Data Cleaning Steps

Step 1: Address the Completeness Issue

• Concatenate df_treatments.csv and df_treatments_cut.csv to ensure all treatment data is included.

Step 2: Address the Structural Issues

- 1. Fix Email and Phone Number Formatting (patients.csv)
 - Use regular expressions to separate the combined email and phone number into two distinct columns.
- 2. Reformat Treatment Data (treatments.csv)
 - Split "Novodra" and "Auralin" into two separate columns: treatment and oral dose.
 - Use the melt() function to restructure the data properly.
 - Extract Dose Values:
 - Split dose data into low dose and high dose values.
 - Remove the letter 'u' from the values.
 - Create three new columns: dose_low, dose_high, and treatment.
- 3. Integrate Adverse Reactions Data (adverse_reactions.csv & treatments.csv)
 - Add a new column df_adverse_condition to df_treatments.csv.
 - Merge both datasets based on common patient names using the merge() function.

Step 3: Address Quality Issues

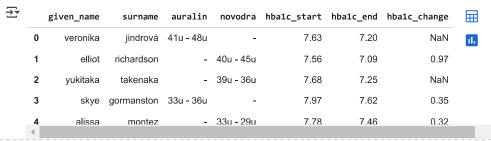
- Handle Improper Data
- · Remove Duplicate Records
- · Handle Missing Values
- · Detect and Fix Corrupt Data
- Fix Incorrect Data

#Colab link:- https://colab.research.google.com/drive/1aoM0eimDwgHJUPxT_aGnLDKbhJW8ePLj?usp=sharing

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

df_treatments = pd.read_csv('/content/treatments.csv')
df_treatments_cut = pd.read_csv('/content/treatments_cut.csv')
df_patients = pd.read_csv('/content/patients.csv')
df_adverse_conditions = pd.read_csv('/content/adverse_reactions.csv')
```

df_treatments.head()



Next steps: (Generate code with df_treatments)

View recommended plots

New interactive sheet

df_treatments.shape



df_adverse_conditions.head()



df_patients.isna().sum()



df_patients[df_patients.isna().any(axis=1)]

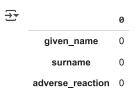
→	patient_id	assigned_sex	given_name	surname	address	city	state	zip_code	country	contact	birthdate	weight	height	bmi
2)9 210	female	Lalita	Eldarkhanov	NaN	NaN	NaN	NaN	NaN	NaN	8/14/1950	143.4	62	26.2
2	19 220	male	Mỹ	Quynh	NaN	NaN	NaN	NaN	NaN	NaN	4/9/1978	237.8	69	35.1
2	30 231	female	Elisabeth	Knudsen	NaN	NaN	NaN	NaN	NaN	NaN	9/23/1976	165.9	63	29.4
2	34 235	female	Martina	Tománková	NaN	NaN	NaN	NaN	NaN	NaN	4/7/1936	199.5	65	33.2
2	12 243	male	John	O'Brian	NaN	NaN	NaN	NaN	NaN	NaN	2/25/1957	205.3	74	26.4
2	19 250	male	Benjamin	Mehler	NaN	NaN	NaN	NaN	NaN	NaN	10/30/1951	146.5	69	21.6
2	57 258	male	Jin	Kung	NaN	NaN	NaN	NaN	NaN	NaN	5/17/1995	231.7	69	34.2
2	34 265	female	Wafiyyah	Asfour	NaN	NaN	NaN	NaN	NaN	NaN	11/3/1989	158.6	63	28.1
2	69 270	female	Flavia	Fiorentino	NaN	NaN	NaN	NaN	NaN	NaN	10/9/1937	175.2	61	33.1
2	78 279	female	Generosa	Cabán	NaN	NaN	NaN	NaN	NaN	NaN	12/16/1962	124.3	69	18.4
2	36 287	male	Lewis	Webb	NaN	NaN	NaN	NaN	NaN	NaN	4/1/1979	155.3	68	23.6
29	297	female	Chỉ	Lâm	NaN	NaN	NaN	NaN	NaN	NaN	5/14/1990	181.1	63	32.1

df_patients= df_patients.dropna()

df_patients

1:58 PW										
patien	t_id a	assigned_sex	given_name	surname	address	city	state	zip_code	country	
0	1	female	Zoe	Wellish	576 Brown Bear Drive	Rancho California	California	92390.0	United States	9170ZoeWellish@su
1	2	female	Pamela	Hill	2370 University Hill Road	Armstrong	Illinois	61812.0	United States	PamelaSHill@cuvox.de+
2	3	male	Jae	Debord	1493 Poling Farm Road	York	Nebraska	68467.0	United States	402-363-6804JaeMDebord
3	4	male	Liêm	Phan	2335 Webster Street	Woodbridge	NJ	7095.0	United States	PhanBaLiem@jourrapide.c
4	5	male	Tim	Neudorf	1428 Turkey Pen Lane	Dothan	AL	36303.0	United States	334-515-7487TimNeudori
498	499	male	Mustafa	Lindström	2530 Victoria Court	Milton Mills	ME	3852.0	United States	0579MustafaLindstrom@jou
			D	Ricliny	494 Clarksburg	Sedona	AZ	86341.0	United	928-284-4492RumanBisliev
499	500	male	Ruman	Disliev	Park Road	Codona			States	
500	501	female	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park	MO	64110.0	United States	6007JinkedeKeizer@t
500	501	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland		64110.0	United	
500	501	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
500	501 erate co	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
500 ✓ steps: Gene tients.isna(501 erate co).sum(0	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
500 steps: Gene tients.isna(patient_id	501 erate co).sum(0 0	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
steps: Genetic tients.isna(patient_id assigned_se:	501 erate co).sum(0 0	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
steps: Gene tients.isna(patient_id assigned_se; given_name	501 erate co).sum(0 0 0	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
steps: General stients.isna(patient_id assigned_se; given_name surname	501 erate co) . sum(0 0 0 0	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
steps: Gene tients.isna(patient_id assigned_se; given_name surname address	501 erate co).sum(0 0 0 0	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
steps: General strents.isna(patient_id assigned_se; given_name surname address city	501 erate co).sum(0 0 0 0 0	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
patient_id assigned_se; given_name surname address city state	501 erate co).sum(0 0 0 0 0 0	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
steps: General strength of the	501 erate co).sum(0 0 0 0 0 0 0 0 0 0 0 0 0	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
patient_id assigned_se; given_name address city state zip_code country contact birthdate	501 erate co).sum(0 0 0 0 0 0 0 0 0 0 0	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	
steps: General steps:	501 erate co).sum(0 0 0 0 0 0 0 0 0 0 0 0 0	female de with df_pat	Jinke	de Keizer	Park Road 649 Nutter Street	Overland Park		64110.0	United	

df_adverse_conditions.isna().sum()



df_treatments.isna().sum() 0 given_name 0 surname 0 auralin 0 0 novodra hba1c_start 0 hba1c_end 0 hba1c_change 109 df_treatments_cut.isna().sum() $\overline{\pm}$ 0 given_name 0 0 surname auralin 0 novodra 0 hba1c_start hba1c_end 0 hba1c_change 28 df_treatments = pd.concat([df_treatments, df_treatments_cut], ignore_index=True) df_treatments **₹** novodra hba1c_start hba1c_end hba1c_change given_name $\overline{\Pi}$ surname auralin 0 veronika jindrová 41u - 48u 7.63 7.20 NaN 1 - 40u - 45u 7.56 7.09 0.97 elliot richardson 2 yukitaka takenaka - 39u - 36u 7.68 7.25 NaN 3 skye gormanston 33u - 36u 7.97 7.62 0.35 4 alissa montez - 33u - 29u 7.78 7.46 0.32 345 rovzan kishiev 32u - 37u 7.75 7.41 0.34 346 jakob jakobsen - 28u - 26u 7.96 7.51 0.95 347 schneider 48u - 56u 7.74 7.44 0.30 bernd 348 - 42u - 44u 7.21 berta napolitani 7.68 NaN sauvé 36u - 46u 349 7.86 7.40 NaN armina 350 rows × 7 columns Next steps: (Generate code with df_treatments) View recommended plots New interactive sheet df_treatments.isna().sum()

```
∓
                        0
       given_name
                        0
         surname
                        0
          auralin
         novodra
       hba1c_start
        hba1c_end
                        0
      hba1c_change 137
df_treatments['hba1c_change']= df_treatments['hba1c_start'] - df_treatments['hba1c_end']
df_treatments.head()
₹
         given_name
                         surname
                                   auralin
                                             novodra hba1c_start hba1c_end hba1c_change
                                                                                                 ☶
      0
                                                               7.63
                                                                           7.20
             veronika
                         jindrová 41u - 48u
                                                                                          0.43
                                                                                                 11.
                                                               7.56
                                                                           7.09
                                                                                          0.47
                elliot
                       richardson
                                            40u - 45u
      2
             yukitaka
                                            39u - 36u
                                                               7.68
                                                                           7.25
                                                                                          0.43
                         takenaka
      3
                skye
                      gormanston 33u - 36u
                                                               7.97
                                                                           7.62
                                                                                          0.35
               alissa
                                            33u - 29u
                                                               7.78
                                                                           7.46
                                                                                          0.32
                          montez
 Next steps:
                                                  View recommended plots
              Generate code with df_treatments
                                                                                New interactive sheet
df_treatments.isna().sum()
∓
                      0
       given_name
                      0
         surname
                      0
                      0
          auralin
                      0
         novodra
                      0
       hba1c_start
                      0
        hba1c_end
      hba1c_change 0
df_patients.head()
\overline{z}
         patient_id assigned_sex given_name
                                                  surname
                                                            address
                                                                            city
                                                                                     state zip_code country
                                                                                                                                      contact birthdat
                                                                 576
                                                              Brown
                                                                         Rancho
                                                                                                         United
                                                                                                                                      951-719-
      0
                             female
                                            Zoe
                                                   Wellish
                                                                                   California
                                                                                              92390.0
                                                                                                                                                7/10/197
                                                               Bear
                                                                        California
                                                                                                         States
                                                                                                                 9170ZoeWellish@superrito.com
                                                               Drive
                                                               2370
                                                                                                                 PamelaSHill@cuvox.de+1 (217)
                                                                                                         United
                   2
                                                                                              61812.0
                                                                                                                                                 4/3/196
      1
                             female
                                         Pamela
                                                           University
                                                                       Armstrong
                                                                                     Illinois
                                                                                                         States
                                                                                                                                     569-3204
                                                            Hill Road
 Next steps:
              Generate code with df_patients
                                                View recommended plots
                                                                              New interactive sheet
df_patients['phone_number'] = df_patients['contact'].str.extract(r'(\+?[0-9 \(\\)-]+)')
 df_patients['email'] = df_patients['contact'].str.extract(r'([a-zA-Z.\_%+-]+@[a-zA-Z.-]+\.[a-zA-Z]\{2,\})') ) 
df_patients = df_patients.drop(columns=['contact'])
```

 $df_patients$

₹		patient_id	assigned_sex	given_name	surname	address	city	state	zip_code	country	birthdate	weight	height	bm:
	0	1	female	Zoe	Wellish	576 Brown Bear Drive	Rancho California	California	92390.0	United States	7/10/1976	121.7	66	19.€
	1	2	female	Pamela	Hill	2370 University Hill Road	Armstrong	Illinois	61812.0	United States	4/3/1967	118.8	66	19.2
	2	3	male	Jae	Debord	1493 Poling Farm Road	York	Nebraska	68467.0	United States	2/19/1980	177.8	71	24.{
	3	4	male	Liêm	Phan	2335 Webster Street	Woodbridge	NJ	7095.0	United States	7/26/1951	220.9	70	31.7
	4	5	male	Tim	Neudorf	1428 Turkey Pen Lane	Dothan	AL	36303.0	United States	2/18/1928	192.3	27	26.′
	498	499	male	Mustafa	Lindström	2530 Victoria Court	Milton Mills	ME	3852.0	United States	4/10/1959	181.1	72	24.6
	499	500	male	Ruman	Bisliev	494 Clarksburg Park Road	Sedona	AZ	86341.0	United States	3/26/1948	239.6	70	34.4
	500	501	female	Jinke	de Keizer	649 Nutter Street	Overland Park	МО	64110.0	United States	1/13/1971	171.2	67	26.{
	1													>

Next steps: Generate code with df_patients View recommended plots New interactive sheet

df_treatments

₹		given_name	surname	auralin	novodra	hba1c_start	hba1c_end	hba1c_change	\blacksquare
	0	veronika	jindrová	41u - 48u	-	7.63	7.20	0.43	ılı
	1	elliot	richardson	-	40u - 45u	7.56	7.09	0.47	+/
	2	yukitaka	takenaka	-	39u - 36u	7.68	7.25	0.43	-
	3	skye	gormanston	33u - 36u	-	7.97	7.62	0.35	
	4	alissa	montez	-	33u - 29u	7.78	7.46	0.32	
	345	rovzan	kishiev	32u - 37u	-	7.75	7.41	0.34	
	346	jakob	jakobsen	=	28u - 26u	7.96	7.51	0.45	
	347	bernd	schneider	48u - 56u	-	7.74	7.44	0.30	
	348	berta	napolitani	=	42u - 44u	7.68	7.21	0.47	
	349	armina	sauvé	36u - 46u	-	7.86	7.40	0.46	
	350 rc	ws × 7 column	ıs						

Next steps: Generate code with df_treatments

• View recommended plots

New interactive sheet

Double-click (or enter) to edit

Melted dataframe

[#] Filter out rows where dose is missing

١

```
df_treatments = df_treatments[df_treatments['dose'] != '-']
```

```
# Display the resulting DataFrame
print(df_treatments)
```

_		given_name	surname	hba1c_start	hba1c_end	hba1c_change	١
	0	veronika	jindrová	7.63	7.20	0.43	
	3	skye	gormanston	7.97	7.62	0.35	
	6	sophia	haugen	7.65	7.27	0.38	
	7	eddie	archer	7.89	7.55	0.34	
	9	asia	woźniak	7.76	7.37	0.39	
	688	christopher	woodward	7.51	7.06	0.45	
	690	maret	sultygov	7.67	7.30	0.37	
	694	lixue	hsueh	9.21	8.80	0.41	
	696	jakob	jakobsen	7.96	7.51	0.45	
	698	berta	napolitani	7.68	7.21	0.47	

```
treatment_type
          auralin 41u - 48u
0
          auralin 33u - 36u
3
           auralin 37u - 42u
7
          auralin 31u - 38u
9
          auralin 30u - 36u
          novodra 55u - 51u
690
          novodra 26u - 23u
694
          novodra 22u - 23u
696
          novodra 28u - 26u
          novodra 42u - 44u
698
```

[350 rows x 7 columns]

Double-click (or enter) to edit

df_treatments.head()

_		given_name	surname	hba1c_start	hba1c_end	hba1c_change	treatment_type	dose	
	0	veronika	jindrová	7.63	7.20	0.43	auralin	41u - 48u	11.
	3	skye	gormanston	7.97	7.62	0.35	auralin	33u - 36u	
	6	sophia	haugen	7.65	7.27	0.38	auralin	37u - 42u	
	7	eddie	archer	7.89	7.55	0.34	auralin	31u - 38u	
	9	asia	woźniak	7.76	7.37	0.39	auralin	30u - 36u	

Next steps: Generate code with df_treatments

View recommended plots

New interactive sheet

df_patients.tail()

₹		patient_id	assigned_sex	given_name	surname	address	city	state	zip_code	country	birthdate	weight	height	bmi
	498	499	male	Mustafa	Lindström	2530 Victoria Court	Milton Mills	ME	3852.0	United States	4/10/1959	181.1	72	24.6
	499	500	male	Ruman	Bisliev	494 Clarksburg Park Road	Sedona	AZ	86341.0	United States	3/26/1948	239.6	70	34.4
	500	501	female	Jinke	de Keizer	649 Nutter Street	Overland Park	МО	64110.0	United States	1/13/1971	171.2	67	26.8
	501	502	female	Chidalu	Onyekaozulu	3652 Boone Crockett	Seattle	WA	98109.0	United States	2/13/1952	176.9	67	27.7

```
# Split the 'dose' column into low and high dose values
df_treatments[['low_dose', 'high_dose']] = df_treatments['dose'].str.split(' - ', expand=True)
```

[#] Remove the 'u' character and convert to integers
df_treatments['low_dose'] = df_treatments['low_dose'].str.replace('u', '').astype(int)
df_treatments['high_dose'] = df_treatments['high_dose'].str.replace('u', '').astype(int)

[#] Keep only relevant columns

final_df = df_treatments[['given_name', 'surname', 'treatment_type', 'low_dose', 'high_dose']]

Display the resulting DataFrame
print(final_df)

$\overline{\Rightarrow}$		given_name		treatment_type	low_dose	high_dose
	0	veronika	jindrová	auralin	41	48
	3	skye	gormanston	auralin	33	36
	6	sophia	haugen	auralin	37	42
	7	eddie	archer	auralin	31	38
	9	asia	woźniak	auralin	30	36
				• • •		
	688	christopher	woodward	novodra	55	51
	690	maret	sultygov	novodra	26	23
	694	lixue	hsueh	novodra	22	23
	696	jakob	jakobsen	novodra	28	26
	698	berta	napolitani	novodra	42	44

[350 rows x 5 columns]

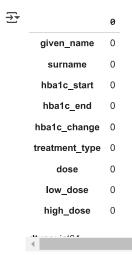
df_treatments.head()

_ →	g	given_name	surname	hba1c_start	hba1c_end	hba1c_change	treatment_type	dose	low_dose	high_dose	
	0	veronika	jindrová	7.63	7.20	0.43	auralin	41u - 48u	41	48	ıl.
	3	skye	gormanston	7.97	7.62	0.35	auralin	33u - 36u	33	36	
	6	sophia	haugen	7.65	7.27	0.38	auralin	37u - 42u	37	42	
	7	eddie	archer	7.89	7.55	0.34	auralin	31u - 38u	31	38	
	9	asia	woźniak	7.76	7.37	0.39	auralin	30u - 36u	30	36	
Next	step	s: Generat	e code with df	f treatments	€ View r	recommended plo	ots New interac	tive sheet			

df_treatments.tail()

-	given_name	surname	hba1c_start	hba1c_end	hba1c_change	treatment_type	dose	low_dose	high_dose	
688	christopher	woodward	7.51	7.06	0.45	novodra	55u - 51u	55	51	
690	maret	sultygov	7.67	7.30	0.37	novodra	26u - 23u	26	23	
694	lixue	hsueh	9.21	8.80	0.41	novodra	22u - 23u	22	23	
696	jakob	jakobsen	7.96	7.51	0.45	novodra	28u - 26u	28	26	
698	berta	napolitani	7.68	7.21	0.47	novodra	42u - 44u	42	44	
4										

df_treatments.isnull().sum()



df_treatments

1:561	•				v_assignmen		<u>_</u> ,p	, J		
	given_name	surname	hba1c_start	hba1c_end	hba1c_change	treatment_type	dose	low_dose	high_dose	
0	veronika	jindrová	7.63	7.20	0.43	auralin	41u - 48u	41	48	11.
3	skye	gormanston	7.97	7.62	0.35	auralin	33u - 36u	33	36	+/
6	sophia	haugen	7.65	7.27	0.38	auralin	37u - 42u	37	42	_
7	eddie	archer	7.89	7.55	0.34	auralin	31u - 38u	31	38	
9	asia	woźniak	7.76	7.37	0.39	auralin	30u - 36u	30	36	
688	christopher	woodward	7.51	7.06	0.45	novodra	55u - 51u	55	51	
690	maret	sultygov	7.67	7.30	0.37	novodra	26u - 23u	26	23	
694	lixue	hsueh	9.21	8.80	0.41	novodra	22u - 23u	22	23	
696	jakob	jakobsen	7.96	7.51	0.45	novodra	28u - 26u	28	26	
698	berta	napolitani	7.68	7.21	0.47	novodra	42u - 44u	42	44	

→ 1

df_treatments = df_treatments.applymap(lambda x: x.lower() if isinstance(x, str) else x)

<ipython-input-38-810883b0f417>:1: FutureWarning: DataFrame.applymap has been deprecated. Use DataFrame.map instead.

df_treatments = df_treatments.applymap(lambda x: x.lower() if isinstance(x, str) else x)

df_treatments.duplicated().sum()

→ 1

 $df_treatments$

	aivon nomo	cummama.	bbole stont	bbala and	hhale chance	+noo+mon+ +uno	dasa	lau daca	high doco
	given_name	Surname	IDATC_Start	nbarc_enu	nbarc_change	treatment_type	uose	TOW_GOSE	high_dose
0	veronika	jindrová	7.63	7.20	0.43	auralin	41u - 48u	41	48
3	skye	gormanston	7.97	7.62	0.35	auralin	33u - 36u	33	36
6	sophia	haugen	7.65	7.27	0.38	auralin	37u - 42u	37	42
7	eddie	archer	7.89	7.55	0.34	auralin	31u - 38u	31	38
9	asia	woźniak	7.76	7.37	0.39	auralin	30u - 36u	30	36
688	christopher	woodward	7.51	7.06	0.45	novodra	55u - 51u	55	51
690	maret	sultygov	7.67	7.30	0.37	novodra	26u - 23u	26	23
694	lixue	hsueh	9.21	8.80	0.41	novodra	22u - 23u	22	23
696	jakob	jakobsen	7.96	7.51	0.45	novodra	28u - 26u	28	26
698	berta	napolitani	7.68	7.21	0.47	novodra	42u - 44u	42	44

df_treatments = df_treatments.reset_index(drop=True)

df_treatments = df_treatments.drop(columns=['dose'])

 ${\tt df_treatments}$

0	given_name	surname	hinada adama							
0		341.114	nbaic_start	hba1c_end	hba1c_change	treatment_type	low_dose	high_dose		
	veronika	jindrová	7.63	7.20	0.43	auralin	41	48	11.	
1	skye	gormanston	7.97	7.62	0.35	auralin	33	36	*/	
2	sophia	haugen	7.65	7.27	0.38	auralin	37	42	-	
3	eddie	archer	7.89	7.55	0.34	auralin	31	38		
4	asia	woźniak	7.76	7.37	0.39	auralin	30	36		
345	christopher	woodward	7.51	7.06	0.45	novodra	55	51		
346	maret	sultygov	7.67	7.30	0.37	novodra	26	23		
347	lixue	hsueh	9.21	8.80	0.41	novodra	22	23		
348	jakob	jakobsen	7.96	7.51	0.45	novodra	28	26		
349	berta	napolitani	7.68	7.21	0.47	novodra	42	44		
steps		code with df_			ommended plots	New interactive				
Jnnamed inal_d	d: 0" in fi df = final_ tasets on g	iven_name an	ns: mns=["Unnamed d surname		s, on=["given_r	name", "surname"], how="le	ft")		
Innamedinal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_de	<pre>d: 0" in fi df = final_ tasets on g nts = df_tr sing values nts["advers updated da</pre>	nal_df.colum df.drop(colu iven_name an eatments.mer in adverse_ e_reaction"]	ns: mns=["Unnamed d surname ge(df_adverse reaction colu = df_treatme	e_conditions mn with "No nts["advers	Reaction" se_reaction"].f	Fillna("No Reacti		ft")		
nnamedinal_deltage dareatmened at misseatmened ethe eatmened at menedatmened at menedatmened at manage at menedatmened at manage at menedatmened at manage a	<pre>d: 0" in fi df = final_ tasets on g nts = df_tr sing values nts["advers updated da</pre>	nal_df.colum df.drop(colu iven_name an eatments.mer in adverse_ e_reaction"] taframe "final_dataf	ns: mns=["Unnamed d surname ge(df_adverse reaction colu = df_treatme	e_conditions mn with "No nts["advers	Reaction"	Fillna("No Reacti		ft")		
nnamedinal_decorption ge dareatmen l misseatmen e the eatmen play (df_tr	<pre>d: 0" in fi df = final_ tasets on g nts = df_tr sing values nts["advers updated da nts.to_csv(the first f</pre>	mal_df.colum df.drop(colu iven_name an eatments.mer in adverse_ e_reaction"] taframe "final_dataf ew rows ead())	ns: mns=["Unnamed d surname ge(df_adverse reaction colu = df_treatme rame_with_rea	e_conditions mn with "No ents["advers	o Reaction" se_reaction"].f	Fillna("No Reacti		ft")		
Innamedinal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_de	<pre>d: 0" in fi df = final_ tasets on g nts = df_tr sing values nts["advers updated da nts.to_csv(the first f reatments.h en_name eronika</pre>	nal_df.colum df.drop(colu iven_name an eatments.mer in adverse_ e_reaction"] taframe "final_dataf ew rows ead()) surname hb jindrová	ns: mns=["Unnamed d surname ge(df_adverse reaction colu = df_treatme rame_with_rea a1c_start ht 7.63	e_conditions mn with "No ents["advers ections.csv'	o Reaction" se_reaction"].f ', index=False) pa1c_change tre 0.43	Fillna("No Reacti) eatment_type \ auralin		ft")		
Innamedinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_de	<pre>d: 0" in fi df = final_ tasets on g nts = df_tr sing values nts["advers updated da nts.to_csv(the first f reatments.h en_name eronika skye go</pre>	nal_df.colum df.drop(colu iven_name an eatments.mer in adverse_ e_reaction"] taframe "final_dataf ew rows ead()) surname hb jindrová rmanston	ns: mns=["Unnamed d surname ge(df_adverse reaction colu = df_treatme rame_with_rea a1c_start ht 7.63 7.97	e_conditions mn with "No ints["advers actions.csv' halc_end ht 7.20 7.62	o Reaction" se_reaction"].f ', index=False) palc_change tre 0.43 0.35	Fillna("No Reacti) eatment_type \ auralin auralin		ft")		
Innamedinal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_definal_de	<pre>d: 0" in fi df = final_ tasets on g nts = df_tr sing values nts["advers updated da nts.to_csv(the first f reatments.h en_name eronika</pre>	nal_df.colum df.drop(colu iven_name an eatments.mer in adverse_ e_reaction"] taframe "final_dataf ew rows ead()) surname hb jindrová	ns: mns=["Unnamed d surname ge(df_adverse reaction colu = df_treatme rame_with_rea a1c_start ht 7.63	e_conditions mn with "No ents["advers ections.csv'	o Reaction" se_reaction"].f ', index=False) pa1c_change tre 0.43	Fillna("No Reacti) eatment_type \ auralin		ft")		
ge dareatmen l miss eatmen et the eatmen play (df_ti give 1 1 2 3	<pre>d: 0" in fi df = final_ tasets on g nts = df_tr sing values nts["advers updated da nts.to_csv(the first f reatments.h en_name eronika skye go sophia</pre>	nal_df.colum df.drop(colu iven_name an eatments.mer in adverse_ e_reaction"] taframe "final_dataf ew rows ead()) surname hb jindrová rmanston haugen	ns: mns=["Unnamed d surname ge(df_adverse reaction colu	e_conditions mn with "No ents["advers actions.csv' palc_end ht 7.20 7.62 7.27	Particular Reaction R	Fillna("No Reacti) eatment_type \ auralin auralin auralin auralin		ft")		
rnnamed rinal real real real real real real real re	d: 0" in fi df = final_ tasets on g nts = df_tr sing values nts["advers updated da nts.to_csv(the first f reatments.h en_name eronika skye go sophia eddie asia w_dose hig	mal_df.colum df.drop(colu iven_name an eatments.mer in adverse_ e_reaction"] taframe "final_dataf ew rows ead()) surname ht jindrová rmanston haugen archer woźniak h_dose adver	ns: mns=["Unnamed d surname ge(df_adverse reaction colu = df_treatme rame_with_rea a1c_start ht 7.63 7.97 7.65 7.89 7.76 se_reaction	e_conditions mn with "No ents["advers actions.csv' palc_end ht 7.20 7.62 7.27 7.55	palc_change tre 0.43 0.35 0.38 0.34	Fillna("No Reacti) eatment_type \ auralin auralin auralin auralin		ft")		
Innamedinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_derinal_de	d: 0" in fi df = final_ tasets on g nts = df_tr sing values nts["advers updated da nts.to_csv(the first f reatments.h en_name eronika skye go sophia eddie asia w_dose hig 41	mal_df.colum df.drop(colu iven_name an eatments.mer in adverse_ e_reaction"] taframe "final_dataf ew rows ead()) surname ht jindrová rmanston haugen archer woźniak h_dose adver 48	ns: mns=["Unnamed d surname ge(df_adverse reaction colu = df_treatme rame_with_rea a1c_start ht 7.63 7.97 7.65 7.89 7.76 se_reaction No Reaction	e_conditions mn with "No ents["advers actions.csv' palc_end ht 7.20 7.62 7.27 7.55	palc_change tre 0.43 0.35 0.38 0.34	Fillna("No Reacti) eatment_type \ auralin auralin auralin auralin		ft")		
Innamed in al_correction al_co	d: 0" in fi df = final_ tasets on g nts = df_tr sing values nts["advers updated da nts.to_csv(the first f reatments.h en_name eronika skye go sophia eddie asia w_dose hig	mal_df.colum df.drop(colu iven_name an eatments.mer in adverse_ e_reaction"] taframe "final_dataf ew rows ead()) surname ht jindrová rmanston haugen archer woźniak h_dose adver	ns: mns=["Unnamed d surname ge(df_adverse reaction colu = df_treatme rame_with_rea a1c_start ht 7.63 7.97 7.65 7.89 7.76 se_reaction	e_conditions mn with "No ents["advers actions.csv' palc_end ht 7.20 7.62 7.27 7.55	palc_change tre 0.43 0.35 0.38 0.34	Fillna("No Reacti) eatment_type \ auralin auralin auralin auralin		ft")		
Innamed in al_derivation of the control of the cont	d: 0" in fi df = final_ tasets on g nts = df_tr sing values nts["advers updated da nts.to_csv(the first freatments.h en_name eronika skye go sophia eddie asia w_dose hig 41 33	nal_df.colum df.drop(colu iven_name an eatments.mer in adverse_ e_reaction"] taframe "final_dataf ew rows ead()) surname hb jindrová rmanston haugen archer woźniak h_dose adver 48 36	ns: mns=["Unnamed d surname ge(df_adverse reaction colu = df_treatme rame_with_rea alc_start ht 7.63 7.97 7.65 7.89 7.76 se_reaction No Reaction No Reaction	e_conditions mn with "No ents["advers actions.csv' palc_end ht 7.20 7.62 7.27 7.55	palc_change tre 0.43 0.35 0.38 0.34	Fillna("No Reacti) eatment_type \ auralin auralin auralin auralin		ft")		

https://colab.research.google.com/drive/1aoM0eimDwgHJUPxT_aGnLDKbhJW8ePLj#scrollTo=VdOZZDKgX_hU

adverse_reaction 348 non-null

dtypes: float64(3), int64(2), object(4)

memory usage: 27.2+ KB

```
₹
                          surname hba1c_start hba1c_end hba1c_change treatment_type low_dose high_dose
           given_name
                                                                                                                   adverse_reaction
 Next 0
              CYCRENIKE code WINDY OYA treatment 3 63 C View 20 commended 9 163s
                                                                                                                         No Reaction
                                                                             New in the radio of the sheet 41
                                           7.97
                                                                                                            36
                 skye gormanston
                                                      7.62
                                                                                   auralin
                                                                                                                         No Reaction
# Check for outliers using IQR method for numerical columns
Q1 = df_treatments[["low_dose", "high_dose"]].quantile(0.25)
Q3 = df_treatments[["low_dose", "high_dose"]].quantile(0.75)
IQR = Q3 - Q1
# Define outlier boundaries
lower\_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR
# Filter out rows with outliers
df_treatments = df_treatments[
    ~((df_treatments["low_dose"] < lower_bound["low_dose"]) | (df_treatments["low_dose"] > upper_bound["low_dose"])) &
     \sim ((df\_treatments["high\_dose"] < lower\_bound["high\_dose"]) \mid (df\_treatments["high\_dose"] > upper\_bound["high\_dose"])) 
]
     SOUTOWS A 9 COMMUNS
df_treatments.info()
    <class 'pandas.core.frame.DataFrame'>
     Index: 348 entries, 0 to 349
     Data columns (total 9 columns):
      # Column
                             Non-Null Count Dtype
      0
          given_name
                             348 non-null
                                             object
                             348 non-null
      1
          surname
                                             object
      2
          hba1c_start
                             348 non-null
                                             float64
          hba1c_end
                             348 non-null
                                             float64
          hba1c_change
                             348 non-null
                                             float64
          treatment_type
                             348 non-null
                                             object
                             348 non-null
          low_dose
                                             int64
                             348 non-null
          high dose
                                             int64
```

object