

Working with tabular data

Importing data

By far the most ubiquitous data format is the “comma-separated values” file format – or simply the csv file format. The csv file format can be used to store a matrix. Here is how:

1. Each row of the file contains numbers separated by commas. There must be as many entries in a row as matrix columns.
 2. If you start a row with “#” it will be ignored when reading the file.
- In the directory, there is file named `data.csv`, we can read it using pandas library
 - The data is a thermodynamic table with saturated steam properties.

Reading data with pandas

```
import pandas as pd
data = pd.read_csv('data.csv', delimiter=',')
data.head(3)
```

	Absolute Pres- sure (bar)	Boil- ing Point (°C)	Spe- cific Vol- (m3/ kg)	Den- sity (steam) (kg m3)	Spe- cific of Liq- uid Wa- ter (sen- sible heat) (kJ/ kg)	Spe- cific of Liq- uid Wa- ter (sen- sible heat) (kcal/ kg)	Spe- cific of Steam (total heat) (kJ/ kg)	Spe- cific of Steam (total heat) (kcal/ kg)	La- tent heat of Va- por- iza- tion (kJ/ kg)	La- tent heat of Va- por- iza- tion (kcal/ kg)	Spe- cific Heat (kJ/ kg K)
0	0.02	17.51	67.006	0.015	73.45	17.54	2533.64	605.15	2460.19	587.61	1.8644
1	0.03	24.10	45.667	0.022	101.00	24.12	2545.64	608.02	2444.65	583.89	1.8694
2	0.04	28.98	34.802	0.029	121.41	29.00	2554.51	610.13	2433.10	581.14	1.8736

Exploratory Analysis

Exploratory Analysis

Counting missing data

- The first thing to check in a tabulated data set is if there's missing data.
- That will indicate either a problem in the data set reading or just say that these records were actually missing!

```
data.isna().sum()
```

```
Absolute Pressure (bar)          0
Boiling Point (oC)               0
Specific Volume (steam) (m3/kg)  1
Density (steam) (kg/m3)         1
Specific Enthalpy of Liquid Water (sensible heat) (kJ/kg)  1
Specific Enthalpy of Liquid Water (sensible heat) (kcal/kg) 1
Specific Enthalpy of Steam (total heat) (kJ/kg)           1
Specific Enthalpy of Steam (total heat) (kcal/kg)          1
Latent heat of Vaporization (kJ/kg)                        1
Latent heat of Vaporization (kcal/kg)                      1
Specific Heat (kJ/kg K)                                    2
dtype: int64
```

Summary Statistics

- When data is missing, pandas will remove the any lines with missing information when you try to calculate anything.

```
summary = data.describe()
summary
```

	Absolute Pres- sure (bar)	Boil- ing Point (°C)	Spe- cific Vol- (m3/ kg)	Den- sity (steam) (kg/ m3)	Spe- cific Enthalpy of Liq- uid Wa- ter (sen- sible heat) (kJ/ kg)	Spe- cific Enthalpy of Liq- uid Wa- ter (sen- sible heat) (kJ/ kg)	Spe- cific Enthalpy of Steam (total heat) (kJ/ kg)	Spe- cific Enthalpy of Steam (total heat) (k- cal/ kg)	La- tent heat of Va- por- iza- tion (kJ/ kg)	La- tent heat of Va- por- iza- tion (k- cal/ kg)	Spe- cific Heat (kJ/ kg K)
count	68.000000	68.000000	67.000000	67.000000	67.000000	67.000000	67.000000	67.000000	67.000000	67.000000	66.000000
mean	9.043235	137.210800	0.850761	14.638073	2.540117	79.799851	84.580809	89.525270	202.041390	351.478402	2277
std	12.242954	14.173370	0.165488	0.272021	3.605472	6.961576	3.145057	2.922357	7.318156	9.958146	476031
min	0.020000	0.255000	0.067000	0.015000	0.345000	0.540000	0.603700	0.651500	0.960000	0.469000	1.864400
25%	0.975000	95.910000	0.146000	0.562500	0.847000	0.825000	0.683100	0.684700	0.614850	0.661050	0.18375
50%	3.750000	136.205000	0.606000	1.908000	1.144000	139.550000	124.660000	124.440000	17.350000	12.890000	215600
75%	14.250000	192.460000	0.978000	0.351000	22.365000	0.000000	0.866050	0.661000	261.785000	40.220000	765700
max	76.000000	203.840000	0.006000	0.093000	0.008330	2704.840000	2802.270000	2855.530000	190.680000	67.610000	406900

Exploratory Analysis

Actively handling missing data

- pandas can handle automatically missing data if you do calculations in it.
- However, if you want to do calculations using another package, errors will appear as they might not handle missing values like pandas

Exploratory Analysis

Actively handling missing data - Code

```
data_no_null_rows = data.dropna(axis = 0)

#Dropna drops all rows that contain at least 1 missing value
```

```
data_no_null_rows.isna().sum()
```

```

Absolute Pressure (bar)                0
Boiling Point (oC)                    0
Specific Volume (steam) (m3/kg)       0
Density (steam) (kg/m3)               0
Specific Enthalpy of Liquid Water (sensible heat) (kJ/kg)  0
Specific Enthalpy of Liquid Water (sensible heat) (kcal/kg) 0
Specific Enthalpy of Steam (total heat) (kJ/kg)            0
Specific Enthalpy of Steam (total heat) (kcal/kg)          0
Latent heat of Vaporization (kJ/kg)    0
Latent heat of Vaporization (kcal/kg)  0
Specific Heat (kJ/kg K)                0
dtype: int64

```

More options in dropna

Manipulating Data Frames

Given our clean data frame `data_no_null_rows`, we may be interested in answering some questions with it:

1. How many rows and columns my dataframe has?
2. How many data points I have for a temperature above 80 degrees celsius and less than 120?
3. What is the temperature and pressure of the record at 75% of the data set size?

Manipulating Data Frames

Question 1

- How many rows and columns my dataframe has?

```
data_no_null_rows.shape
```

```
(66, 11)
```

Question 2

- How many data points I have for a temperature above 80 degrees celsius and less than 120?

```

above_80 = data_no_null_rows["Boiling Point (oC)"] > 80.0
below_120 = data_no_null_rows["Boiling Point (oC)"] < 120.0
below_120.head(3)

```

```

0    True
1    True
2    True
Name: Boiling Point (oC), dtype: bool

```

```
a80_b120 = above_80 & below_120
data_no_null_rows[a80_b120].shape
```

```
(15, 11)
```

Question 3

- What is the temperature and pressure of the record at 75% of the data set size?

```
idx_75 = int(data_no_null_rows.shape[0]*.75)
data_no_null_rows.loc[idx_75, ["Boiling Point (oC)", "Absolute Pressure (bar)"]]
```

```
Boiling Point (oC)      187.96
Absolute Pressure (bar)  12.00
Name: 49, dtype: float64
```

Data visualization

Python has many visualization libraries. Here are some popular ones for data science:

- Matplotlib
- seaborn
- plotly

Data Visualization

2D - scatterplots

```
import matplotlib.pyplot as plt
plt.rcParams["figure.figsize"] = (5.0, 2.5)
data_no_null_rows.plot(x = "Boiling Point (oC)", y = "Absolute Pressure (bar)",
kind = "scatter")
```

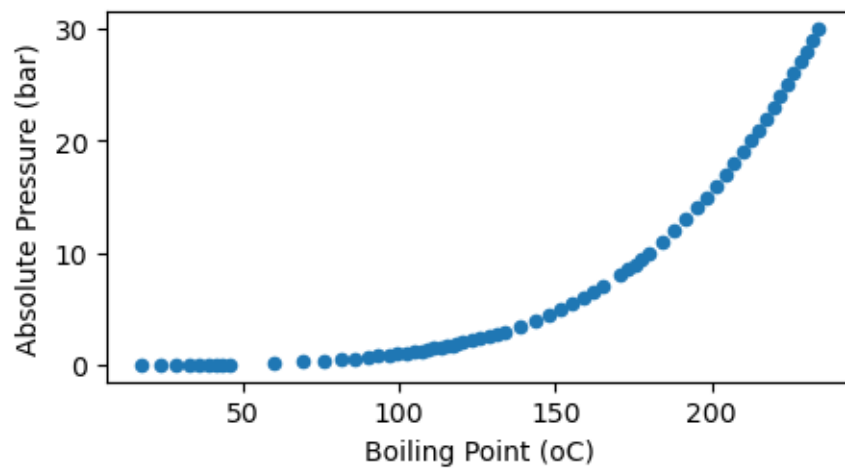
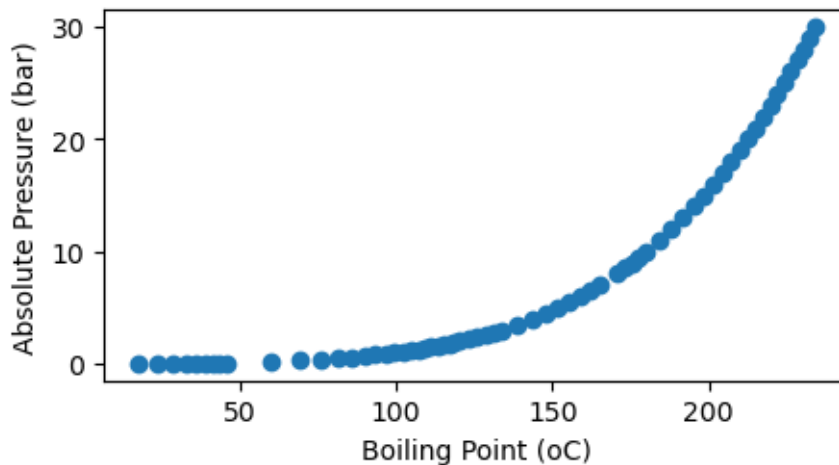


Figure 1: Temperature vs Pressure - Saturated Steam

- You can also plot without using pandas

```
plt.scatter(data_no_null_rows["Boiling Point (oC)"],
data_no_null_rows["Absolute Pressure (bar)"])
plt.xlabel("Boiling Point (oC)")
plt.ylabel("Absolute Pressure (bar)")
```

```
Text(0, 0.5, 'Absolute Pressure (bar)')
```

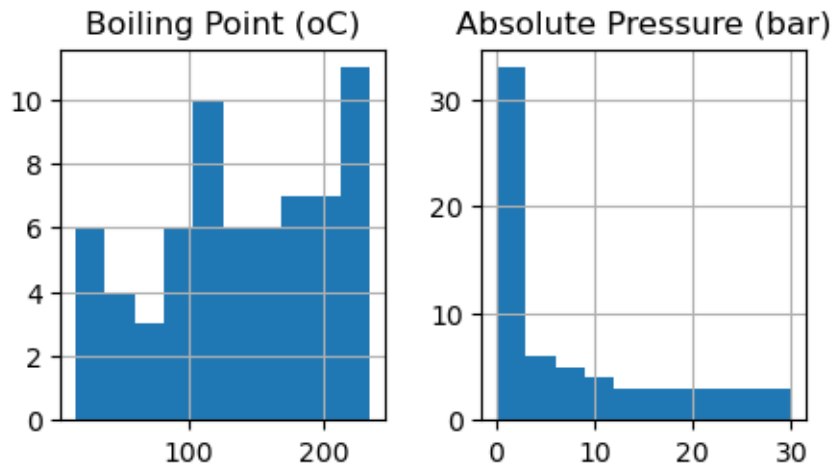


2D - histograms

```
data_no_null_rows.hist(["Boiling Point (oC)", "Absolute Pressure (bar)"])
```

```
array([[<Axes: title={'center': 'Boiling Point (oC)'>,
      <Axes: title={'center': 'Absolute Pressure (bar)'>]],
      dtype=object)
```

Histograms



2D - Box plot

```
data_no_null_rows.boxplot(["Boiling Point (oC)", "Absolute Pressure (bar)"])
```

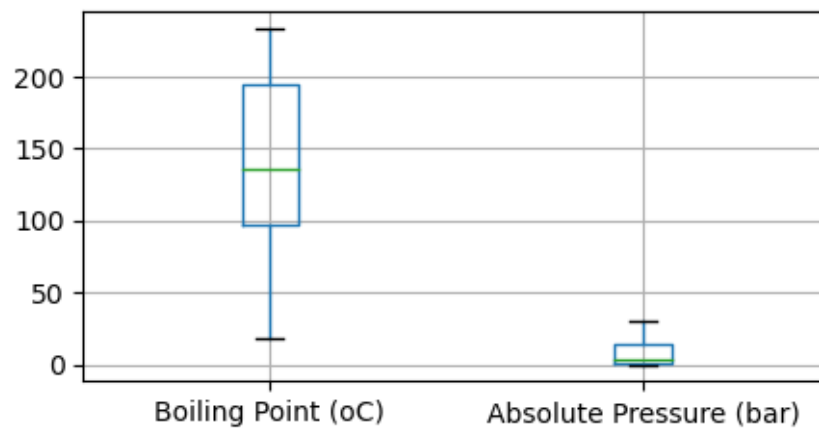


Figure 2: Boxplot