v1.20.4 CS-E4002 Python for Engineers ▼

<

« Calculations Course materials Lambda in Python »

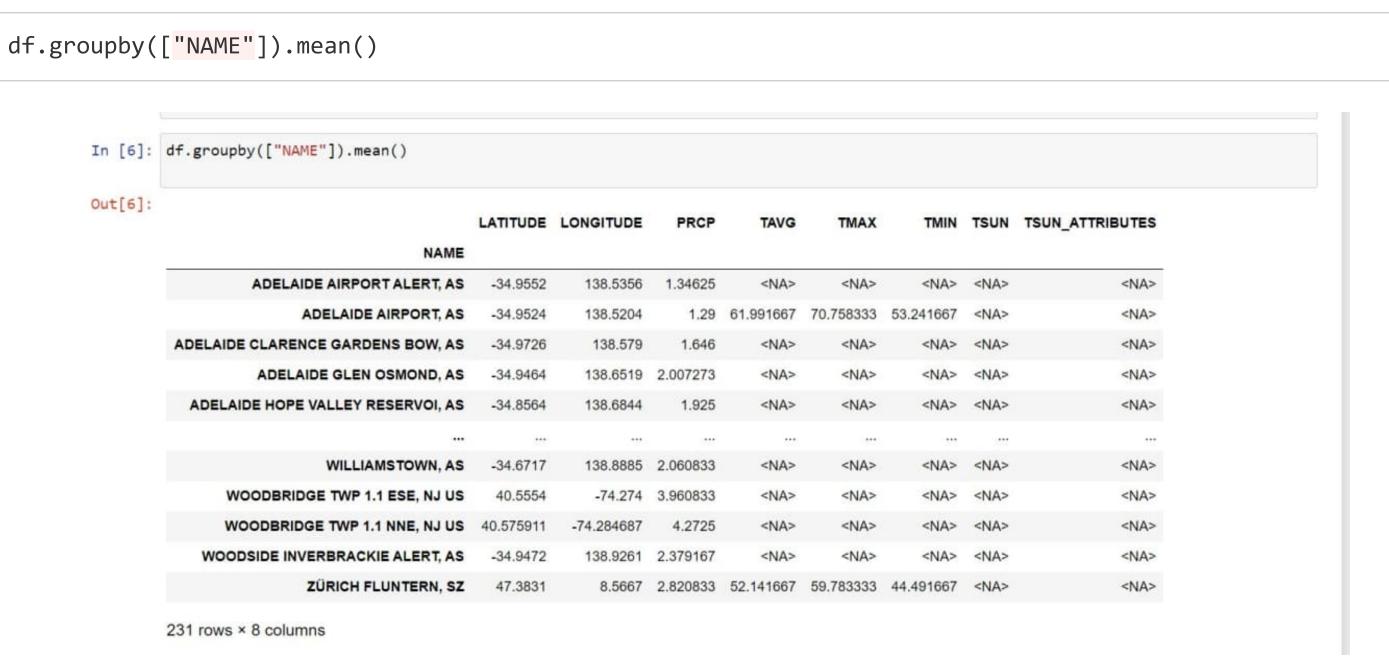
## Data Analysis

CS-E4002 / Calculations / Data Analysis

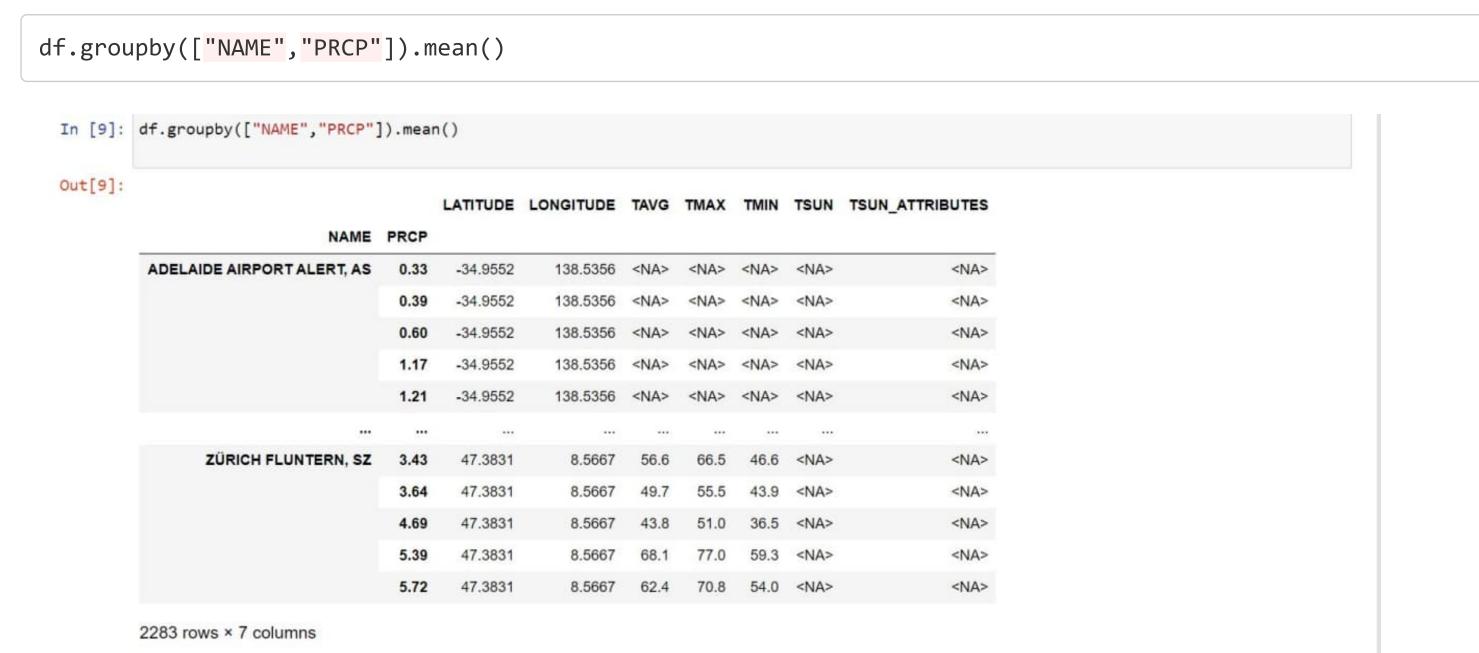
Pandas has many useful built-in functions to ease data analysis. We will present some here, but take your time to discover them.

## **Overlay Structure**

A very convenient one is the possibility to group the dataframe according to one of the columns. As an example if you have weather data, and you are interested in the average temperature per measuring station. You can group the dataframe by the name of the station (in our example the weather station name is stored in the column "NAME") and get the corresponding average temperature.



If you are interested in a subgroup of a subgroup, Pandas offers the possibility to group the dataframe by a list of columns - this will generate subgroups accordingly



The order of the elements in the list corresponds to the substructure Pandas applies This means that df.groupby(["PRCP","NAME"]).mean() is not the same as df.groupby(["NAME","PRCP"]).mean()

## Adding Columns and Rows

Sometimes we want to add an additional column to a dataframe. This can be done by just assigning a series as a new column to the dataframe. Use a name that was not used so far and pandas adds it as a new column. Make sure the series has the same size as the rest of the dataframe. Lets say we want to add as an additional information the variation in temperature. We want to store in an additional column the difference between the minimum and the maximum temperature. We can first calculate the difference between the maximum and the minimum temperature measured:

```
df_diffT = df["TMAX"]-df["TMIN"]
```

This results in a pandas series with the same number of rows as the original dataframe.

We can now just generate a new dataframe column by assigning this series to the dataframe by choosing the corresponding columnname

```
df["DIFFT"] = df["TMAX"]-df["TMIN"]
```

We could also add as an information if the station of an entry is on the north hemisphere by adding a column to the dataframe, which will have the value "TRUE" for all stations lying in the northern hemisphere and FALSE otherwise

```
df["N_HEM"] = (df["LATITUDE"] > 0)
```

If you want to add another row to your dataframe you can make use of the "append" functionality. As "append" does not change the dataframe (is immutable) do not forget to reassign the new dataframe. We will enter a copy of the fourth row to avoid all the typing for our weather data, but you can define your new row in the form of a dictionary completely from scratch (as indicated in the code below in a shortened version) For our weather data this could look like that:

```
new_row = df.iloc[3]
#new_row = {'STATION:'US1NJMD0086', 'NAME':'HIGHLAND', 'LATITUDE':40.50542, 'LONGITUDE':-74.42978}
#append row to the dataframe
df = df.append(new_row, ignore_index=True)
```

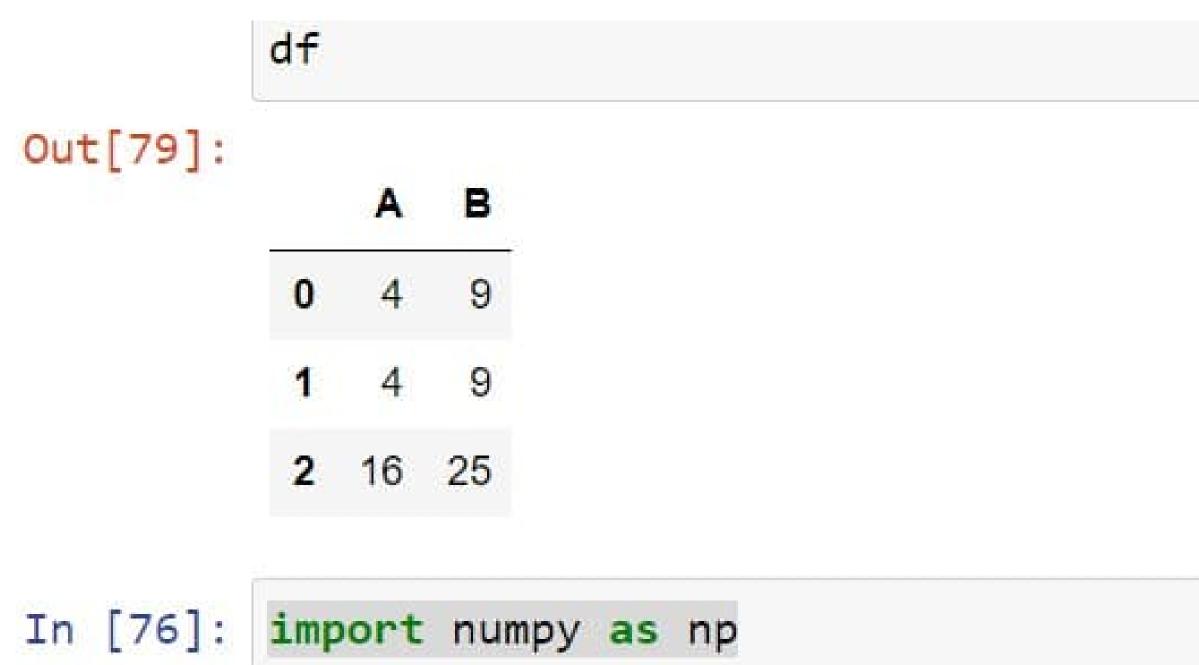
## Changing parts of the dataframe

Sometimes we want to apply changes to a whole column. If we for example want to change the unit of the temperature from Fahrenheit to Celsius, we can do that directly to a whole column. We first distract 32 and then multiply it by 5/9.0 and then you have your average data stored in Celcius.

```
df_new["TAVG"] = df_new["TAVG"]-32
df_new["TAVG"] = df_new["TAVG"]*(5/9.0)
```

Maybe you want to do changes to the whole dataframe, and not only a colum. For that pandas offers a higher order function called apply. A higher order function is a function that takes as an argument another function. With the apply method pandas applies the function to each data entry. If you have rows and columns filled with numbers and you want to have the square root of these numbers instead you can apply the np.sqrt function to each element like this

```
import numpy as np
df.apply(np.sqrt)
```



```
In [76]: import numpy as np
In [80]: df.apply(np.sqrt)
Out[80]:
```

Feedback 🕑

A+ v1.20.4