# Checkpoint #1

Skladišta i rudarenje podataka

#### Dataset

- Koristi se dataset: "Election data 2019."
- Predstavlja rezultate parlamentarnih izbora u Portugalu, 2019. godine
- Prati promjene rezultata po teritorijima kroz vrijeme te ostale metrike koje će biti navedene u nastavku

#### Prvih 5 redaka dataseta

```
0 2019-10-06 20:10:02 Território Nacional
            0 2019-10-06 20:10:02 Território Nacional
  availableMandates numParishes numParishesApproved blankVotes \
                                                           9652
                                               1081
               226
                           3092
                                               1081
                                                          9652
               226
                           3092
                                               1081
                                                          9652
                                               1081
                                                           9652
  blankVotesPercentage nullVotes ... pre.votersPercentage \
                            8874 ...
                                                     52.66
                  2.5
                                                     52.66
                            8874 ...
                                                     52.66
                  2.5
                            8874 ...
                                                     52.66
  pre.subscribedVoters pre.totalVoters
                                         Party Mandates Percentage \
                                            PS
                                                              33.28
               813743
                                428546
                                       PPD/PSD
               813743
                                                               6.81
               813743
                                428546
                                        CDS-PP
                                                               4.90
                                                               4.59
               813743
  validVotesPercentage Votes Hondt FinalMandates
                40.22 147993
                                                77
                                                19
                        17757
                                                12
[5 rows x 28 columns]
```

## Stupci dataseta

```
# Ispis imena stupaca
print(data.columns.values)

['TimeElapsed' 'time' 'territoryName' 'totalMandates' 'availableMandates'
  'numParishes' 'numParishesApproved' 'blankVotes' 'blankVotesPercentage'
  'nullVotes' 'nullVotesPercentage' 'votersPercentage' 'subscribedVoters'
  'totalVoters' 'pre.blankVotes' 'pre.blankVotesPercentage' 'pre.nullVotes'
  'pre.nullVotesPercentage' 'pre.votersPercentage' 'pre.subscribedVoters'
  'pre.totalVoters' 'Party' 'Mandates' 'Percentage' 'validVotesPercentage'
  'Votes' 'Hondt' 'FinalMandates']
```

#### Dataset se sastoji od stupaca:

- Time elapsed, time, territoryName, totalMandates,
- availableMandates, numParishes,
- numParishesApproved, blankVotes,
- blankVotesPercentage, nullVotes,
- nullVotesPercentage, votersPercentage,
- subscribedVoters, totalVoters, pre.blankVotes,
- pre.blankVotesPercentage, pre.nullVotes,
- pre.nullVotesPercentage, pre.votersPercentage,
- pre.subscribedVoters, pre.totalVoters, Party,
- Mandates, Percentage, validVotesPercentage,
- Votes, Hondt, FinalMandates

#### Veličina dataseta

• Dataset se sastoji od 21643 retka i 28 stupaca

```
(21643, 28)
```

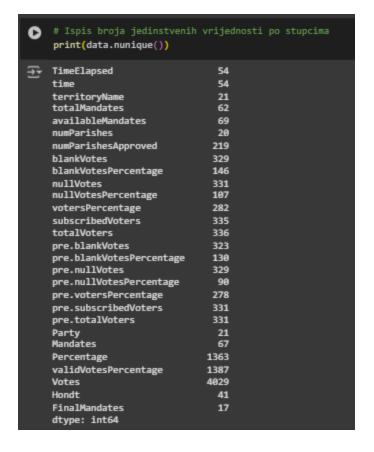
## Null vrijednosti

Dataset nema niti 1 null vrijednost



## Jedinstvene vrijednosti

Svaki stupac ima između 17 i 4029 unique vrijednosti



## Tipovi podataka

• Tipovi podataka u datasetu su: int64, float64 i object

TimeElapsed int64 time object territoryName object totalMandates int64 availableMandates int64 numParishes int64 numParishesApproved int64 blankVotes int64 blankVotesPercentage float64 nullVotes int64 nullVotesPercentage float64 subscribedVoters int64 totalVoters int64 totalVotes int64 pre.blankVotesPercentage float64 pre.nullVotes pre.nullVotes int64 pre.votersPercentage float64 pre.votalVoters int64 Party object Mandates int64 Percentage float64 Votes int64 Percentage float64 Votes int64 Hondt int64 FinalMandates int64 FinalMandates int64 FinalMandates int64

#### Broj jedinstvenih vrijednosti po stupcima

	lapsed
135	403
140	403
150	403
155	403
160	403
165	403
170	403
175	403
180	403
185	403
190	403
195	403
200	403
205	403
210	403
215	403
220	403
225	403
230	403
235	403
240	403
245	403
250	403
255	403
260	403
145	403
265	403
70	403
75	403
35	403
40	403
45	403
50	403
55	403
60	403
65	403
130	403
80	403
85	403
90 95	403 403
100	403 403
100	403 403
110	403
115	403
120	403
125	403
20	386
10	386
15	386
5	386
25	386
30	386
9	386
Name:	count, dtype: int64

```
2019-10-06 22:25:01
2019-10-06 22:30:01
                        403
2019-10-06 22:40:02
                        403
2019-10-06 22:45:01
                        403
2019-10-06 22:50:01
                        403
2019-10-06 22:55:02
                        483
2019-10-06 23:00:01
                        403
2019-10-06 23:05:01
                        403
2019-10-06 23:10:02
                        403
2019-10-06 23:15:01
                        403
2019-10-06 23:20:01
2019-10-06 23:25:02
                        403
2019-10-06 23:30:01
2019-10-06 23:35:02
                        403
2019-10-06 23:40:01
                        493
2019-10-06 23:45:02
                        403
2019-10-06 23:50:02
2019-10-06 23:55:01
                        403
403
2019-10-07 00:00:02
                        483
2019-10-07 00:05:01
                        403
2019-10-07 00:10:02
                        403
2019-10-07 00:15:01
2019-10-07 00:25:01
2019-10-07 00:30:02
                        403
2019-10-06 22:35:01
                        403
2019-10-07 00:35:01
                        403
2019-10-06 21:20:02
                        483
2019-10-06 21:25:02
                        403
2019-10-06 20:45:02
                        403
2019-10-06 20:50:01
                        483
2019-10-06 20:55:01
                        403
2019-10-06 21:00:01
                        403
2019-10-06 21:05:01
                        403
2019-10-06 21:10:01
                        403
2019-10-06 21:15:02
                        403
2019-10-06 22:20:02
                        403
2019-10-06 21:30:01
                        403
2019-10-06 21:35:01
                        403
2019-10-06 21:40:02
                        493
2019-10-06 21:45:02
                        483
2019-10-06 21:50:01 2019-10-06 21:55:01
                        403
2019-10-06 22:00:02
                        403
2019-10-06 22:05:02
2019-10-06 22:10:01
                        403
2019-10-06 22:15:01
2019-10-06 20:30:02
2019-10-06 20:20:02
2019-10-06 20:25:02
2019-10-06 20:15:02
                        386
2019-10-06 20:35:02
2019-10-06 20:40:02
                        386
2019-10-06 20:10:02
Name: count, dtype: int64
```

```
Name: count, dtype: int64
territorvName
Território Nacional
Braga
                       1134
Porto
                       1134
Leiria
                       1134
Coimbra
Lisboa
Viseu
Viana do Castelo
Setúbal
Santarém
                       1026
Êvora
                       1026
Faro
Guarda
                       1026
                        972
Castelo Branco
                        918
Portalegre
                        918
                        864
Vila Real
Name: count, dtype: int64
totalMandates
       1726
       1419
129
         21
132
         21
144
         21
149
         21
Name: count, Length: 62, dtype: int64
```

```
lame: count, Length: 62, dtype: int64
       9686
       1873
       1201
       1150
       1089
         21
         21
         21
Name: count, Length: 69, dtype: int64
3892
        1134
110
        1134
347
        1134
155
        1080
147
        1080
134
        1080
208
        1080
277
        1080
        1026
242
        1026
141
        1026
120
         972
         918
197
         864
156
         799
Name: count, dtype: int64
numParishesApproved
       1220
242
        680
        679
208
        648
170
171
Wame: count, Length: 219, dtype: int64
```

```
blankVotes
1371
        716
2026
        665
3617
697
        620
1688
        608
         16
1870
2241
         16
      count, Length: 329, dtype: int64
blankVotesPercentage
2.27
        859
2.97
        798
        740
0.54
2.64
        686
5.46
5.01
4.87
         17
1.98
      count, Length: 146, dtype: int64
nullVotes
1094
2057
        665
2292
        620
1047
        608
988
        595
        ...
1985
         16
618
Name: count, Length: 331, dtype: int64
nullVotesPercentage
1.77
        1117
1.93
         791
1.78
2.30
         701
2.29
           18
2.67
2.46
          16
Name: count, Length: 107, dtype: int64
```

#### Broj jedinstvenih vrijednosti po stupcima #2

	Percent	tage			
52.26	714				
50.58	665				
50.60	640				
50.29	620				
55.09	612				
45.74	16				
44.20	16				
44.82	16				
44.21	16				
44.28	16				
	count,	Length:	282,	atype:	1NT64
	ibedVot	tone			
122987					
151535					
240917					
258144					
136696					
130030		•			
161553	10	5			
216808					
186681					
53405	16				
181546	16	5			
		Length:	335,	dtype:	int64
totalV	oters				
64269	686	•			
76649	669				
129821					
74026	608				
51212	599	5			
73425	16				
77266	16				
81364	16				
22931	16				
84261	.10		224		
	count,	Length:	336,	atype:	10164
	ankVote				
1230	680	25			
1689	665				
3228	640				
4235	620				
1018	620				
1587	16				
1377	16				
1534	16				
393	16				
1451	16				
Name:	count,	Length:	323,	dtype:	int64

```
re.blankVotesPercentage
       1272
         904
1.92
        760
          17
17
     count, Length: 130, dtype: int64
       665
       640
2141
       620
1724
 ame: count, Length: 329, dtype: int64
re.nullVotesPercentage
1.67
       1215
        971
        903
          19
 ame: count, Length: 90, dtype: int64
57.77
        640
50.74
        630
46.72
48.30
 me: count, Length: 278, dtype: int64
```

```
pre.subscribedVoters
163462
          665
253219
255821
          628
           620
           16
197712
           16
57540
           16
           16
Name: count, Length: 331, dtype
          665
          640
128488
          620
199712
          628
77716
           16
93570
           16
81567
           16
25973
           16
90346
           16
Name: count, Length: 331, dtype
Party
              1127
PCTP/MRPP
              1127
              1127
PDR
              1127
              1127
PPD/PSD
              1127
              1127
              1127
              1127
PAN
              1127
PCP-PEV
              1127
CDS_PP
              1127
              1127
              1073
R.I.R.
              1026
MPT
              1019
PTP
              1019
              972
PURP
               972
JPP
               486
MAS
               425
Name: count, dtype: int64
```

```
Mandates
      19338
        489
        393
Name: count, Length: 67, dtype: int64
0.20
0.18
         516
0.17
         503
0.21
         446
0.16
         427
39.61
36.35
35.71
Name: count, Length: 1363, dtype: int64
validVotesPercentage
0.19
         556
         499
         441
0.12
0.22
         423
4.06
19.05
12.35
2.93
Name: count, Length: 1387, dtype: int64
```

```
17754
      1137
       723
       696
       316
       212
       107
        77
15
20
        61
12
        47
10
        39
17
        37
        37
23
        28
14
        21
91
18
        16
90
        16
11
        15
22
        15
69
        14
16
        11
76
70
92
82
13
19
74
21
93
68
81
Name: count, dtype: int64
```

Final	Mandates	i.	
0	17823		
1	1134		
2	695		
3	641		
4 5	378		
5	270		
12	108		
8	108		
106	54		
15	54		
17	54		
20	54		
6	54		
7	54		
77	54		
19	54		
9	54		
Name:	count,	dtype:	int64

## Dimenzije dataseta

- Neke od dimenzija po kojima bi se ovaj dataset mogao promatrati su:
- Vrijeme
- Teritorj
- Stranke
- Mandati
- Glasovi
- •