
You are currently looking at **version 1.1** of this notebook. To download notebooks and datafiles, as well as get help on Jupyter notebooks in the Coursera platform, visit the [Jupyter Notebook FAQ \(https://www.coursera.org/learn/python-text-mining/resources/d9pwm\)](https://www.coursera.org/learn/python-text-mining/resources/d9pwm), course resource.

Assignment 1

In this assignment, you'll be working with messy medical data and using regex to extract relevant information from the data.

Each line of the `dates.txt` file corresponds to a medical note. Each note has a date that needs to be extracted, but each date is encoded in one of many formats.

The goal of this assignment is to correctly identify all of the different date variants encoded in this dataset and to properly normalize and sort the dates.

Here is a list of some of the variants you might encounter in this dataset:

- 04/20/2009; 04/20/09; 4/20/09; 4/3/09
- Mar-20-2009; Mar 20, 2009; March 20, 2009; Mar. 20, 2009; Mar 20 2009;
- 20 Mar 2009; 20 March 2009; 20 Mar. 2009; 20 March, 2009
- Mar 20th, 2009; Mar 21st, 2009; Mar 22nd, 2009
- Feb 2009; Sep 2009; Oct 2010
- 6/2008; 12/2009
- 2009; 2010

Once you have extracted these date patterns from the text, the next step is to sort them in ascending chronological order according to the following rules:

- Assume all dates in `xx/xx/xx` format are `mm/dd/yy`
- Assume all dates where year is encoded in only two digits are years from the 1900's (e.g. 1/5/89 is January 5th, 1989)
- If the day is missing (e.g. 9/2009), assume it is the first day of the month (e.g. September 1, 2009).
- If the month is missing (e.g. 2010), assume it is the first of January of that year (e.g. January 1, 2010).
- Watch out for potential typos as this is a raw, real-life derived dataset.

With these rules in mind, find the correct date in each note and return a pandas Series in chronological order of the original Series' indices.

For example if the original series was this:

```
0    1999
1    2010
2    1978
3    2015
4    1985
```

Your function should return this:

```
0    2
1    4
2    0
3    1
4    3
```

Your score will be calculated using [Kendall's tau \(https://en.wikipedia.org/wiki/Kendall_rank_correlation_coefficient\)](https://en.wikipedia.org/wiki/Kendall_rank_correlation_coefficient), a correlation measure for ordinal data.

This function should return a Series of length 500 and dtype int.

```
In [1]: # Hide warnings
import warnings
warnings.filterwarnings('ignore')

import pandas as pd

with open('dates.txt') as reader:
    data = pd.Series(reader.readlines())

data.head()
```

```
Out[1]: 0      03/25/93 Total time of visit (in minutes):\n
1      6/18/85 Primary Care Doctor:\n
2      sshe plans to move as of 7/8/71 In-Home Servic...\n
3      7 on 9/27/75 Audit C Score Current:\n
4      2/6/96 sleep studyPain Treatment Pain Level (N...\n
dtype: object
```

```
In [2]: data.describe()
```

```
Out[2]: count      500
unique      500
top      12/1975 Primary Care Doctor:\n
freq      1
dtype: object
```

```
In [3]: def date_sorter():

    # Extract different date formats
    format1 = '(\d{1,2}[/|\\-]\d{1,2}[/|\\-]\d{2,4})'
    format2 = '(\d{1,2}[/|\\-][1|2]\d{3})'
    format3 = '([1|2]\d{3})'

    format4 = '((?:Jan|Feb|Mar|Apr|May|Jun|Jul|Aug|Sep|Oct|Nov|Dec)[\S]*[+\\s]\d{1,2},[,{0,1}[+\\s]\d{4})'
    format5 = '((?:Jan|Feb|Mar|Apr|May|Jun|Jul|Aug|Sep|Oct|Nov|Dec)[\S]*[+\\s]\d{4})'
    format6 = '(\d{1,2}[+\\s](?:Jan|Feb|Mar|Apr|May|Jun|Jul|Aug|Sep|Oct|Nov|Dec)[\S]*[+\\s]\d{4})'

    all_format = '%s|s|s|s|s|s|s' % (format1, format2, format3, format4, format5, format6)
    extracted_date = data.str.extract(all_format)

    # Correct typos
    extracted_date = extracted_date.iloc[:,0].str.replace('Janaury', 'January').str.replace('Decemeber', 'December')

    # Convert argument to datetime.
    extracted_date = pd.Series(pd.to_datetime(extracted_date))

    # Sort ascending
    extracted_date = extracted_date.sort_values(ascending=True).index

    index = pd.Series(extracted_date.values)

    return index

# date_sorter().count()
# date_sorter().head()
```

```
In [4]: date_sorter().count()
```

```
Out[4]: 500
```

```
In [5]: date_sorter().head()
```

```
Out[5]: 0      9
1     84
2      2
3     53
4     28
dtype: int64
```

```
In [ ]:
```