



# Galaxy Data Technologies

where data meets technology

[Home](#)[About](#)[Data Analysis](#)[Data Visualization](#)[Machine Learning](#)[Deployment](#)[Projects](#)

## Fuzzy String Matching With Pandas and FuzzyWuzzy

[Data Analysis](#)[Pandas](#)

sammy ongaya Fuzzy, FuzzyWuzzy, string matching

One of the big challenges with data is that it is unstructured and incomplete. How do we tell that UK is the same as United Kingdom and Unite King also same as United Kingdom? Does Jpan mean Japan? Or is guogle same as google? Well, this is what we need to understand and correct before we can start analyzing and predicting our data else we risk ignoring very important data. The data that is being generated daily is unstructured and the objective of a

### Recent posts

- [Apriori Algorithm](#)
- [Singular Value Decomposition](#)
- [Principal Component Analysis](#)
- [Hierarchical Agglomerative Clustering](#)
- [K-Means Clustering Algorithm](#)
- [Unsupervised Machine Learning](#)
- [Multi-Layer Perceptron Model](#)
- [Logistic Regression Algorithm](#)
- [Linear Regression Algorithm](#)
- [Random Forest Algorithm](#)

[Deixe uma](#)

data analyst is to standardize it into a uniform format for easier and accurate analysis. In this post we are going to learn about fuzzy string matching with Pandas and FuzzyWuzzy.

## Fuzzy String Matching With Pandas and FuzzyWuzzy

Fuzzy string matching or searching is a process of approximating strings that matches a particular pattern. Note that it gives an approximate and there is no guarantee that the string can be exact, however, sometimes the string accurately matches the pattern. How close the string is to a given match is measured by the *edit distance*. FuzzyWuzzy uses *Levenshtein Distance* to calculate the *edit distance*.

### String Matching Problem Definition

From Wikipedia here is how string matching problem is defined. *Given a pattern string  $P = p_1 p_2 \dots p_m$  and a text string  $T = t_1 t_2 \dots t_n$ , find a substring  $T_{j' : j}$ ,  $j = t_{j'} \dots t_j$  in  $T$ , which, of all substrings of  $T$ , has the smallest edit distance to the pattern  $P$ . A brute-force approach would be to compute the edit distance to  $P$  for all substrings of  $T$ , and then choose the substring with the minimum distance.*

### Application of String Matching

String matching has a wide range of applications from spell checking, spam filtering, and plagiarism detection among other uses.

### Installing FuzzyWuzzy

FuzzyWuzzy is an open-source library developed by the **SeatGeek**. SeatGeek is the Web's largest event ticket search engine. The library was created due to the problem faced by having many different ways of referring to the same event, adding and hiding location on events and concerts. Nowadays, the problem is being experienced in different domains.

## Topics

- Data Analysis (24)
  - NumPy (8)
  - Pandas (16)
- Data Visualization (20)
  - Matplotlib (14)
  - Superset (6)
- Deployment (8)
  - Deploying Data Analysis With Flask (6)
  - Jupyter Notebook Presentation (2)
- Machine Learning (22)
  - Supervised Learning (9)
  - Unsupervised Learning (6)
- Projects (8)
  - Web Scraping (6)
  - WhatsApp Data Analysis (2)

## Categories

- Data Analysis
  - NumPy
  - Pandas
- Data Visualization
  - Matplotlib
  - Superset
- Deployment
  - Deploying Data Analysis With Flask
  - Jupyter Notebook Presentation
- Machine Learning
  - Supervised Learning
  - Unsupervised Learning

Deixe uma

To install FuzzyWuzzy you can use the *pip* command as follows;

```
Pip install fuzzywuzzy
```

You can also use the conda to install FuzzyWuzzy. The following command will install the library;

```
conda install -c conda-forge fuzzywuzzy
```

## Importing The Library

To use FuzzyWuzzy we need to import the required libraries as follows.

```
1 from fuzzywuzzy import fuzz
2 from fuzzywuzzy import process
```

## Calculating Simple Ratio

```
1 from fuzzywuzzy import fuzz
2
3 simple_ratio=fuzz.ratio("Hello world", "Hello world!")
4
5 print(simple_ratio)
```

Output

96

## Calculating Partial Ratio

```
1 from fuzzywuzzy import fuzz
2
3 partial=fuzz.partial_ratio("Hello world", "Hello world!")
4
5 print(partial)
```

Output

100

## Calculating Token Sort Ration

Projects

Web Scraping

WhatsApp Data Analysis

## Tag Cloud

Anaconda analysis Apache Superset BI  
Cognos cryptography csv dashboard data  
data analysis database data exploration  
dataframe data science data  
visualization dimensionality  
reduction Facebook Graph API  
jupyter k-means logo machine  
learning matplotlib matrix  
MySQL ndarray notebook NumPy  
overview pandas pip python seaborn  
singular value decomposition split SQL  
statistics superset Tableau Twitter  
unsupervised machine learning virtualenv  
visualization visualization  
scraping

Deixe uma

## Downloads

```
1 from fuzzywuzzy import fuzz
2
3 sort_ratio=fuzz.token_sort_ratio("Hello world", "World Hello!")
4
5 print(sort_ratio)
```

Output

100

## Calculating Token Set Ratio

```
1 from fuzzywuzzy import fuzz
2
3 set_ratio=fuzz.token_set_ratio("Hello", "World Hello!")
4
5 print(set_ratio)
```

Output

100

## Extracting Strings From List of Choices

When we have a list of options and we want to match with our string we can use the *extract* or *extractOne* methods. The *limit* argument in extract method specifies the number of matches to return. The extractOne method returns exactly one match with the highest ratio.

### Extract

```
1 from fuzzywuzzy import process
2
3 choices = ["Hello world is an introductory phrase in programming",
4           "Game of throne is a movie",
5           "Albert Einstein developed the theory of relativity",
6           "Asia is the largest continent in the world"]
7
8 extract_all=process.extract("Hello world", choices, limit=2)
9
10 print(extract_all)
```

Output

NumPy sample csv data (30 downloads)  
births-and-deaths (30 downloads)  
pima-indians-diabetes (30 downloads)  
iris data set (31 downloads)  
boston-house-price-dataset (31  
downloads)  
live-graph-test (32 downloads)  
American Housing Affordability (35  
downloads)  
FMEL-Dataset (36 downloads)  
NumPy sample text data (40 downloads)  
string matched country names (41  
downloads)

## Archives

July 2018 (5)  
June 2018 (10)  
May 2018 (11)  
April 2018 (8)  
March 2018 (16)  
February 2018 (5)  
January 2018 (5)  
December 2017 (12)  
November 2017 (8)  
October 2017 (2)

Deixe uma

```
[('Hello world is an introductory phrase in programming', 90), ('Asia is the largest continent in the world', 80)]
```

### extractOne

```
1 from fuzzywuzzy import process
2
3 choices = ["Hello world is an introductory phrase in programming",
4           "Game of throne is a movie",
5           "Albert Einstein developed the theory of relativity",
6           "Asia is the largest continent in the world"]
7
8 extract_one=process.extractOne("Hello world", choices)
9
10 print(extract_one)
```

### Output

```
('Hello world is an introductory phrase in programming', 90)
```

## Fuzzy String Matching With Pandas.

We have seen how to install FuzzyWuzzy and learned how to use it. Now we are going to apply what we have learned in both Pandas and string matching to clean our sample data. If you are not familiar with Pandas I recommend to visit my previous posts on Pandas. In this project we have one data set which contains a list of countries, their country codes, area code, and region among other details in the wrong country names csv file. The country names are not correctly typed. The objective is to match the wrong country names and return the correct name for each wrongly spelled country name.

To solve the above problem we need a list of correct country names. This is provided in a country names csv file. We will string match the wrong country names with the correct country names and return a match with a ratio for each match. In our Pandas DataFrame we will create two new columns one for correct\_country\_name and another for country\_names\_ratio. We will save the output to a string\_matched\_country\_names csv file.

Deixe uma

Here are the csv files we need; wrong country names (64 downloads) , country names (59 downloads)

```
1 from fuzzywuzzy import process
2 import pandas as pd
3
4
5 names_array=[]
6 ratio_array=[]
7 def match_names(wrong_names,correct_names):
8     for row in wrong_names:
9         x=process.extractOne(row, correct_names)
10        names_array.append(x[0])
11        ratio_array.append(x[1])
12    return names_array,ratio_array
13
14
15 #Wrong country names dataset
16 df=pd.read_csv("wrong-country-names.csv",encoding="ISO-8859-1")
17 wrong_names=df['name'].dropna().values
18
19 #Correct country names dataset
20 choices_df=pd.read_csv("country-names.csv",encoding="ISO-8859-1")
21 correct_names=choices_df['name'].values
22
23 name_match,ratio_match=match_names(wrong_names,correct_names)
24
25 df['correct_country_name']=pd.Series(name_match)
26 df['country_names_ratio']=pd.Series(ratio_match)
27
28 df.to_csv("string_matched_country_names.csv")
29
30 print(df[['name','correct_country_name','country_names_ratio']].head(10))
```

Output

	name	correct_country_name	country_names_ratio
0	Ålend Islends	Åland Islands	83
1	elbenie	Niue	77
2	endorre	Andorra	71
3	eustrie	Austria	71
4	Belerus	Belarus	86
5	Belgium	Belgium	100
6	Bosnie end Herzegovine	Bosnia and Herzegovina	86
7	Bulgerie	Bulgaria	75
8	Croetie	Croatia	71
9	Czech Republic	Czech Republic	100

Deixe uma

In the above example you notice something about the matching. The second name is not well matched, *elbenie* is supposed to be *albania*. However, our string matching code has managed to get most of the matching correct. A point to note also is about the ratio, the ratio is generated by the edit distance between pattern and the string to be matched. The lower ratio does not imply that the matching is inaccurate also the higher ratio does not mean that the matching is accurate, apart from the ratio being 100 which will always be accurate unless it's a partial ratio . This can be seen on the second country name whose ratio is high but yields inaccurate results and the third country name whose ratio is low but yields accurate matching.

When working with Pandas and FuzzyWuzzy ensure that the DataFrame of strings to be matched does not have NaN else the code will throw an exception. Here is the output csv after string matching; string matched country names (41 downloads)

## Conclusion

In this post you have seen how to use string matching to compare and similar data. FuzzyWuzzy is a great library for string approximation. There are other techniques for string matching like NLP techniques that we'll see in other post. Fuzzy string approximation uses a brute-force approach in comparing string with a given pattern. The challenges of this approach is that it does not semantically work well with synonyms. In fuzzy string matching United Kingdom and England are two different strings yet semantically we know they refer to the same thing. Fuzzy string matching only applies to text data and not to other data such as sound or images. However, fuzzy string matching remains to be a very important process for data scientist in matching similar data.

Deixe uma

## About Galaxy Data Technologies

**Galaxy data technologies** is an online learning platform that focuses on technology and data. We leverage data science, software engineering and computer science concepts to learn how to solve complex problems using data.

## Galaxy Data Technologies



## Menu

- [Home](#)
- [About](#)
- [Data Analysis](#)
- [Data Visualization](#)
- [Machine Learning](#)
- [Deployment](#)
- [Projects](#)

© 2018 Galaxy Data Technologies

Powered by Massive Data Technologies

Deixe uma