

To NumPy

Understand how DataFrames can be converted to 2-D NumPy arrays.

Chapter Goals:

- Learn how to convert a DataFrame to a NumPy matrix
- Write code to modify an MLB dataset and convert it to a NumPy matrix

A. Machine learning

The DataFrame object is great for storing a dataset and performing data analysis in Python. However, most machine learning frameworks (e.g. TensorFlow), work directly with NumPy data. Furthermore, the NumPy data used as input to machine learning models must solely contain quantitative values.

Therefore, to use a DataFrame's data with a machine learning model, we need to convert the DataFrame to a NumPy matrix of quantitative data. So even the categorical features of a DataFrame, such as gender and birthplace, must be converted to quantitative values.

B. Indicator features

When converting a DataFrame to a NumPy matrix of quantitative data, we need to find a way to modify the categorical features in the DataFrame.

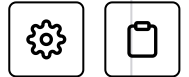
The easiest way to do this is to convert each categorical feature into a set of *indicator features* for each of its categories. The indicator feature for a specific category represents whether or not a given data sample belongs to that category.

The code below shows a DataFrame with indicator features.

```
1 # predefined non-indicator DataFrame
2 print('{}\n'.format(df))
```



```
3
4 # predefined indicator DataFrame
5 print('{}\n'.format(indicator_df))
```



Output

1.314s

```
      color
r1    red
r2   blue
r3  green
r4    red
r5    red
r6   blue

blue green red
```



In the code above, the DataFrame `df` has a single categorical feature called `color`. The corresponding indicator features for `color` are shown in `indicator_df`.

Note that an indicator feature contains `1` when the row has that particular category, and `0` if the row does not.

C. Converting to indicators

In pandas, we convert each categorical feature of a DataFrame to indicator features with the `get_dummies` (https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.get_dummies.html) function. The function takes in a DataFrame as its required argument, and returns the DataFrame with each of its categorical features converted to indicator features.

The code below demonstrates how to use the `get_dummies` function.

```
1 # predefined df
2 print('{}\n'.format(df))
3
```



```

4 converted = pd.get_dummies(df)
5 print('{}\n'.format(converted.columns))
6
7 print('{}\n'.format(converted[['teamID_BOS',
8                               'teamID_PIT']]))
9 print('{}\n'.format(converted[['lgID_AL',
10                               'lgID_NL']]))

```



Output

1.340s

	lgID	teamID
playerID		
bettsmo01	AL	BOS
martest01	NL	PIT
pedrodu01	AL	BOS
polangr01	NL	PIT

Index(['lgID_AL', 'lgID_NL', 'teamID_BOS', 'teamID_PIT'], dtype='object')

Note that the indicator features have the original categorical feature's label as a prefix. This makes it easy to see where each indicator feature originally came from.

D. Converting to NumPy

After converting all the categorical features to indicator features, the DataFrame should have all quantitative data. We can then convert to a NumPy matrix using the `values` (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.values.html#pandas.DataFrame.values>) function.

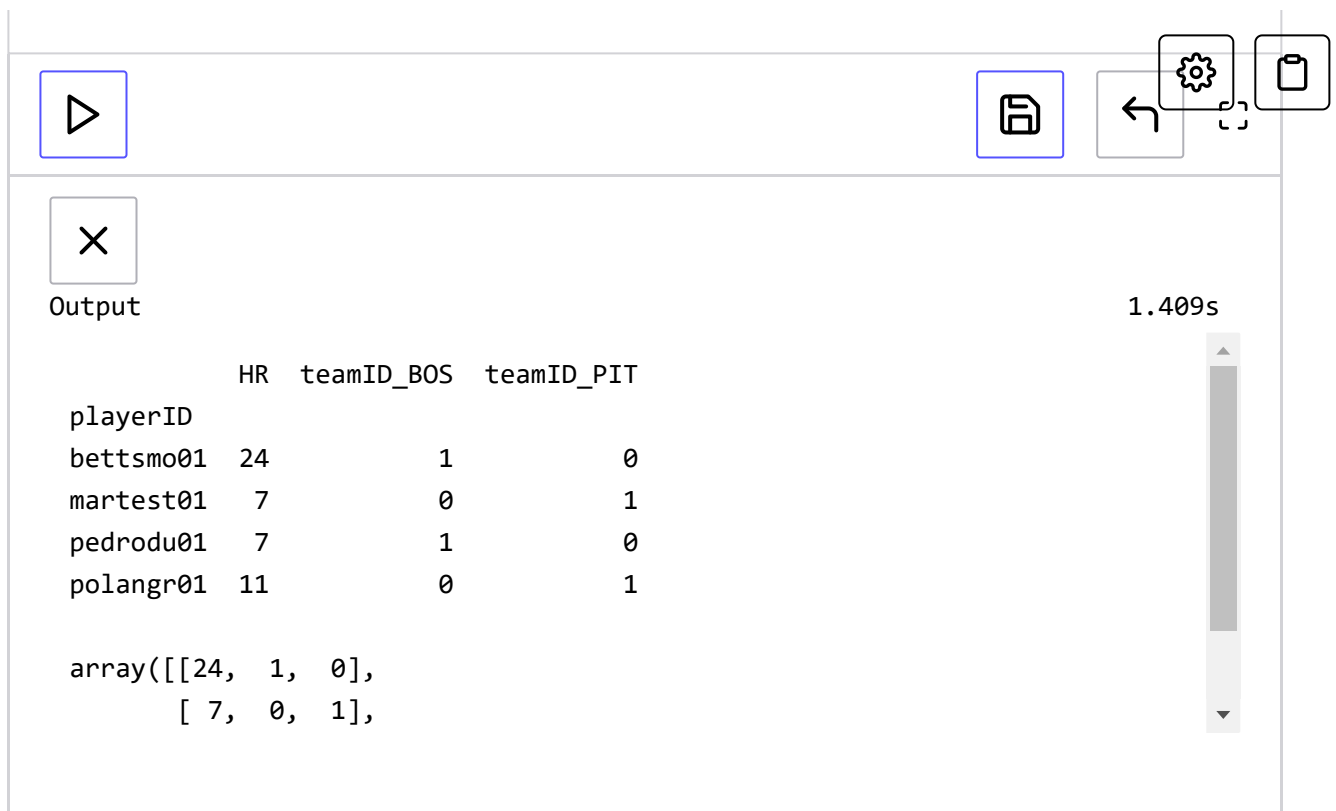
The code below converts a DataFrame, `df` into a NumPy matrix.

```

1 # predefined indicator df
2 print('{}\n'.format(df))
3
4 n_matrix = df.values
5 print(repr(n_matrix))

```





The image shows a Jupyter Notebook interface. At the top, there are icons for running (play), saving (disk), undo (left arrow), settings (gear), and copy (document). Below these is a code cell with a close button (X) and the word "Output". The output area displays a table of data and a NumPy array. The table has columns: playerID, HR, teamID_BOS, and teamID_PIT. The data rows are: bettsmo01 (24 HR, 1 BOS, 0 PIT), martest01 (7 HR, 0 BOS, 1 PIT), pedrodu01 (7 HR, 1 BOS, 0 PIT), and polangr01 (11 HR, 0 BOS, 1 PIT). Below the table, the NumPy array is shown: `array([[24, 1, 0], [7, 0, 1],`. A vertical scrollbar is on the right, and the execution time "1.409s" is shown in the top right corner of the output area.

playerID	HR	teamID_BOS	teamID_PIT
bettsmo01	24	1	0
martest01	7	0	1
pedrodu01	7	1	0
polangr01	11	0	1

```
array([[24, 1, 0],
       [ 7, 0, 1],
```

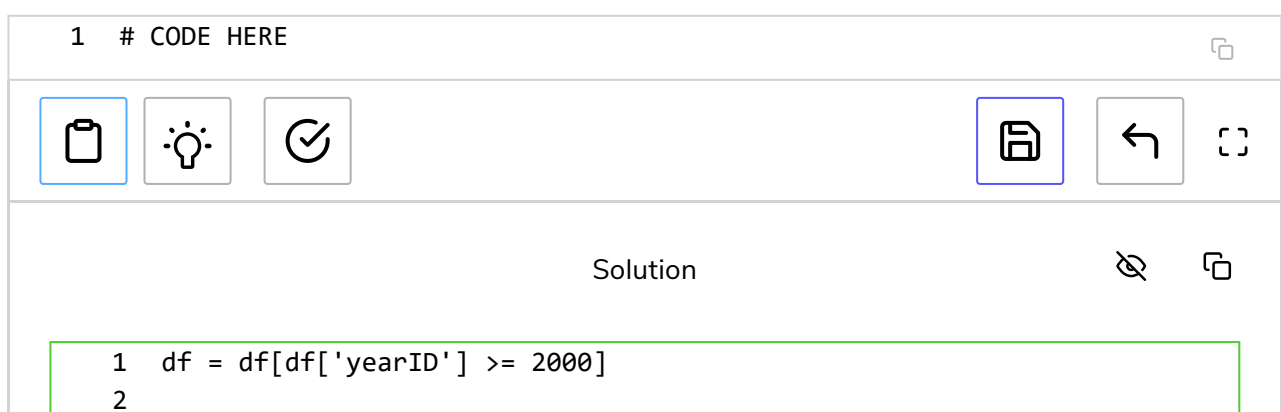
The rows and columns of the output matrix correspond to the rows and columns of the same position in the DataFrame. In the code above, the first column of the NumPy matrix represents `HR`, the second column represents `teamID_BOS`, and the third column represents `teamID_PIT`.

Time to Code!

The code exercise for this chapter will be to convert a DataFrame of MLB statistics (`df`) into a NumPy matrix.

We only want the data in `df` to be from the current century, so we need to first apply a filter.

Filter `df` for rows where `'yearID'` is at least 2000, then reset `df` equal to the filtered output.



The image shows a Jupyter Notebook interface. At the top, there is a code cell with the text "1 # CODE HERE". Below the code cell are icons for saving (disk), undo (left arrow), redo (right arrow), and copy (document). Below these is a solution area with the word "Solution" and a close button (X). The solution area contains a code cell with the following code:






```
1 df = df[df['yearID'] >= 2000]
2
```



We also don't want any of the NaN values in our data. We can filter those out using the special `dropna` (<https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.dropna.html>) function.

Set `df` equal to `df.dropna` applied with no arguments.

```
1 # CODE HERE
```



Solution

```
1 df = df.dropna()  
2  
3
```






Finally, we want to convert each categorical feature into a set of indicator features for each of its categories.

Then we can convert `df` into a NumPy matrix.

Set `df` equal to `pd.get_dummies` with `df` as the only argument.

Set `matrix` equal to `df.values`.

```
1 # CODE HERE
```



Solution

```
1 df = pd.get_dummies(df)  
2 matrix = df.values  
3
```

← Back

Plotting

Next



Quiz

☒ Mark as Completed

27% completed, meet the criteria and claim your course certificate!

Buy Certificate



Report an
Issue



Ask a Question

(https://discuss.educative.io/tag/to-numpy__data-analysis-with-pandas__machine-learning-for-software-engineers)