

# Part 2 — Workshop 4

TECH2: Introduction to Programming, Data, and Information Technology

Richard Foltyn  
*Norwegian School of Economics (NHH)*

October 25, 2024

See GitHub repository for notebooks and data:

<https://github.com/richardfoltyn/TECH2-H24>

## Contents

1	Exercise: Business cycle correlations	1
2	Exercise: Loading many data files	4
3	Exercise: Decade averages of macro time series	5
4	Exercise: Merging the Titanic data	10

## 1 Exercise: Business cycle correlations

For this exercise, you'll be using macroeconomic data from the folder `../data/FRED`.

1. There are seven decade-specific files named `FRED_monthly_19X0.csv` where `X` identifies the decade (`X` takes on the values 5, 6, 7, 8, 9, 0, 1). Write a loop that reads in all seven files as DataFrames and store them in a list.

*Hint:* Recall from the lecture that you should use `pd.read_csv(..., parse_dates=['DATE'])` to automatically parse strings stored in the `DATE` column as dates.

2. Use `pd.concat()` to concatenate these data sets into a single DataFrame and set the `DATE` column as the index.
3. You realize that your data does not include GDP since this variable is only reported at quarterly frequency. Load the GDP data from the file `GDP.csv` and merge it with your monthly data using an *inner join*.
4. You want to compute how (percent) changes of the variables in your data correlate with percent changes in GDP.
  1. Create a *new* DataFrame which contains the percent changes in CPI and GDP (using `pct_change()`, see also the last exercise in workshop 3), and the absolute changes for the remaining variables (using `diff()`).
  2. Compute the correlation of the percent changes in GDP with the (percent) changes of all other variables (using `corr()`). What does the sign and magnitude of the correlation coefficient tell you?

---

*Solution.*

## Part (1)

```
[1]: # Uncomment this to use files in the local data/ directory
DATA_PATH = '../data/FRED'

# Uncomment this to load data directly from GitHub
# DATA_PATH = 'https://raw.githubusercontent.com/richardfoltyn/TECH2-H24/main/data/FRED'
```

There are many ways to load the seven files we need. One possibility is to loop over the decades 1950, 1960, ..., construct the decade-specific file name and load the decade-specific file.

```
[2]: import numpy as np
import pandas as pd
import os.path

# Create years representing decades: 1950, 1960, ....
years = np.arange(1950, 2011, 10)

data = []
for year in years:
    # File name for current decade
    fn = f'FRED_monthly_{year}.csv'

    # Join data folder + filename to get path to CSV file
    path = os.path.join(DATA_PATH, fn)

    # Load decade data
    df = pd.read_csv(path, parse_dates=['DATE'])

    data.append(df)
```

## Part (2)

```
[3]: # Concatenate decade data along the row axis
df = pd.concat(data, axis=0)

# Print first 5 observations
df.head(5)
```

```
[3]:
```

	DATE	CPI	UNRATE	FEDFUNDS	REALRATE	LFPART
0	1950-01-01	23.5	6.5	NaN	NaN	58.9
1	1950-02-01	23.6	6.4	NaN	NaN	58.9
2	1950-03-01	23.6	6.3	NaN	NaN	58.8
3	1950-04-01	23.6	5.8	NaN	NaN	59.2
4	1950-05-01	23.8	5.5	NaN	NaN	59.1

```
[4]: # Print last 5 observations
df.tail(5)
```

```
[4]:
```

	DATE	CPI	UNRATE	FEDFUNDS	REALRATE	LFPART
115	2019-08-01	256.0	3.6	2.1	0.6	63.1
116	2019-09-01	256.4	3.5	2.0	0.3	63.2
117	2019-10-01	257.2	3.6	1.8	-0.0	63.3
118	2019-11-01	257.9	3.6	1.6	-0.2	63.3
119	2019-12-01	258.6	3.6	1.6	-0.3	63.3

The index in the concatenated data is not unique, as you can easily verify:

```
[5]: # Selecting the obs with label 0 returns 7 rows!
df.loc[0]
```

```
[5]:
```

	DATE	CPI	UNRATE	FEDFUNDS	REALRATE	LFPART
0	1950-01-01	23.5	6.5	NaN	NaN	58.9
0	1960-01-01	29.4	5.2	4.0	NaN	59.1
0	1970-01-01	37.9	3.9	9.0	NaN	60.4
0	1980-01-01	78.0	6.3	13.8	NaN	64.0
0	1990-01-01	127.5	5.4	8.2	3.8	66.8
0	2000-01-01	169.3	4.0	5.4	2.7	67.3
0	2010-01-01	217.5	9.8	0.1	-0.8	64.8

It is advisable to always work with a unique index in pandas, and for this data set the most natural unique index is the date.

```
[6]: # Set DATE column as index
df = df.set_index('DATE')
```

### Part (3)

```
[7]: # Path to GDP data
fn = os.path.join(DATA_PATH, 'GDP.csv')

# Load GDP data
gdp = pd.read_csv(fn, parse_dates=['DATE'], index_col='DATE')

gdp.head(5)
```

```
[7]:
```

	GDP
DATE	
1947-01-01	2182.7
1947-04-01	2176.9
1947-07-01	2172.4
1947-10-01	2206.5
1948-01-01	2239.7

We merge the GDP using an *inner join*, which discards all months where GDP is not reported.

```
[8]: # Merge the GDP data using an inner join
df = df.join(gdp, how='inner')

df.head(5)
```

```
[8]:
```

	CPI	UNRATE	FEDFUNDS	REALRATE	LFPART	GDP
DATE						
1950-01-01	23.5	6.5	NaN	NaN	58.9	2346.1
1950-04-01	23.6	5.8	NaN	NaN	59.2	2417.7
1950-07-01	24.1	5.0	NaN	NaN	59.1	2511.1
1950-10-01	24.5	4.2	NaN	NaN	59.4	2559.2
1951-01-01	25.4	3.7	NaN	NaN	59.1	2594.0

### Part (4.1)

We can compute (percent) changes for multiple columns at once, so there is no need to even loop over variables:

```
[9]: # Compute percent changes for CPI and GDP
df_changes = df[['CPI', 'GDP']].pct_change() * 100

# Other variables for which to compute absolute changes
variables = ['UNRATE', 'FEDFUNDS', 'REALRATE', 'LFPART']

# Compute absolute changes, add to DataFrame
```

```
df_changes[variables] = df[variables].diff()

df_changes.head(5)
```

```
[9]:
```

	CPI	GDP	UNRATE	FEDFUNDS	REALRATE	LFPART
DATE						
1950-01-01	NaN	NaN	NaN	NaN	NaN	NaN
1950-04-01	0.425532	3.051873	-0.7	NaN	NaN	0.3
1950-07-01	2.118644	3.863176	-0.8	NaN	NaN	-0.1
1950-10-01	1.659751	1.915495	-0.8	NaN	NaN	0.3
1951-01-01	3.673469	1.359800	-0.5	NaN	NaN	-0.3

## Part (4.2)

The `corr()` method returns the whole (symmetric) correlation matrix. We are only interested in the correlations with GDP changes, so we can select that particular row.

```
[10]: # Compute correlation matrix, keep only GDP row
df_changes.corr().loc['GDP']
```

```
[10]: CPI          -0.113091
      GDP           1.000000
      UNRATE       -0.564872
      FEDFUNDS      0.206370
      REALRATE      0.074500
      LFPART        0.019639
      Name: GDP, dtype: float64
```

As we can see, some (changes in) variables are more highly correlated with GDP changes than others. For example, the unemployment rate is highly negatively correlated with GDP growth, i.e., in good times (large positive GDP changes), the unemployment rate drops.

## 2 Exercise: Loading many data files

In the previous exercise, you loaded the individual files by specifying an explicit list of file names. This can become tedious or infeasible if your data is spread across many files with varying file name patterns. Python offers the possibility to iterate over all files in a directory (for example, using `os.listdir()`), or to iterate over files that match a pattern, for example using `glob.glob()`.

Repeat parts (1) and (2) from the previous exercise, but now iterate over the input files using `glob.glob()`. You'll need to use a wildcard `*` and make sure to match only the relevant files in `../data/FRED`, i.e., those that start with `FRED_monthly`.

### Solution.

```
[11]: # Uncomment this to use files in the local data/ directory
      DATA_PATH = '../data/FRED'

      # Uncomment this to load data directly from GitHub
      # DATA_PATH = 'https://raw.githubusercontent.com/richardfoltn/TECH2-H24/main/data/FRED'
```

```
[12]: import pandas as pd
      import glob

      # Pattern to match only the desired files in data/FRED. The wildcard *
      # matches anything.
```

```

pattern = f'{DATA_PATH}/FRED_monthly*.csv'

data = []

for file in glob.glob(pattern):
    print(f'Loading file {file}')
    d = pd.read_csv(file, parse_dates=['DATE'], index_col='DATE')
    data.append(d)

# Concatenate all DataFrames
df = pd.concat(data, axis=0)

# Sort index in case files have been loaded in unexpected order
df = df.sort_index()

# Print first 12 rows
df.head(12)

```

```

Loading file
/home/richard/repos/teaching/TECH2-H24/data/FRED/FRED_monthly_1950.csv
Loading file
/home/richard/repos/teaching/TECH2-H24/data/FRED/FRED_monthly_1960.csv
Loading file
/home/richard/repos/teaching/TECH2-H24/data/FRED/FRED_monthly_1970.csv
Loading file
/home/richard/repos/teaching/TECH2-H24/data/FRED/FRED_monthly_1980.csv
Loading file
/home/richard/repos/teaching/TECH2-H24/data/FRED/FRED_monthly_1990.csv
Loading file
/home/richard/repos/teaching/TECH2-H24/data/FRED/FRED_monthly_2000.csv
Loading file
/home/richard/repos/teaching/TECH2-H24/data/FRED/FRED_monthly_2010.csv

```

```

[12]:

```

	CPI	UNRATE	FEDFUNDS	REALRATE	LFPART
DATE					
1950-01-01	23.5	6.5	NaN	NaN	58.9
1950-02-01	23.6	6.4	NaN	NaN	58.9
1950-03-01	23.6	6.3	NaN	NaN	58.8
1950-04-01	23.6	5.8	NaN	NaN	59.2
1950-05-01	23.8	5.5	NaN	NaN	59.1
1950-06-01	23.9	5.4	NaN	NaN	59.4
1950-07-01	24.1	5.0	NaN	NaN	59.1
1950-08-01	24.2	4.5	NaN	NaN	59.5
1950-09-01	24.3	4.4	NaN	NaN	59.2
1950-10-01	24.5	4.2	NaN	NaN	59.4
1950-11-01	24.6	4.2	NaN	NaN	59.3
1950-12-01	25.0	4.3	NaN	NaN	59.2

---

### 3 Exercise: Decade averages of macro time series

For this exercise, you'll be using macroeconomic data from the folder `../data/FRED`.

1. There are five files containing monthly observations on annual inflation (INFLATION), the Fed Funds rate (FEDFUNDS), the labor force participation rate (LFPART), the 1-year real interest rate (REALRATE) and the unemployment rate (UNRATE). Write a loop to import these and merge them on DATE into a single DataFrame using *outer joins* (recall that `merge()` and `join()` operate on only two DataFrames at a time).

*Hint:* Recall from the lecture that you should use `pd.read_csv(..., parse_dates=['DATE'])` to automatically parse strings stored in the DATE column as dates.

2. Your friend is a pandas guru and tells you that you don't need to iteratively merge many files but can instead directly use `pd.concat()` for merging many DataFrames in a single step. Repeat the previous part using `pd.concat()` instead, and verify that you get the same result (you can do this using `compare()`).
3. You want to compute the average value of each variable by decade, but you want to include only decades without *any* missing values for *all* variables.
  1. Create a variable Decade which stores the decade (1940, 1950, ...) for each observation.

*Hint:* You should have set the DATE as the DataFrame index. Then you can access the calendar year using the attribute `df.index.year` which can be used to compute the decade.
  2. Write a function `num_missing(x)` which takes as argument `x` a Series and returns the number of missing values in this Series.
  3. Compute the number of missing values by decade for each variable using a `groupby()` operation and the function `num_missing` you wrote.
  4. Aggregate this data across all variables to create an indicator for each decade whether there are any missing values. This can be done in many ways but will require aggregation across columns, e.g., with `sum(..., axis=1)`.
  5. Merge this decade-level indicator data back into the original DataFrame (*many-to-one* merge).
4. Using this indicator, drop all observations which are in a decade with missing values.
5. Compute the decade average for each variable.

### Challenge

- Your pandas guru friend claims that all the steps in 3.2 to 3.5 can be done with a single one-liner using `transform()`. Can you come up with a solution?

---

### Solution.

#### Part (1)

```
[13]: # Uncomment this to use files in the local data/ directory
DATA_PATH = '../data/FRED'

# Uncomment this to load data directly from GitHub
# DATA_PATH = 'https://raw.githubusercontent.com/richardfoltn/TECH2-H24/main/data/FRED'

[14]: import pandas as pd
import os.path

# Variables to be imported
variables = ['INFLATION', 'FEDFUNDS', 'LFPART', 'REALRATE', 'UNRATE']

df = None
for var in variables:
    # File path for current variable
    fn = os.path.join(DATA_PATH, f'{var}.csv')

    # Load data for current variable, parse DATE and set DATE as index
    d = pd.read_csv(fn, parse_dates=['DATE'], index_col='DATE')

    # Merge to existing DataFrame if it is defined, otherwise this is the
    # first variable in the loop.
```

```

if df is None:
    df = d
else:
    # Merges on common index by default
    df = df.join(d, how='outer')

# Print first 5 rows
df.head(5)

```

```

[14]:
      INFLATION  FEDFUNDS  LFPART  REALRATE  UNRATE
DATE
1948-01-01    10.2      NaN    58.6      NaN    3.4
1948-02-01     9.7      NaN    58.9      NaN    3.8
1948-03-01     6.8      NaN    58.5      NaN    4.0
1948-04-01     8.2      NaN    59.0      NaN    3.9
1948-05-01     9.1      NaN    58.3      NaN    3.5

```

## Part (2)

When using `pd.concat()` to merge files, we first import all the DataFrames and store them in a list, and concatenate all at the end.

```

[15]: # List to store individual DataFrames
data = []

for var in variables:
    # File path for current variable
    fn = os.path.join(DATA_PATH, f'{var}.csv')

    # Load data for current variable, parse DATE and set DATE as index
    d = pd.read_csv(fn, parse_dates=['DATE'], index_col='DATE')

    data.append(d)

# Concatenate all DataFrames along the column axis. We specify join='outer'
# as we want to perform an outer join along the OTHER (row) axis.
df2 = pd.concat(data, axis=1, join='outer')

# Print first 12 rows
df2.head(5)

```

```

[15]:
      INFLATION  FEDFUNDS  LFPART  REALRATE  UNRATE
DATE
1948-01-01    10.2      NaN    58.6      NaN    3.4
1948-02-01     9.7      NaN    58.9      NaN    3.8
1948-03-01     6.8      NaN    58.5      NaN    4.0
1948-04-01     8.2      NaN    59.0      NaN    3.9
1948-05-01     9.1      NaN    58.3      NaN    3.5

```

We can compare the second DataFrame to the first one using `compare()`. This method only prints the differences, so if it prints an empty DataFrame, the DataFrames are identical.

```

[16]: # Compare both DataFrames. compare() only displays differences!
df.compare(df2)

```

```

[16]: Empty DataFrame
Columns: []
Index: []

```

### Part (3.1)

We first need to create a new column `Decade` which stores the decade 1950, 1960, etc., corresponding to each observation.

```
[17]: # Extract calendar year
year = df.index.year

# Create decade from calendar year using truncated division (could also use np.floor())
decade = (year // 10) * 10

# Assign to new column
df['Decade'] = decade

# Verify that Decade variable looks as expected
df['Decade'].value_counts()
```

```
[17]: Decade
1950    120
1960    120
1970    120
1980    120
1990    120
2000    120
2010    120
2020     58
1940     24
Name: count, dtype: int64
```

### Part (3.2)

```
[18]: # Define function that returns the number of missing observations
def num_missing(x):
    n = x.isna().sum()
    return n
```

```
[19]: # Test num_missing() on one column
num_missing(df['INFLATION'])
```

```
[19]: 2
```

### Part (3.3)

```
[20]: # Compute number of missing by decade for each variable
df_miss = df.groupby('Decade').agg(num_missing)

# Print table of missing observations
df_miss
```

```
[20]:
```

	INFLATION	FEDFUNDS	LFPART	REALRATE	UNRATE
Decade					
1940	0	24	0	24	0
1950	0	54	0	120	0
1960	0	0	0	120	0
1970	0	0	0	120	0
1980	0	0	0	24	0
1990	0	0	0	0	0
2000	0	0	0	0	0
2010	0	0	0	0	0



2020	2	1	1	0	1
------	---	---	---	---	---

### Part (3.4)

We sum across columns to get the total number of missing values by decade and check whether this number is zero.

```
[21]: # Indicator whether none of the observations are missing:
no_miss = (df_miss.sum(axis=1) == 0)

no_miss
```

```
[21]: Decade
1940    False
1950    False
1960    False
1970    False
1980    False
1990     True
2000     True
2010     True
2020    False
dtype: bool
```

Before merging the missing indicator with the original data, we convert it to a DataFrame using `to_frame()`. This allows us to assign a meaningful column name.

```
[22]: # Convert to DataFrame, assign meaningful column name
no_miss = no_miss.to_frame('NotMissing')

# Reset index to move Decade back into columns
no_miss = no_miss.reset_index()

no_miss
```

```
[22]:   Decade  NotMissing
0   1940      False
1   1950      False
2   1960      False
3   1970      False
4   1980      False
5   1990       True
6   2000       True
7   2010       True
8   2020      False
```

### Part (3.5)

```
[23]: # Merge back into original data set
df = df.merge(no_miss, how='left', on='Decade')

df.head(5)
```

```
[23]:   INFLATION  FEDFUNDS  LFPART  REALRATE  UNRATE  Decade  NotMissing
0      10.2        NaN    58.6        NaN     3.4   1940      False
1       9.7        NaN    58.9        NaN     3.8   1940      False
2       6.8        NaN    58.5        NaN     4.0   1940      False
3       8.2        NaN    59.0        NaN     3.9   1940      False
4       9.1        NaN    58.3        NaN     3.5   1940      False
```

#### Part (4)

```
[24]: # Keep only decades without any missing observations
df_no_miss = df.loc[df['NotMissing']].copy()

print(f'Final number of observations: {len(df_no_miss)}')
```

Final number of observations: 360

#### Part (5)

```
[25]: # drop NotMissing, don't want averages of these
df_no_miss = df_no_miss.drop(columns=['NotMissing'])

# Compute decade means
df_no_miss.groupby('Decade').mean()
```

```
[25]:
```

	INFLATION	FEDFUNDS	LFPART	REALRATE	UNRATE
Decade					
1990	3.006667	5.140000	66.668333	2.206667	5.762500
2000	2.569167	2.952500	66.236667	1.023333	5.541667
2010	1.775833	0.618333	63.295000	-0.732500	6.220833

#### Challenge

The pandas guru friend was right, we can use `transform()` and a lambda expression to create the indicator by decade and assign it to each observation in one line:

```
[26]: df['NotMissing2'] = df.groupby('Decade').transform(lambda x: x.isna().sum()).sum(axis=1)
      → == 0
```

```
[27]: # Compare original and new indicator to ensure they are the same
(df['NotMissing'] == df['NotMissing2']).all()
```

```
[27]: True
```

---

## 4 Exercise: Merging the Titanic data

In this exercise, you'll be working with the the original Titanic data set in `titanic.csv` and additional (partly fictitious) information on passengers stored in `titanic-additional.csv`, both located in the `data/` folder.

The goal of the exercise is to calculate the survival rates by country of residence (for this exercise we restrict ourselves to the UK, so these will be England, Scotland, etc.).

1. Load the `titanic.csv` and `titanic-additional.csv` into two DataFrames.

Inspect the columns contained in both data sets. As you can see, the original data contains the full name including the title and potentially maiden name (for married women) in a single column. The additional data contains this information in separate columns. You want to merge these data sets, but you first need to create common keys in both DataFrames.

2. Since the only common information is the name, you'll need to extract the individual name components from the original DataFrame and use these as merge keys.

Focusing only on men (who have names that are much easier to parse), split the Name column into the tokens Title, FirstName and LastName, just like the columns in the second DataFrame.

*Hint:* This is the same task as in the last exercise in Workshop 2. You can just use your solution here.

3. Merge the two data sets based on the columns Title, FirstName and LastName you just created using a *left join* (*one-to-one* merge). Tabulate the columns and the number of non-missing observations to make sure that merging worked.

*Note:* The additional data set contains address information only for passengers from the UK, so some of these fields will be missing.

4. You are now in a position to merge the country of residence (*many-to-one* merge). Load the country data from UK\_post\_codes.csv which contains the UK post code prefix (which you can ignore), the corresponding city, and the corresponding country.

Merge this data with your passenger data set using a *left join* (what is the correct merge key?).

5. Tabulate the number of observations by Country, including the number of observations with missing Country (these are passengers residing outside the UK).

Finally, compute the mean survival rate by country.

---

### Solution.

#### Part (1)

```
[28]: # Uncomment this to use files in the local data/ directory
DATA_PATH = '../data'

# Uncomment this to load data directly from GitHub
# DATA_PATH = 'https://raw.githubusercontent.com/richardfoltyn/TECH2-H24/main/data'
```

Import the original Titanic data:

```
[29]: import pandas as pd
import os.path

# Path to original data
fn1 = os.path.join(DATA_PATH, 'titanic.csv')

# Read in original data set
df1 = pd.read_csv(fn1, index_col='PassengerId')

# Inspect first 5 rows of the original data set
df1.head(5)
```

```
[29]:
```

	Survived	Pclass	\
PassengerId			
1	0	3	
2	1	1	
3	1	3	
4	1	1	
5	0	3	

  

	Name	Sex	Age	\
PassengerId				
1	Braund, Mr. Owen Harris	male	22.0	
2	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	

3		Heikkinen, Miss Laina	female	26.0
4		Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0
5		Allen, Mr. William Henry	male	35.0

	Ticket	Fare	Cabin	Embarked
PassengerId				
1	A/5 21171	7.2500	NaN	S
2	PC 17599	71.2833	C85	C
3	STON/O2. 3101282	7.9250	NaN	S
4	113803	53.1000	C123	S
5	373450	8.0500	NaN	S

Import the additional data:

```
[30]: # Path to additional data
fn2 = os.path.join(DATA_PATH, 'titanic-additional.csv')

# Read in additional data. Note: This one does not have PassengerId
df2 = pd.read_csv(fn2)

# Inspect 5 rows of the additional data
df2.head(5)
```

```
[30]: Title      LastName      FirstName MaidenName      City Postcode \
0  Mr.      Christmann      Emil      NaN      Chester CH6 34H
1  Miss     Heikkinen      Laina      NaN      Bolton BL0 1XG
2  Lady.    Duff Gordon      Lucille Christiana Sutherland      NaN      NaN
3  Miss     Pettersson      Ellen Natalia      NaN      Northampton NN0 H5R
4  Mr.      Odahl      Nils Martin      NaN      Derby DE7 QZ7

      Address
0      3 Graham ways
1      0 Griffin wells
2      NaN
3      889 Murray glen
4  Studio 2, Long courts
```

## Part (2)

```
[31]: # Restrict sample to men
df1 = df1.query('Sex == "male"')
```

```
[32]: # Create DataFrame of name tokens: first column contains the last name,
# second column contains the title and first name
names = df1['Name'].str.split(',', expand=True)

# Trim any residual spaces at the beginning and end
for col in names.columns:
    names[col] = names[col].str.strip()
```

```
[33]: # Extract the last name from the first column
last_name = names.loc[:, 0]

last_name.head(5)
```

```
[33]: PassengerId
1      Braund
5      Allen
6      Moran
7      McCarthy
8      Palsson
```

Name: 0, dtype: object

```
[34]: # Title and first name (potentially multiple) are separated by space.
# partition() splits on the first occurrence of the given character and returns
# three columns.
title_first = names[1].str.strip().str.partition(' ')

title_first.head(5)
```

```
[34]:
```

	0	1	2
PassengerId			
1	Mr.	Owen	Harris
5	Mr.	William	Henry
6	Mr.	James	
7	Mr.	Timothy	J
8	Master	Gosta	Leonard

```
[35]: # Extract title from 1st column, strip any remaining white space
title = title_first[0].str.strip()
title.head(5)
```

```
[35]: PassengerId
1      Mr.
5      Mr.
6      Mr.
7      Mr.
8  Master
Name: 0, dtype: object
```

```
[36]: # Extract first name(s) from 3rd column, strip any remaining white space
first_name = title_first[2].str.strip()

# Print first 5 observations
first_name.head(5)
```

```
[36]: PassengerId
1      Owen Harris
5  William Henry
6      James
7  Timothy J
8  Gosta Leonard
Name: 2, dtype: object
```

```
[37]: # Add all name components back into original DataFrame
df1['FirstName'] = first_name
df1['LastName'] = last_name
df1['Title'] = title

# Delete Name column
del df1['Name']
```

```
[38]: df1.head(5)
```

```
[38]:
```

	Survived	Pclass	Sex	Age	Ticket	Fare	Cabin	Embarked	\
PassengerId									
1	0	3	male	22.0	A/5 21171	7.2500	NaN	S	
5	0	3	male	35.0	373450	8.0500	NaN	S	
6	0	3	male	NaN	330877	8.4583	NaN	Q	
7	0	1	male	54.0	17463	51.8625	E46	S	
8	0	3	male	2.0	349909	21.0750	NaN	S	

  

	FirstName	LastName	Title
--	-----------	----------	-------

PassengerId			
1	Owen Harris	Braund	Mr.
5	William Henry	Allen	Mr.
6	James	Moran	Mr.
7	Timothy J	McCarthy	Mr.
8	Gosta Leonard	Palsson	Master

### Part (3)

```
[39]: # Merge on Title, First name and Last name, keep only left data
keys = ['Title', 'FirstName', 'LastName']
df_merged = df1.merge(df2, on=keys, how='left')
```

```
[40]: # Print missing statistics for merged data
df_merged.info(show_counts=True)
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 577 entries, 0 to 576
Data columns (total 15 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Survived    577 non-null    int64
1   Pclass      577 non-null    int64
2   Sex         577 non-null    object
3   Age         453 non-null    float64
4   Ticket      577 non-null    object
5   Fare        577 non-null    float64
6   Cabin       107 non-null    object
7   Embarked    577 non-null    object
8   FirstName   577 non-null    object
9   LastName    577 non-null    object
10  Title       577 non-null    object
11  MaidenName   0 non-null      object
12  City        471 non-null    object
13  Postcode     471 non-null    object
14  Address      471 non-null    object
dtypes: float64(2), int64(2), object(11)
memory usage: 67.7+ KB
```

### Part (4)

```
[41]: # Path to UK post code data
fn = os.path.join(DATA_PATH, 'UK_post_codes.csv')

# Load UK post code data
df_codes = pd.read_csv(fn)

# Drop the Prefix column, we don't need it for this analysis
del df_codes['Prefix']

df_codes.head(5)
```

```
[41]:      City  Country
0  Aberdeen  Scotland
1   St Albans  England
2 Birmingham  England
3      Bath    England
4  Blackburn  England
```

```
[42]: # Merge in Country data using City as the merge key
df_merged = df_merged.merge(df_codes, on='City', how='left')

# Confirm that merging worked
df_merged.info(show_counts=True)
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 780 entries, 0 to 779
Data columns (total 16 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Survived    780 non-null    int64
1   Pclass      780 non-null    int64
2   Sex         780 non-null    object
3   Age         586 non-null    float64
4   Ticket      780 non-null    object
5   Fare        780 non-null    float64
6   Cabin       142 non-null    object
7   Embarked    780 non-null    object
8   FirstName   780 non-null    object
9   LastName    780 non-null    object
10  Title       780 non-null    object
11  MaidenName   0 non-null      object
12  City         674 non-null    object
13  Postcode     674 non-null    object
14  Address      674 non-null    object
15  Country      674 non-null    object
dtypes: float64(2), int64(2), object(12)
memory usage: 97.6+ KB
```

## Part (5)

```
[43]: # Number of observations by country, including missing
df_merged['Country'].value_counts(dropna=False)
```

```
[43]: Country
England      601
NaN          106
Scotland      67
Northern Ireland    6
Name: count, dtype: int64
```

```
[44]: # Compute survival rate by country of residence
df_merged.groupby('Country')['Survived'].mean()
```

```
[44]: Country
England      0.114809
Northern Ireland  0.000000
Scotland      0.179104
Name: Survived, dtype: float64
```

---