This notebook is an exercise in the **Data Cleaning** course. You can reference the tutorial at [this link](#).

---

In this exercise, you'll apply what you learned in the **Inconsistent data entry** tutorial.

## Setup

The questions below will give you feedback on your work. Run the following cell to set up the feedback system.

```python
from learntools.core import binder
binder.bind(globals())
from learntools.data_cleaning.ex5 import *
print("Setup Complete")
```

## Get our environment set up

The first thing we'll need to do is load in the libraries and dataset we'll be using. We use the same dataset from the tutorial.

```python
# modules we'll use
import pandas as pd
import numpy as np

# helpful modules
import fuzzywuzzy
from fuzzywuzzy import process
import chardet

# read in all our data
professors = pd.read_csv("../input/pakistan-intellectual-capital/pakistan_intellectual_capital.csv")

# set seed for reproducibility
np.random.seed(0)
```

Next, we'll redo all of the work that we did in the tutorial.

```python
# convert to lower case
professors['Country'] = professors['Country'].str.lower()
# remove trailing white spaces
professors['Country'] = professors['Country'].str.strip()

# get the top 10 closest matches to "south korea"
countries = professors['Country'].unique()
matches = fuzzywuzzy.process.extract("south korea", countries, limit=10, scorer=fuzzywuzzy.fuzz.token_sort_ratio)

def replace_matches_in_column(df, column, string_to_match, min_ratio = 47):
    # get a list of unique strings
    strings = df[column].unique()

    # get the top 10 closest matches to our input string
    matches = fuzzywuzzy.process.extract(string_to_match, strings,
                                         limit=10, scorer=fuzzywuzzy.fuzz.token_sort_ratio)

    # only get matches with a ratio > 90
    close_matches = [matches[0] for matches in matches if matches[1] >= min_ratio]

    # get the rows of all the close matches in our dataframe
    rows_with_matches = df[column].isin(close_matches)

    # replace all rows with close matches with the input matches
    df.loc[rows_with_matches, column] = string_to_match

    # let us know the function's done
    print("All done!")

replace_matches_in_column(df=professors, column='Country', string_to_match="south korea")
countries = professors['Country'].unique()
```

## 1) Examine another column

Write code below to take a look at all the unique values in the "Graduated from" column.

```python
# TODO: Your code here
countries = professors['Country'].unique()
```

Do you notice any inconsistencies in the data? Can any of the inconsistencies in the data be fixed by removing white spaces at the beginning and end of cells?

Once you have answered these questions, run the code cell below to get credit for your work.

```python
# Check your answer (Run this code cell to receive credit!)
q1.check()
```

```python
# Line below will give you a hint
#q1.hint()
```

## 2) Do some text pre-processing

Convert every entry in the "Graduated from" column in the `professors` DataFrame to remove white spaces at the beginning and end of cells.

```python
professors['Graduated from'] = professors['Graduated from'].str.strip()


q2.check()
```

```python
# Lines below will give you a hint or solution code
#q2.hint()
#q2.solution()
```

## 3) Continue working with countries

In the tutorial, we focused on cleaning up inconsistencies in the "Country" column. Run the code cell below to view the list of unique values that we ended with.

Take another look at the "Country" column and see if there's any more data cleaning we need to do.

It looks like 'usa' and 'usofa' should be the same country. Correct the "Country" column in the dataframe so that 'usofa' appears instead as 'usa'.

**Use the most recent version of the DataFrame (with the whitespaces at the beginning and end of cells removed) from question 2.**

```python
# get all the unique values in the 'City' column
countries = professors['Country'].unique()

# sort them alphabetically and then take a closer look
countries.sort()
countries
```

```python
# TODO: Your code here!
matches = fuzzywuzzy.process.extract("usa", countries, limit=10, scorer=fuzzywuzzy.fuzz.token_sort_ratio)
replace_matches_in_column(df=professors, column='Country', string_to_match="usa", min_ratio=70)

# Check your answer
q3.check()
```

```python
# Lines below will give you a hint or solution code
#q3.hint()
#q3.solution()
```

## Congratulations!

Congratulations for completing the **Data Cleaning** course on Kaggle Learn!

To practice your new skills, you're encouraged to download and investigate some of [Kaggle's Datasets](#).

---