# Data Analysis and Machine Learning Education in the Physics Department at the University of Illinois

## Mark Neubauer

### *University of Illinois at Urbana-Champaign*

# Origins of my DS/ML course in Physics

The idea for developing a course in data analysis and ML applications for our undergraduates grew out of townhall discussions ~6 years ago in our department

- At the time, there were no such courses in our Physics department (now there are several)
- Data science was on the rise, especially driven by the revolution around machine learning
- Students recognize the high value of training at the intersection of data science, AI and physics

I pitched this new course in 2017 at our PAB meeting & taught in Fall 2018, Spring 2019, Fall 2019, Spring 2022
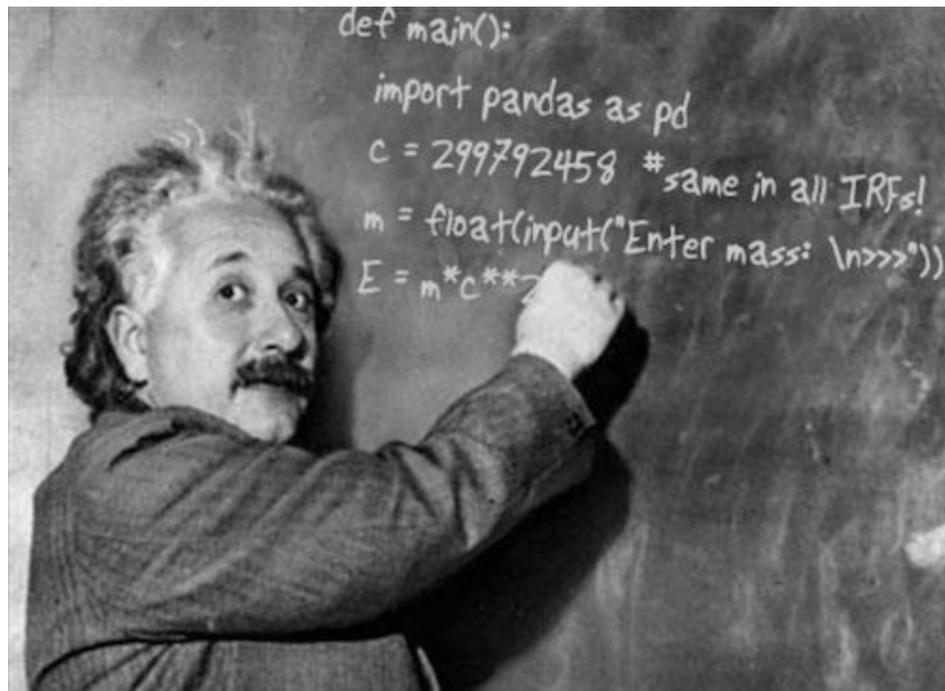
# PHY 398MLA

## *Data Analysis & ML Applications for Physicists*

Topical focus:

- Techniques for analysis and interpretation of scientific data
- Machine learning principles and applications to physics
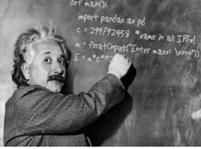


Technical approach:

- Open everything
- Minimize coding focus: Data Science to advance Physics is the goal

# PHY 398MLA

## *Data Analysis & ML Applications for Physicists*

**Motivation**

❖ We live in a data-centric world, with people and machines learning from vast amounts of data.

❖ Early-career physicists need a solid understanding in the basics of data analysis, data-driven inference and machine learning, and a working knowledge of modern tools and techniques from data science
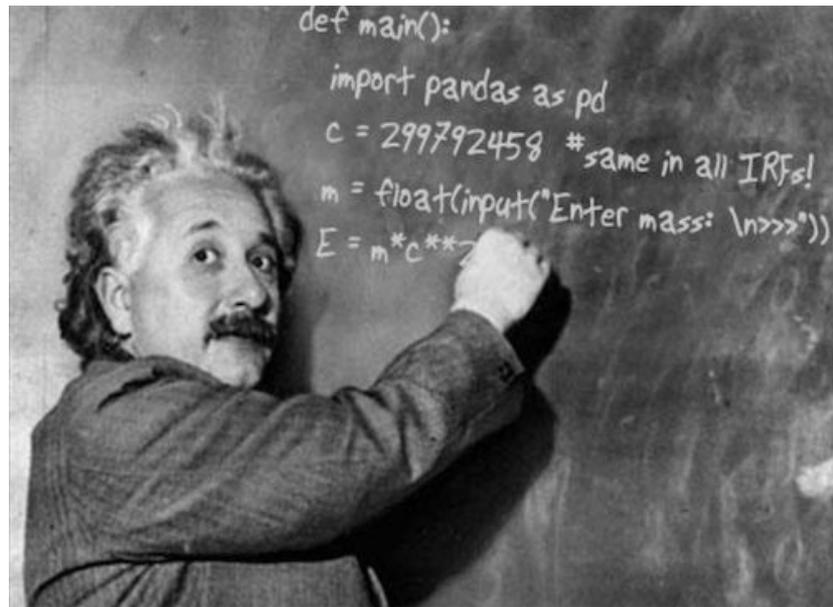
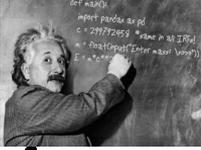# PHY 398MLA
# *Data Analysis & ML Applications for Physicists*

## *Why do this in Physics?*



- Data science is an exploding area both academia and industry
  - E.g. we find that many problems in physics map to vision problems amenable to automation via ML methods
- Many of the future (and current) jobs will be in this area
- Physics students need a solid foundation in data analysis & interpretation to do research. They need this education and training to be productive in science and thrive as the next generation of leaders

# PHY 398MLA

## *Data Analysis & ML Applications for Physicists*

---

**Structure**

- ❖ Two credit-hour course
- ❖ One 2 hour lecture each week
  - o Attendance and class participation are mandatory
- ❖ 8 Homework assignments
  - o New problems posted ~weekly, due within 1 week
- ❖ A final project based on analysis of open science data

---

# PHY 398MLA

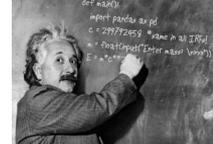## *Data Analysis & ML Applications for Physicists*

**Pedagogy in Motion**

❖ At the center of the course is *scientific data*
  o Students learn how to analyze it, simulate it, gain insight from it, and get machines to learn from it
❖ Students acquire a working knowledge of DS tools
  o All lecture materials, homework and projects are python-based within in Jupyter notebooks (.ipynb)
  o Materials are managed within a Github organization
    o Students submit homework (.ipynb) to their private Github repo. Homework is managed by nbgrader
  o In-class time is a mixture of collective listening and active learning with *in-situ* coding exercises

# PHY 398MLA

## *Data Analysis & ML Applications for Physicists*



**PrairieLearn** is **an online problem-driven learning system for creating homeworks and tests**
- Developed by **the University of Illinois in 2014**

# PHY 398MLA

# *Data Analysis & ML Applications for Physicists*

**PrairieLearn Workspaces** allow students to work in persistent remote containers via in-browser frontends such as VS Code and JupyterLab.

Workspace questions are integrated with the standard PrairieLearn autograding pipeline.

The remote containers are configured by instructors to provide *custom, uniform environments per question*.

# PHY 398MLA

## *Data Analysis & ML Applications for Physicists*



10

# PHY 398MLA

# *Data Analysis & ML Applications for Physicists*

**Topics**

1) Handling and Visualizing Data
2) Finding structure in data
3) Measuring and reducing dimensionality
4) Adapting linear methods to nonlinear problems
5) Estimating probability density
6) Probability theory
7) Statistical methods
8) Bayesian statistics
9) Markov-chain Monte Carlo in practice
10) Stochastic processes and Markov-chain theory

11) Variational inference
12) Optimization
13) Computational graphs
14) Probabilistic programming
15) Bayesian model selection
16) Learning in a probabilistic context
17) Supervised learning in Scikit–Learn
18) Cross validation
19) Neural networks
20) Deep learning

# PHY 398MLA

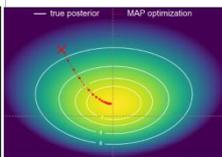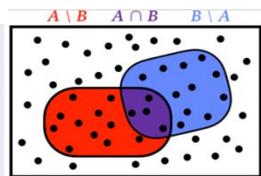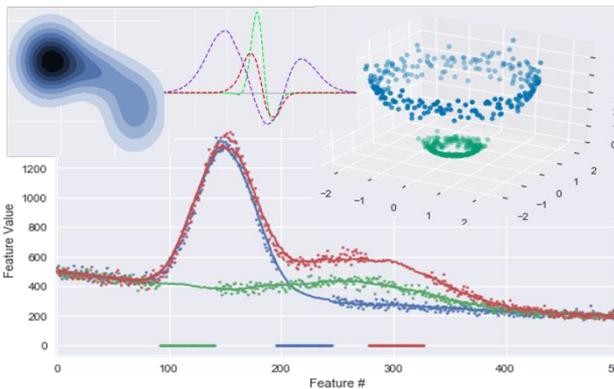## *Data Analysis & ML Applications for Physicists*

*Some challenges*:

- Coding is new for many students. Many are unfamiliar with Python, Git, Jupyter notebooks, etc … → provide ample resources/examples
- Challenge to keep notebooks working over years (e.g. TF API changes, …) and the software packages/tools current
- Modules for physics applications (thanks to DSECOP Fellows!)
- Access to GPU resources for training deep neural network models

*Work in progress:*

- Developing an advanced version for MEng in Instrumental Physics @ Illinois

https://a3d3.ai