

pandas_intro

August 23, 2017

0.1 Pandas

Pandas je modul namenjen brzom i efikasnom radu sa podacima.

```
In[8]: import pandas as pd
```

Za predstavljanje podataka se koriste dve strukture: **Series** i **DataFrame**

Series su strukture koje odgovaraju jednodimenzionim nizovima (čiji elementi mogu biti raznorodni) i koje imaju pridruženi niz obeležja tzv. indeks.

DataFrame su strukture koje odgovaraju tabelarnim podacima (eng. spreadsheet-like data structures) i koje imaju indekse i vrsta i kolona.

NaN označava nedostajuću vrednost. Sa **isnull** i **notnull** se može proveriti da li je neka vrednost nedostajuća ili ne.

```
In[2]: # stampa
       #\?
```

```
In[9]: # kreiranje serije
p = pd.Series([1, 1], [2, 2], [3, 3])
p
import numpy as np
p = pd.Series([1, 1], [2, 2], [3, 3])
p[1]
```

```
<Series>
0    (1, 1)
1    (2, 2)
2    (3, 3)
dtype: object
<Series>
```

```
In[10]: # ocitavanje elemenata serije
        #print(points.values)
        print(p.values)
```



```
In[52]: # pristup pojedinacnim elementima
data[1]['age']
```

```
Out[52]: 'Jh'
```

```
In[53]: data
```

```
Out[53]: pandas.DataFrame
```

```
In[54]: data['age']
```

```
Out[54]: pandas.Series
```

```
In[55]: data[1]['age']
```

```
Out[55]: 4
```

Mogu se koristiti različiti metodi filtriranja podataka:

```
In[63]: # izdvajanje korisnika koji imaju vise od 25 godina
data[data['age'] > 25]
print("\n")
data[data['age'] > 25]
```

```
age  name
1    30  New York Jh
2    33  Bed Pe
```

```
0    Fa
1    Te
2    Te
3    Fa
Name: age, dtype: object
```

Postoji mogućnost promene indeksa bilo za serije bilo za tabelarne podatke. Ukoliko se dodaju nove vrste/kolone, njihove vrednosti su NaN (konstanta np.nan)

```
In[68]: data = data[['age', 'name', 'city', 'country']]
data
data = data[['age', 'name', 'city', 'country']]
data
```

```
age  name  city  country
0    24    B  An    NaN
1    30  New York Jh    NaN
2    33  Bed Pe    NaN
3    21    Ld Ld    NaN
```

```

      age  num_m1
0   24.0    NaN
1   30.0    NaN
x    NaN    NaN    NaN    NaN
y    NaN    NaN    NaN    NaN
z    NaN    NaN    NaN    NaN
w    NaN    NaN    NaN    NaN

```

Provera nedostajucih vrednosti je moguca kroz **isnull** ili **notnull**

```
In[80]: pd.isnull()
```

```

Out[80]: 0    True
         1    True
         2    True
         3    True
         Name: num_m1, dtype: bool

```

Primer korišćenja raspoloživih podataka

```

In[15]: # učitavanje podataka iz CSV fajla
        # podaci su preuzeti sa https://www.kaggle.com/the-guardian/extinct-languages
        # i predstavljaju informacije o jezicima koji su ugrozeni ili u izumiranju
        # Kategorije koje se razmatraju su:
        # vulnerable - jezik se uci u nekom uzem kontekstu
        # definitely endangered - jezik se ne uci kao maternji jezik
        # severely endangered - jezik govore pripadnici starijih generacija, dok g
        # critically endangered - jezik govore samo pripadnici starijih generacija
        # extinct - jezici ciji govornici vise ne postoje

```

```
data = pd.read_csv('data/agg.csv')
```

```

In[70]: # dodatne informacije o read_csv funkciji
        # stampa \?

        # na raspolaganju su i funkcije za analiziranje Excel formata, JSON formata

```

```

In[71]: # rezultat citanje je DataFrame struktura
        data

```

```
Out[71]: pandas.DataFrame
```

```

In[72]: # ispis broja vrsti ucitanog skupa podataka
        data

```

```
Out[72]: 2722
```

```

In[1]: # ispis prvih 5 redova
        # stampa languages.head(5)

```



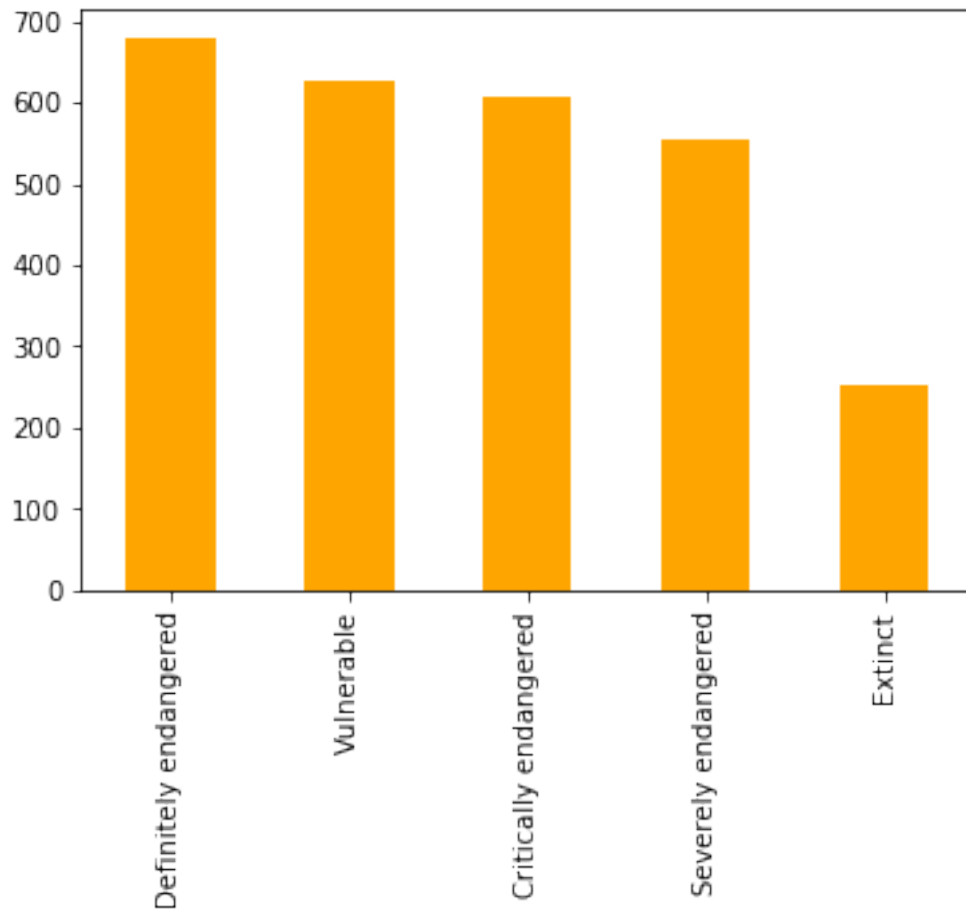
```
In[80]: # value_counts prebrojava razlicite vrednosti u kolonama
# stampa \?
```

```
In[81]: # izdvajamo informaciju o broju jezika iz svake kategorije
# kategorije koje su nam na raspolaganju su sadrzane u koloni Degree of en
print('Degree of en')
```

```
Out[81]: Degree of en
Definitely endangered    680
Vulnerable                628
Critically endangered    607
Severely endangered      554
Extinct                   253
Name: Degree of en, dtype: int64
```

```
In[82]: # iscrtavamo histogram na osnovu izracunatih podataka
print('Degree of en')
```

```
Out[82]: <matplotlib.figure.Figure at 0x8b48d0>
```



**** Dodatni materijali: ****

- <https://jupyter.org>
- **knjiga:** Python for Data Analysis, Wes McKinney

In[]: