# Data Versioner

Castle Leonard
March 8, 2022
DATA 515 - Project Presentation
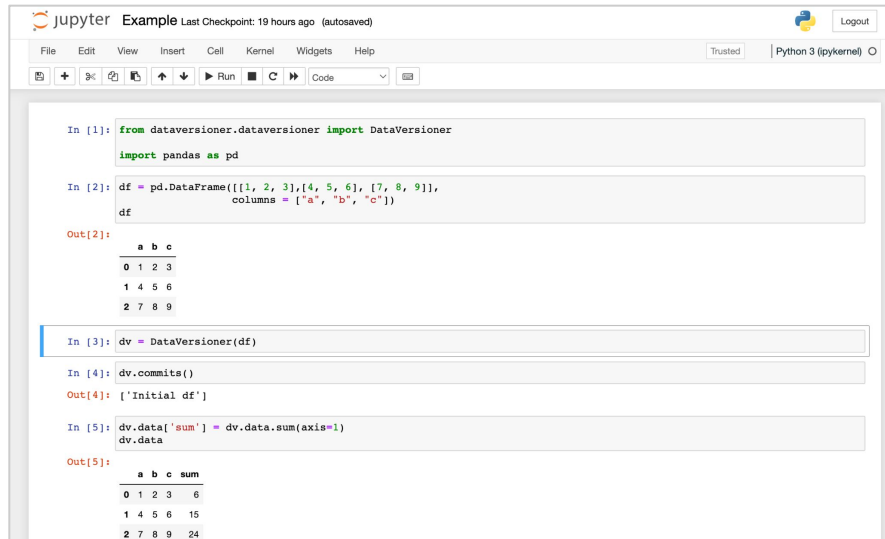
# Background & Use Cases

**Problem:** How can we avoid interrupting and stemming our creativity in experimental and exploratory phases of data science without losing track of variations and findings? Especially in the wild west that is Jupyter notebooks?

**Solution:** Data versioner is a python package that offers git-like tracking of your pandas dataframe so that your experimental variations and EDA findings are preserved, even as your notebook evolves.

# Design & Project Structure



| | |
|---|---|
| 📁 | tests |
| 📄 | __init__.py |
| 📄 | committree.py |
| 📄 | datacommit.py |
| 📄 | dataversioner.py |

# Software Engineering Lessons & Future Work

## Lessons

1. **Trade-offs everywhere** - In many choices I encountered trade offs between user experience, code complexity, development time, and constraints to the future of the package.

2. **Start with MVP** - It is so easy to get swept up with all the aspirations and anticipations of needs for sophisticated functionalities, but premature optimization

## Future

1. **Extensible multi-inheritance** - My dream implementation of this package is a wrapper class of the initializing dataframe type. This would be much less obtrusive during data manipulation (data_versioner.data is no longer necessary) and would enable the package to be extended relatively easily to provide seamless user experience for all supported dataframe types.

2. **Automatic dataframe modification logging and diffing**