# Improving your notebook workflow

By Corrie Bartelheimer

# Messy Notebooks

# Messy?

- Easy to end up with a messy notebook pile

- Rapid prototyping leads to less documentation

- Reproducibility of results

# Collaboration?

- Will future-you still understand your notebooks?

- Or your colleague?

- Will they be able to run the notebooks?



Jupyter **PyTorch Prototyp** Last Checkpoint: 06/21/2018 (autosaved)

File  Edit  View  Insert  Cell  Kernel  Widgets  Help

```python
from random import random
from numpy import array
from numpy import sum
from os import listdir
from os.path import join, isdir
import random
import sys


from IPython.display import clear_output, display
from sklearn.metrics import accuracy_score, f1_score
from sklearn.model_selection import cross_validate, cross_val_predict
from scipy.sparse import vstack, hstack
from sklearn.utils import shuffle
from sklearn.calibration import CalibratedClassifierCV
from sklearn.svm import LinearSVC
from keras.models import Sequential
from keras.layers import Bidirectional, Embedding, Dropout, Masking, TimeDistributed, Dense, LSTM, Input
from keras import optimizers
from sklearn.feature_extraction.text import TfidfVectorizer
import numpy as np
import pandas as pd
```

**Where is the PyTorch import?**

```python
def read_data(base_dir, class_dir):
    print('read {}'.format(class_dir))
    dir = join(base_dir, class_dir)
    doc_ids = [f for f in listdir(dir) if isdir(join(dir, f))]
    folder_names = [join(dir, l) for l in doc_ids]
    page_content = []
    file_names = []
    for folder_name in folder_names:
        page_file_names = [f for f in listdir(folder_name) if f.endswith('.txt')]
        for page_file_name in page_file_names:
            file = open(join(folder_name, page_file_name), errors='ignore').read()
            if len(_file) > 0:
                page_content.append(_file)
                file_names.append(folder_name)
    return {'label': [class_dir] * len(page_content), 'file_name': file_names,
            'content': page_content}
```

```python
base_dir = '../omc_out'
#base_dir = './pku_single_page_ocr'
print('Read data from FS')
labels = [f for f in listdir(base_dir) if isdir(join(ba
#labels = [f for f in listdir(base_dir) if isdir(join(b
labels.sort()
print('#labels: {}'.format(len(labels)))
```

**Which data was used here?**

```python
base_dir_pku = '.'
#base_dir = './pku_single_page_ocr'
print('Read data from FS')
labels = [f for f in listdir(base_dir_pku) if isdir(join(base_dir_pku,f)) and not f == '.ipynb_checkpoints' and not
#labels = [f for f in listdir(base_dir) if isdir(join(base_dir,f)) and not f == 'LEGITIMATION']
#labels = labels_pku
labels.sort()
print('#labels: {}'.format(len(labels)))
```

```
Read data from FS
#labels: 9
```

4

# Version Control?

- Diff too large

- Or too cryptic

# Version Control!

- Convert notebooks to markdown

- Using nbconvert*

- Can be automated via Jupyter SaveHook

```
jupyter nbconvert --to md notebook.ipynb
```

*Other options are for example: jupytext, reviewnb

# Version Control!

- Diff of both code and results

- Rich Markdown diff is also possible

# Collaboration!

- Use a code (and analysis) review process

- In our team, we use the GitHub workflow

- A review encourages clean-ups and documentation

# Reproducibility!

- Document where your data comes from

- Or better: access data programmatically

- Run the whole notebook before commiting

# Some more tidying

- Introduce a naming convention for notebooks, e.g. `1.0-cba-initial-data-exploration`

- Include a TL;DR summary of the question you're trying to solve and your conclusion

- Use a default folder structure

```
├── LICENSE
├── Makefile           <- Makefile with commands like `make data` or `make train`
├── README.md          <- The top-level README for developers using this project.
├── data
│   ├── external       <- Data from third party sources.
│   ├── interim        <- Intermediate data that has been transformed.
│   ├── processed      <- The final, canonical data sets for modeling.
│   └── raw            <- The original, immutable data dump.
│
├── docs               <- A default Sphinx project; see sphinx-doc.org for details
│
├── models             <- Trained and serialized models, model predictions, or model summaries
│
├── notebooks          <- Jupyter notebooks. Naming convention is a number (for ordering),
│                         the creator's initials, and a short `-` delimited description, e.g.
│                         `1.0-jqp-initial-data-exploration`.
│
├── references         <- Data dictionaries, manuals, and all other explanatory materials.
│
├── reports            <- Generated analysis as HTML, PDF, LaTeX, etc.
│   └── figures        <- Generated graphics and figures to be used in reporting
│
├── requirements.txt   <- The requirements file for reproducing the analysis environment, e.g.
│                         generated with `pip freeze > requirements.txt`
│
├── setup.py           <- Make this project pip installable with `pip install -e`
├── src                <- Source code for use in this project.
│   ├── __init__.py    <- Makes src a Python module
│   │
│   ├── data           <- Scripts to download or generate data
│   │   └── make_dataset.py
│   │
│   ├── features       <- Scripts to turn raw data into features for modeling
│   │   └── build_features.py
│   │
│   ├── models         <- Scripts to train models and then use trained models to make
│   │   │                 predictions
│   │   ├── predict_model.py
│   │   └── train_model.py
│   │
│   └── visualization  <- Scripts to create exploratory and results oriented visualizations
│       └── visualize.py
│
└── tox.ini            <- tox file with settings for running tox; see tox.testrun.org
```
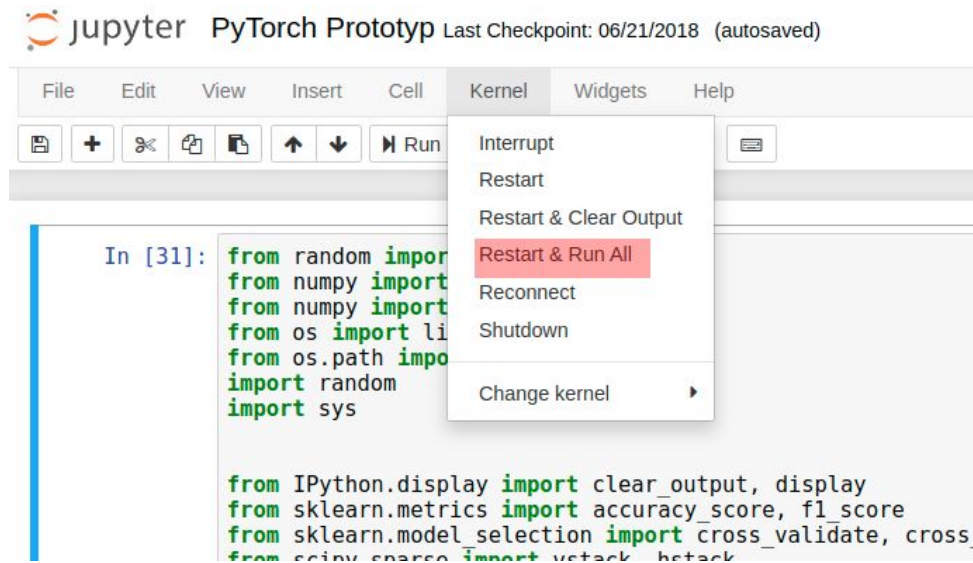
# Summary

- Version Control with converted notebooks

- Collaborate and use Code Reviews

- Make notebooks reproducible

- Clean up your notebooks

# Thanks!

**Any questions?**

@corrieaar    corriebar