

```
In [1]: import pandas as pd
import numpy as np

In [2]: df1 = pd.read_csv("../heart.csv")
df2 = pd.read_csv("../AIQuality.csv", sep=";")
```

In [3]: df1.head()

	Age	Sex	ChestPainType	RestingBP	Cholesterol	FastingBS	RestingECG	MaxHR	ExerciseAngina	Oldpeak	ST_Slope	HeartDisease
0	40	M	ATA	140	289	0	Normal	172	N	0.0	Up	0
1	49	F	NAP	160	180	0	Normal	156	N	1.0	Flat	1
2	37	M	ATA	130	283	0	ST	98	N	0.0	Up	0
3	48	F	ASY	138	214	0	Normal	108	Y	1.5	Flat	1
4	54	M	NAP	150	195	0	Normal	122	N	0.0	Up	0

In [4]: df2.head()

	Date	Time	CO(CT)	PT08.S1(CO)	NMHC(CT)	CH6(CT)	PT08.S2(NMHC)	NOx(CT)	PT08.S3(NOX)	NO2(CT)	PT08.S4(NO2)	PT08.S5(O3)	T	RH	AH	Unnamed: 15	Unnamed: 16
0	1003/2004	18.00.00	2.6	1360.0	150.0	11.9	1046.0	166.0	1056.0	113.0	1692.0	1268.0	13.6	48.9	0.7578	NaN	NaN
1	1003/2004	19.00.00	2	1292.0	112.0	9.4	955.0	103.0	1174.0	92.0	1559.0	972.0	13.3	47.7	0.7255	NaN	NaN
2	1003/2004	20.00.00	2.2	1402.0	88.0	8.0	939.0	131.0	1140.0	114.0	1555.0	1074.0	11.0	54.0	0.7502	NaN	NaN
3	1003/2004	21.00.00	2.2	1376.0	80.0	8.2	948.0	172.0	1092.0	122.0	1584.0	1203.0	11.0	60.0	0.7867	NaN	NaN
4	1003/2004	22.00.00	1.6	1272.0	51.0	6.5	836.0	131.0	1205.0	116.0	1490.0	1110.0	11.2	59.6	0.7888	NaN	NaN

## DATA CLEANING

In [5]: df1.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 918 entries, 0 to 917
Data columns (total 12 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   Age         918 non-null    int64
 1   Sex         918 non-null    object
 2   ChestPainType  918 non-null    object
 3   RestingBP   918 non-null    int64
 4   Cholesterol  918 non-null    int64
 5   FastingBS   918 non-null    int64
 6   RestingECG  918 non-null    object
 7   MaxHR       918 non-null    int64
 8   ExerciseAngina  918 non-null    object
 9   Oldpeak     918 non-null    float64
10  ST_Slope    918 non-null    object
11  HeartDisease  918 non-null    int64
dtypes: float64(1), int64(6), object(5)
memory usage: 86.2+ KB
```

In [6]: df2.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9471 entries, 0 to 9470
Data columns (total 17 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   Date        9357 non-null    object
 1   Time        9357 non-null    object
 2   CO(CT)      9357 non-null    object
 3   PT08_S1(CO) 9357 non-null    float64
 4   NMHC(CT)    9357 non-null    float64
 5   CH6(CT)     9357 non-null    object
 6   PT08_S2(NMHC) 9357 non-null    float64
 7   NOx(CT)     9357 non-null    float64
 8   PT08_S3(NOX) 9357 non-null    float64
 9   NO2(CT)     9357 non-null    float64
10  PT08_S4(NO2) 9357 non-null    float64
11  PT08_S5(O3) 9357 non-null    float64
12  T            9357 non-null    object
13  RH           9357 non-null    object
14  AH           9357 non-null    object
15  Unnamed: 15  0 non-null      float64
16  Unnamed: 16  0 non-null      float64
dtypes: float64(19), object(7)
memory usage: 1.2+ MB
```

In [7]: df1.isnull().sum()

	Age	Sex	ChestPainType	RestingBP	Cholesterol	FastingBS	RestingECG	MaxHR	ExerciseAngina	Oldpeak	ST_Slope	HeartDisease
	0	0	0	0	0	0	0	0	0	0	0	0

In [8]: df2.isnull().sum()

	Date	Time	CO(CT)	PT08_S1(CO)	NMHC(CT)	CH6(CT)	PT08_S2(NMHC)	NOx(CT)	PT08_S3(NOX)	NO2(CT)	PT08_S4(NO2)	PT08_S5(O3)	T	RH	AH	Unnamed: 15	Unnamed: 16
	114	114	114	114	114	114	114	114	114	114	114	114	114	114	114	114	114
	9471	9471	9471	9471	9471	9471	9471	9471	9471	9471	9471	9471	9471	9471	9471	9471	9471
	16	16	16	16	16	16	16	16	16	16	16	16	16	16	16	16	16
	dtype: int64																

In [9]: df2.drop(["Unnamed: 15", "Unnamed: 16"], axis=1, inplace=True)

In [10]: df2.sample(10)

	Date	Time	CO(CT)	PT08_S1(CO)	NMHC(CT)	CH6(CT)	PT08_S2(NMHC)	NOx(CT)	PT08_S3(NOX)	NO2(CT)	PT08_S4(NO2)	PT08_S5(O3)	T	RH	AH
4381	301/2004	18.00.00	6.3	1408.0	-200.0	38.0	1502.0	692.0	640.0	178.0	1843.0	1598.0	11.6	64.3	0.8789
4638	286/2004	18.00.00	4.1	1236.0	-200.0	21.3	1332.0	641.0	601.0	152.0	1683.0	1590.0	17.8	45.7	0.8224
1408	1408/2004	18.00.00	0.9	964.0	-200.0	5.1	769.0	40.0	873.0	56.0	1546.0	1572.0	35.2	32.1	1.7888
3262	2407/2004	16.00.00	1.1	1012.0	-200.0	5.2	775.0	49.0	899.0	53.0	1547.0	562.0	36.6	28.2	1.6970
698	06/04/2004	20.00.00	4.3	1319.0	544.0	15.8	1172.0	232.0	746.0	136.0	1699.0	1425.0	15.8	39.8	0.7096
5421	22/03/2004	18.00.00	3.7	1476.0	-200.0	22.0	1352.0	365.0	474.0	112.0	1996.0	1568.0	26.4	48.8	1.6928
575	03/04/2004	17.00.00	-200	1398.0	-200.0	15.2	1156.0	-200.0	812.0	-200.0	1722.0	1076.0	25.1	27.1	0.8496
8411	24/02/2005	05.00.00	0.6	883.0	-200.0	1.2	518.0	87.0	1135.0	81.0	862.0	606.0	3.3	84.5	0.6612
8646	06/03/2005	00.00.00	2.4	1201.0	-200.0	9.2	949.0	358.0	665.0	170.0	1220.0	1260.0	6.8	72.1	0.7143
3077	16/07/2004	23.00.00	2.2	1090.0	-200.0	13.6	1103.0	142.0	695.0	119.0	1686.0	1235.0	25.7	33.5	1.0918

In [11]: df2.droptna(inplace=True)

In [12]: df2.tail()

	Date	Time	CO(CT)	PT08_S1(CO)	NMHC(CT)	CH6(CT)	PT08_S2(NMHC)	NOx(CT)	PT08_S3(NOX)	NO2(CT)	PT08_S4(NO2)	PT08_S5(O3)	T	RH	AH
9352	04/04/2005	10.00.00	3.1	1314.0	-200.0	13.5	1101.0	477.0	539.0	190.0	1374.0	1729.0	21.9	29.3	0.7568
9353	04/04/2005	11.00.00	2.4	1163.0	-200.0	11.4	1027.0	353.0	604.0	179.0	1264.0	1268.0	24.3	23.7	0.7119
9354	04/04/2005	12.00.00	2.4	1142.0	-200.0	12.4	1063.0	293.0	603.0	175.0	1241.0	1092.0	26.9	18.3	0.6406
9355	04/04/2005	13.00.00	2.1	1003.0	-200.0	9.5	961.0	235.0	702.0	156.0	1041.0	770.0	28.3	13.5	0.5139
9356	04/04/2005	14.00.00	2.2	1071.0	-200.0	11.9	1047.0	265.0	654.0	168.0	1129.0	816.0	28.5	13.1	0.5028

In [13]: df2.drop\_duplicates()

	Date	Time	CO(CT)	PT08_S1(CO)	NMHC(CT)	CH6(CT)	PT08_S2(NMHC)	NOx(CT)	PT08_S3(NOX)	NO2(CT)	PT08_S4(NO2)	PT08_S5(O3)	T	RH	AH	
0	1003/2004	18.00.00	2.6	1360.0	150.0	11.9	1046.0	166.0	1056.0	113.0	1692.0	1268.0	13.6	48.9	0.7578	
1	1003/2004	19.00.00	2	1292.0	112.0	9.4	955.0	103.0	1174.0	92.0	1559.0	972.0	13.3	47.7	0.7255	
2	1003/2004	20.00.00	2.2	1402.0	88.0	8.0	939.0	131.0	1140.0	114.0	1555.0	1074.0	11.0	54.0	0.7502	
3	1003/2004	21.00.00	2.2	1376.0	80.0	8.2	948.0	172.0	1092.0	122.0	1584.0	1203.0	11.0	60.0	0.7867	
4	1003/2004	22.00.00	1.6	1272.0	51.0	6.5	836.0	131.0	1205.0	116.0	1490.0	1110.0	11.2	59.6	0.7888	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
9352	04/04/2005	10.00.00	3.1	1314.0	-200.0	13.5	1101.0	472.0	539.0	190.0	1374.0	1729.0	21.9	29.3	0.7568	
9353	04/04/2005	11.00.00	2.4	1163.0	-200.0	11.4	1027.0	353.0	604.0	179.0	1264.0	1268.0	24.3	23.7	0.7119	
9354	04/04/2005	12.00.00	2.4	1142.0	-200.0	12.4	1063.0	293.0	603.0	175.0	1241.0	1092.0	26.9	18.3	0.6406	
9355	04/04/2005	13.00.00	2.1	1003.0	-200.0	9.5	961.0	235.0	702.0	156.0	1041.0	770.0	28.3	13.5	0.5139	
9356	04/04/2005	14.00.00	2.2	1071.0	-200.0	11.9	1047.0	265.0	654.0	168.0	1129.0	816.0	28.5	13.1	0.5028	

9357 rows × 15 columns

In [14]: df1.drop\_duplicates()

	Age	Sex	ChestPainType	RestingBP	Cholesterol	FastingBS	RestingECG	MaxHR	ExerciseAngina	Oldpeak	ST_Slope	HeartDisease
0	40	M	ATA	140	289	0	Normal	172	N	0.0	Up	0
1	49	F	NAP	160	180	0	Normal	156	N	1.0	Flat	1
2	37	M	ATA	130	283	0	ST	98	N	0.0	Up	0
3	48	F	ASY	138	214	0	Normal	108	Y	1.5	Flat	1
4	54	M	NAP	150	195	0	Normal	122	N	0.0	Up	0
...	...	...	...	...	...	...	...	...	...	...	...	...
913	45	M	TA	110	264	0	Normal	132	N	1.2	Flat	1
914	68	M	ASY	144	193	1	Normal	141	N	3.4	Flat	1
915	57	M	ASY	130	131	0	Normal	115	Y	1.2	Flat	1
916	57	F	ATA	130	236	0	LVH	174	N	0.0	Flat	1
917	38	M	NAP	138	175	0	Normal	173	N	0.0	Up	0

918 rows × 12 columns

In [15]: df1.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 918 entries, 0 to 917
Data columns (total 12 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   Age         918 non-null    int64
 1   Sex         918 non-null    object
 2   ChestPainType  918 non-null    object
 3   RestingBP   918 non-null    int64
 4   Cholesterol  918 non-null    int64
 5   FastingBS   918 non-null    int64
 6   RestingECG  918 non-null    object
 7   MaxHR       918 non-null    int64
 8   ExerciseAngina  918 non-null    object
 9   Oldpeak     918 non-null    float64
10  ST_Slope    918 non-null    object
11  HeartDisease  918 non-null    int64
dtypes: float64(1), int64(6), object(5)
memory usage: 86.2+ KB
```

In [16]: df2.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 9357 entries, 0 to 9356
Data columns (total 15 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   Date        9357 non-null    object
 1   Time        9357 non-null    object
 2   CO(CT)      9357 non-null    object
 3   PT08_S1(CO) 9357 non-null    float64
 4   NMHC(CT)    9357 non-null    float64
 5   CH6(CT)     9357 non-null    object
 6   PT08_S2(NMHC) 9357 non-null    float64
 7   NOx(CT)     9357 non-null    float64
 8   PT08_S3(NOX) 9357 non-null    float64
 9   NO2(CT)     9357 non-null    float64
10  PT08_S4(NO2) 9357 non-null    float64
11  PT08_S5(O3) 9357 non-null    float64
12  T            9357 non-null    object
13  RH           9357 non-null    object
14  AH           9357 non-null    object
dtypes: float64(8), object(7)
memory usage: 1.1+ MB
```

In [17]: df2.isnull().sum()

	Date	Time	CO(CT)	PT08_S1(CO)	NMHC(CT)	CH6(CT)	PT08_S2(NMHC)	NOx(CT)	PT08_S3(NOX)	NO2(CT)	PT08_S4(NO2)	PT08_S5(O3)	T	RH	AH
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

## Data transformation

In [18]: df2.replace(".", "", regex=True, inplace=True)

In [19]: df2.drop(["Date", "Time"], axis=1, inplace=True)

In [20]: df2.head()

	CO(CT)	PT08_S1(CO)	NMHC(CT)	CH6(CT)	PT08_S2(NMHC)	NOx(CT)	PT08_S3(NOX)	NO2(CT)	PT08_S4(NO2)	PT08_S5(O3)	T	RH	AH
0	2.6	1360.0	150.0	11.9	1046.0	166.0	1056.0	113.0	1692.0	1268.0	13.6	48.9	0.7578
1	2	1292.0	112.0	9.4	955.0	103.0	1174.0	92.0	1559.0	972.0	13.3	47.7	0.7255
2	2.2	1402.0	88.0	8.0	939.0	131.0	1140.0	114.0	1555.0	1074.0	11.0	54.0	0.7502
3	2.2	1376.0	80.0	8.2	948.0	172.0	1092.0	122.0	1584.0	1203.0	11.0	60.0	0.7867
4	1.6	1272.0	51.0	6.5	836.0	131.0	1205.0	116.0	1490.0	1110.0	11.2	59.6	0.7888

In [21]: df2.dtypes

	CO(CT)	PT08_S1(CO)	NMHC(CT)	CH6(CT)	PT08_S2(NMHC)	NOx(CT)	PT08_S3(NOX)	NO2(CT)	PT08_S4(NO2)	PT08_S5(O3)	T	RH	AH
	object	float64	float64	object	float64	float64	float64	float64	float64	float64	object	object	object

In [22]: df2.dtypes

	CO(CT)	PT08_S1(CO)	NMHC(CT)	CH6(CT)	PT08_S2(NMHC)	NOx(CT)	PT08_S3(NOX)	NO2(CT)	PT08_S4(NO2)	PT08_S5(O3)	T	RH	AH
	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64	float64

In [23]: df1.head()

	Age	Sex	ChestPainType	RestingBP	Cholesterol	FastingBS	RestingECG	MaxHR	ExerciseAngina	Oldpeak	ST_Slope	HeartDisease
0	40	M	ATA									