

Import dataset

```
In [183].. from pathlib import Path
import pandas as pd
filepath = Path('data.csv')
```

```
In [184].. average_temp = pd.read_csv(filepath,
header=4,
parse_dates=True,
usecols=[ 'Date',
'Value',
'Anomaly'])

average_temp.head()
```

Out[184]..

	Date	Value	Anomaly
0	194001	55.8	-0.5
1	194101	56.4	0.1
2	194201	57.7	1.4
3	194301	56.3	0.0
4	194401	56.1	-0.2

Locate missing values and change them to nan

```
In [185].. import numpy as np
average_temp = average_temp.replace(-99: np.nan)
print(average_temp.isna())
```

Out[185]..

	Date	Value	Anomaly
0	False	False	False
1	False	False	False
2	False	False	False
3	False	False	False
4	False	False	False
..
78	201801	59.4	3.1
79	201901	57.8	1.5
80	202001	57.9	1.6
81	202101	58.4	2.1
82	202201	58.6	2.3

Use the interpolate function to put a value in the Nan's place

```
In [186].. average_temp.interpolate(inplace = True)
print(average_temp)
```

Out[186]..

	Date	Value	Anomaly
0	194001	55.8	-0.5
1	194101	56.4	0.1
2	194201	57.7	1.4
3	194301	56.3	0.0
4	194401	56.1	-0.2
..
78	201801	59.4	3.1
79	201901	57.8	1.5
80	202001	57.9	1.6
81	202101	58.4	2.1
82	202201	58.6	2.3

Convert the index to datetime format

```
In [187].. average_temp['Date'] = pd.to_datetime(average_temp['Date'], format = '%Y%m')
average_temp.set_index('Date', inplace=True)
average_temp.head()
```

Out[187]..

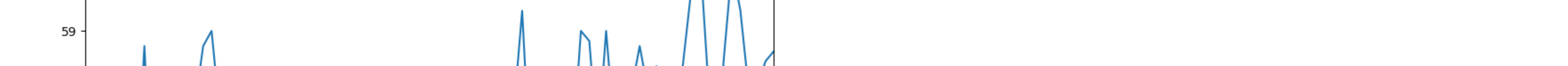
	Date	Value	Anomaly
1940-01-01	1940-01-01	55.8	-0.5
1941-01-01	1941-01-01	56.4	0.1
1942-01-01	1942-01-01	57.7	1.4
1943-01-01	1943-01-01	56.3	0.0
1944-01-01	1944-01-01	56.1	-0.2

Plot the average temperature time series, the corresponding histogram, and kernel density plot

```
In [188].. average_temp['Value'].plot(figsize=(10, 6), title = 'Average temperature time series')
```



```
In [189].. average_temp['Value'].plot.hist(figsize=(10, 6), title = 'Average temperature histogram')
```



```
In [190].. average_temp['Value'].plot.kde(figsize=(10, 6), title = 'Average temperature kernel density plot')
```



Generate descriptive statistics

```
In [191].. average_temp.describe()
```

Out[191]..

	Value	Anomaly
count	83.000000	83.000000
mean	56.766265	0.466265
std	1.536331	1.536331
min	53.600000	-2.700000
25%	55.600000	-0.700000
50%	56.600000	0.300000
75%	57.800000	1.500000
max	61.000000	4.700000

Task 1. Vacation search results Set datetime index for each dataframe.

```
In [192].. poland_df = pd.read_csv('poland.csv', parse_dates=[0], header=1)
uk_df = pd.read_csv('uk.csv', parse_dates=[0], header=1)
usa_df = pd.read_csv('usa.csv', parse_dates=[0], header=1)
```

```
In [193].. poland_df.set_index('Month', inplace=True)
uk_df.set_index('Month', inplace=True)
usa_df.set_index('Month', inplace=True)
```

```
In [194].. usa_df.head()
```

Out[194]..

vacation: (United States)	
Month	
2004-01-01	97
2004-02-01	87
2004-03-01	83
2004-04-01	79
2004-05-01	82

Rename the columns to the country name

```
In [195].. usa_df.rename(columns={"vacation: (United States)": "US"}, inplace=True)
poland_df.rename(columns={"vacation: (Poland)": "PL"}, inplace=True)
uk_df.rename(columns={"vacation: (United Kingdom)": "UK"}, inplace=True)
```

```
In [196].. usa_df.head()
```

Out[196]..

US	
Month	
2004-01-01	97
2004-02-01	87
2004-03-01	83
2004-04-01	79
2004-05-01	82

Combine the search counts in one dataframe

```
In [197].. combined_df = pd.concat([usa_df, poland_df, uk_df], axis=1)
combined_df.head()
```

Out[197]..

	US	PL	UK
Month			
2004-01-01	97	0	40
2004-02-01	87	100	37
2004-03-01	83	0	47
2004-04-01	79	73	36
2004-05-01	82	0	38

Present the time series for all countries in one plot

```
In [198].. combined_df.plot(figsize=(10, 6), title='Vacation results time series')
```



Generate descriptive statistics

```
In [199].. combined_df.describe()
```

Out[199]..

	US	PL	UK
count	243.000000	243.000000	243.000000
mean	57.395062	42.510288	27.160494
std	15.221515	17.128914	14.386739
min	26.000000	0.000000	15.000000
25%	47.000000	30.500000	20.000000
50%	55.000000	40.000000	23.000000
75%	65.000000	51.500000	27.000000
max	100.000000	100.000000	100.000000

Show three histograms in one plot

```
In [200].. combined_df.plot.hist(figsize=(10, 6), title='Vacation results histograms')
```



Show three kernel densities in one plot

```
In [201].. combined_df.plot.kde(figsize=(10, 6), title='Vacation results kernel densities')
```

