                                                     0
5 rews = 13 calumns
data.tail()
                                                                                                       lag int64 | missing_value int64 | missing_stderr int64 | missing_sample_size int64 | value
    geo_value object signal object
                                               time_value datetime64[ns]
                                                                           issue datetime64[ns]
                                                                           2021-11-15T00:00:00.00000
                                               2021-11-12T00:00:00.00000
553 06107
                     confirmed_incidence_num
                                               2021-11-12T00:00:00.00000
                                                                           2021-11-15T00:00:00.00000
554 06109
                     confirmed_incidence_num
555 06111
                     confirmed incidence num
                                                                                                                                      5
                                                                                                                                                            5
                                               2021-11-12T00:00:00.00000
                                                                          2021-11-15T00:00:00.00000
556 06113
                     confirmed_incidence_num
                                                                                                                                                                                       0
                                               2021-11-12T00:00:00.00000
                                                                          2021-11-15T00:00:00.00000
557 06115
                     confirmed_incidence_num
S rows x 12 columns
labels = data['value']
# number of observations
labels.size
9480
# looking for NA values for value column
data.isna().sum()
geo value
signal
time_value
issue
lag
missing_value
missing_stderr
missing_sample_size
value
stderr
                       9480
sample_size
                       9480
geo_type
                          О
data source
                          0
dtvpe: int64
features we are using:
```

```
smoothed_outpatient_cli source name : chng
```

smoothed_cli source: doctor-visits

smoothed_covid19_from_claims source: hospital-admissions

sum_anosmia_ageusia_raw_search, ageusia_raw_search, anosmia_raw_search source: google-symptoms

```
# many missing dates in month of november
chng = covidcast.signal("chng", "smoothed_outpatient_cli",
geo_values=google_sum_fips, start_day=date(2020, 2, 20), end_day=date(2021, 11, 12))
chng.head()
```

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211003 for geography 'county' NoDataWarning)

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211004 for geography 'county' NoDataWarning)

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211005 for geography 'county' NoDataWarning)

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211006 for geography 'county' NoDataWarning)

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211007 for geography 'county' NoDataWarning)

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211008 for geography 'county' NoDataWarning)

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211009 for geography 'county' NoDataWarning)

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211010 for geography 'county' NoDataWarning)

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211011 for geography 'county' NoDataWarning)

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211012 for geography 'county' NoDataWarning)

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211013 for geography 'county' NoDataWarning)

/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211014 for geography 'county' NoDataWarning)
/root/venv/lib/python3.7/site-packages/covidcast.py:425: NoDataWarning: No chng smoothed outpatient cli data found on 20211015 for geography 'county'

NoDataWarning)
/root/venv/lib/python3.7/site-packages/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211016 for geography 'county'

NoDataWarning)
/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:425: NoDataWarning: No chng smoothed_outpatient_cli data found on 20211017 for geography 'county'

NoDataWarning)
/root/venv/lib/pvthon3.7/site-packages/covidcast/covidcast.pv:425: NoDataWarning: No chng smoothed outpatient cli data found on 20211018 for geography 'county'

geo_value object	signal object	<pre>time_value datetime64[ns]</pre>	issue datetime64[ns]	lag int64	missing_value int64	missing_stderr int64	missing_sample_size int64	value
@ 06001	smoothed_outpatient_cli	2020-02-20T00:00:00.00000 0	2021-02-21T00:00:00.00000 0	367	0	5	5	0.0220
11 06013	smoothed_outpatient_cli	2020-02-20T00:00:00.00000 0	2021-02-21T00:00:00.00000 0	367	0	5	5	0.0203
2 2 06019	smoothed_outpatient_cli	2020-02-20T00:00:00.00000 0	2021-02-21T00:00:00.00000 0	367	0	5	5	0.0370
38 06029	smoothed_outpatient_cli	2020-02-20T00:00:00.00000 0	2021-02-21T00:00:00.00000 0	367	0	5	5	0.0101
41 06037	smoothed_outpatient_cli	2020-02-20T00:00:00.00000 0	2021-02-21T00:00:00.00000 0	367	0	5	5	0.0041

S rows x 1% columns

chr	ng.tail()								
	geo_value object	signal object	time_value datetime64[ns]	issue datetime64[ns]	lag int64	missing_value int64	missing_stderr int64	missing_sample_size int64	valu
1102	06075	smoothed_outpatient_cli	2021-10-02T00:00:00.00000 0	2021-10-07T00:00:00.00000 0	5	0	5	5	0.67
1111	06077	smoothed_outpatient_cli	2021-10-02T00:00:00.00000 0	2021-10-07T00:00:00.00000 0	5	0	5	5	3.49
1122	06081	smoothed_outpatient_cli	2021-10-02T00:00:00.00000 0	2021-10-07T00:00:00.00000 0	5	0	5	5	1.62
113	06085	smoothed_outpatient_cli	2021-10-02T00:00:00.00000 0	2021-10-07T00:00:00.00000 0	5	0	5	5	4.89
1141	06111	smoothed_outpatient_cli	2021-10-02T00:00:00.00000 0	2021-10-07T00:00:00.00000 0	5	0	5	5	1.67
5 1	ows = 13 columns								

chng.shape

(8865, 13)

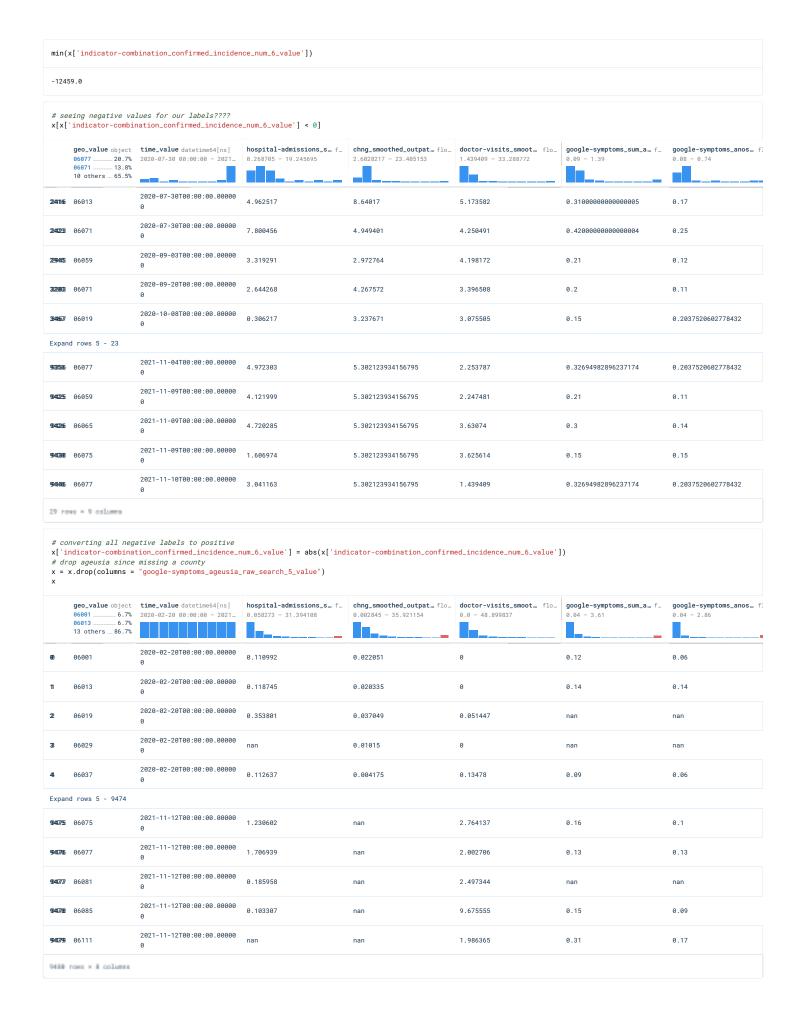
```
chnq.isna().sum()
geo_value
signal
time_value
lag
missing_stderr
missing_sample_size
                          А
value
stderr
                       8865
                       8865
sample_size
geo_type
data_source
dtype: int64
hosp = covidcast.signal("hospital-admissions", "smoothed_covid19_from_claims",
geo_values=google_sum_fips, start_day=date(2020, 2, 20), end_day=date(2021, 11, 12))
hosp.head()
   geo_value object signal object
                                                                                                       lag int64 missing_value int64 missing_stderr int64 missing_sample_size int64 val
                                                time_value datetime64[ns]
                                                                           issue datetime64[ns]
                    smoothed_covid19_from_cla
                                               2020-02-20T00:00:00.00000
                                                                           2020-06-21T00:00:00.00000
@ 06001
                                                                                                        122
                     smoothed_covid19_from_cla
                                                                                                        122
11 06013
                                                                                                                                                                                      0.1
                    smoothed_covid19_from_cla
                                               2020-02-20T00:00:00.00000
                                                                           2020-06-21T00:00:00.00000
2 06019
                                                                                                        122
                                                                                                                                                            5
                                                                                                                                                                                      и з
                    smoothed_covid19_from_cla
                                               2020-02-20T00:00:00.00000
                                                                           2020-06-21T00:00:00.00000
3 06037
                                                                                                        122
                    ims
                     smoothed_covid19_from_cla
                                                2020-02-20T00:00:00.00000
                                                                           2020-06-21T00:00:00.00000
                                                                                                        122
                                                                                                                                      5
                                                                                                                                                            5
                                                                                                                                                                                      0.0
44 06059
S rows x 12 columns
hosp.tail()
    geo_value object | signal object
                                                 time_value datetime64[ns]
                                                                                                        lag int64 missing_value int64 missing_stderr int64 missing_sample_size int64 va
                      smoothed_covid19_from_cla
                                                 2021-11-12T00:00:00.00000
                                                                             2021-11-28T00:00:00.00000
   06073
                                                                                                                                                             5
                                                                                                                                                                                       2.
                     smoothed covid19 from cla
                                                2021-11-12T00:00:00.00000
                                                                            2021-11-28T00:00:00.00000
    06075
                      smoothed_covid19_from_cla
                                                2021-11-12T00:00:00.00000
                                                                            2021-11-26T00:00:00.00000
1111 06077
                                                2021-11-12T00:00:00.00000
                                                                             2021-11-28T00:00:00.00000
                      smoothed_covid19_from_cla
1122 06081
                                                                                                                                                             5
                                                                                                                                                                                       0.
                                                                            2021-11-28T00:00:00.00000
                     smoothed covid19 from cla 2021-11-12T00:00:00.00000
113 06085
                     ims
                                                 0
                                                                            0
5 rews = 13 columns
hosp.shape
(9117, 13)
google_sum = covidcast.signal("google-symptoms", "sum_anosmia_ageusia_raw_search"
geo_values=google_sum_fips, start_day=date(2020, 2, 20), end_day=date(2021, 11, 12))
google_sum.head()
   geo_value object | signal object
                                                time value datetime64[ns]
                                                                           issue datetime64[ns]
                                                                                                       lag int64 missing_value int64 missing_stderr int64 missing_sample_size int64 val
                    sum_anosmia_ageusia_raw_s
                                                                                                                                                            5
@ 06001
                                                                                                       329
                                                                                                                                                                                      0.1
                                                2020-02-20T00:00:00.00000
                                                                           2021-01-14T00:00:00.00000
                    sum_anosmia_aqeusia_raw_s
11 06013
                                                                                                       329
                                                                                                                                                            5
                                                                                                                                                                                      0.1
                    earch
                    sum_anosmia_ageusia_raw_s
                                               2020-02-20T00:00:00.00000
                                                                           2021-01-14T00:00:00.00000
                    earch
                                                                           0
                                               2020-02-20T00:00:00.00000
                                                                           2021-01-14T00:00:00.00000
38 06059
                                                                                                       329
                                                                                                                                                                                      0.1
                                                                           2021-01-14T00:00:00.00000
                    sum_anosmia_aqeusia_raw_s
                                               2020-02-20T00:00:00.00000
4 96965
                                                                                                       329
                                                                                                                                                                                      0.1
                    earch
```

```
google_ageusia = covidcast.signal("google-symptoms", "ageusia_raw_search",
geo_values=google_sum_fips, start_day=date(2020, 2, 20), end_day=date(2021, 11, 12))
google_ageusia.head()
   geo_value object signal object
                                       time_value datetime64[ns]
                                                                  issue datetime64[ns]
                                                                                               lag int64 missing_value int64 missing_stderr int64 missing_sample_size int64 value floate
                                        2020-02-20T00:00:00.00000
                                                                   2021-01-14T00:00:00.00000
68 96991
                    ageusia_raw_search
                                                                                               329
                                                                                                                                                   5
                                                                                                                                                                             0.06
                                        2020-02-20100:00:00.00000
                                                                   2021-01-14T00:00:00.00000
  06037
                    ageusia_raw_search
                                                                                               329
                                                                                                                                                                             0.03
                                        2020-02-20T00:00:00.00000
                                                                   2021-01-14T00:00:00.00000
                                                                                               329
22 06059
                    ageusia_raw_search
                                                                                                                                                   5
                                                                                                                                                                             0.08
                                        2020-02-20T00:00:00.00000
                                                                   2021-01-14T00:00:00.00000
38 96965
                    ageusia_raw_search
                                                                                               329
                                                                                                                                                   5
                                                                                                                                                                             0.04
                                        2020-02-20T00:00:00.00000
                                                                   2021-01-14T00:00:00.00000
                    ageusia_raw_search
                                                                                                                                                                             0.04
5 rows x 13 calumns
google_anosmia = covidcast.signal("google-symptoms", "anosmia_raw_search"
geo_values=google_sum_fips, start_day=date(2020, 2, 20), end_day=date(2021, 11, 12))
google_anosmia.head()
/root/yeny/lib/python3.7/site-packages/covidcast/covidcast.py:429: RuntimeWarning: Problem obtaining google-symptoms anosmia raw search data on 20210528 for geography 'county': error:
  RuntimeWarning)
/root/venv/lib/python3.7/site-packages/covidcast.py:429: RuntimeWarning: Problem obtaining google-symptoms anosmia_raw_search data on 20210530 for geography 'county': error:
  RuntimeWarning)
/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:429: RuntimeWarning: Problem obtaining google-symptoms anosmia_raw_search data on 20210531 for geography 'county': error:
  RuntimeWarning)
/root/venv/lib/python3.7/site-packages/covidcast.py:429: RuntimeWarning: Problem obtaining google-symptoms anosmia_raw_search data on 20210607 for geography 'county': error:
  RuntimeWarning)
/root/venv/lib/python3.7/site-packages/covidcast.py:429: RuntimeWarning: Problem obtaining google-symptoms anosmia_raw_search data on 20210614 for geography 'county': error:
  RuntimeWarning)
/root/venv/lib/python3.7/site-packages/covidcast.py:429: RuntimeWarning: Problem obtaining google-symptoms anosmia_raw_search data on 20210616 for geography 'county': error:
  RuntimeWarning)
/root/venv/lib/python3.7/site-packages/covidcast/covidcast.py:429: RuntimeWarning: Problem obtaining google-symptoms anosmia_raw_search data on 20210618 for geography 'county': error:
  RuntimeWarning)
   geo_value object | signal object
                                       time_value datetime64[ns]
                                                                                               lag int64 missing_value int64 missing_stderr int64 missing_sample_size int64 value float6
                                                                  issue datetime64[ns]
                                        2020-02-20T00:00:00.00000
                                                                   2021-01-14T00:00:00.00000
6 06001
                    anosmia raw search
                                                                                               329
                                                                                                                                                                             0.06
                                        2020-02-20100:00:00.00000
                                                                   2021-01-14T00:00:00.00000
                    anosmia_raw_search
                                                                   2021-01-14T00:00:00.00000
                                        2020-02-20T00:00:00.00000
                                                                                              329
                                                                                                                                                   5
22 06037
                    anosmia raw search
                                                                                                                                                                             0.06
                                        2020-02-20T00:00:00.00000
                                                                   2021-01-14T00:00:00.00000
38 06059
                    anosmia raw search
                                                                                               329
                                                                                                                             5
                                                                                                                                                   5
                                                                                                                                                                             0.09
                                        2020-02-20T00:00:00.00000
                                                                   2021-01-14T00:00:00.00000
                    anosmia_raw_search
5 rows = 13 calumns
# checking to see if counties match up in different signals
fips = []
for i in ageusia_fips:
    if i in hosp_fips:
         fips.append(i)
doc = covidcast.signal("doctor-visits", "smoothed_cli"
geo_values=google_sum_fips, start_day=date(2020, 2, 20), end_day=date(2021, 11, 12))
doc.head()
   geo_value object | signal object | time_value datetime64[ns]
                                                            issue datetime64[ns]
                                                                                        lag int64 missing_value int64 missing_stderr int64 missing_sample_size int64 value float64 st
                                  2020-02-20T00:00:00.00000
                                                             2020-06-09T00:00:00.00000
6 06001
                    smoothed cli
                                                                                         110
                                                                                                                                                                                     No
                                  2020-02-20T00:00:00.00000
                                                             2020-06-09T00:00:00.00000
1 96913
                    smoothed_cli
                                                                                         110
                                                                                                                                                                                     No
                                  2020-02-20T00:00:00.00000
                                                             2020-06-09700:00:00.00000
                    smoothed_cli
                                                                                                                                                                       0.051447
                                                                                                                                                                                     No
                                  2020-02-20T00:00:00.00000
                                                             2020-06-09T00:00:00.00000
                    smoothed cli
                                                                                         110
                                                                                                                                             5
38 06029
                                                                                                                                                                       0
                                                                                                                                                                                     No
                                  2020-02-20T00:00:00.00000
                                                             2020-06-09T00:00:00.00000
4 96937
                    smoothed_cli
                                                                                         110
                                                                                                                                             5
                                                                                                                                                                       0.13478
                                                                                                                                                                                     No
```

400.	tail()												
g	geo_value object	signal object	time_value date	etime64[ns]	issue datetime64	[ns]	lag int64	missing_value int64	missing_stderr	int64 mis s	sing_sample_size int6	4 value float6	54
3896 0	96103	smoothed_cli	2021-11-12T00: 0	00:00.00000	2021-11-27T00:0	0:00.00000	15	0	5	5		2.29704	
410E 0	96107	smoothed_cli	2021-11-12T00:	00:00.00000	2021-11-27T00:0	0:00.00000	15	0	5	5		3.001187	
4111 0	96111	smoothed_cli	2021-11-12T00:	00:00.00000	2021-11-27T00:0	0:00.00000	15	0	5	5		1.986365	
4422 0	n6113	smoothed_cli	0 2021-11-12T00:	00:00.00000	0 2021-11-27T00:0	0:00.00000	15	0	5	5		0.976299	
			0 2021-11-12T00:	00:00.00000	0 2021-11-27T00:0	0:00.00000							
413 0		smoothed_cli	0		0		15	0	5	5		1.265493	
5 rew	s = 13 calumns												
np.aı	rray(doc).shape												
(9480	9, 13)												
doc.:	isna().sum()												
geo_v		0											
signa time_	-value	0 0											
issue	e	0											
lag missi	ing_value	0 0											
missi	ing_stderr	0											
missi value	ing_sample_size	0 0											
stder		28440											
	le_size	28440											
geo_t	type	28440 0											
geo_t data_ dtype		28440 0 0	gnals([hosp, ch	ing, doc, god	ogle_sum, google	e_anosmia,	google_age	usia, data])					
geo_t data_ dtype	type _source e: int64 ed = covidcast.	28440 0 0	gnals([hosp, ch	ing, doc, god	ogle_sum, googl	e_anosmia,	google_age	usia, data])					
geo_t data_ dtype	type _source =: int64 ed = covidcast. ed geo_value object	28440 0 0 aggregate_sig	datetime64[ns]	hospital-add	missions… date…	hospital-ac	missions_s…	f hospital-admis			issions_s… f… hospi		:_s
geo_t data_ dtype	type _source e: int64 ed = covidcast.	28440 0 0 aggregate_sig	datetime64[ns]	hospital-adı	missions… date…		missions_s…			spital-admi) - 5.0	issions_s_ f_ hospi 5.0 -		:_s
geo_t data_ dtype merge	type _source e: int64 ed = covidcast. ed geo_value object 06081	28440 0 0 aggregate_sig	datetime64[ns]	hospital-add	missions date 0:00:00 - 2021	hospital-ac	missions_s…	f hospital-admis					:_\$
geo_t data_ dtype merge merge	type _source e: int64 ed = covidcast. ed geo_value object 060816.7% 060816.7%	28440 0 0 aggregate_sig	datetime64[ns] 0:00:00 - 2021	hospital-adi 2020-06-21 0	missions date 0:00:00 - 2021 00:00:00	hospital-ac 14.0 - 122.6	missions_s…	f hospital-admis	5.6		5.0 -		;_S
geo_t data_ dtype merge merge	type _source =: int64 ed = covidcast. ed geo_value object 06001	28440 0 0 aggregate_sig	datetime64[ns] 10:80:80 - 2021	hospital-add 2020-06-21 01 2020-06-21 0 2020-06-21 0	missions date 0:00:00 - 2021 00:00:00 00:00:00	hospital-ac 14.0 - 122.6	missions_s…	f hospital-admis 0.0 - 0.0	5.6		5.0 -		:_S
geo_t data_ dtype merge merge 11 22	type _source a: int64 ed = covidcast. ed geo_value object	28440 0 0 aggregate_sig time_value 0 2020-02-20 0 2020-02-20 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0	datetime64[ns] 10:80:80 - 2021	hospital-add 2020-06-21 01 2020-06-21 (2020-06-21 (2020-06-21 (missions date 0:00:00 - 2021 00:00:00 00:00:00	hospital-ac 14.0 - 122.6 122	missions_s…	f	5.6) - 5.0	5.0 -		;_S
geo_t data_dtype merge merge 11 22 38	type _source a: int64 ed = covidcast. ed geo_value object 060016.7% 0601386.7% 06001 060018601380001 06001860001	28440 0 0 aggregate_sig	datetime64[ns] 10:00:00 - 2021_ 10:00:00:00.00000 100:00:00.00000	hospital-add 2020-06-21 01 2020-06-21 0 2020-06-21 0 2020-06-21 0	missions date 0:00:00 - 2021 00:00:00 00:00:00 00:00:00	hospital-ac 14.0 - 122.6 122 122 122	missions_s…	f hospital-admis 0.0 - 0.0 0 0 nan	5.6 5 5) - 5.0	5.0 - 5 5		:_ \$
geo_t data dtype merge merge 11 22 33	type _source a: int64 ed = covidcast. ed geo_value object 06001	28440 0 0 aggregate_sig	datetime64[ns] 10:00:00 - 2021_ 10:00:00:00.00000 100:00:00.00000 100:00:00.00000	hospital-add 2020-06-21 01 2020-06-21 (2020-06-21 (2020-06-21 (missions date 0:00:00 - 2021 00:00:00 00:00:00 00:00:00	hospital-ac 14.0 - 122.6 122	missions_s…	f	5.6) - 5.0	5.0 -		::_s
geo_t data dtype merge merge 60 11 22 33 44 Expansion	type _source a: int64 ed = covidcast. ed geo_value object 060016.7% 0601386.7% 06001 060018601380001 06001860001	28446 0 0 aggregate_sig time_value (2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0	datetime64[ns] 10:00:00 - 2021_ 10:00:00:00.00000 100:00:00.00000 100:00:00.00000	hospital-add 2020-06-21 01 2020-06-21 0 2020-06-21 0 2020-06-21 0	missions date 0:00:00 - 2021 00:00:00 00:00:00 00:00:00	hospital-ac 14.0 - 122.6 122 122 122	missions_s…	f hospital-admis 0.0 - 0.0 0 0 nan	5.6 5 5) - 5.0	5.0 - 5 5		:_s
geo_t data dtype merge merge 89 11 22 33 44 Expan 9947/75	type source source s: int64 ed = covidcast. ed geo_value object 66001	28440 0 0 aggregate_sig time_value c 2020-02-20 0 2020-02-20 0 2020-02-20 0 2020-02-20 0 2020-02-20 0 2020-02-20 0 0 2020-02-20 0 2020-02-20 0 2020-02-20 0 2020-02-20 0 2020-02-20 0	datetime64[ns] 0:00:00 - 2021	hospital-add 2020-06-21 (0) 2020-06-21 (1) 2020-06-21 (1) 2020-06-21 (1) NaT	missions date 9:00:00 - 2021 00:00:00 00:00:00 00:00:00 00:00:00	hospital-ac 14.0 - 122.6 122 122 122 122 122	missions_s…	f hospital-admis	5.6 5 5 5 nar) - 5.0	5.0 - 5 5 5 nan 5		:_s
geo_t data dtype merge merge 11 22 33 44 Expan 9947/75	type source a: int64 ed = covidcast. ed geo_value object 86601	28440 0 0 aggregate_sig time_value of 2020-02-20 0 2020-02-20 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0	datetime64[ns] 0:00:00 - 2021	hospital-add 2020-06-21 (0) 2020-06-21 (1) 2020-06-21 (1) 2020-06-21 (1) NaT 2020-06-21 (1) 2021-11-28 (1)	missions date	hospital-ac 14.0 - 122.6 122 122 122 122 122 14	missions_s…	f hospital-admis	5) - 5.0	5.0 - 5 5 5 nan 5		:_s
geo_t data dtype merge merge 99 11 22 33 44 Expan: 9944775 9944777	type source a: int64 ed = covidcast. ed geo_value object 06001	28440 0 0 aggregate_sig time_value of 2020-02-20 0 2020-02-20 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2021-11-12T 0 2021-11-12T 0	datetime64[ns] 10:00:00 - 2021_ 10:00:00:00 - 2021_ 100:00:00:00.00000 100:00:00.000000 100:00:00.000000 100:00:00.000000 100:00:00.0000000000	hospital-add 2020-06-21 (0) 2020-06-21 (1) 2020-06-21 (1) 2020-06-21 (1) NAT 2020-06-21 (1) 2021-11-28 (1) 2021-11-28 (1)	missions date 0:00:00 - 2021 00:00:00 00:00:00 00:00:00 00:00:00 00:00:00 00:00:00	hospital-ac 14.0 - 122.6 122 122 122 122 16 14	missions_s…	f hospital-admis 0.0 - 0.0 0 0 nan 0	5) - 5.0	5.0 - 5 5 5 nan 5		::_\$
geo_t data dtype merge merge 99 1 2 2 3 4 Expan: 994775 9947776	type source a: int64 ed = covidcast. ed geo_value object 86601	28440 0 0 aggregate_sig time_value 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2020-02-20T 0 2021-11-12T 0 2021-11-12T 0 2021-11-12T 0	datetime64[ns] 00:00:00 - 2021	hospital-add 2020-06-21 (0) 2020-06-21 (1) 2020-06-21 (1) 2020-06-21 (1) NaT 2020-06-21 (1) 2021-11-28 (1)	missions date 0:00:00 - 2021 00:00:00 00:00:00 00:00:00 00:00:00 00:00:00 00:00:00	hospital-ac 14.0 - 122.6 122 122 122 122 122 14	missions_s…	f hospital-admis	5) - 5.0	5.0 - 5 5 5 nan 5		:_S.

```
merged.columns
Index(['geo_value', 'time_value',
        hospital-admissions_smoothed_covid19_from_claims_0_issue',
        hospital-admissions_smoothed_covid19_from_claims_0_lag
        hospital-admissions\_smoothed\_covid19\_from\_claims\_0\_missing\_value',
        'hospital-admissions_smoothed_covid19_from_claims_0_missing_stderr'
        hospital-admissions\_smoothed\_covid19\_from\_claims\_0\_missing\_sample\_size',
        'hospital-admissions_smoothed_covid19_from_claims_0_value'
        hospital-admissions_smoothed_covid19_from_claims_0_stderr'
        'hospital-admissions_smoothed_covid19_from_claims_0_sample_size',
        chng_smoothed_outpatient_cli_1_issue'
        chnq_smoothed_outpatient_cli_1_lag',
        chng_smoothed_outpatient_cli_1_missing_value',
        chng_smoothed_outpatient_cli_1_missing_stderr'
        chng_smoothed_outpatient_cli_1_missing_sample_size',
        chng_smoothed_outpatient_cli_1_value
        chng_smoothed_outpatient_cli_1_stderr
        'chng_smoothed_outpatient_cli_1_sample_size',
        'doctor-visits_smoothed_cli_2_issue',
        'doctor-visits_smoothed_cli_2_lag'
        'doctor-visits_smoothed_cli_2_missing_value',
        doctor-visits smoothed cli 2 missing stderr'
        doctor-visits_smoothed_cli_2_missing_sample_size',
        doctor-visits_smoothed_cli_2_value',
        doctor-visits_smoothed_cli_2_stderr'
        doctor-visits_smoothed_cli_2_sample_size',
        google-symptoms_sum_anosmia_ageusia_raw_search_3_issue',
        google-symptoms_sum_anosmia_ageusia_raw_search_3_lag',
        google-symptoms_sum_anosmia_ageusia_raw_search_3_missing_value',
        google-symptoms\_sum\_anosmia\_ageusia\_raw\_search\_3\_missing\_stderr'
        google-symptoms_sum_anosmia_ageusia_raw_search_3_missing_sample_size',
        google-symptoms_sum_anosmia_ageusia_raw_search_3_value'
        |qooqle-symptoms_sum_anosmia_ageusia_raw_search_3_stderr',
x = merged.drop(columns = [
         hospital-admissions_smoothed_covid19_from_claims_0_issue',
        'hospital-admissions_smoothed_covid19_from_claims_0_lag'
        'hospital-admissions_smoothed_covid19_from_claims_0_missing_value',
'hospital-admissions_smoothed_covid19_from_claims_0_missing_stderr'
        hospital-admissions_smoothed_covid19_from_claims_0_missing_sample_size',
        'hospital-admissions_smoothed_covid19_from_claims_0_stderr',
'hospital-admissions_smoothed_covid19_from_claims_0_sample_size',
         chng_smoothed_outpatient_cli_1_issue',
        chng_smoothed_outpatient_cli_1_lag',
        'chng_smoothed_outpatient_cli_1_missing_value',
'chng_smoothed_outpatient_cli_1_missing_stderr'
        chng_smoothed_outpatient_cli_1_missing_sample_size',
        chng_smoothed_outpatient_cli_1_stderr
        chng_smoothed_outpatient_cli_1_sample_size',
        doctor-visits_smoothed_cli_2_issue',
        doctor-visits_smoothed_cli_2_lag'
        doctor-visits_smoothed_cli_2_missing_value',
        doctor-visits_smoothed_cli_2_missing_stderr
         doctor-visits_smoothed_cli_2_missing_sample_size',
        doctor-visits_smoothed_cli_2_stderr
        'doctor-visits_smoothed_cli_2_sample_size'
         google-symptoms_sum_anosmia_ageusia_raw_search_3_issue',
         google-symptoms_sum_anosmia_ageusia_raw_search_3_lag',
         google-symptoms_sum_anosmia_ageusia_raw_search_3_missing_stderr',
google-symptoms_sum_anosmia_ageusia_raw_search_3_missing_sample_size',
         google-symptoms_sum_anosmia_ageusia_raw_search_3_stderr
        google-symptoms_sum_anosmia_ageusia_raw_search_3_sample_size',
google-symptoms_anosmia_raw_search_4_issue',
         google-symptoms_anosmia_raw_search_4_lag',
         google-symptoms_anosmia_raw_search_4_missing_value',
        'google-symptoms_anosmia_raw_search_4_missing_stderr'
         google-symptoms_anosmia_raw_search_4_missing_sample_size',
         google-symptoms_anosmia_raw_search_4_stderr
         \verb|google-symptoms_anosmia_raw_search_4_sample_size'|,\\
         google-symptoms_ageusia_raw_search_5_issue',
         google-symptoms_ageusia_raw_search_5_lag',
         google-symptoms_ageusia_raw_search_5_missing_value',
         google-symptoms_ageusia_raw_search_5_missing_stderr'
         google-symptoms_ageusia_raw_search_5_missing_sample_size',
         google-symptoms_ageusia_raw_search_5_stderr
         {\tt google-symptoms\_ageusia\_raw\_search\_5\_sample\_size'}
        'indicator-combination confirmed incidence num 6 issue'.
         indicator-combination\_confirmed\_incidence\_num\_6\_lag'
        'indicator-combination_confirmed_incidence_num_6_missing_value',
        'indicator-combination_confirmed_incidence_num_6_missing_stderr'
         indicator-combination\_confirmed\_incidence\_num\_6\_missing\_sample\_size',
        'indicator-combination_confirmed_incidence_num_6_stderr
        'indicator-combination_confirmed_incidence_num_6_sample_size',
        'geo_type'])
      geo_value object time_value datetime64[ns]
                                                    hospital-admissions_s... f...
                                                                                 chng_smoothed_outpat... flo... doctor-visits_smoot... flo... google-symptoms_sum_a... f... google-symptoms_anos... fl
             6.7% 2020-02-20 00:00:00 - 2021.
                                                                                                                                               - 3.61
                                                                                                                                                                       0.04 - 2.86
                ... 6.7%
```

```
13 others ... 86.7%
                        2020-02-20T00:00:00.00000
                                                     0.110992
                                                                                                              0
      06001
                                                                                 0.022051
                                                                                                                                           0.12
                                                                                                                                                                        0.06
                        2020-02-20T00:00:00.00000
      06013
                                                     0.118745
                                                                                  0.020335
                                                                                                                                           0.14
                                                                                                                                                                        0.14
                        2020-02-20T00:00:00.00000
      06019
                                                     0.353801
                                                                                  0.037049
                                                                                                               0.051447
22
                                                                                                                                           nan
                                                                                                                                                                        nan
                        2020-02-20T00:00:00.00000
                                                                                  0.01015
33
      06029
                                                                                                                                           nan
                                                                                                                                                                        nan
                        2020-02-20T00:00:00.00000
      06037
                                                     0.112637
                                                                                  0.004175
                                                                                                               0.13478
                                                                                                                                           0.09
                                                                                                                                                                        0.06
Expand rows 5 - 9474
                        2021-11-12T00:00:00.00000
9941775 96975
                                                     1.239692
                                                                                  nan
                                                                                                               2.764137
                                                                                                                                           0.16
                                                                                                                                                                        0.1
                        2021-11-12T00:00:00.00000
     96977
                                                     1.706939
                                                                                                               2.002706
                                                                                                                                           0.13
                                                                                                                                                                        0.13
                        0
                        2021-11-12T00:00:00.00000
9941777 06081
                                                     0.185958
                                                                                                               2.497344
                                                                                                                                           nan
                                                                                                                                                                        nan
                        2021-11-12T00:00:00.00000
                                                     0.103307
99417798 06085
                                                                                  nan
                                                                                                               9.675555
                                                                                                                                           0.15
                                                                                                                                                                        0.09
                        2021-11-12T00:00:00.00000
99417799 96111
                                                                                                               1.986365
                                                                                                                                           0.31
                                                                                                                                                                        0.17
9488 rows = 9 columns
x.columns
Index(['geo_value', 'time_value',
        hospital-admissions_smoothed_covid19_from_claims_0_value',
        chng_smoothed_outpatient_cli_1_value',
        'doctor-visits_smoothed_cli_2_value',
        'google-symptoms_sum_anosmia_ageusia_raw_search_3_value',
        'google-symptoms_anosmia_raw_search_4_value',
        'google-symptoms_ageusia_raw_search_5_value',
        'indicator-combination_confirmed_incidence_num_6_value'],
      dtype='object')
x.isna().sum()
geo_value
hospital-admissions\_smoothed\_covid19\_from\_claims\_0\_value
                                                               363
chng_smoothed_outpatient_cli_1_value
                                                               615
{\tt doctor-visits\_smoothed\_cli\_2\_value}
                                                                 а
{\tt google-symptoms\_sum\_anosmia\_ageusia\_raw\_search\_3\_value}
                                                               710
google-symptoms_anosmia_raw_search_4_value
                                                               986
qooqle-symptoms_aqeusia_raw_search_5_value
                                                              2277
indicator-combination_confirmed_incidence_num_6_value
dtype: int64
np.unique(x[x['google-symptoms\_ageusia\_raw\_search\_5\_value'].isna()]['geo\_value'])\\
# these are the unique counties that the ageusia column in the merged dataset is na for ?
{\tt array(['06001', '06013', '06019', '06029', '06065', '06067', '06071',}
        '06075', '06077', '06081', '06085', '06111'], dtype=object)
x.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 9480 entries, 0 to 9479
Data columns (total 9 columns):
 # Column
                                                                 Non-Null Count Dtype
 0
     geo_value
                                                                  9480 non-null
                                                                                  object
                                                                  9480 non-null
                                                                                  datetime64[ns]
     hospital-admissions\_smoothed\_covid19\_from\_claims\_0\_value
                                                                 9117 non-null
                                                                                  float64
     chng_smoothed_outpatient_cli_1_value
                                                                  8865 non-null
                                                                                  float64
     doctor-visits_smoothed_cli_2_value
                                                                  9480 non-null
                                                                                  float64
     google-symptoms_sum_anosmia_ageusia_raw_search_3_value
                                                                 8770 non-null
                                                                                  float64
     {\tt google-symptoms\_anosmia\_raw\_search\_4\_value}
                                                                 8494 non-null
                                                                                  float64
     google-symptoms_ageusia_raw_search_5_value
                                                                 7203 non-null
                                                                                  float64
     indicator-combination\_confirmed\_incidence\_num\_6\_value
                                                                 9480 non-null
                                                                                  float64
dtypes: datetime64[ns](1), float64(7), object(1)
memory usage: 998.7+ KB
```



```
x.isna().sum()
geo_value
                                                                    0
time_value
hospital-admissions_smoothed_covid19_from_claims_0_value
                                                                  363
chng_smoothed_outpatient_cli_1_value
                                                                  615
doctor-visits_smoothed_cli_2_value
                                                                    0
google-symptoms_sum_anosmia_ageusia_raw_search_3_value
                                                                  710
google-symptoms_anosmia_raw_search_4_value
                                                                  986
indicator\hbox{-}combination\_confirmed\_incidence\_num\_6\_value
                                                                    А
dtype: int64
# imputing by forward filling based on previous observation in each county
updated_x = x
updated_x['hospital-admissions_smoothed_covid19_from_claims_0_value'] = x.groupby('geo_value')['hospital-admissions_smoothed_covid19_from_claims_0_value'].fillna(method='f
updated_x['chng_smoothed_outpatient_cli_1_value'] = x.groupby('geo_value')['chng_smoothed_outpatient_cli_1_value'].fillna(method='ffill')
updated_x['google-symptoms_sum_anosmia_ageusia_raw_search_3_value'] = x.groupby('geo_value')['google-symptoms_sum_anosmia_ageusia_raw_search_3_value'].fillna(method='ffill')
updated_x['google-symptoms_anosmia_raw_search_4_value'] = x.groupby('geo_value')['google-symptoms_anosmia_raw_search_4_value'].fillna(method='ffill')
      geo_value object
06001 ......6.7%
                         time_value datetime64[ns]
2020-02-20 00:00:00 - 2021.
                                                        hospital-admissions_s... f...
                                                                                      {\tt chng\_smoothed\_outpat}... ~{\tt flo}...
                                                                                                                    doctor-visits_smoot... flo..
                                                                                                                                                   google-symptoms_sum_a... f...
                                                                                                                                                                                 google-symptoms_anos... fl
                                                        0.058273 - 31.394108
                                                                                      0.002845 - 35.921154
                                                                                                                     0.0 - 48.899837
                                                                                                                                                        - 3.61
                                                                                                                                                                                  0.04 - 2.86
       06013
                   6.7%
       13 others ... 86.7%
                          2020-02-20T00:00:00.00000
                                                                                                                     0
      06001
                                                        0.110992
                                                                                      0.022051
                                                                                                                                                   0.12
                                                                                                                                                                                  0.06
                         2020-02-20T00:00:00.00000
      06013
                                                        0.118745
                                                                                      0.020335
                                                                                                                     а
                                                                                                                                                   0.14
                                                                                                                                                                                  0.14
                         2020-02-20T00:00:00.00000
                                                        0.353801
                                                                                      0.037049
                                                                                                                     0.051447
                          2020-02-20T00:00:00.00000
      06029
                                                                                      0.01015
                                                        nan
                                                                                                                                                   nan
                                                                                                                                                                                  nan
                         2020-02-20T00:00:00.00000
      96937
                                                        0.112637
                                                                                      0.004175
                                                                                                                     0.13478
                                                                                                                                                   0.09
                                                                                                                                                                                 0.06
Expand rows 5 - 9474
                          2021-11-12T00:00:00.00000
                                                        1.230602
9941775 06075
                                                                                      0.6704379
                                                                                                                    2.764137
                                                                                                                                                   0.16
                                                                                                                                                                                 0.1
                         2021-11-12T00:00:00.00000
      96977
                                                        1.706939
                                                                                      3.4975448
                                                                                                                     2.002706
                                                                                                                                                                                  0.13
                         2021-11-12T00:00:00.00000
                                                        0.185958
                                                                                      1.6268549
                                                                                                                     2.497344
                          2021-11-12T00:00:00.00000
                                                        0.103307
99417788 06085
                                                                                      4.8956309
                                                                                                                     9.675555
                                                                                                                                                   0.15
                                                                                                                                                                                  0.09
                         2021-11-12T00:00:00.00000
99417799 96111
                                                        18.908364
                                                                                      1.6790429
                                                                                                                     1.986365
                                                                                                                                                   0.31
                                                                                                                                                                                 0.17
9488 rows = 8 columns
updated_x.isna().sum()
qeo_value
hospital\text{-}admissions\_smoothed\_covid19\_from\_claims\_0\_value
chng_smoothed_outpatient_cli_1_value
                                                                  0
doctor-visits_smoothed_cli_2_value
google-symptoms_sum_anosmia_ageusia_raw_search_3_value
                                                                  11
google-symptoms_anosmia_raw_search_4_value
                                                                  29
indicator-combination_confirmed_incidence_num_6_value
dtvpe: int64
# drop remaining NA values
updated_x = updated_x.dropna(axis = 0)
updated_x.isna().sum()
geo_value
time value
hospital-admissions_smoothed_covid19_from_claims_0_value
chng_smoothed_outpatient_cli_1_value
doctor-visits_smoothed_cli_2_value
google-symptoms_sum_anosmia_ageusia_raw_search_3_value
google-symptoms_anosmia_raw_search_4_value
indicator-combination_confirmed_incidence_num_6_value
dtype: int64
updated_x
```

