
▼ Everything is Better with Friends

Using SAS in Python Applications with SASPy and Open-Source Tooling (Beyond the Basics)

▼ Setup for Part 3

Getting setup to use Google Colab with SAS OnDemand for Academics (ODA)

1. To execute code cells, you'll need credentials for the following accounts:

- Google. (If you're not already signed in, you should see a **Sign In** button in the upper right corner. You can also visit <https://accounts.google.com/signup> to create an account for free.)
- SAS OnDemand for Academics. (You can create an account for free at <https://welcome.oda.sas.com/> using an existing SAS Profile account. If you don't already have a SAS Profile account, you can create one for free using the "Don't have a SAS Profile?" link on the ODA login page.)

2. We recommend enabling line numbers using the Tools menu: **Tools -> Settings -> Editor -> Show line numbers -> Save**

3. We also recommend enabling the Table of Contents using the View menu: **View -> Table of contents**

4. To save a copy of this notebook, along with any edits you make, please use the File menu: **File -> Save a copy in Drive**

5. Looking for "extra credit"? Please let us know if you spot any typos!

▼ Connect to SAS OnDemand for Academics (ODA) and start a SAS session

Instructions:

1. Determine the Region for your ODA account by logging into <https://welcome.oda.sas.com/>. You should see a value like `Asia Pacific 1`, `Asia Pacific 2`, `Europe 1`, `United States 1`, or `United States 2` next to your username in the upper-right corner. (For more information about Regions and using Python in Jupyter Notebooks, please see the ODA documentation at https://support.sas.com/ondemand/cag_new.html#region and <https://support.sas.com/ondemand/saspy.html>.)
2. If your ODA account is associated with a Region other than `United States 1`, comment out Line 11 by adding a number sign (`#`) at the beginning of the line, and then uncomment the list of servers corresponding to your Region.

Note: As of the time of creation of this Notebook, only the Regions listed below were available. If your SAS ODA account is associated with a Region that's not listed, you will need to manually add the appropriate servers.
3. Click anywhere in the code cell, and run the cell using Shift-Enter.
4. At the prompt `Please enter the OMR user id`, enter either your SAS ODA user ID or the email address associated with your ODA account.
5. At the prompt `Please enter the password for OMR user`, enter the password for your SAS ODA account.

```
1 !pip install saspy
2
3 import saspy
4
5 sas = saspy.SASsession(
6     java='/usr/bin/java',
7     iomport=8591,
8     encoding='utf-8',
9
10     # For Region "United States 1", uncomment the line below.
11     iomhost = ['odaws01-usw2.oda.sas.com', 'odaws02-usw2.oda.sas.com', 'odaws03-usw2.oda.sas.com', 'odaws04-usw2.
12
13     # For Region "United States 2", uncomment the line below.
14     #iomhost = ['odaws01-usw2-2.oda.sas.com', 'odaws02-usw2-2.oda.sas.com'],
15
16     # For Region "Europe 1", uncomment the line below.
17     #iomhost = ['odaws01-euw1.oda.sas.com', 'odaws02-euw1.oda.sas.com'],
18
19     # For Region "Asia Pacific 1", uncomment the line below.
20     #iomhost = ['odaws01-apsel.oda.sas.com', 'odaws02-apsel.oda.sas.com'],
```

```

21
22     # For Region "Asia Pacific 2", uncomment the line below.
23     #iomhost = ['odaws01-apsel-2.oda.sas.com', 'odaws02-apsel-2.oda.sas.com'],
24
25 )
26 print(sas)

```

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>
 Collecting saspy
 Downloading saspy-4.3.2.tar.gz (9.9 MB)
 |██| 9.9 MB 4.7 MB/s
 Building wheels for collected packages: saspy
 Building wheel for saspy (setup.py) ... done
 Created wheel for saspy: filename=saspy-4.3.2-py3-none-any.whl size=9929529 sha256=ae5fb36c680ee683cc9c11be
 Stored in directory: /root/.cache/pip/wheels/55/a2/91/45db2a8ca68bcb2ee02c28366dc2c36d40e1af670a9ba96e12
 Successfully built saspy
 Installing collected packages: saspy
 Successfully installed saspy-4.3.2
 Using SAS Config named: default
 Please enter the OMR user id: isaiah.lankham@ucop.edu
 Please enter the password for OMR user :
 SAS Connection established. Subprocess id is 133

```

Access Method           = IOM
SAS Config name         = default
SAS Config file         = /usr/local/lib/python3.7/dist-packages/saspy/sascfg.py
WORK Path               = /saswork/SAS_workABFB00000530_odaws01-usw2.oda.sas.com/SAS_workE01900000530_odaws01-u
SAS Version             = 9.04.01M6P11072018
SASPy Version          = 4.3.2
Teach me SAS            = False
Batch                   = False
Results                 = Pandas
SAS Session Encoding    = utf-8
Python Encoding value   = utf-8
SAS process Pid value   = 1328

```

Note: This establishes a connection from Python in Google Colab to a SAS session running in SAS ODA.

▼ Install and import additional packages

```
1 # Install the faker module for generating fake data
2 !pip install faker
3
4 # Install the rich module for colorful printing
5 !pip install rich
6
7 # Initialize a Faker session called fake
8 from faker import Faker
9 fake = Faker()
10
11 # We'll use IPython to display DataFrames or HTML content
12 from IPython.display import display, HTML
13
14 # We'll use the pandas package to create and manipulate DataFrame objects
15 import pandas
16
17 # We're overwriting the default print function with rich.print
18 from rich import print
```

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>
Collecting faker

Downloading Faker-14.2.0-py3-none-any.whl (1.6 MB)

|██| 1.6 MB 5.1 MB/s

Requirement already satisfied: python-dateutil>=2.4 in /usr/local/lib/python3.7/dist-packages (from faker) (2

Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.7/dist-packages (from fak

Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/dist-packages (from python-dateutil>=2.4-

Installing collected packages: faker

Successfully installed faker-14.2.0

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>

Collecting rich

Downloading rich-12.5.1-py3-none-any.whl (235 kB)

|██| 235 kB 4.9 MB/s


Requirement already satisfied: pygments<3.0.0,>=2.6.0 in /usr/local/lib/python3.7/dist-packages (from rich) (

Collecting commonmark<0.10.0,>=0.9.0

Downloading commonmark-0.9.1-py2.py3-none-any.whl (51 kB)

|██| 51 kB 5.9 MB/s

```
Requirement already satisfied: typing-extensions<5.0,>=4.0.0 in /usr/local/lib/python3.7/dist-packages (from  
Installing collected packages: commonmark, rich  
Successfully installed commonmark-0.9.1 rich-12.5.1
```



▼ Part 3. Merging and appending datasets in SAS and Python

▼ Section 3.1. Create class_df

```
1 # Let's start by importing and displaying the dataset sashelp.class from SAS ODA.  
2  
3 # Use the sas.sasdata2dataframe method to copy the contents of sashelp.class into DataFrame  
4 # class_df, and use dataset options to get only the first few rows and specific columns.  
5 class_df = sas.sasdata2dataframe(  
6     table='class',  
7     libref='sashelp',  
8     dsopts={  
9         'keep': ['Name', 'Age', 'Height'],  
10        'obs': 7  
11    },  
12 )  
13  
14 # Display the resulting DataFrame.  
15 display(class_df)
```

	Name	Age	Height
0	Alfred	14.0	69.0
1	Alice	13.0	56.5

Concept Check 3.1

- Short Answer: What are some other dataset options we might consider including?
- Fun Fact: The DataFrame `class_df` is a data structure kept in memory in our local Google Colab session, having copied the contents of the dataset `sashelp.class` over the wire.

6 Jane 12.0 50.8

Solution: An overview of the dataset options available in SASPy can be found at

<https://sassoftware.github.io/saspy/api.html#saspy.SASsession.sasdata>

▼ Section 3.2. Create `additional_students_df`

```

1 # Now imagine we want to create additional example student records to append to sashelp.class, but
2 # we don't feel like being creative.
3
4 # Instead, let's use a "standard recipe" to have Python make a DataFrame with fake values for us!
5
6 # Set the number of rows of mock data to generate.
7 number_of_rows = 5
8
9 # Define a function that returns a fake first name that's not already in the Name column of the
10 # DataFrame class_df defined above.
11 def create_distinct_first_name():
12     random_first_name = fake.unique.first_name()
13     while random_first_name in class_df['Name'].unique():
14         random_first_name = fake.unique.first_name()
15     return random_first_name
16
17 # Generate and display a DataFrame called additional_students_df with the specified number_of_rows

```

```

18 # Note: This example uses more advanced Python features, including list comprehensions. There's no
19 # need to focus on anything other than the overall intent of creating a DataFrame with mock data.
20 additional_students_df = pandas.DataFrame(
21     [
22         {
23             'Name': create_distinct_first_name(),
24             'Age': fake.pyint(min_value=10,max_value=19),
25             'Height': fake.pyfloat(right_digits=1,min_value=50,max_value=75),
26         }
27         for _
28         in range(number_of_rows)
29     ]
30 )
31
32 display(additional_students_df)

```

	Name	Age	Height
0	Christopher	14	51.4
1	Antonio	19	59.6
2	Mackenzie	15	68.2
3	Eduardo	10	50.5
4	Dustin	17	57.8

Concept Check 3.2

- Try this, and see what happens: Change the underscore (`_`) on Line 27 to any valid Python variable name, and then rerun the code cell above.
- True or False: The code runs just as well, whether an underscore (`_`) or an actual variable name is used.
- Fun Facts:
 - It's standard convention in the Python community to use an underscore (`_`) for a variable whose actual value isn't important.

- You might remember the "Fun Fact" from Part 2 about DataFrame operations being slow when embedded in a for-loop. Here, we've effectively turned this on its head by using a list comprehension to embed a for-loop inside of a DataFrame operation, which is a fairly standard (and highly efficient) Python practice.

Solution: True! The single underscore (`_`) is commonly used as a general-purpose "throwaway" variable in Python.

▼ Section 3.3. Create `appended_df`

```
1 # Now that we have two DataFrames with identical columns, we can vertically combine (aka append or
2 # union) them to create a new, taller dataset.
3
4 # Starting with class_df, append each row from additional_students_df, and display the result.
5 appended_df = class_df.append(additional_students_df, ignore_index=True)
6 display(appended_df)
```



```
19         ;
20     quit;
21     proc print data=appended_sds;
22     run;
23     '''
24 )
25
26 # But what type of object exactly do we get back from sas.submit?
27 print(f'The type of sas_submit_return_value: {type(sas_submit_return_value)}')
28
29 # And what are the keys in this dict?
30 print(f'\n\nThe keys in sas_submit_return_value: {list(sas_submit_return_value.keys())}')
31
32 # Output the SAS log, which corresponds to the key 'LOG' in the dict returned by sas.submit.
33 print('\n\n\nThe LOG component of sas_submit_return_value:')
34 sas_submit_log = sas_submit_return_value['LOG']
35 print(sas_submit_log)
36
37 # Render and display the SAS HTML results, which corresponds to the key 'LST' in the dict returned
38 # by sas.submit.
39 print('\n\n\nThe LST (aka results) component of sas_submit_return_value:')
40 sas_submit_results = sas_submit_return_value['LST']
41 display(HTML(sas_submit_results))
```

The type of sas_submit_return_value: <class 'dict'>

The keys in sas_submit_return_value: ['LOG', 'LST']

The LOG component of sas_submit_return_value:

72

The SAS System

Monday, September 12, 2022 11:26:00 PM

2022 11:26:00 PM

```
586      ods listing close;ods html5 (id=saspy_internal) file=_tomods1 options(bitmap_mode='inline') device=
style=HTMLBlue;
586      ! ods graphics on / outputfmt=png;
587
588
589      proc sql;
590          create table appended_sds as
591              select Name, Age, Height from sashelp.class(obs=7)
592              union all corr
593              select * from additional_students_sds
594          ;
595      quit;
596      proc print data=appended_sds;
597      run;
598
599
600
601      ods html5 (id=saspy_internal) close;ods listing;
602
```

73

The SAS System

Monday, September 12, 2022 11:26:00 PM

2022 11:26:00 PM

603

The LST (aka results) component of sas_submit_return_value:

The SAS System

Obs	Name	Age	Height
1	Alfred	14	69.0
2	Alice	13	56.5
3	Barbara	13	65.3

Concept Check 3.4

- Short Answer: List some others ways to vertically combine datasets in SAS.
- Fun Facts:
 - Because the SAS dataset `work.appended_sds` is created inside of the `submit` method, our Python session has no access to it. To use the contents of `work.appended_sds` in our Python session, create a `SASdata` object as follows:


```
sas.sasdata(table='appended_sds')
```
 - The `submit` method always returns a dictionary with the same two keys, `LOG` and `LST`.

Solution: Options include PROC APPEND and A SET statement in a DATA step.

▼ Section 3.5. Create age_ranges_df

```

1 # Now suppose we want to add a column describing the age range for each student in our full
2 # DataFrame.
3
4 # We'll start by building a lookup table:
5
6 # As a first step, let's make a DataFrame with a single column called Age, using the built-in range
7 # function to populate this column with the integer values 10 through 19. (Note: The range function
8 # always stops at one less than the specified upper bound.)
9 age_ranges_df = pandas.DataFrame(
10     {
11         'Age': range(10,20)
12     }
13 )
14
15 # Now add a column to age_ranges_df by "binning" integers in age_ranges_df['Age'] into the
16 # categories 10-12 (tween) and 13-19 (teen).
17
18 # In other words, use the pandas.cut method with these arguments:
19 # * bins=[10,12,19], which creates the values ranges [10,12] and (12,19]
20 # * labels=['tween','teen'], which specifies the label for each range, in order
21 age_ranges_df['Age_Range'] = pandas.cut(

```

```

22     age_ranges_df['Age'],
23     bins=[10,12,19],
24     labels=['tween','teen'],
25     include_lowest=True,
26 )
27 display(age_ranges_df)

```

	Age	Age_Range
0	10	tween
1	11	tween
2	12	tween
3	13	teen
4	14	teen
5	15	teen
6	16	teen
7	17	teen
8	18	teen
9	19	teen

Concept Check 3.5

- Try this, and see what happens: Comment out the option `include_lowest=True` on Line 25, and then rerun the code cell above.
- True or False: Commenting out the option `include_lowest=True` on Line 25 won't affect the contents of the resulting DataFrame.
- Fun Fact: The `pandas.cut` method could have instead been used directly on DataFrame `appended_df` to bin values of `Age`.

Solution: False! Deleting `include_lowest=True` will cause the value 10 to not have a corresponding label.

▼ Section 3.6. Create merged_df

```
1 # Given the lookup table age_ranges_df, which has some (but not all) of its columns in common with
2 # appended_df, we can horizontally combine (aka merge or join) them to create a new, wider dataset.
3
4 # Specifically, starting with appended_df, we'll add an Age_Range column whose values are determined
5 # by matching values of Age in age_ranges_df as part of a left join, resulting in a DataFrame having
6 # the same height as appended_df.
7 merged_df = appended_df.merge(
8     age_ranges_df,
9     on='Age',
10    how='left',
11 )
12
13 # Since we're not using the Height column, drop it from merged_df, and display the result.
14 merged_df.drop(columns=['Height'], inplace=True)
15 display(merged_df)
```

	Name	Age	Age_Range
0	Alfred	14.0	teen

Concept Check 3.6

- Try this, and see what happens: Delete the option `inplace=True` on Line 14, and then rerun the code cell above.
- True or False: Deleting the option `inplace=True` on Line 14 won't affect the contents of the resulting DataFrame.
- Fun Fact: Just like many SQL dialects, `pandas.merge` supports left, right, (full) outer, inner, and cross joins.

Solution: False! The column will not be dropped successfully unless you use `inplace=True`, or unless you assign the result `merged_df.drop(columns=['Height']` to a variable.

```
0      Alfred    14.0      teen
```

▼ Section 3.7. Create merged_sds

```
1 # TANGENT/ASIDE: Since the DataFrame additional_students_df was already copied over to SAS ODA
2 # above, we could have instead copied over age_ranges_df and used PROC SQL to perform the left join.
3 # Let's see how these two methods compare!
4
5 # Make a SAS dataset from age_ranges_df.
6 sas.dataframe2sasdata(
7     age_ranges_df,
8     table="age_ranges_sds",
9     libref="Work"
10 )
11
12 # Use the sas.submit method to submit a PROC SQL step directly to ODA, and capture the resulting
13 # dict in sas_submit_return_value.
14 sas_submit_return_value = sas.submit(
15     '''
16         proc sql;
17             create table merged_sds as
18                 select
```

```

19             A.Name
20             ,A.Age
21             ,B.Age_Range
22         from
23             appended_sds as A
24             left join
25             age_ranges_sds as B
26             on A.Age = B.Age
27         ;
28     quit;
29     proc print data=merged_sds;
30     run;
31     '''
32 )
33
34 # Output the SAS log, which corresponds to the key 'LOG' in the dict returned by sas.submit.
35 sas_submit_log = sas_submit_return_value['LOG']
36 print(sas_submit_log)
37
38 # Render and display the SAS results HTML, which corresponds to the key 'LST' in the dict returned
39 # by sas.submit.
40 sas_submit_results = sas_submit_return_value['LST']
41 display(HTML(sas_submit_results))

```


84
2022 11:26:00 PM

The SAS System

Monday, September 13, 2022 11:26:00 PM

```
695      ods listing close;ods html5 (id=saspy_internal) file=_tomods1 options(bitmap_mode='inline') device=
style=HTMLBlue;
695      ! ods graphics on / outputfmt=png;
696
697
698          proc sql;
699              create table merged_sds as
700                  select
701                      A.Name
702                      ,A.Age
703                      ,B.Age_Range
704                  from
705                      appended_sds as A
706                      left join
707                      age_ranges_sds as B
708                      on A.Age = B.Age
709              ;
710          quit;
711          proc print data=merged_sds;
712          run;
713
714
715
716      ods html5 (id=saspy_internal) close;ods listing;
717
```

85
2022 11:26:00 PM

The SAS System

Monday, September 13, 2022 11:26:00 PM

718

The SAS System

Obs	Name	Age	Age_Range
1	Eduardo	10	tween
2	James	12	tween
3	Jane	12	tween
4	Barbara	13	teen

Concept Check 3.7

- Short Answer: List some other ways to get the same result in SAS without a join.
- Fun Fact: Just like in Section 3.4, because the SAS dataset `Work.merged_sds` is created inside of the `submit` method, our Python session has no access to it unless we create a `SASdata` object as follows: `sas.sasdata(table='merged_sds')`

Solution: Options include a `MERGE` or `UPDATE` statements in a `DATA` step, or creating age ranges using `PROC FORMAT`.

▼ Section 3.8. Additional Exercises

For practice, we recommend the following:

1. Run the code cell below to convert the `Height` column of `class_df` to integer values and create a reference table called `height_ranges_df`.
2. Then add new code, repeating the steps in Sections 3.6 to merge `height_ranges_df` with `class_df` on `Height`.

For additional practice, you might also try applying `pandas.cut` directly to `class_df`.

```
1 # Make sure the values of height are formatted as integers.
2 class_df['Height'] = class_df['Height'].apply(int)
3 display(class_df)
4
5 # Create a range of possible heights (in inches, per the contents of sashelp.class).
6 height_ranges_df = pandas.DataFrame(
7     {
8         'Height': range(50,75)
9     }
10 )
11
12 # Bin the possible values of height in the reference dataset.
13 height_ranges_df['Height_Range'] = pandas.cut(
14     height_ranges_df['Height'],
15     bins=[50,62,75],
16     labels=['short', 'tall'],
17     include_lowest=True,
```

```
18 )  
19 display(height_ranges_df)
```

	Name	Age	Height
0	Alfred	14.0	69
1	Alice	13.0	56
2	Barbara	13.0	65
3	Carol	14.0	62
4	Henry	14.0	63
5	James	12.0	57
6	Jane	12.0	59

	Height	Height_Range
0	50	short
1	51	short

```

1 # Merge in the binned values of height from the reference dataset.
2 merged_height_df = class_df.merge(
3     height_ranges_df,
4     on='Height',
5     how='left',
6 )
7
8 # Since we're not using the Age column, drop it from merged_df, and display the result.
9 merged_height_df.drop(columns=['Age'], inplace=True)
10 display(merged_height_df)
11
12 # Or just apply pandas.cut directly to class_df (without dropping the Age column).
13 class_df['Height_Range'] = pandas.cut(
14     class_df['Height'],
15     bins=[50,62,75],
16     labels=['short','tall'],
17     include_lowest=True,
18 )
19 display(class_df)

```

	Name	Height	Height_Range
0	Alfred	69	tall
1	Alice	56	short
2	Barbara	65	tall
3	Carol	62	short
4	Henry	63	tall
5	James	57	short
6	Jane	59	short

	Name	Age	Height	Height_Range
0	Alfred	14.0	69	tall
1	Alice	13.0	56	short
2	Barbara	13.0	65	tall
3	Carol	14.0	62	short
4	Henry	14.0	63	tall
5	James	12.0	57	short
6	Jane	12.0	59	short

▼ Notes and Resources

Want some ideas for what to do next? Here are our suggestions:

1. For more about the `faker` package, which can generate many other types of fake data, see <https://faker.readthedocs.io/>
2. For more about the `pandas` package, including the methods used above, see the following:
 - <https://pandas.pydata.org/docs/reference/api/pandas.cut.html>

- <https://pandas.pydata.org/docs/reference/api/pandas.unique.html>
- <https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.append.html>
- <https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.drop.html>
- <https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.merge.html>

3. For more about the `rich` package, see <https://rich.readthedocs.io/>

4. For more about the `saspy` package, including the methods used above, see the following:

- <https://sassoftware.github.io/saspy/api.html#saspy.SASsession.dataframe2sasdata>
- <https://sassoftware.github.io/saspy/api.html#saspy.SASsession.sasdata2dataframe>
- <https://sassoftware.github.io/saspy/api.html#saspy.SASsession.submit>

5. For more about some of the Python features used, such as functions and list comprehensions, we recommend the following chapters of [A Whirlwind Tour of Python](#):

- <https://jakevdp.github.io/WhirlwindTourOfPython/08-defining-functions.html>
- <https://jakevdp.github.io/WhirlwindTourOfPython/11-list-comprehensions.html>

6. We welcome follow-up conversations. You can connect with us on LinkedIn or email us at isaiah.lankham@gmail.com and matthew.t.slaughter@gmail.com

7. If you have a GitHub account (or don't mind creating one), you can also chat with us on Gitter at <https://gitter.im/saspy-bffs/community>