

Python and content analysis: Lessons learned

Jason Kiley
Oklahoma State University

Overview

- Python: data analysis is not (necessarily) programming.
- Workflows: move research forward faster.
- Manuscripts: common areas for improvement.

Python Fluency

Software Development

Good-enough Programming

Data preparation

Basics

Basics

Skills	<ul style="list-style-type: none">• Software: Anaconda, Python interpreter, Jupyter Notebooks• Variable types: strings, ints, floats• Objects and methods: lists, dictionaries• Packages: importing and installing• Documentation: official and community
Time	2-4 hours
Necessity	Largely unavoidable

Data Preparation

Skills	<ul style="list-style-type: none">• Software: pandas• Reading data formats (built-in)• Slicing, views, <code>df.loc[]</code>• Operations on columns and rows• Reshaping• Merging and querying
Time	1-2 days and ongoing
Necessity	Needed and high ROI

Good-enough Programming

Skills	<ul style="list-style-type: none">• Loops• Writing functions• Reading and writing files (the hard way)• Throwing and handling exceptions• Using additional packages• End point: working, reusable script
Time	1 week and ongoing; divisible
Necessity	Helpful and good ROI

Software Development

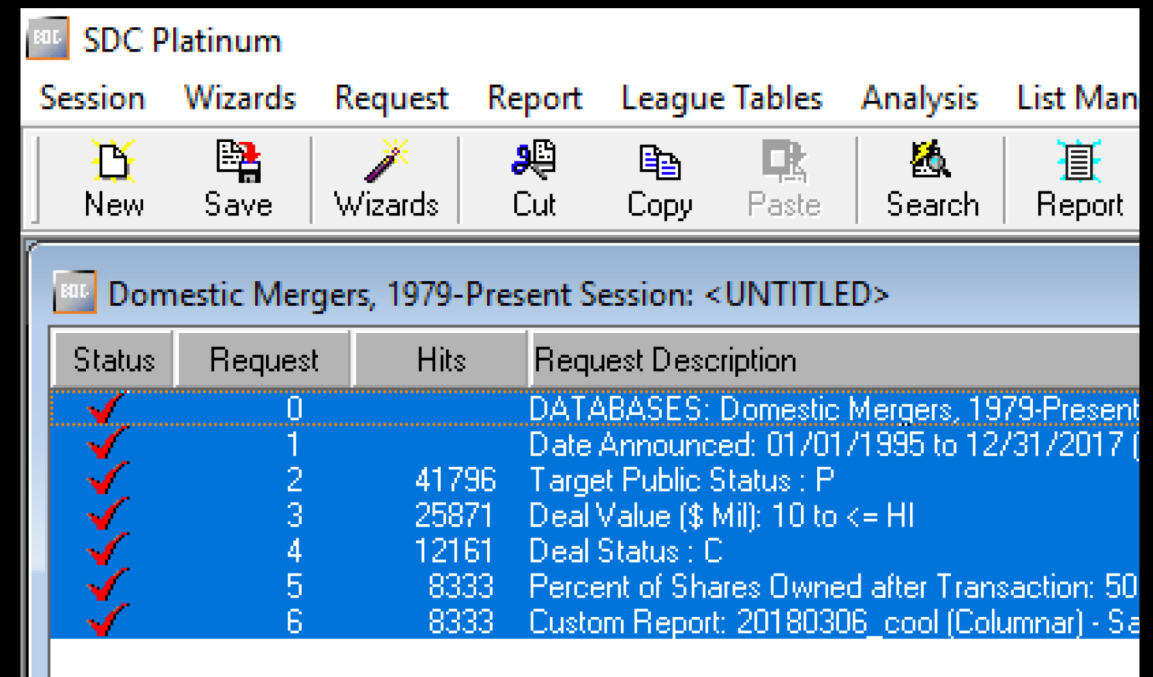
Skills	<ul style="list-style-type: none">• Classes and inheritance• Package development• Version control• Unit testing and continuous integration• Cross-version support• Open source contributions
Time	A lot
Necessity	Not at all; good for the field

Workflows

- Documentation: can I authoritatively show my work?
- Reproducibility: can someone else reproduce my data and results?
- Efficiency: can I easily make changes in the middle of my data pipeline?

Documentation

- Data sources and queries with enough specificity to recreate it.
- Not perfect; our databases are not versioned.
- Helps mitigate the “one more variable” problem.



The screenshot shows the SDC Platinum software interface. The title bar reads 'SDC Platinum'. The menu bar includes 'Session', 'Wizards', 'Request', 'Report', 'League Tables', 'Analysis', and 'List Man'. The toolbar contains icons for 'New', 'Save', 'Wizards', 'Cut', 'Copy', 'Paste', 'Search', and 'Report'. The main window title is 'Domestic Mergers, 1979-Present Session: <UNTITLED>'. Below the title bar is a table with the following data:

Status	Request	Hits	Request Description
✓	0		DATABASES: Domestic Mergers, 1979-Present
✓	1		Date Announced: 01/01/1995 to 12/31/2017
✓	2	41796	Target Public Status : P
✓	3	25871	Deal Value (\$ Mil): 10 to <= HI
✓	4	12161	Deal Status : C
✓	5	8333	Percent of Shares Owned after Transaction: 50
✓	6	8333	Custom Report: 20180306_cool (Columnar) - Se

Fix versions

- Fix the software versions that you use, so updates do not affect your study.
- Conda environments make this fairly easy.
- Update as needed, but check your results.

81 lines (129 sloc) | 3.65 KB

```
1  name: cool
2  channels:
3    - defaults
4    - conda-forge
5  dependencies:
6    - textblob=0.15.1=py_0
7    - appdirs=1.4.3=py36h28b3542_0
8    - appnope=0.1.0=py36hf537a9a_0
9    - arrow-cpp=0.9.0=py36ha51b053_7
10   - asn1crypto=0.24.0=py36_0
11   - attrs=18.1.0=py36_0
12   - automat=0.7.0=py36_0
13   - backcall=0.1.0=py36_0
14   - blas=1.0=mkl
15   - bleach=2.1.3=py36_0
16   - boost-cpp=1.65.1=h1de35cc_4
17   - bzip2=1.0.6=h1de35cc_5
```

Code that runs

- Make all changes to data in code that runs cleanly on the original data.
- Reproducibility is a big favor to future you.
- Rerunning it proves that it is authoritative.

Prep

```
In [39]: _WRDS_COLUMNS = {'ticker': 'id_ticker',  
                        'evtdate': 'deal_date_ann',  
                        'cret': 'ret_cret_m1p1',  
                        'car': 'ret_car_m1p1',  
                        'bhar': 'ret_bhar_m1p1'}  
  
wrdssevent = pd.read_stata('./source/20180309_wrdssevent.dta')  
  
wrdssevent = wrdssevent.loc[:, ['ticker', 'evtdate', 'cret', 'car', 'bhar']]  
wrdssevent.rename(columns=_WRDS_COLUMNS, inplace=True)  
wrdssevent['id_ticker'] = wrdssevent['id_ticker'].str.lower()  
wrdssevent['deal_date_ann'] = wrdssevent['deal_date_ann'].dt.date
```

```
In [40]: wrdssevent.head(15)
```

```
Out[40]:
```

	id_ticker	deal_date_ann	ret_cret_m1p1	ret_car_m1p1	ret_bhar_m1p1
0	sunw	2005-06-02	-0.039370	-0.045034	-0.045667
1	sunw	2000-09-18	0.001596	0.005746	0.006149
2	sunw	1999-08-23	0.008518	-0.077037	-0.079339
3	sunw	2005-06-27	-0.002667	-0.001747	-0.001603
4	orcl	2016-05-02	-0.016117	-0.008230	-0.008284
5	orcl	2006-11-02	-0.037358	-0.033634	-0.033339
6	orcl	2013-02-04	-0.000845	-0.010991	-0.011517
7	orcl	2003-06-06	-0.053093	-0.053275	-0.051736
8	orcl	2005-03-08	0.005271	0.019543	0.019060

Github: jtkiley

**Data Curation Workshop with Tim
Hannigan, Hovig Tchalian, and Laura
Nelson**

Manuscripts

- There is no substitute for a great theoretical question.
- Validate your measures. Can you demonstrate that the measure captures the construct?
- Know your data. When your measure performs badly, do you know why?
- Don't chase advanced methods without demonstrating why they are superior.

Questions