

# Inconsistent Data Entry

Efficiently fix typos in your data.

Tutorial   Data

---



Learn Tutorial

Data Cleaning

---

In this notebook, we're going to learn how to clean up inconsistent text entries.

Let's get started!

## Get our environment set up

The first thing we'll need to do is load in the libraries and dataset we'll be using.

```
# 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)
```

## Do some preliminary text pre-processing

We'll begin by taking a quick look at the first few rows of the data.

```
professors.head()
```

	Unnamed: 0	S#	Teacher Name	University Currently Teaching	Department	Province University Located	Designation	Terminal Degree	Graduated from	Country	Year	Area of Specialization/Research Interests
0	2	3	Dr. Abdul Basit	University of Balochistan	Computer Science & IT	Balochistan	Assistant Professor	PhD	Asian Institute of Technology	Thailand	NaN	Software Engineering, DBMS
1	4	5	Dr. Waheed Noor	University of Balochistan	Computer Science & IT	Balochistan	Assistant Professor	PhD	Asian Institute of Technology	Thailand	NaN	DBMS
2	5	6	Dr. Junaid Baber	University of Balochistan	Computer Science & IT	Balochistan	Assistant Professor	PhD	Asian Institute of Technology	Thailand	NaN	Information processing, Multimedia mining
3	6	7	Dr. Maheen Bakhtyar	University of Balochistan	Computer Science & IT	Balochistan	Assistant Professor	PhD	Asian Institute of Technology	Thailand	NaN	NLP, Information Retrieval, Question Answering...
4	24	25	Samina Aftab	Sardar Bahadur Khan	Computer Science	Balochistan	Lecturer	BS	Balochistan University of	Pakistan	2005.0	VLSI Electronics

Say we're interested in cleaning up the "Country" column to make sure there's no data entry inconsistencies in it. We could go through and check each row by hand, of course, and hand-correct inconsistencies when we find them. There's a more efficient way to do this, though!

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

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

array([' Germany', ' New Zealand', ' Sweden', ' USA', 'Australia',
      'Austria', 'Canada', 'China', 'Finland', 'France', 'Greece',
      'HongKong', 'Ireland', 'Italy', 'Japan', 'Macau', 'Malaysia',
      'Mauritius', 'Netherland', 'New Zealand', 'Norway', 'Pakistan',
      'Portugal', 'Russian Federation', 'Saudi Arabia', 'Scotland',
      'Singapore', 'South Korea', 'SouthKorea', 'Spain', 'Sweden',
      'Thailand', 'Turkey', 'UK', 'USA', 'USofA', 'Urbana', 'germany'],
      dtype=object)
```

Just looking at this, I can see some problems due to inconsistent data entry: ' Germany', and 'germany', for example, or ' New Zealand' and 'New Zealand'.

The first thing I'm going to do is make everything lower case (I can change it back at the end if I like) and remove any white spaces at the beginning and end of cells. Inconsistencies in capitalizations and trailing white spaces are very common in text data and you can fix a good 80% of your text data entry inconsistencies by doing this.

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

Next we're going to tackle more difficult inconsistencies.

## Use fuzzy matching to correct inconsistent data entry

Alright, let's take another look at the 'Country' column and see if there's any more data cleaning we need to do.

```
:  
# get all the unique values in the 'Country' column  
countries = professors['Country'].unique()  
  
# sort them alphabetically and then take a closer look  
countries.sort()  
countries  
  
:  
array(['australia', 'austria', 'canada', 'china', 'finland', 'france',  
      'germany', 'greece', 'hongkong', 'ireland', 'italy', 'japan',  
      'macau', 'malaysia', 'mauritius', 'netherlands', 'new zealand',  
      'norway', 'pakistan', 'portugal', 'russian federation',  
      'saudi arabia', 'scotland', 'singapore', 'south korea',  
      'southkorea', 'spain', 'sweden', 'thailand', 'turkey', 'uk',  
      'urbana', 'usa', 'usofa'], dtype=object)
```

It does look like there is another inconsistency: 'southkorea' and 'south korea' should be the same.

We're going to use the [fuzzywuzzy](#) package to help identify which strings are closest to each other. This dataset is small enough that we could probably correct errors by hand, but that approach doesn't scale well. (Would you want to correct a thousand errors by hand? What about ten thousand? Automating things as early as possible is generally a good idea. Plus, it's fun!)

**Fuzzy matching:** The process of automatically finding text strings that are very similar to the target string. In general, a string is considered "closer" to another one the fewer characters you'd need to change if you were transforming one string into another. So "apple" and "snapple" are two changes away from each other (add "s" and "n") while "in" and "on" are one change away (replace "i" with "o"). You won't always be able to rely on fuzzy matching 100%, but it will usually end up saving you at least a little time.

Fuzzywuzzy returns a ratio given two strings. The closer the ratio is to 100, the smaller the edit distance between the two strings. Here, we're going to get the ten strings from our list of cities that have the closest distance to "south korea".

```
# get the top 10 closest matches to "south korea"
matches = fuzzywuzzy.process.extract("south korea", countries, limit=10, scorer=fuzzywuzzy.fuzz.token_sort_ratio)

# take a look at them
matches
```

```
[('south korea', 100),
 ('southkorea', 48),
 ('saudi arabia', 43),
 ('norway', 35),
 ('austria', 33),
 ('ireland', 33),
 ('pakistan', 32),
 ('portugal', 32),
 ('scotland', 32),
 ('australia', 30)]
```



We can see that two of the items in the cities are very close to "south korea": "south korea" and "southkorea". Let's replace all rows in our "Country" column that have a ratio of > 47 with "south korea".

To do this, I'm going to write a function. (It's a good idea to write a general purpose function you can reuse if you think you might have to do a specific task more than once or twice. This keeps you from having to copy and paste code too often, which saves time and can help prevent mistakes.)

```
# function to replace rows in the provided column of the provided dataframe
# that match the provided string above the provided ratio with the provided string
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!")
```

Now that we have a function, we can put it to the test!

:

```
# use the function we just wrote to replace close matches to "south korea" with "south korea"  
replace_matches_in_column(df=professors, column='Country', string_to_match="south korea")
```

All done!

And now let's check the unique values in our "Country" column again and make sure we've tidied up "south korea" correctly.

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

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

```
array(['australia', 'austria', 'canada', 'china', 'finland', 'france',
      'germany', 'greece', 'hongkong', 'ireland', 'italy', 'japan',
      'macau', 'malaysia', 'mauritius', 'netherland', 'new zealand',
      'norway', 'pakistan', 'portugal', 'russian federation',
      'saudi arabia', 'scotland', 'singapore', 'south korea', 'spain',
      'sweden', 'thailand', 'turkey', 'uk', 'urbana', 'usa', 'usofa'],
      dtype=object)
```