

```
[In (1)]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
[In (2)]: patients_df = pd.read_csv(r"C:\Users\Jai Shree Shyam\Desktop\data cleaning files\patients.csv")
treatments_df = pd.read_csv(r"C:\Users\Jai Shree Shyam\Desktop\data cleaning files\treatments.csv")
adverse_reaction_df = pd.read_csv(r"C:\Users\Jai Shree Shyam\Desktop\data cleaning files\adverse_reactions.csv")
treatments_cut_df = pd.read_csv(r"C:\Users\Jai Shree Shyam\Desktop\data cleaning files\treatments_cut.csv")
```

```
[In (3)]: print(patients.shape)
print(treatments.shape)
print(adverse_reaction.shape)
print(treatments_cut.shape)
```

```
(503, 14)
(289, 7)
(289, 7)
(279, 7)
```

```
[In (4)]: # Export data for manual assessment
with pd.ExcelWriter("clinical_trials.xlsx") as writer:
    patients.to_excel(writer, sheet_name="patients")
    treatments.to_excel(writer, sheet_name="treatments")
    adverse_reaction.to_excel(writer, sheet_name="adverse_reactions")
    treatments_cut.to_excel(writer, sheet_name="adverse_reaction")
```

Data Description

This is a dataset about 503patients of which 350 patients participated in a clinical trial. None of the patients were using Novodra (a popular injectable insulin) or auralin (the oral insulin) being researched as their primary source of insulin before. All were experiencing elevated hba1c levels.

All 350 patients were treated with Novodra to establish a baseline. Hba1c level and insulin dose. After 4 weeks which is enough time to capture all the change in hba1c that can be attributed by the switch to auralin or novodra.

175 patients switched to auralin at 24 weeks: 175 patients switched or remain to novodra for 24 weeks.

Data about patients feeling some adverse effect is also recorded.

Description about columns of the data.

Table- patients

- patient_id- patient id of all the patient is given | 1 to 503 total 503 patients are there.
- assigned_sex- patient sex is given ["Male", "Female"]
- given_name- name of the patient is given.
- surname- surname of each patient is given.
- address- Address is given starting from a digit and followed by alphabets.
- city- city- name is given of particular person living in that locality.
- State- State- name is given of particular person belongs to that state.
- zip_code- Zip code of his or her locality is given and must be of 5 digit.
- country- country name is given, United states is given because all the data regarding patients is of united states.
- birthdate- date of birth of the patients is given, in which year they were born.
- weight- weight of the patients is given in (pounds)
- height- height of the patients is given in (inches)
- bmi- body mass index is given of every patient.

Table- treatments and treatment_cut

- given_name-volunteer name is given who take part in clinical trials.
- surname- surname of particular volunteer is given who take part in clinical trials.
- auralin- auralin is the name of oral insulin drug which is given to volunteer consisting usage of auralin (41u - 48u) is given starting usage before 24 weeks and after 24 weeks.
- novodra- novodra is the name of the injectable insulin drug which is given to volunteer consisting usage of novodra (41u - 48u) is given starting usage before 24 weeks and after 24 weeks.
- hba1c_start- it measures the level of sugar present in hemoglobin. hba1c_start means the rate before the volunteer has.
- hba1c_end- it measures the level of sugar present in hemoglobin. hba1c_end means the rate after the volunteer has.
- hba1c_change- it measures the change in hba1c_start and hba1c_end after 24 weeks after clinical trials is completed.

Table- adverse_reaction

- given_name- volunteers name is given who participated in the clinical trials.
- surname- surname is given of the volunteers who participated in the clinical trials.
- adverse_reaction- adverse_reaction is given to certain volunteer after taking trials of the doses.

Data Assessment

Dirty Data [content issue]

Table [Patient]

- patient_name column has row number 9 has david must be correct as David. [Accuracy]
- state column has some name which is fully and some are in abbrevations. [Consistency]
- zip_code column has 6 values 5 digit but have 4 digit is given. [Validity]
- some columns have no values in address, city, state, zip_code, country and contact consisting of 12 patients [completeness]
- john Doe patients is duplicated 5 times in the data have to delete it [validity]
- patient_id 211 female Camila Zateva has weight around 48 pounds which is not possible [Accuracy]
- patient_id 5 male Tim Neufort has height of 27 inches which is not possible [Accuracy]

Table [treatments and treatments_cut]

- given_name and surname having all small letter have to change in capitalize first word [consistency]
- novodra and novodra having additional unit have to delete it [validity]
- hba1c change column have some Missing value [completeness]
- patient_id 139 joseph day is repeated or duplicated in the data [validity]
- treatments and treatments_cut column have some misvalue return in hba1c column [Accuracy]

Table [adverse_reaction]

- column Name given, surname must be in capitalize first word. [consistency].

Messy Data [structural issue]

Table [Patient]

- Column contact has mixed phone no. and email id have to separate one column into two column phone no. and email id.

Table [treatments and treatments_cut]

- Column auralin and novodra having starting and ending both dosage have to split into two different columns.
- auralin and Novodra both have to in one column.

Table [adverse_reaction]

- adverse_reaction doesn't have seprate table must have to merge with treatment table.

```
[In (5)]: # Making copy of all the original files.
patients_df = patients.copy()
treatments_df = treatments.copy()
adverse_reactions_df = adverse_reactions.copy()
```

Patients_df

```
[In (6)]: # Patients
patients_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 503 entries, 0 to 502
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
--  --
0   patient_id  503 non-null     int64
1   assigned_sex  503 non-null     object
2   given_name   503 non-null     object
3   surname      503 non-null     object
4   address      491 non-null     object
5   city         491 non-null     object
6   state        491 non-null     object
7   zip_code     491 non-null     float64
8   country      491 non-null     object
9   contact      491 non-null     object
10  birthdate    503 non-null     object
11  weight       503 non-null     float64
12  height       503 non-null     float64
13  bmi          503 non-null     float64
dtypes: float64(3), int64(2), object(9)
memory usage: 55.1+ KB
```

```
[In (7)]: # patients_df[patients_df["address"] !=na]]
# These 12 patients has missing values in all (address, city,state,zip_code, country, contact)
```

	patient_id	assigned_sex	given_name	surname	address	city	state	zip_code	country	contact	birthdate	weight	height	bmi
209	210	female	Lata	Elderharov	NaN	NaN	NaN	NaN	NaN	NaN	8/14/1950	143.4	62	26.2
219	220	male	My	Qaym	NaN	NaN	NaN	NaN	NaN	NaN	4/9/1978	237.8	69	35.1
230	231	female	Elisabeth	Koudsten	NaN	NaN	NaN	NaN	NaN	NaN	9/23/1976	165.9	63	29.4
234	235	female	Martina	Tomáriková	NaN	NaN	NaN	NaN	NaN	NaN	4/7/1936	195.5	65	33.2
242	243	male	John	O'Brian	NaN	NaN	NaN	NaN	NaN	NaN	2/25/1957	205.3	74	26.4
249	250	male	Benjamin	Mehler	NaN	NaN	NaN	NaN	NaN	NaN	10/30/1951	146.5	69	21.6
257	258	male	Jin	Hung	NaN	NaN	NaN	NaN	NaN	NaN	5/17/1995	231.7	69	34.2
264	265	female	Wafiyah	Asfour	NaN	NaN	NaN	NaN	NaN	NaN	1/10/1989	150.9	63	28.1
269	270	female	Flavia	Florentino	NaN	NaN	NaN	NaN	NaN	NaN	1/10/1937	175.2	61	33.1
278	279	female	Genessa	Cabán	NaN	NaN	NaN	NaN	NaN	NaN	12/16/1962	124.3	69	18.4
286	287	male	Lewis	Veebo	NaN	NaN	NaN	NaN	NaN	NaN	4/1/1979	155.3	68	23.6
296	297	female	Chi	Lâm	NaN	NaN	NaN	NaN	NaN	NaN	5/14/1990	181.1	63	32.1

```
[In (8)]: # patient_duplicated
patients_df[patients_df.duplicated(subset=["given_name", "surname"])]
# John Doe patient is duplicated 5 times in the data which we can't do anything to lean names.
```

```
[In (8)]: patient_id assigned_sex given_name surname address city state zip_code country contact birthdate weight height bmi
229 230 male John Doe 123 Main Street New York NY 12345.0 United States johndoe@gmail.com1234567890 1/1/1975 180.0 72 24.4
237 238 male John Doe 123 Main Street New York NY 12345.0 United States johndoe@gmail.com1234567890 1/1/1975 180.0 72 24.4
244 245 male John Doe 123 Main Street New York NY 12345.0 United States johndoe@gmail.com1234567890 1/1/1975 180.0 72 24.4
251 252 male John Doe 123 Main Street New York NY 12345.0 United States johndoe@gmail.com1234567890 1/1/1975 180.0 72 24.4
277 278 male John Doe 123 Main Street New York NY 12345.0 United States johndoe@gmail.com1234567890 1/1/1975 180.0 72 24.4
```

```
[In (9)]: # patient_df.describe()
patients_df.describe()
```

	patient_id	zip_code	weight	height	bmi
count	503.000000	491.000000	503.000000	503.000000	503.000000
mean	250.000000	49884.11828	173.434990	66.634195	27.488997
std	146.347859	30265.807442	33.916741	4.411297	5.278438
min	1.000000	1002.000000	48.000000	27.000000	17.100000
25%	126.500000	21340.500000	148.300000	63.000000	23.300000
50%	262.000000	48957.000000	175.300000	67.000000	27.200000
75%	377.500000	75679.000000	195.500000	70.000000	31.750000
max	503.000000	97810.000000	255.900000	79.000000	37.700000

```
[In (10)]: # patients_df table consists weight as min of 48 pounds.
```

```
patients_df[patients_df["weight"] !=49,000000]]
# weight is not correctly given, should be change as it affects the distribution of the data and mean.
```

	patient_id	assigned_sex	given_name	surname	address	city	state	zip_code	country	contact	birthdate	weight	height	bmi
210	211	female	Camila	Zateva	4689 Briarleaf Lane	Woonosr	OH	44691.0	United States	330-202-2145CamilaZateva@spermo.com	11/20/1938	48.8	63	19.1

```
[In (11)]: # patients table consists height as min of 27 inches
patients_df[patients_df["height"] !=27.000000]]
# height is 27 inches which is not possible as comparison to weight which is 182.3 pounds
```

	patient_id	assigned_sex	given_name	surname	address	city	state	zip_code	country	contact	birthdate	weight	height	bmi
4	5	male	Tim	Neufort	1428 Turkey Pen Lane	Dodhan	AL	36303.0	United States	334-515-7487TimNeufort@curvex.de	2/18/1928	182.3	27	26.1

treatments_df and treatments_cut_df

```
[In (12)]: print(treatments_df.head())
print(treatments_cut_df.head())
```

```
given_name surname auralin novodra hba1c_start hba1c_end \
0 veronika jindrová 41u - 48u - 7.63 7.20
1 elliot richardson - 48u - 48u 7.06 7.89
2 yukitaka takemaka 30u - 36u - 7.76 7.44
3 skye gormanston 33u - 36u - 7.07 7.42
4 alissa montez - 33u - 28u 7.78 7.46
```

```
hba1c_change
0 NaN
1 8.37
2 NaN
3 8.35
4 8.32
```

```
given_name surname auralin novodra hba1c_start hba1c_end \
0 jada reynolds 22u - 36u - 7.16 7.22
1 iunungga heilmann 67u - 67u - 7.85 7.45
2 alvin svensson 36u - 39u - 7.78 7.34
3 tde luing 41u - 64u 7.64 7.22
4 aranda ribeiro 36u - 44u - 7.85 7.47
```

```
hba1c_change
0 8.34
1 NaN
2 NaN
3 8.32
4 8.38
```

```
[In (13)]: print(treatments_df.tail())
print(treatments_cut_df.tail())
```

```
given_name surname auralin novodra hba1c_start hba1c_end \
275 albina zeticic 45u - 55u - 7.93 7.73
276 john reicheimann 49u - 49u - 9.04 8.67
277 matthea lillibee 23u - 36u - 7.04 6.67
278 valille prince 31u - 36u - 7.64 7.28
279 samuel gubrandsson 53u - 56u - 8.00 7.64
```

```
hba1c_change
275 0.28
276 NaN
277 0.37
278 0.36
279 0.38
```

```
given_name surname auralin novodra hba1c_start hba1c_end \
65 rowan kistner 32u - 37u - 7.75 7.41
66 jakob jakobsen 28u - 28u - 7.96 7.51
67 bernd schneider 48u - 56u - 7.74 7.44
68 betta napollitan 42u - 42u - 7.68 7.23
69 armina saure 36u - 46u - 7.86 7.40
```

```
hba1c_change
65 8.34
66 8.95
67 8.38
68 NaN
69 NaN
```

```
[In (14)]: treatments_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 289 entries, 0 to 278
Data columns (total 7 columns):
#   Column      Non-Null Count  Dtype
--  --
0   given_name  289 non-null     object
1   surname     289 non-null     object
2   auralin     289 non-null     object
3   novodra     289 non-null     object
4   hba1c_start 289 non-null     float64
5   hba1c_end   289 non-null     float64
6   hba1c_change 174 non-null     float64
dtypes: float64(3), object(4)
memory usage: 15.4+ KB
```

```
[In (15)]: treatments_cut_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8 entries, 0 to 7
Data columns (total 7 columns):
#   Column      Non-Null Count  Dtype
--  --
0   given_name  78 non-null     object
1   surname     78 non-null     object
2   auralin     78 non-null     object
3   novodra     78 non-null     object
4   hba1c_start 78 non-null     float64
5   hba1c_end   78 non-null     float64
6   hba1c_change 42 non-null     float64
dtypes: float64(3), object(4)
memory usage: 4.4+ KB
```

```
[In (16)]: treatments_df[treatments_df.duplicated()]
# joseph day volunteer is repeated one line have to delete it.
```

```
[In (16)]: given_name surname auralin novodra hba1c_start hba1c_change
136 joseph day 29u-36u - 7.7 7.19 NaN
```

```
[In (17)]: treatments_df[treatments_df.duplicated(subset=["given_name","surname"])]
given_name surname auralin novodra hba1c_start hba1c_change
136 joseph day 29u-36u - 7.7 7.19 NaN
```

```
[In (18)]: treatments_cut_df[treatments_cut_df.duplicated()]
```

```
[In (18)]: given_name surname auralin novodra hba1c_start hba1c_change
```

```
[In (19)]: treatments_cut_df[treatments_cut_df.duplicated(subset=["given_name","surname"])]
```

```
[In (20)]: # treatment [describe()]
treatments_df.describe()
```

	hba1c_start	hba1c_end	hba1c_change
count	280.000000	280.000000	171.000000
mean	7.895929	7.588286	0.540623
std	0.566838	0.560672	0.279555
min	7.500000	7.010000	0.200000
25%	7.600000	7.270000	0.300000
50%	7.600000	7.400000	0.300000
75%	7.970000	7.570000	0.870000
max	9.950000	9.800000	0.900000

```
[In (21)]: treatments_df[treatments_df["hba1c_start"] !=9.950000]]
```

	given_name	surname	auralin	novodra	hba1c_start	hba1c_end	hba1c_change
166	ernie	allen	36u - 42u	-	9.95	9.98	0.37

```
[In (22)]: treatments_df.sort_values("hba1c_start")
# if some normal as many volunteer is having hba1c_start in 8 something range.
```

	given_name	surname	auralin	novodra	hba1c_start	hba1c_end	hba1c_change
270	mila	mattsson	30u - 43u	-	7.50	7.17	0.33
113	kari	kallanen	36u - 43u	-	7.50	7.11	NaN
126	jenna	vitelnska	- 22u - 23u	-	7.50	7.08	0.92
53	nasser	mansour	- 33u - 33u	-	7.51	7.06	0.95
105	fray	sheppard	- 31u - 30u	-	7.51	7.17	0.34
--	--	--	--	--	--	--	--
25	berndt	konami	- 44u - 43u	-	9.82	9.40	0.82
171	juzyna	kowalczyk	24u - 34u	-	9.84	9.44	NaN
81	robert	wagner	43u - 49u	-	9.84	9.52	0.32
75	machene	mosky	- 44u - 45u	-	9.87	9.48	0.99
166	ernie	allen	36u - 42u	-	9.95	9.98	0.37

280 rows x 7 columns

```
[In (23)]: treatments_df.sort_values("hba1c_change",na_position="first")
# there is incorrect values of hba1c_change that is difference between hba1c_start and hba1c_end
```

	given_name	surname	auralin	novodra	hba1c_start	hba1c_end	hba1c_change
0	veronika	jindrová	41u - 48u	-	7.63	7.20	NaN
2	yukitaka	takemaka	- 39u - 36u	-	7.68	7.25	NaN
9	essa	wozniak	30u-36u	-	7.76	7.37	NaN
10	joseph	day	29u-36u	-	7.70	7.19	NaN
--	--	--	--	--	--	--	--
49	jackson	addison	- 43u-42u	-	7.99	7.51	0.99
17	ghis	can	- 36u-36u	-	7.68	7.40	0.98
32	laura	erichmann	- 43u - 48u	-	7.95	7.46	0.99
265	wu	shung	- 47u-48u	-	7.61	7.12	0.99
138	governa	rocha	- 23u-23u	-	7.87	7.38	0.99

280 rows x 7 columns

```
[In (24)]: # treatments_cut_df
treatments_cut_df.describe()
```

```
# Same problem with treatments_cut_df
```

	hba1c_start	hba1c_end	hba1c_change
count	70.000000	70.000000	42.000000
mean	7.838800	7.443143	0.518810
std	0.423007	0.418706	0.270719
min	7.510000	7.020000	0.280000
25%	7.640000	7.232500	0.340000
50%	7.730000	7.342500	0.370000
75%	7.860000	7.467500	0.970000
max	9.910000	9.460000	0.970000

Data Cleaning

Patients table

```
[In (25)]: patients_df[patients_df["address"] !=na]]
# filling nan values in address not possible as, we do it by our own
# we can't delete the regarding row missing data of which we can't do anything to lean names.
```

2	slve	gornaston	7.95	7.62	0.35	auralin	33a-35u	33	36
3	sophia	hugen	7.67	7.27	0.38	auralin	37u-42u	37	42
4	eddie	arther	7.89	7.55	0.34	auralin	31u-38u	31	38
5	asia	wodrak	7.76	7.37	0.39	auralin	30u-35u	30	35
...
346	christoph	woodward	7.51	7.06	0.45	novotna	55u-51u	55	51
347	maart	subtorg	7.67	7.30	0.37	novotna	26u-23u	26	23
348	lurec	luseh	9.21	8.80	0.41	novotna	22u-23u	22	23
349	jacob	jakobson	7.96	7.51	0.45	novotna	28u-26u	28	26
350	berth	napstani	7.68	7.21	0.47	novotna	42u-44u	42	44

350 rows x 9 columns

```
In [4]: treatments_df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 350 entries, 1 to 358
```