# Communication

## Presenting Results and Conclusions

Data 102 - Fernando Pérez
Slides credit: Lindsey Heagy

# Communication

## Presenting Results and Conclusions

Data 102 - Fernando Pérez
Slides credit: Lindsey Heagy

*presentation matters

# Outline

- Course so far:
  - Techniques for decision making
  - Understanding of their foundations and assumptions
  - (little) impact of those
- You have RESULTS! CONCLUSIONS!!
- Now What?
- What are results, and conclusions in this context?
- Models? Data?
  - Hydro example: a geophysicist draws a model of the ground and draws a line, and hands this model to a hydrologist.
  - The hydrologist then makes decisions, runs simulations, MCMC, generates confidence intervals.
  - Whose 'model', whose 'data', whose 'truth'?
  - Moritz's point: model -> data -> truth -> action. The model enforces reality
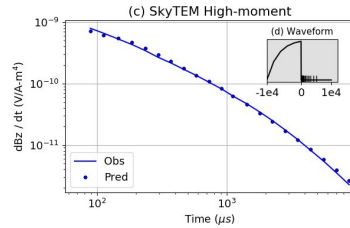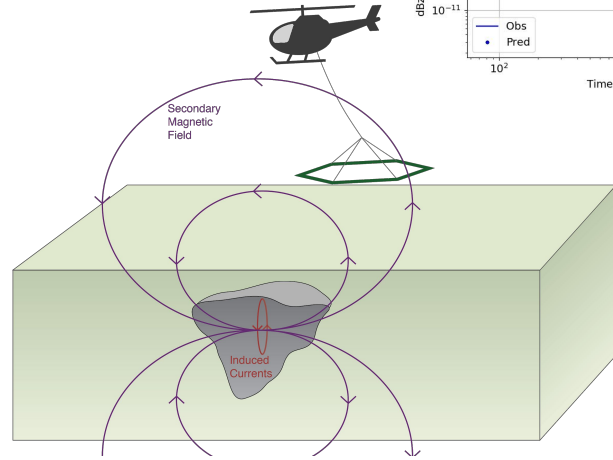- How do we open this?

# Course themes

- Techniques for decision making
- Understanding of their foundations and assumptions
- Impact of these
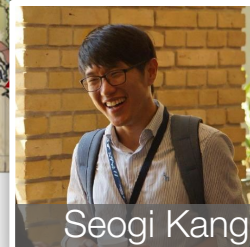
# You have RESULTS! CONCLUSIONS!!

## Now What?
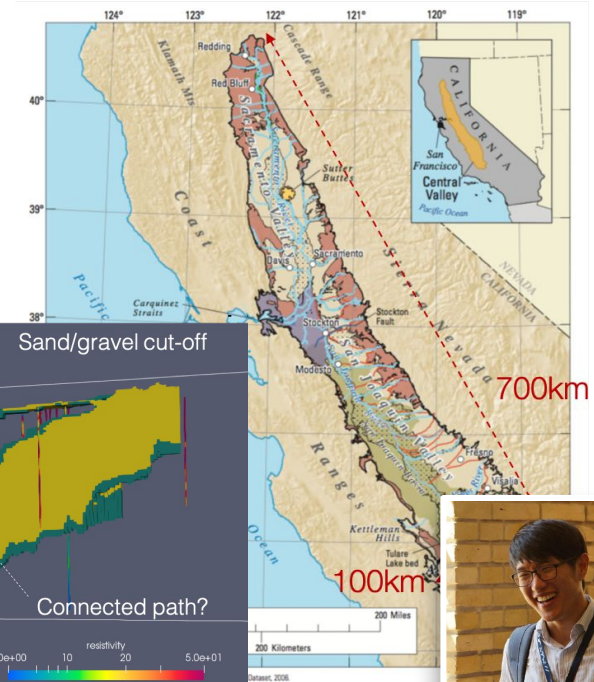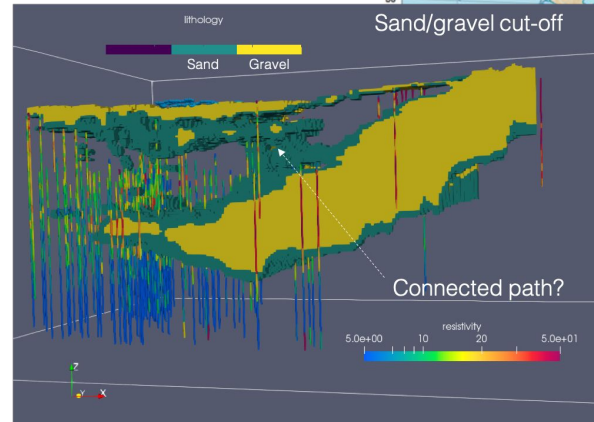
# What are results, and conclusions in this context?

Data?

Models?



(c) SkyTEM High-moment

(d) Waveform

Obs
Pred

Time (μs)

Sand/gravel cut-off

lithology

Sand    Gravel

Connected path?

resistivity

Seogi Kang

700km

100km

# what are "models" made of?

- Algorithmic ideas
- Mathematical structure (choices of features, etc.)
- Data to feed them!

Today, "model" often refers to an "embodied model" that has been "fed data".

# models → data → truth?

communication?

publication, teaching, ...

results

# models → data → truth?

dialog = ?

results

# models → data → truth?

dialog = ?

results

instructions, environment

# the science more than the paper

An article about computational science in a scientific publication is not the scholarship itself, it is merely advertising of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures.

-- Buckheit and Donoho (paraphrasing Claerbout)
  WaveLab and Reproducible Research, 1995

# the science more than the paper

An article about computational science in a scientific publication is not the scholarship itself, it is merely advertising of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures.

-- Buckheit and Donoho (paraphrasing Claerbout)
   WaveLab and Reproducible Research, 1995

^ (and a place to run the code?)

# Core skills

- Version control: Git and GitHub
- Programming: Python
- Process automation: Make
- Data analysis: Numpy, Pandas, Matplotlib, NLTK, Scikit-Learn, …
- Documentation: Sphinx
- Software testing: PyTest
- Continuous Integration: Travis
- Reproducible containers: Binder

# Git and Python workflow everywhere

# Computational hygiene: a daily habit

# Explicit dependency management



```yaml
! environment.yml  ✕

Fernando Perez, 4 months ago | 1 author (Fernando Perez)
1   name: s159-sphinx
2   channels:
3   - conda-forge
4   - defaults
5   dependencies:
6   - ghp-import=0.5.5
7   - ipython=6.1.0
8   - jupyter_client=5.1.0
9   - jupyter_core=4.3.0
10  - nbconvert=5.3.1
11  - nbformat=4.4.0
12  - pandoc=1.19.2.1
13  - python=3.6.3
14  - sphinx=1.6.3
15  - pip:
16    - commonmark==0.5.4
17    - nbsphinx==0.2.17
18    - recommonmark==0.4.0
19
        Fernando Perez, 4 months ago • Add environment.yml
```

**ANACONDA.**

```
(master)alpamayo[stat159]> conda env create -f environment.yml
Solving environment: done

Downloading and Extracting Packages
xz 5.2.3: ############################################################## | 100%
alabaster 0.7.10: ################################################################ | 100%
entrypoints 0.2.3: ############################################################### | 100%
pytz 2018.3: ################################################################### | 100%
nbconvert 5.3.1: ################################################################ | 100%
```

```
#
# To activate this environment, use:
# > source activate s159-sphinx
#
# To deactivate an active environment, use:
# > source deactivate
#

(master)alpamayo[stat159]>
```

# Automation and Testing: SW Carpentry

# Continuous Integration with Travis

shareable, interactive, reproducible
environments from your public git repository

mybinder.org

# Black holes! LIGO, Sept 14, 2015



FIG. 1. The gravitational-wave event GW150914 observed by the LIGO Hanford (H1, left column panels) and Livingston (L1, right column panels) detectors. Times are shown relative to September 14, 2015 at 09:50:45 UTC. For visualization, all time series are filtered with a 35–350 Hz bandpass filter to suppress large fluctuations outside the detectors' most sensitive frequency band, and band-reject

**Make sound files**

Make wav (sound) files from the filtered, downsampled data, +-2s around the event.

```python
# make wav (sound) files from the whitened data, +-2s around the event.
from glob import glob
from IPython.display import display, Audio

from scipy.io import wavfile

# function to keep the data within integer limits, and write to wavfile:
def write_wavfile(filename,fs,data):
    d = np.int16(data/np.max(np.abs(data)) * 32767 * 0.9)
    wavfile.write(filename,int(fs), d)

tevent = 1126259462.422          # Mon Sep 14 09:50:45 GMT 2015
deltat = 2.                      # seconds around the event

# index into the strain time series for this time interval:
indxt = np.where((time >= tevent-deltat) & (time < tevent+deltat))

# write the files:
write_wavfile("GW150914_H1_whitenbp.wav",int(fs), strain_H1_whitenbp[indxt])
write_wavfile("GW150914_L1_whitenbp.wav",int(fs), strain_L1_whitenbp[indxt])
write_wavfile("GW150914_NR_whitenbp.wav",int(fs), NR_H1_whitenbp)

for wav in glob('*whitenbp.wav'):
    display(wav)
    display(Audio(filename=wav))
```

```
'GW150914_H1_whitenbp.wav'
```



http://bit.ly/black-holes-woop

# complete set of instructions

*capture the steps: what is a notebook?*

A document + An interface + An environment

# repo2docker



repo2docker deterministically build a docker image
from a repository with documented dependencies

# complete development environment

*define dependencies following community standards of practice*

**python™**

requirements.txt

requirements.txt ×

```
1  numpy==1.16.*
2  matplotlib==3.*
3  seaborn==0.8.1
4  pandas
5
```

Line 1, Column 1

environment.yml

environment.yml ×

```
1   name: example-environment
2   channels:
3     - conda-forge
4   dependencies:
5     - numpy
6     - psutil
7     - toolz
8     - matplotlib
9     - dill
10    - pandas
11    - partd
12    - bokeh
13    - dask
14
```

Line 1, Column 1                    Spaces: 2

**R**

runtime.txt

runtime.txt ×

```
1  r-2019-04-10
```

install.R

install.R ×

```
1  install.packages("tidyverse")
2  install.packages("rmarkdown")
3  install.packages("httr")
4  install.packages("shinydashboard")
5  install.packages('leaflet')
6
```

Line 1, Column 1                    Spaces: 4

Spaces: 4

# complete development environment

*repo2docker*



dependencies        repo2docker        container file

# Example 1: real-world replication

## Are there statistical links between the direction of European weather systems and ENSO, the solar cycle or stratospheric aerosols?

Benjamin A. Laken and Frode Stordal

Section for Meteorology and Oceanography, Department of Geosciences, University of Oslo, Oslo, Norway

BAL, 0000-0003-2021-6258; FS, 0000-0002-5190-6473

# Are there statistical links between the direction of European weather systems and ENSO, the solar cycle or stratospheric aerosols?

Benjamin A. Laken and Frode Stordal

Section for Meteorology and Oceanography, Department of Geosciences, University of Oslo, Oslo, Norway

BAL, 0000-0003-2021-6258; FS, 0000-0002-5190-6473

benlaken / European_wind

Watch ▾ 1 | ★ Star 4 | ⑂ Fork 1

<> Code | ⓘ Issues 0 | ⇄ Pull requests 1 | ▦ Projects 0 | 𝄞 Wiki | �📊 Insights

Repo relating to a study of European synoptic weather types.

weather-systems    climate-science    enso    aod

| 🕓 26 commits | ⑂ 1 branch | 🏷 0 releases | 👥 1 contributor |

Branch: master ▾ | New pull request | Create new file | Upload files | Find file | Clone or download ▾

benlaken update readme info | Latest commit 24c0b05 on Feb 23, 2016

| 📁 Data | added Sato index to main dataframe | 2 years ago |
| 📁 Figs | changed nomanclature | 2 years ago |
| 📄 .gitignore | pre-modication sync | 2 years ago |
| 📄 HBGWL_analysis.ipynb | update readme info | 2 years ago |
| 📄 HBGWL_functions.py | changed nomanclature | 2 years ago |
| 📄 README.md | update readme info | 2 years ago |
| 📄 sato_pandas.py | added Sato index to main dataframe | 2 years ago |

📖 README.md

# Are there statistical links between the direction of European weather systems and ENSO, the solar cycle or stratospheric aerosols?

Repo relating to a study of European synoptic weather types published in the Royal Society Journal Open Science. Published 17 February 2016, DOI: 10.1098/rsos.150320.

# ROYAL SOCIETY OPEN SCIENCE

Home   Content   Information for   About us   Sign up   Submit

## Are there statistical links between the direction of European weather systems and ENSO, the solar cycle or stratospheric aerosols?

Benjamin A. Laken, Frode Stordal

Published 17 February 2016. DOI: 10.1098/rsos.150320



weather system direction by season

**Figure 3.** Violin plots showing the frequency (days/month) with which weather systems come from cardinal compass directions, grouped by season. Standard error of the mean values were on average 0.19 days/month and did not exceed 0.36 days/month. The violins, like box plots, show the first and third quartiles and median values on horizontal lines, in addition to kernel density estimations (KDEs) reflected around the centre of the categorical sample.

other positive relationships between adjacent compass directions may suggest that these data are biased towards the cardinal compass directions: i.e. the fact that more positive associations between closely related flow directions may indicate a bias towards selecting weather-types corresponding to cardinal directions.

Before any analysis of changes in the direction of weather systems associated with given forcings, seasonal variability is removed from these data. This is achieved by subtracting monthly climatological means from the dataset. All resulting data are described as an anomaly, denoted by $\delta$. We note that following deseasonalization, these frequency data continue to show significant correlations between directions as described in figure 4.

---

benlaken / European_wind

Watch 1   Star 4   Fork 1

<> Code   Issues 0   Pull requests 1   Projects 0   Wiki   Insights

Branch: master ▾   European_wind / HBGWL_analysis.ipynb   Find file   Copy path

benlaken update readme info   24c0b05 on Feb 23, 2016

1 contributor

2.76 MB   Download   History

## Are there statistical links between the direction of European weather systems and ENSO, the solar cycle or stratospheric aerosols?

Code by Benjamin A. Laken, from work published in the journal Royal Society Open Science. Published 17 February 2016. DOI: 10.1098/rsos.150320.

```
For SON
|------>N    4.22µ, 0.24sem
|------>NE   0.74µ, 0.09sem
|------>E    1.91µ, 0.16sem
|------>SE   1.33µ, 0.14sem
|------>S    2.99µ, 0.21sem
|------>SW   1.73µ, 0.16sem
|------>W    8.55µ, 0.32sem
|------>NW   2.02µ, 0.16sem
```

In [8]:  `hbgwl.figure_seasons(data=monthlywind)`

```
/Users/Ben/anaconda/lib/python3.4/site-packages/matplotlib/__init__.py:892: UserWarni
ng: axes.color_cycle is deprecated and replaced with axes.prop_cycle; please use the
latter.
  warnings.warn(self.msg_depr % (key, alt_key))
/Users/Ben/anaconda/lib/python3.4/site-packages/seaborn/categorical.py:1791: UserWarn
ing: The violinplot API has been changed. Attempting to adjust your arguments for the
new API (which might not work). Please update your code. See the version 0.6 release
notes for more info.
  warnings.warn(msg, UserWarning)
/Users/Ben/anaconda/lib/python3.4/site-packages/matplotlib/figure.py:397: UserWarnin
g: matplotlib is currently using a non-GUI backend, so cannot show the figure
  "matplotlib is currently using a non-GUI backend, "
```



Weather system direction by season

This can be shown another way below, however we have used the violin plot in the manuscript as the polar plot may give the false impression that these metrics relate to surface winds with a specific direction as viewed by an observing site (while they actually relate to the origin of regional-scale weather systems estimated over a large area).

# Stat 159/259 - Reproducible and Collaborative Data Science

All materials for this course are available on GitHub.

The class syllabus will be updated over the course of the first couple of weeks of class.

## MNT: update to more reliable method of creating legends

**Open** tacaswell wants to merge 1 commit into `benlaken:master` from `tacaswell:fix_legend`

💬 Conversation 4    Commits 1    Files changed 1

**tacaswell** commented 16 days ago

This fixes a bug identified by **@fperez**

When calling `ax.legend` with one arg the list of strings is zipped with the available artists in the Axes. This is brittle because it assumes the contents and order of this list. In Matplotlib 1.5.1 we added a legend handler for the artists that are used to draw the fill_between regions which caused them to be included in the list of artists which will go in the legend and due to internal details about how Matplotlib store's the children artists of the Axes the fill_between artists are listed before the errorbar artists and thus the legends end up shifted.

The primary change in this commit is to pass the correct label into the plotting calls via the `label=` kwarg and to call `legend` with no args.

**fperez** commented 11 days ago

**@benlaken**, thanks so much for making your research openly available, and in terms of reproducibility you're already doing better than the vast majority of the scientific community!

For context, the reason we found out about this, was b/c I used two of your papers as a homework and class project in my course on Reproducible and Collaborative Data Science at UC Berkeley. I hope you don't mind :)

For the first homework, the students had to practice with replicating your monsoon rainfall notebook. This meant downloading and being able to run it again, via github, working as a team.

Then for the project, they worked with this (European_wind) repo, and there they had to pretty much figure out the things you mention above: wrap the project in a Makefile along with an `environment.yml` (we're using conda envs, but same idea as `pip freeze`), while figuring out the versions you'd used, etc.

It was great to show this very page today during class, while we were discussing the project (their deadline was last night), and for them to see how the author of the paper they were working on was responding so kindly and openly, while identifying the same issues they were working on. I couldn't have timed it better if I'd tried :)

Once we're done grading, happy to send your way the `environment.yml` and `Makefile`, if you'd like to add them to the repo to make it a bit easier in the future...

Many thanks again! Open science rocks :)

**benlaken** commented 11 days ago    Owner

Thanks **@fperez** for the kind words - means a lot coming from you. I am also very happy to hear that my work has been useful on your course. You have some lucky students: I wish I had a similar course when I was studying!

Please do send a PR with the fixes and I will merge to Master - if it is useful for you, I can also leave a branch in its current state?

And indeed, rock-on Open Science! 🎸 💥

# Standard workflow: Makefile and environment.yml

# "Atomic unit" of communicable results

- Data: included in repo or linked if too large.

- Clean, tested code.

- Analysis notebooks and supporting code

  - Break down your analysis into as many notebooks as is reasonable for convenient reading and execution.

- Main narrative notebook: summarizes and discusses results.

- Reproducibility support: Makefile and environment.yml

- Good repository practices: README.md, LICENSE, .gitignore.

  - Use Victoria Stodden's ENABLING REPRODUCIBLE RESEARCH: LICENSING SCIENTIFIC INNOVATION.

A "Standard Playbook"

# Brief Analysis on the Marginal Effects of Studying

`build passing` `launch binder`

As students, we have often wondered what effect an extra hour of studying will have on our grades. When trying to determine whether staying up an extra hour to study for that final exam is truly worth it, we usually are limited by imperfect information and our own superstitions. In this project, we attempt to estimate the "true" marginal effect of studying on students' grades. We try to model the effects of studying first using OLS and then various instruments and 2 stage least squares. This repository is also meant to serve as an example of what a reproducible econometric analysis would look like.

## Required Installations

The only installation needed to run this repo is Anaconda. Click here to learn about how to install Anaconda. Once installed, you should be good to go!

## Using Binder

We've enabled Binder for this project which allows you to view jupyter notebooks in an executable environment. Feel free to click the link at the top of this README to launch the binder.

## Getting Started

Download the repo onto your local machine and open your command prompt. Simply type in the following commands to run the analysis:

```
make clean
make env
source activate study
make run
```

After all your notebooks have run you should see new files in the results, fig, and data directories. Read about our approach and results in main.ipynb. All the figures from our analysis are saved in the fig directory and our regressions are saved in the results directory as dataframes. You can load in these dataframes and work with them as regression instances (i.e. you can call `.summary()` , `.params()` etc. click here for OLS documentation and here for 2SLS documentation)

## Licensing

In an effort to enable reproducible, collaborative reserach our project is subject to the MIT License which allows you to modify and distribute the above code for both private and commercial usage. See LICENSE to learn more.

---

nadavtadelis Merge pull request #27 from berkeley-stat159-f17/nadav_actual_final ... Latest

| | |
|---|---|
| data | added reproducibility aspects, split model fitting into 2 ntbks; NOTE... |
| fig | Fix typos in data_exploration.ipynb |
| results | Fix typos in model_fitting_2.ipynb |
| .gitignore | Add caches to .gitignore |
| .mailmap | adding mailmap to account for config issues |
| .travis.yml | Add pandas install to Travis |
| LICENSE | Added LICENSE |
| Makefile | added reproducibility aspects, split model fitting into 2 ntbks; NOTE... |
| README.md | added reproducibility sentence to README |
| data_exploration.ipynb | Minor add to data_exploration.ipynb |
| environment.yml | added reproducibility aspects, split model fitting into 2 ntbks; NOTE... |
| instructions.md | Add note about grades in team work |
| main.ipynb | correction to instruments justification |
| model_fitting_1.ipynb | Fix typos in model_fitting_1.ipynb |
| model_fitting_2.ipynb | Fix typos in model_fitting_2.ipynb |
| p3functions.py | Add two_way function |
| tests.py | Add two_way function |

## Analysis notebooks

## Code and tests

While these histograms give us some information about the distributions of these individual variables, they don't help with understanding how these variables interact with our dependent variable G3. So lets look at some violin plots to visualize some of these interactions.

For the violin plots we split G3 into 5 bins to more clearly visualize the interactions. We also show the distributions relative to which school the students come from to determine whether there is a difference in the two schools.
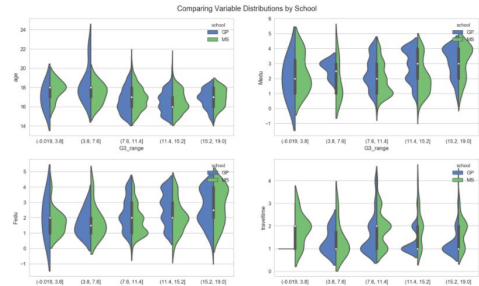
```python
In [6]: # Splitting G3 into ranges to get a cleaner visual
student_perf['G3_range'] = pd.cut(student_perf.G3, 5, retbins = True)[0]

# Creating the plots
plt.figure(figsize=(16, 36))
sns.set(style="whitegrid", palette="muted", color_codes=True)

plt.subplots_adjust(top=0.97)
plt.suptitle('Comparing Variable Distributions by School')

sns.despine()
for column_index, column in enumerate(['age', 'Medu', 'Fedu', 'traveltime', 'studytime',
'freetime', 'failures',
                                       'absences', 'famrel', 'goout', 'Dalc', 'Walc', 'he
alth', 'G1', 'G2']):
    if column == 'G3_range':
        continue
    plt.subplot(8, 2, column_index + 1)
    sns.violinplot(x='G3_range', y=column, hue = 'school', split = True, data=student_per
f)

plt.savefig('fig/distrbyschool.png');
```

---

Branch: master ▾ **project-3-p2-ka-jo-ta** / p3functions.py

s-johnson Add two_way function

1 contributor

41 lines (34 sloc) | 1.38 KB          Raw

```python
1  import pandas as pd
2  import numpy as np
3
4  def make_indicators(df, names):
5      """Make indicator columns in dataframe df of whether existing columns are
6      equal to given values.
7
8      Args:
9          df (pandas.DataFrame): Dataframe to be modified.
10         names (dict)       : Dictionary containing:
11                              - Keys: Desired indicator column names
12                              - Values: Two item tuple containing:
13                                  - Original dataframe column
14                                  - Value to compare to column
15
16         Returns:
17             void: Dataframe df is modified in place.
18     """
19     for k,v in names.items():
20         df[k] = 1*(df[v[0]] == v[1])
```

Branch: master ▾ **project-3-p2-ka-jo-ta** / tests.py

s-johnson Add two_way function

1 contributor

31 lines (24 sloc) | 906 Bytes          Raw

```python
1  import pandas as pd
2  import numpy as np
3  import numpy.testing as npt
4
5  from p3functions import *
6
7
8  def test_make_indicators():
9      d = {'col1': [1, 2], 'col2': [3, 4]}
10     df = pd.DataFrame(data=d)
11     names = {'ind1': ('col1', 2), 'ind2': ('col2', 3)}
12     make_indicators(df,names)
13     exp_d = {'col1': [1, 2], 'col2': [3, 4], 'ind1': [0, 1], 'ind2': [1,0]}
14     exp = pd.DataFrame(data=exp_d)
15     obs = df
16     assert obs.equals(exp)
```
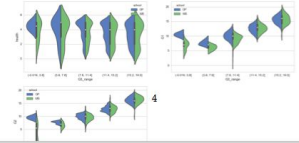
---

## Project1-Main-Narrative

January 4, 2018

### 1 The effects of studying on high school students

**Authors:** Nadav Tadelis, Sarah Johnson, Chitwan Kaudan

#### 1.1 Abstract

As students, a large part of our daily life is take
formed guesses about how much an extra hour
is not ideal; when making allocation decisions,
sue. Specifically, if imperfect information causes
grades, then we make poor decisions about how
(effectively resulting in a loss of utility). If we we
ing on grades, then we could calibrate our inner
naive OLS, then addressing endogenity by using
marginal increase in study time per week can inc

#### 1.2 Exploratory Data Analysis

The data being used are from the public archive
collected by Paulo Cortez of the University of M
Below is a list of all included variables:

in our regression.

Now that we have established our data are clean we can move on to trying to answer our question regarding the marginal effect of studying on grades.

---

1 We need to make the additional assumption that in secondary school (where parents are notified when students are absent), absences are only caused by illnesses and emergencies (which are independent of study time). Without this assumption it would be plausible that students are skipping school because they value leisure over studying, implying a negative correlation between study time and absences.

#### 1.3 Initial Naive OLS fit

The first step is to build a model and make some assumptions to define the relationship between grades and studying. Let an individual's grade be $G_i$ and weekly hours of studying be $S_i$ and their "ability" be $A_i$. Then we can write:

$$G_i = \beta_0 + \beta_1 S_i + \beta_2 A_i + U_i$$

##### 1.5.1 References

Card, D., & Krueger, A. (1992). Does School Quality Matter? Returns to Education and the Characteristics of Public Schools in the United States. *Journal of Political Economy*, 100(1), 1-40.

Cortez, P. and Silva, A. Using Data Mining to Predict Secondary School Student Performance. In A. Brito and J. Teixeira Eds., Proceedings of 5th FUture BUsiness TEChnology Conference (FUBUTEC 2008) pp. 5-12, Porto, Portugal, April, 2008, EUROSIS, ISBN 978-9077381-39-7.

Greene, W. H. (2000). Econometric analysis. Upper Saddle River, N.J.: Prentice Hall.

MacKinnon, J.G. and H. White. (1985), Some heteroskedasticity consistent covariance matrix estimators with improved finite sample properties. *Journal of Econometrics*, 29, 53-57.

##### 1.5.2 Author Contributions

- Nadav Tadelis: Had idea from a project he did in Econ 142, worked to pick right instruments to improve the 2SLS model, wrote analysis in main.ipynb, created visualizations, and wrote/coded model fitting notebooks.
- Sarah Johnson: Helped brainstorm instruments to improve 2SLS, wrote analysis in main.ipynb, created functions and tests, and integrated testing through Travis.
- Chitwan Kaudan: Helped brainstorm instruments to improve 2SLS, wrote analysis in main.ipynb, worked on reproducibility aspects, created environment and makefile, and structured notebooks.

# In closing

The ideas of data analysis ought to survive a look at how data is analyzed.

-- "*The future of data analysis*", 1961
John Tukey (1915-2000)