# Developing a Machine Learning Course for Chemical Engineering

Alex Albaugh

AIChE National Meeting

October 4, 2025

Boston, MA

# A Course For All ChE/MSE Students

- Wayne State Chemical Engineering & Materials Science
  - 26,000 student public R1 university in Detroit, MI
  - 91 B.S. students
  - 19 M.S. students
  - 18 PhD. students
  - 13 faculty
- Motivation
  - Modernize curriculum
  - Expand elective offerings
  - Integrate computational methods
- Designed for B.S., M.S., and Ph.D. students
  - Cross-listed at upper division undergrad/PhD level
  - In-person, shared lectures
  - Students from biology, chemistry, physics
  - Not a programming course!
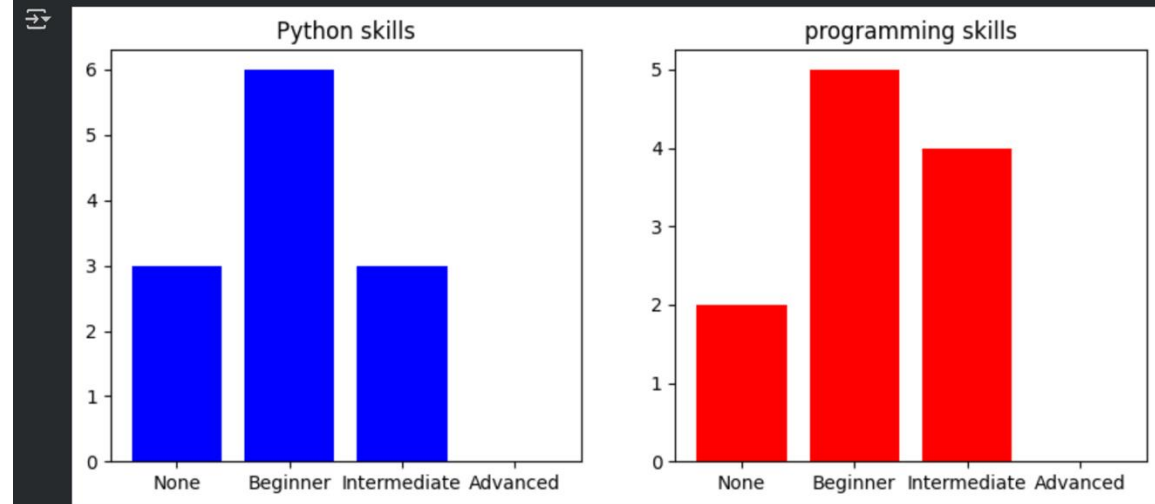
# Python, Our Language of Choice

- Used in intro. programming course

- Widely adopted

- Portable

- User friendly

- Libraries for everything
  - `numpy`: computation
  - `pandas`: data I/O
  - `matplotlib`: plotting
  - `sklearn`: general ML
  - `keras`: neural networks

- Free

```python
import numpy as np
import matplotlib.pyplot as plt

categories = ['None','Beginner','Intermediate','Advanced']
python = [3,6,3,0]
programming = [2,5,4,0]

fig,axes = plt.subplots(1,2,figsize=(10,4))
axes[0].bar(categories,python,color='b')
axes[1].bar(categories,programming,color='r')
axes[0].set_title('Python skills')
axes[1].set_title('programming skills')

plt.show()
```
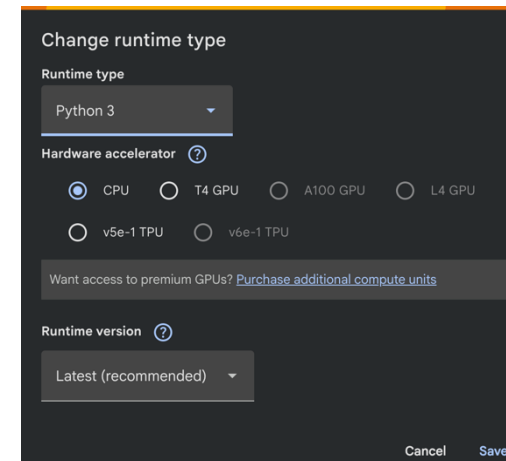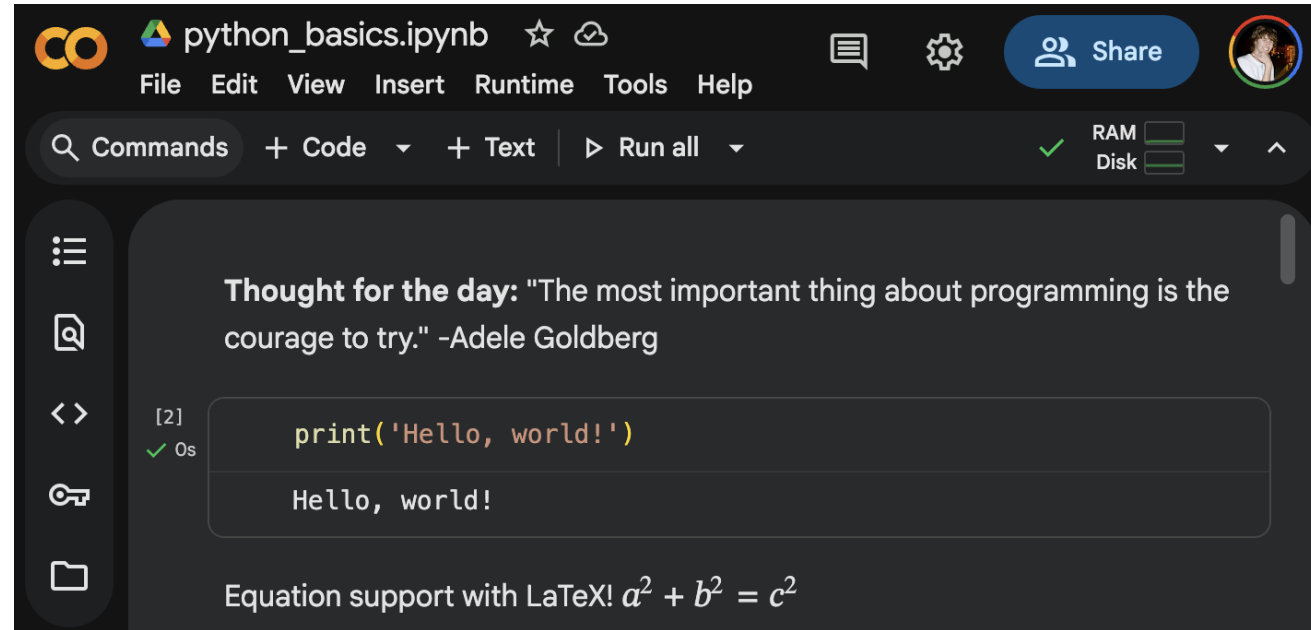
Lesson 2: Python Basics

# Google Colab, Our Environment of Choice

- Supports notebooks (.ipynb)
- No installation
- No version management
- No environment management
- Compatible with any machine
- "Explain error" feature
- Requirements:
  - Browser + internet
  - Google account
- Server-side compute
- Access to advanced hardware
- Free



Lesson 2:
Python Basics

colab.research.google.com

# Tip: Share Data Sets via Github URL

- Upload data sets to public Github

- Share URL

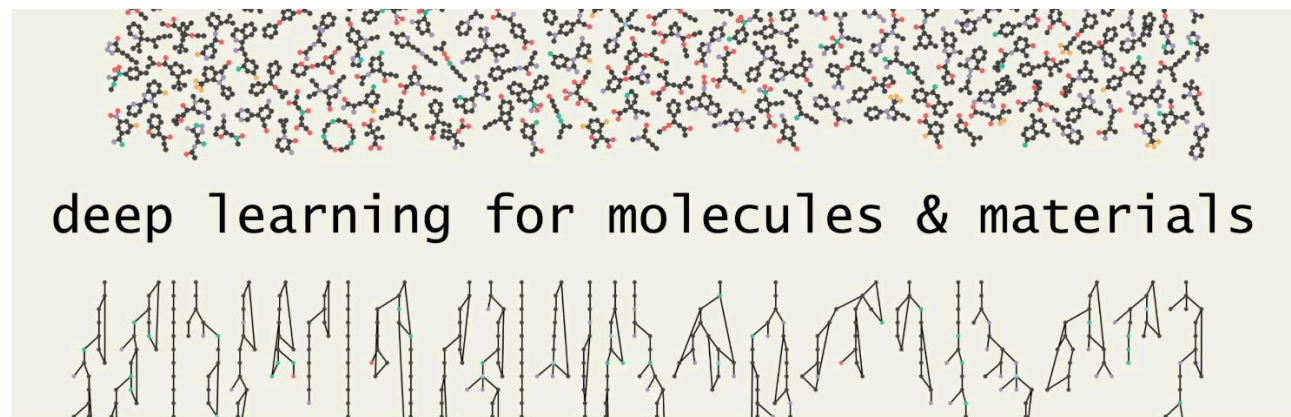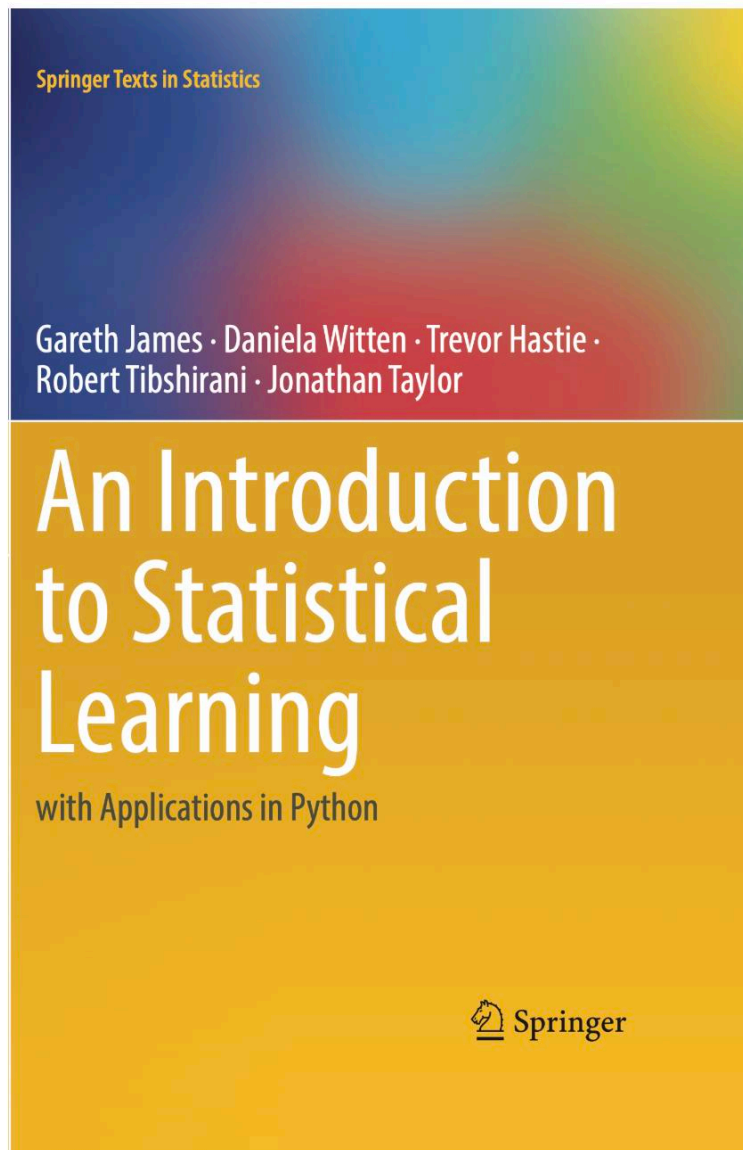- Load directly from URL in Colab

- No download/upload issues

- Free

```python
#read data
df = pd.read_csv('https://raw.githubusercontent.com/albaugh/CHE7507/refs/heads/main/Lecture8/combined_cycle_power_plant.csv')

#set features
X = df[['AT', 'AP', 'RH', 'V']].values

#set target
y = df['PE'].values
```

Lesson 8: Validation

# Bridging Intro. ML and ChE Applications



Andrew White

Great chemical engineering and chemistry examples
Jumps in with advanced methods
Free

Excellent introduction to concepts with broad coverage
No chemical applications
Free

# Lesson Structure: Theory + Practice



Lesson 13: Convolutional Neural Networks

| Concept | Examples |
|---|---|
| Python programming, packages | |
| Linear algebra | Process modeling, Markov models, chemical reaction networks |
| Statistics | Kinetic theory of gases, Maxwell-Boltzmann statistics |



Lesson 4: Linear Algebra



Lesson 21: Statistics

# Content: Core Machine Learning Concepts

| Concept | Examples |
|---|---|
| Regression, classification | Handwriting classification, FDA approval prediction |
| Bias-variance, training/testing/validating, regularization | Sabatier principle, power plant output prediction |
| Featurization | SMILEs strings, molecular descriptors, molecular features |
| Optimization | Optimizing reactor conditions |
| Ethics | LLMs & copyright law, ML in medical decision making, environmental impact |



Lesson 8: Validation



Lesson 15: Featurization



Lesson 9: Optimization

# Content: Basic Machine Learning Methods

| Concept | Examples |
|---|---|
| Regression (linear, multiple, polynomial, ridge, LASSO, logistic) | Control of liquid tank height, Sabatier principle |
| Principal component analysis | Periodic table analysis |
| K-nearest neighbors, K-means | Phase classification |



Lesson 6: Polynomial Regression



Lesson 20: PCA, KNN, K-Means

# Content: Advanced ML Methods

| Concept | Examples |
| --- | --- |
| Auto-differentiation, backprop. | |
| Neural networks (dense, convolutional, recurrent) | MOF $H_2$ loading prediction, small molecule solubility prediction |
| Discriminant analysis | Periodic table classification |
| Decision trees, random forests | Surfactant phase behavior |
| Molecular simulation | Covid-19 protein spike, machine-learned potentials |
| AlphaFold | Protein structure from sequences |

Lesson 10:
Intro to NNs

Lesson 23:
Decision Trees

# Tip: Run AlphaFold Easily with ColabFold

Mirdita, *et al.*, *Nature Methods*, 2022

# Relate Assignments to Real Applications

2. **Solving a Process Flow Diagram with Linear Algebra.** A reactor-separator-recycle section is a common paradigm in chemical processing. Calculations involving recycle streams can be a pain to do by hand. In the following process, we feed a stream of mostly A and a little C into a reactor where some A is converted to B. C is inert and does not participate in the reaction. After the reactor, the stream is fed to a flash drum where the bottom stream is a concentrated stream of our product B. The top stream is split into a purge and a recycle. The recycle is mixed with the feed stream and fed back into the reactor. Our mass flow rates are given as $m_{Xi}$ where $X = A, B, C$ is the species and $i = 1, 2, 3, 4, 5, 6$ is the stream number.



2. **Optimizing a Chemical Reactor.** A batch chemical reactor is a closed reaction vessel that is charged with reactants at time $t = 0$. The reaction proceeds for a certain amount of time before the products are extracted. Using an external jacket, we can control the temperature of the reactor. We are designing a reactor to optimize the depicted reaction scheme. Our reactant is species A, which undergoes a first order reaction with rate constant $k_1$ into species B. Species B can then undergo a reaction into species C with rate constant $k_2$ or species D with rate constant $k_3$. Species B is the desired product and we want to maximize its concentration.



batch reactor

5. **(CHE 7507 only) Polymer Layers, Recurrent Layers.** A block polymer consists of a section made of one type of polymer and a section made of a different type of polymer. These block polymers form alternating layers, called lamellae. A computational study looked at how the composition of a block affects the width (period)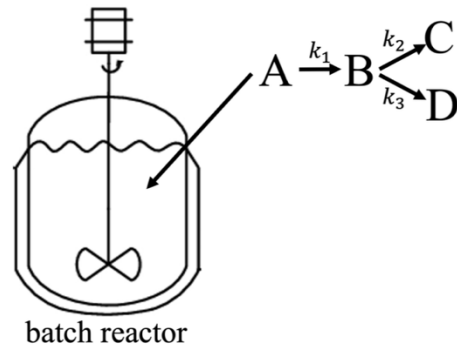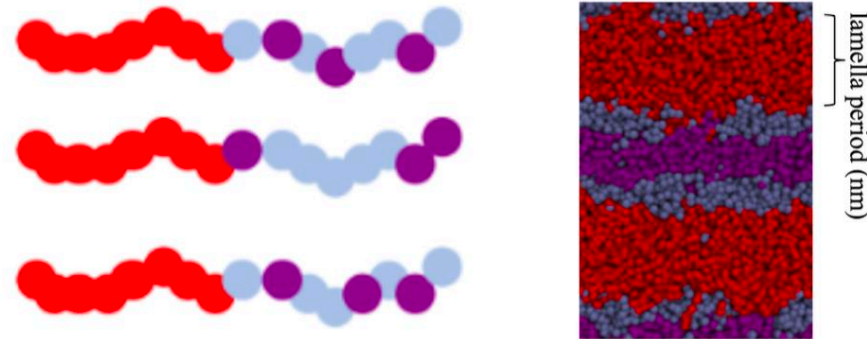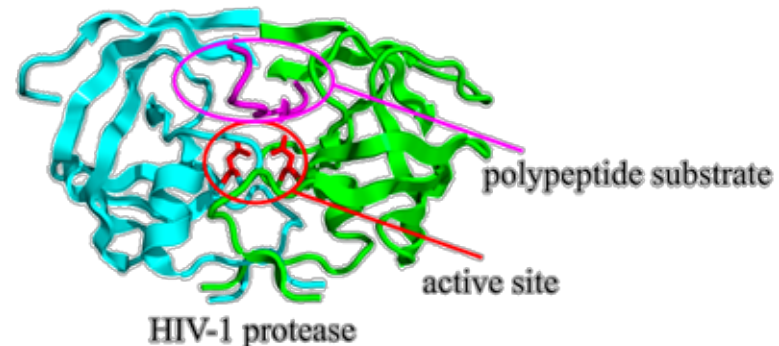 of these layers based on the composition of one of the block polymer sections. You can find this study here: `https://pubs.acs.org/doi/full/10.1021/acs.macromol.3c02401`



lamella period (nm)

2. **Fighting a Virus with Machine Learning.** Pictured below is HIV-1 protease, a critical component of the human immunodeficiency virus (HIV). HIV-1 protease is an enzyme, which is a protein that catalyzes a chemical reaction. Enzymes have active sites where the chemical reactions take place. Specifically, a protease cuts (cleaves) another protein into pieces. HIV-1 protease is critical in the "life cycle" of the virus because it cuts long proteins into functional pieces. Because it is critical to virus reproduction, HIV-1 protease is a drug target. An inhibitor is a molecule that can bind to an enzyme active site and cause it to stop work, essentially clogging up the enzyme. HIV-1 protease inhibitors are a common class of drug for treating HIV infections.



polypeptide substrate

active site

HIV-1 protease

Mysona, *et al.*, *Macromolecules*, 2024; Rögnvaldsson, *et al.*, *Bioinformatics*, 2015

# Tip: UCI ML Repo for Lots of Clean Datasets



**UC Irvine**
**Machine Learning Repository**

### Toxicity
Donated on 5/4/2022

The dataset includes 171 molecules designed for functional domains of a core clock protein, CRY1, responsible for generating circadian rhythm. 56 of the molecules are toxic and the rest are non-toxic.

| Dataset Characteristics | Subject Area | Associated Tasks |
|---|---|---|
| Tabular | Biology | Classification |

| Feature Type | # Instances | # Features |
|---|---|---|
| - | 171 | 1203 |

### Superconductivty Data
Donated on 10/11/2018

Two file s contain data on 21263 superconductors and their relevant features.

| Dataset Characteristics | Subject Area | Associated Tasks |
|---|---|---|
| Multivariate | Physics and Chemistry | Regression |

| Feature Type | # Instances | # Features |
|---|---|---|
| Real | 21263 | 81 |

### HIV-1 protease cleavage
Donated on 4/24/2015

The data contains lists of octamers (8 amino acids) and a flag (-1 or 1) depending on whether HIV-1 protease will cleave in the central position (between amino acids 4 and 5).

| Dataset Characteristics | Subject Area | Associated Tasks |
|---|---|---|
| Multivariate | Biology | Classification |

| Feature Type | # Instances | # Features |
|---|---|---|
| Categorical | 6590 | 1 |

archive.ics.uci.edu

# AI Use In Class

- LLM use (Claude, Copilot, ChatGPT) is permitted for programming
  - Not a programming course!
  - Introduce students to AI tools
  - Asked students to describe AI use on assignments
  - Use in a controlled environment

- Some problems are marked *"No AI"*
  - Honor system
  - How do I enforce this or design around it?

- Explore how LLMs work
  - Relationship between thermodynamic and LLM temperature
  - Shortcomings and legal issues

# Student Feedback

- "I appreciated that we started with learning some basics of python because as a beginner that was interested in the class, **I worried that it would be harder for me to keep up with the coding**."

- "remove most of the introductory stuff"

- "The scope of practical application, based on this courses learning objectives, has the opportunity to extend far beyond foundational chemical engineering practices. I enjoy this course because it introduces students to coding (Python), machine learning/data analytic techniques, and the **opportunity to experiment with large language models and AI platforms**."

- "It is really interesting to see the applications of machine learning in a variety of topics, in addition I got to see the relevant topics **combining information from previous classes** such as Schrodinger's equation of states from Physical Chemistry or reaction kinetics and reactor design."

- "The assignments helped us **experiment and practice** the concepts we learned."

- "I enjoyed the way your slides on each topic sort of **introduced topics from the ground floor–up**."
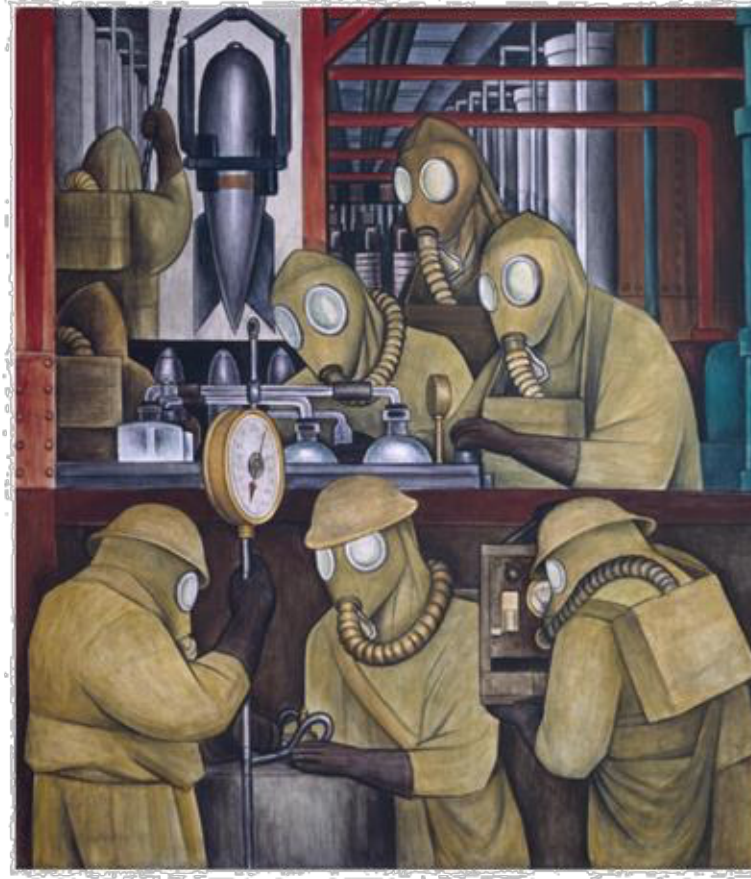
# Acknowledgements

- Jeffrey Potoff, Wayne State
- Camille Bishop, Wayne State
- Ali Seyedi, Wayne State
- Yamil Colón, Notre Dame
- Jindal Shah, Oklahoma State
- Bingqing Cheng, UC Berkeley
- Cory Simon, Oregon State
- Teresa Head-Gordon, UC Berkeley

# AI Tools Are Not Good or Bad, They Are What We Make Them



"Commercial Chemical Operations"



"Manufacture of Poison Gas Bombs"

Questions?

*Detroit Industry Murals*, 1932-1933, Diego Rivera, Detroit Institute of Arts