# A Practical Tutorial for Jupyter Notebooks

Extracted from Logan Ward's notes

Postdoctoral Scholar

University of Chicago

17 October 2018

globus labs
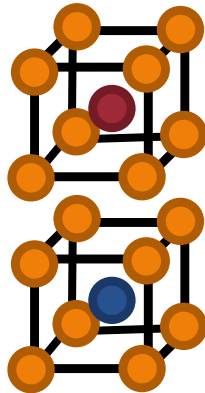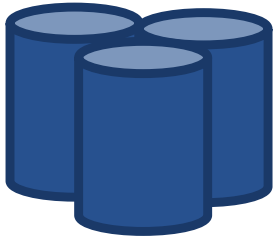
# Materials Informatics Workflow

Collect ⟩ Process ⟩ Represent ⟩ Learn

$\vec{X}$  $\vec{y}$

$\Delta H_f = -1.0$

$\Delta H_f = -0.5$

| $Z_A$ | $Z_B$ | $\Delta H_f$ |
|-------|-------|--------------|
| 3 | 4 | -1.0 |
| 3 | 5 | -0.5 |

$\Delta H_f = f(Z_A, Z_B)$

# Data Representation

Collect ⟩ Process ⟩ **Represent** ⟩ Learn

***Encode domain knowledge into inputs***

$$Property = \boldsymbol{f}(Attributes)$$

LiF   NaPb   $Na_2O$

$x_{Na}$

NaPb

$Na_2O$ LiF

$|\Delta X|$

Representation of material

Ex: $Attributes = \boldsymbol{g}(x_H, x_{He}, \dots)$

**What does a representation need?**

*Completeness:* Differentiate materials

*Efficiency*: Quick to compute

*Accuracy:* Capture important effects

**End product:** Machine-learning compatible inputs

# Fit a model to the data

Collect ＞ Process ＞ Represent ＞ Learn

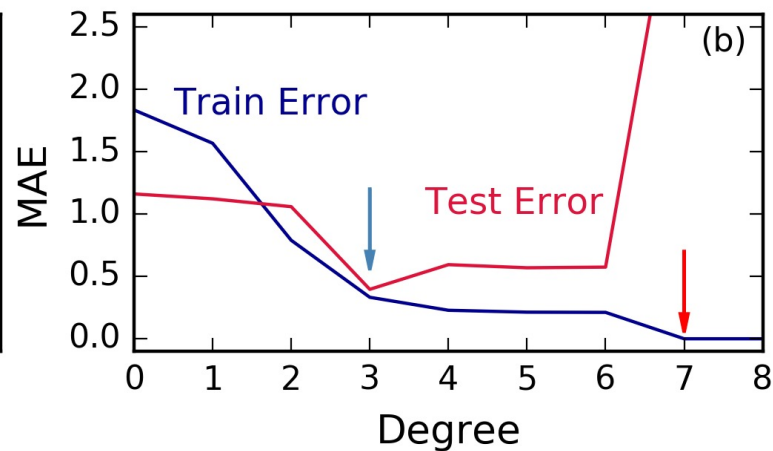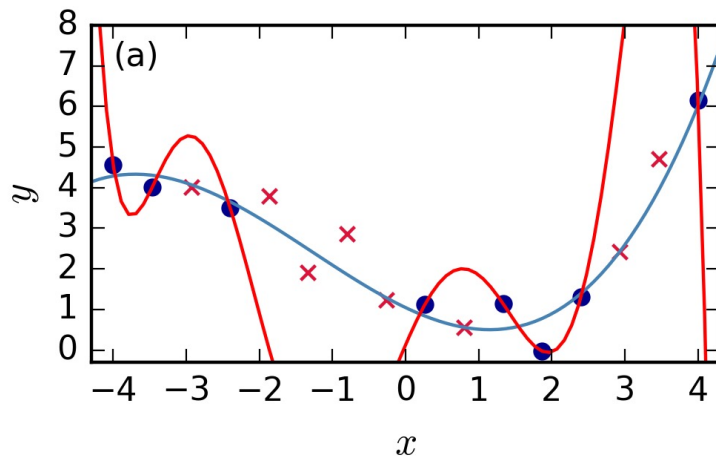## There are many algorithms, how do you decide?

1. Identify algorithms with desired properties.
   Ex: differentiability

2. Use cross-validation to optimize parameters

3. Pick the one with the lowest "loss"

- **Numpy**
- **Scipy**
- **Scikit-learn**
- **Theano**
- **TensorFlow**
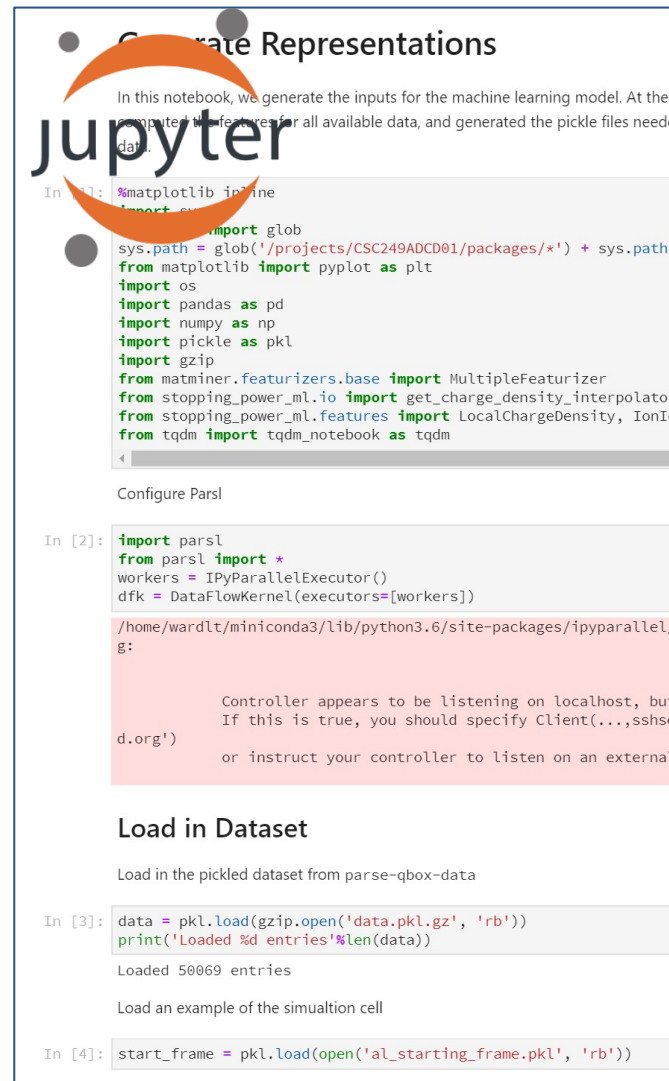- **Keras**
- **PyTorch**
- **Pandas**
- **Matplotlib**

# How to Organize? [Jupyter]

# What is Jupyter? Why care?

- Jupyter is environment designed for reproducible computational science

- Stores code with outputs and documentation in a single notebook
  - Modeled after Mathematica
  - https://www.theatlantic.com/science/archive/2018/04/the-scientific-paper-is-obsolete/556676/

- Why do *I* care? Jupyter lets me…
  - organize my code
  - easily run on a remote system
  - keep track of results and rationale
  - communicate my research better

# How I write a notebook

*One notebook per "experiment" or "idea"*

**Outline:**

1. Title and short abstract
   1. "What am I doing here and why"
2. Load in libraries
   1. Put them up front, so it crashes early
3. Load in data from disk
   1. Ex: training data, results from other notebook
4. Each step in their own block
   1. Introduction
   2. Code <u>and explanation</u>
   3. Figure/visualization
   4. Explanation of finding

# Example Step

## Simple Test: Channel

Introduction

Here, we have a particle traveling forward or backwards along the channel of FCC Al. The path forward and backwards are indentical, so we should get the same stopping power

```python
In [5]: @App('python', dfk)
        def compute_stopping_power(starting_point, direction, traj_computer=traj_computer):
            return traj_computer.compute_stopping_power(starting_point, direction, 1)
```

Set up calculations for the stopping power in the channel

```python
In [6]: forward = compute_stopping_power([0,0.75,0.75], [1,0,0])
```

```python
In [7]: backward = compute_stopping_power([0,0.75,0.75], [-1,0,0])
```

Wait for them to finish        Explanation

```python
In [8]: %%time
        forward = forward.result(); backward = backward.result()

CPU times: user 12 ms, sys: 12 ms, total: 24 ms
Wall time: 1.8 s
```

```python
In [9]: print('Forward stopping power: ', forward[0])
        print('Backward stopping power: ', backward[0])
        print('Difference: ', backward[0]-forward[0])
```

Visualization and conclusion

```
Forward stopping power:  0.2343000208236084
Backward stopping power:  0.23431760358559353
Difference:  1.758276198512987e-05
```

*Finding*: Consistent with my initial expectations, they are indeed the same (within numerical tolerances).

Someone should be able to understand this without knowing Python!

# The Full Narrative

## Assessing Forward/Backward Asymmetry in Stopping Powers

Our stopping power model predicts different stopping powers for particles traveling forward than those traveling backwards on some trajectories. The purpose of this notebook is to showcase why this is valid.

```python
In [1]:  %matplotlib inline
         from matplotlib import pyplot as plt
         from stopping_power_ml.features import LocalChargeDensity
         import pandas as pd
         import numpy as np
         import pickle as pkl
         import os
```

Configure parsl

```python
In [2]:  import parsl
         from parsl import *
         #from parsl_config import config
         workers = IPyParallelExecutor()
         dfk = DataFlowKernel(executors=[workers])
         print("Parsl version : ", parsl.__version__)
```

```
Parsl version :  0.5.0
/home/wardlt/miniconda3/lib/python3.6/site-packages/ipyparallel/client/client.py:458: RuntimeWarning:


        Controller appears to be listening on localhost, but not on this machine.
        If this is true, you should specify Client(...,sshserver='you@js-168-224.jetstream-cloud.org')

        or instruct your controller to listen on an external IP.
```

## Load in the Tools

We'll need the trajectory computer, and the charge density so that we can make illustrative plots.

```python
In [3]:  traj_computer = pkl.load(open('traj_computer.pkl', 'rb'))
```

```python
In [4]:  charge_density = pkl.load(open(os.path.join('..', 'density_interp.pkl'), 'rb'))
```

[Link to GitHub Page](#)

# Pitfalls When Making Notebooks

1. Not including documentation
   **Advice:** Write what you're going to do first

2. Notebook not capturing entire process
   **Problem:** Jupyter lets you execute cells out of order
   **Advice:** Periodically "Restart Kernel and Run All Cells"

3. Duplicate code between notebooks
   **Advice:** Make a separate module for common code

4. Library conflicts
   **Advice:** Run mature projects in container, separate machines
   **Advice:** Make a "requirements.txt" file

5. The "one cell notebooks"
   **Advice:** <10 lines of code per cell

See "I Don't Like Jupyter Notebooks": https://t.co/30peBFwTbv

# Organizing Multiple Notebooks

↑ > single-velocity

| Name ▲ | Last Modified |
|---|---|
| 📁 convergence-detection | a month ago |
| 📁 data | a month ago |
| 📁 feature-analysis | a month ago |
| 📁 figures | a month ago |
| 📁 manifold | 17 hours ago |
| 📁 neural-network | 15 hours ago |
| 📁 runinfo | 15 hours ago |
| 📄 0_collect-subset.ipynb | 16 hours ago |
| 📄 1_build-machine-learni... | 16 hours ago |
| 📄 2_evaluate-direction-d... | 16 hours ago |
| 📄 3_forward-backward-as... | in a few seconds |
| 📄 best_model.pkl | 16 hours ago |
| 📄 best_weight.pkl | a month ago |
| 📄 direction_stopping.pkl | 16 hours ago |
| 📄 run-notebooks.sh | a month ago |
| 📄 stopping_power_result... | 16 hours ago |
| 📄 traj_computer.pkl | 16 hours ago |

Re-usable datasets in separate folder

Figures in separate folder
(with meaningful names!)
(making pub-ready figures!)

Dependent tasks in separate folder

Notebooks numbered by execution order

Results shared between steps using pickled files

A bash script to re-run all notebooks

# Running Notebooks from Terminal

```bash
#! /bin/bash

# This script just runs through all of the current notebooks in the proper order

runJupyter() {
    jupyter nbconvert --execute --inplace --ExecutePreprocessor.timeout=-1 $1
}

# If we get too many files, we should probably create a file containing
#  names of notebooks and the order in which they should be executed
for n in `seq 0 1`; do
    runJupyter ${n}_*nb
done
```

Rerun your project often, put in "assert" statements for key results!

# Keeping Track of Project

Commits on May 21, 2018

Switched to using spherical coordinates for traj optimizer
WardLT committed 3 days ago

Commits on May 20, 2018

Use global optimizer to find best trajectory ...
WardLT committed 4 days ago
d027633 <>

**Use GitHub to track changes**

| Overview | Yours | Active | Stale | All branches | Q Search branches... |

Default branch

master Updated 3 days ago by WardLT

Your branches

**Make branches for dead-ends**

Branch: master ▼    stopping-power-ml / 1_compute-representation.ipynb    Find file | Copy path

WardLT Switched back to using time offsets for the change density    6d7d3f2 7 days ago

1 contributor

459 lines (458 sloc) | 86.5 KB

**GitHub renders your notebooks: Great for Sharing with Collaborators!**

## Generate Representations

In this notebook, we generate the inputs for the machine learning model. At the end of this notebook, we will have computed the features for all available data, and generated the pickle files needed to run these models for any available data.

# Take-Away Points

1. Write your notebooks like papers
   *Explain what you are doing, why, and what you found*
   *Learn Markdown to include links, equations, pictures*

2. Break up complex projects into multiple steps
   *Each notebook should tell a single story*

3. Notebooks should be easy to re-run
   *Use "Restart and Run," bash scripts*

4. Track your changes with Git
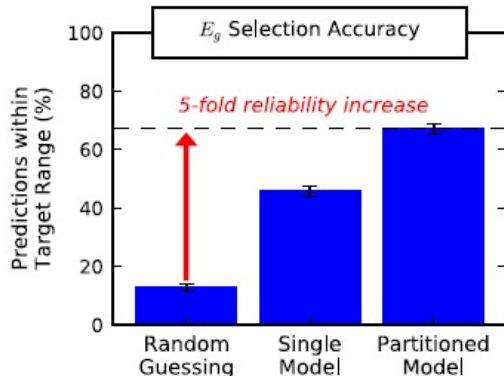   *Explain what changes you made, and why*

# Publication: Not Just Papers

# Is my 2016 paper reproducible?

# I released the code, but not *all of it*
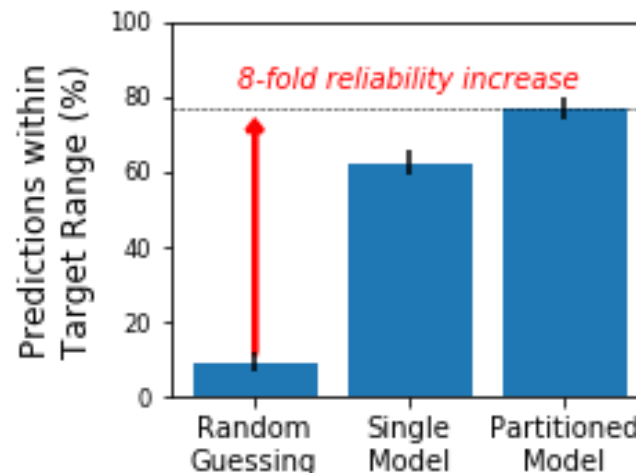
# … and there are differences

**Reported Results in 2016(?)**

**Replication in 2017**



**Without Exact Scripts From <2016, Cannot Verify Results**

Table ... predicted using machine learning to be candidates for solar cell applications

| Composition | E |
| --- | --- |
| ScHg$_4$Cl$_7$ | |
| V$_2$Hg$_3$Cl$_7$ | 1.1 |
| Mn$_6$CCl$_8$ | 1.28 |
| Hf$_4$S$_{11}$Cl$_2$ | 1.11 |
| VCu$_5$Cl$_9$ | |

Abbreviations: DF... materials database...

Out[11]:

| | Entry | bandgap_predicted |
| --- | --- | --- |
| 1037 | CoB2F9 | 1.380256 |
| 3414 | YbAs7Cl6 | 1.156973 |
| 3884 | Tl3OsO3.5 | 1.078918 |
| 1920 | Cs8CoSe5 | 1.074358 |
| | | 1.119002 |

**Publication Should Go Beyond the Paper!**

# How Do We Avoid this Problem?

## What do I need to publish?

- *Datasets*: In a well-described format

- *Scripts*: Not just the core methods

- *Outputs:* Exact version from the paper

- *Models:* In a user-friendly way

Save Anaconda/Python versions and Environments/ library version!

CITRINE
INFORMATICS

GitHub

WholeTale

MATERIALS DATA FACILITY

DLHub