

```
In [1]: import numpy as np
import pandas as pd
import os
from glob import glob
from glob import glob

In [2]: filenames = []
for file in glob('data/*'):
    filenames.append(os.path.join('data', file))

In [3]: filenames[10]

Out [3]: ['data\101391.yaml',
'data\101351.yaml',
'data\1001353.yaml',
'data\1004729.yaml',
'data\1007657.yaml',
'data\1007659.yaml',
'data\101779.yaml',
'data\101981.yaml',
'data\1019983.yaml']

In [4]: final_df = pd.DataFrame()
count = 0
for file in glob(filenames):
    with open(file, 'r') as f:
        data = json.load(f)
        df = pd.DataFrame(data)
        final_df = final_df.append(df)
        count += 1

final_df

104320433 [04:25:00.00, 5.391s]

NotImplementedError Traceback (most recent call last)
<ipython-input-4-23b569e39220> in <module>
      4     with open(file, 'r') as f:
--> 5         data = pd.read_json(normalize_data(f))
      6         df = pd.DataFrame(data)
      7         final_df = final_df.append(df)

c:\users\py\appdata\local\programs\python\python39\lib\site-packages\pandas\core\fr
normalize.py in normalize_data(data, record_path, meta, meta_prefix, record_prefix,
errors, as_p, max_level)
    421     data = list(data)
    422     else:
--> 423         raise NotImplementedError
    424
    425

NotImplementedError:

In [5]: backup = final_df.copy()
del backup

In [6]: final_df

Out [6]:
innings meta.data.version meta.created meta.revision info.dates info.gender info.match.type info
0 (Team: 'Australia', 0.9 2017-02-18 2 [2017-02-17] male T20
[Team: 'Australia', 0.9 2017-02-19 2 [2017-02-19] male T20
[Team: 'Australia', 0.9 2017-02-23 1 [2017-02-22] male T20
[Team: 'Hong Kong', 0.9 2016-09-12 1 [2016-09-05] male T20
[Team: 'Zimbabwe', 0.9 2016-06-19 1 [2016-06-18] male T20
-- -- -- -- -- -- -- --
[Team: 'Sri Lanka', 0.9 2016-03-05 2 [2016-03-04] male T20
[Team: 'Bangladesh', 0.9 2016-03-08 1 [2016-03-06] male T20
[Team: 'Netherlands', 0.9 2016-02-03 1 [2016-02-03] male T20
[Team: 'Australia', 0.9 2016-09-12 1 [2016-09-06] male T20
[Team: 'Sri Lanka', 0.9 2016-09-12 1 [2016-09-09] male T20
1432 rows x 8 columns

In [7]: final_df.drop(columns=['meta.created',
'meta.revision',
'info.outcome.bowl.out',
'info.bowl.out',
'info.supersubs.south.africa',
'info.supersubs.new.zaland',
'info.outcome.eliminator',
'info.outcome.result',
'info.outcome.method',
'info.neutral.venue',
'info.match.type.number',
'info.outcome.by.wickets',
], inplace=True)

In [8]: final_df

Out [8]:
innings info.dates info.gender info.match.type info.outcome.winner info.overs info.player.of.match
0 [Team: 'Australia', 0.9 2017-02-18 2 [2017-02-17] male T20 Sri Lanka 20 [DAS Gunaratne]
[Team: 'Australia', 0.9 2017-02-19 2 [2017-02-19] male T20 Sri Lanka 20 [DAS Gunaratne]
[Team: 'Australia', 0.9 2017-02-23 1 [2017-02-22] male T20 Sri Lanka 20 [DAS Gunaratne]
[Team: 'Hong Kong', 0.9 2016-09-12 1 [2016-09-05] male T20 Hong Kong 20 [Ireland Hong Kong]
[Team: 'Zimbabwe', 0.9 2016-06-19 1 [2016-06-18] male T20 Zimbabwe 20 [E Chigumbura]
-- -- -- -- -- -- -- --
[Team: 'Sri Lanka', 0.9 2016-03-05 2 [2016-03-04] male T20 Pakistan 20 [Umar Akmal]
[Team: 'Bangladesh', 0.9 2016-03-08 1 [2016-03-06] male T20 India 20 [S Dhawan]
[Team: 'Netherlands', 0.9 2016-02-03 1 [2016-02-03] male T20 Netherlands 20 [Mudassar Bukhan]
[Team: 'Australia', 0.9 2016-09-12 1 [2016-09-06] male T20 Australia 20 [G Maxwell]
[Team: 'Sri Lanka', 0.9 2016-09-12 1 [2016-09-09] male T20 Australia 20 [G Maxwell]
1432 rows x 14 columns

In [9]: final_df['info.gender'].value_counts()

Out [9]: male 966
female 466
Name: info.gender, dtype: int64

In [10]: final_df = final_df[final_df['info.gender'] == 'male']
final_df.drop(columns=['info.gender'], inplace=True)
final_df

c:\users\py\appdata\local\programs\python\python39\lib\site-packages\pandas\core\fr
ame.py:490: SettingWithCopyWarning:
A value is being set on a copy of a slice from a DataFrame
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/1
return super().drop()

Out [10]:
innings info.dates info.match.type info.outcome.winner info.overs info.player.of.match info.team
0 [Team: 'Australia', 0.9 2017-02-18 2 [2017-02-17] male T20 Sri Lanka 20 [DAS Gunaratne] [Australia, Sri Lanka]
[Team: 'Australia', 0.9 2017-02-19 2 [2017-02-19] male T20 Sri Lanka 20 [DAS Gunaratne] [Australia, Sri Lanka]
[Team: 'Australia', 0.9 2017-02-23 1 [2017-02-22] male T20 Sri Lanka 20 [DAS Gunaratne] [Australia, Sri Lanka]
[Team: 'Hong Kong', 0.9 2016-09-12 1 [2016-09-05] male T20 Hong Kong 20 [Ireland Hong Kong]
[Team: 'Zimbabwe', 0.9 2016-06-19 1 [2016-06-18] male T20 Zimbabwe 20 [E Chigumbura] [Zimbabwe, India]
-- -- -- -- -- -- -- --
[Team: 'Sri Lanka', 0.9 2016-03-05 2 [2016-03-04] male T20 Pakistan 20 [Umar Akmal] [Pakistan, Sri Lanka]
[Team: 'Bangladesh', 0.9 2016-03-08 1 [2016-03-06] male T20 India 20 [S Dhawan] [Bangladesh, India]
[Team: 'Netherlands', 0.9 2016-02-03 1 [2016-02-03] male T20 Netherlands 20 [Mudassar Bukhan] [United Arab Emirates, Netherlands]
[Team: 'Australia', 0.9 2016-09-12 1 [2016-09-06] male T20 Australia 20 [G Maxwell] [Sri Lanka, Australia]
[Team: 'Sri Lanka', 0.9 2016-09-12 1 [2016-09-09] male T20 Australia 20 [G Maxwell] [Sri Lanka, Australia]
966 rows x 13 columns

In [11]: final_df['info.match.type'].value_counts()

Out [11]: T20 966
Name: info.match.type, dtype: int64

In [12]: final_df['info.overs'].value_counts()

Out [12]: 20 963
Name: info.overs, dtype: int64

In [13]: final_df = final_df[final_df['info.overs'] == 20]
final_df.drop(columns=['info.overs', 'info.match.type'], inplace=True)
final_df

c:\users\py\appdata\local\programs\python\python39\lib\site-packages\pandas\core\fr
ame.py:490: SettingWithCopyWarning:
A value is being set on a copy of a slice from a DataFrame
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/1
return super().drop()

Out [13]:
innings info.dates info.outcome.winner info.player.of.match info.teams info.toss.decision info.tos
0 [Team: 'Australia', 0.9 2017-02-18 2 [2017-02-17] male T20 Sri Lanka [DAS Gunaratne] [Australia, Sri Lanka] field
[Team: 'Australia', 0.9 2017-02-19 2 [2017-02-19] male T20 Sri Lanka [DAS Gunaratne] [Australia, Sri Lanka] field
[Team: 'Australia', 0.9 2017-02-23 1 [2017-02-22] male T20 Sri Lanka [DAS Gunaratne] [Australia, Sri Lanka] field
[Team: 'Hong Kong', 0.9 2016-09-12 1 [2016-09-05] male T20 Hong Kong [Ireland, Hong Kong] bat Hc
[Team: 'Zimbabwe', 0.9 2016-06-19 1 [2016-06-18] male T20 Zimbabwe [E Chigumbura] [Zimbabwe, India] field
-- -- -- -- -- -- -- --
[Team: 'Sri Lanka', 0.9 2016-03-05 2 [2016-03-04] male T20 Pakistan [Umar Akmal] [Pakistan, Sri Lanka] field
[Team: 'Bangladesh', 0.9 2016-03-08 1 [2016-03-06] male T20 India [S Dhawan] [Bangladesh, India] field
[Team: 'Netherlands', 0.9 2016-02-03 1 [2016-02-03] male T20 Netherlands [Mudassar Bukhan] [United Arab Emirates, Netherlands] field Uni
[Team: 'Australia', 0.9 2016-09-12 1 [2016-09-06] male T20 Australia [G Maxwell] [Sri Lanka, Australia] field
[Team: 'Sri Lanka', 0.9 2016-09-12 1 [2016-09-09] male T20 Australia [G Maxwell] [Sri Lanka, Australia] bat
963 rows x 11 columns

In [14]: pickle.dump(final_df, open('dataset_level1', 'wb'))

In [15]: matches = pickle.load(open('dataset_level1', 'rb'))

Out [15]:
innings info.dates info.outcome.winner info.player.of.match info.teams info.toss.decision info.tos
0 [Team: 'Australia', 0.9 2017-02-18 2 [2017-02-17] male T20 Sri Lanka [DAS Gunaratne] [Australia, Sri Lanka] field
[Team: 'Australia', 0.9 2017-02-19 2 [2017-02-19] male T20 Sri Lanka [DAS Gunaratne] [Australia, Sri Lanka] field
[Team: 'Australia', 0.9 2017-02-23 1 [2017-02-22] male T20 Sri Lanka [DAS Gunaratne] [Australia, Sri Lanka] field
[Team: 'Hong Kong', 0.9 2016-09-12 1 [2016-09-05] male T20 Hong Kong [Ireland, Hong Kong] bat Hc
[Team: 'Zimbabwe', 0.9 2016-06-19 1 [2016-06-18] male T20 Zimbabwe [E Chigumbura] [Zimbabwe, India] field
-- -- -- -- -- -- -- --
[Team: 'Sri Lanka', 0.9 2016-03-05 2 [2016-03-04] male T20 Pakistan [Umar Akmal] [Pakistan, Sri Lanka] field
[Team: 'Bangladesh', 0.9 2016-03-08 1 [2016-03-06] male T20 India [S Dhawan] [Bangladesh, India] field
[Team: 'Netherlands', 0.9 2016-02-03 1 [2016-02-03] male T20 Netherlands [Mudassar Bukhan] [United Arab Emirates, Netherlands] field Uni
[Team: 'Australia', 0.9 2016-09-12 1 [2016-09-06] male T20 Australia [G Maxwell] [Sri Lanka, Australia] field
[Team: 'Sri Lanka', 0.9 2016-09-12 1 [2016-09-09] male T20 Australia [G Maxwell] [Sri Lanka, Australia] bat
963 rows x 11 columns

In [16]: matches.loc[0]['innings'][0]['1st innings']['deliveries']

Out [16]: [0.1: {'batsman': 'AJ Finch',
'bowler': 'SI Malinga',
'run': '1',
'wicket': 'none',
'extras': 0,
'total': 1},
0.2: {'batsman': 'AJ Finch',
'bowler': 'SI Malinga',
'run': '1',
'wicket': 'none',
'extras': 0,
'total': 1},
0.3: {'batsman': 'AJ Finch',
'bowler': 'SI Malinga',
'run': '1',
'wicket
```


