# AI in Drug Discovery – An Overview
# Session 1

**September 16, 2024**

# Who we are!

**Pat Walters**
Relay Therapeutics

**Raquel López-Ríos de Castro**
Charité Berlin, MSKCC NYC

**Michael Backenköhler**
Saarland University

**Andrea Volkamer**
Saarland University

# What we will do today

**Session 1 - 1:30 - 2:30 pm**

- An introduction to Artificial Intelligence (AI) and Machine Learning (ML)
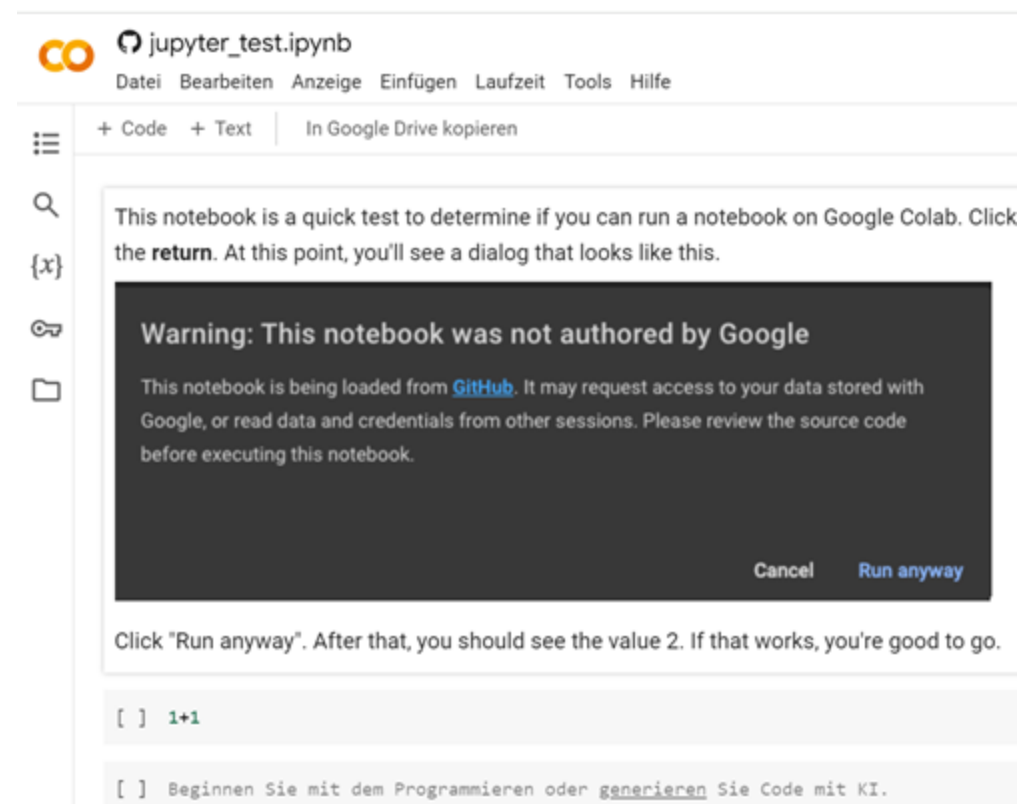- Molecular representations
- AI architectures

**Session 2 - 3:00 - 4:00 pm**

- The importance of data quality for AI/ML
- Exploratory data analysis
- Data preprocessing
- Applicability domains

**Session 3 - 4:30 - 5:30 pm**

- AI in Practice
- Molecule generation
- Active learning

Lectures supported by hands-on sessions …

## Definitions Can Be Tricky

**artificial intelligence (AI),** the ability of a digital computer or computer-controlled robot to perform tasks commonly associated with intelligent beings.

Not a well-defined statement



https://www.britannica.com/technology/artificial-intelligence

# AI and "The Rise of the Machines"



Andrew Chen Retweeted

Mat Velloso @matvelloso · Nov 22
Difference between machine learning and **AI**:

If it is written in Python, it's probably machine learning

If it is written in **PowerPoint**, it's probably **AI**

166    6.6K    19K

Show this thread

# What Is Machine Learning?

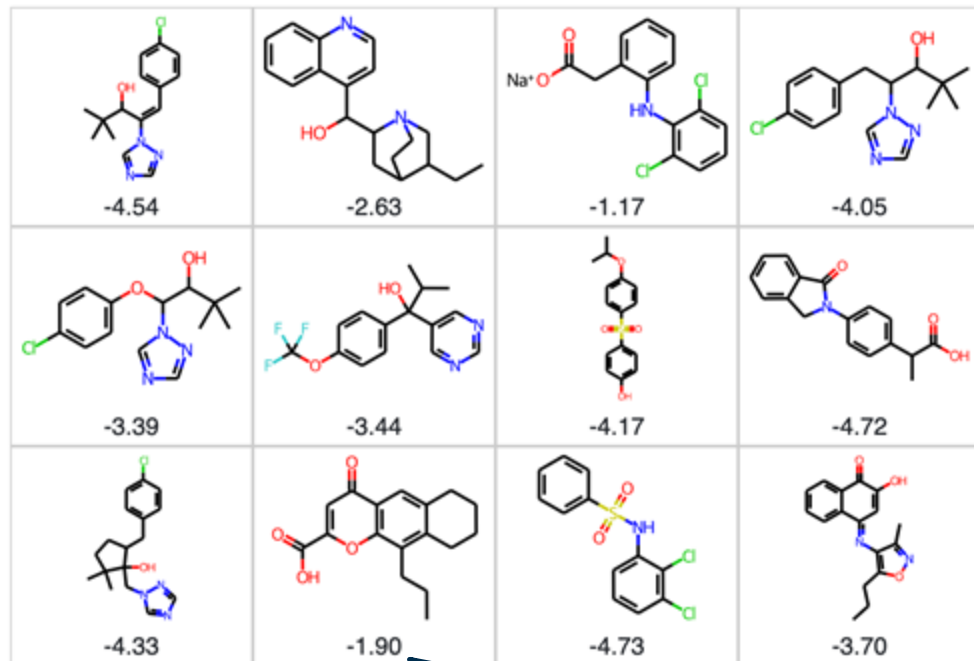Machine learning is all about labeling things using examples

**Cassie Kozyrkov, Google**
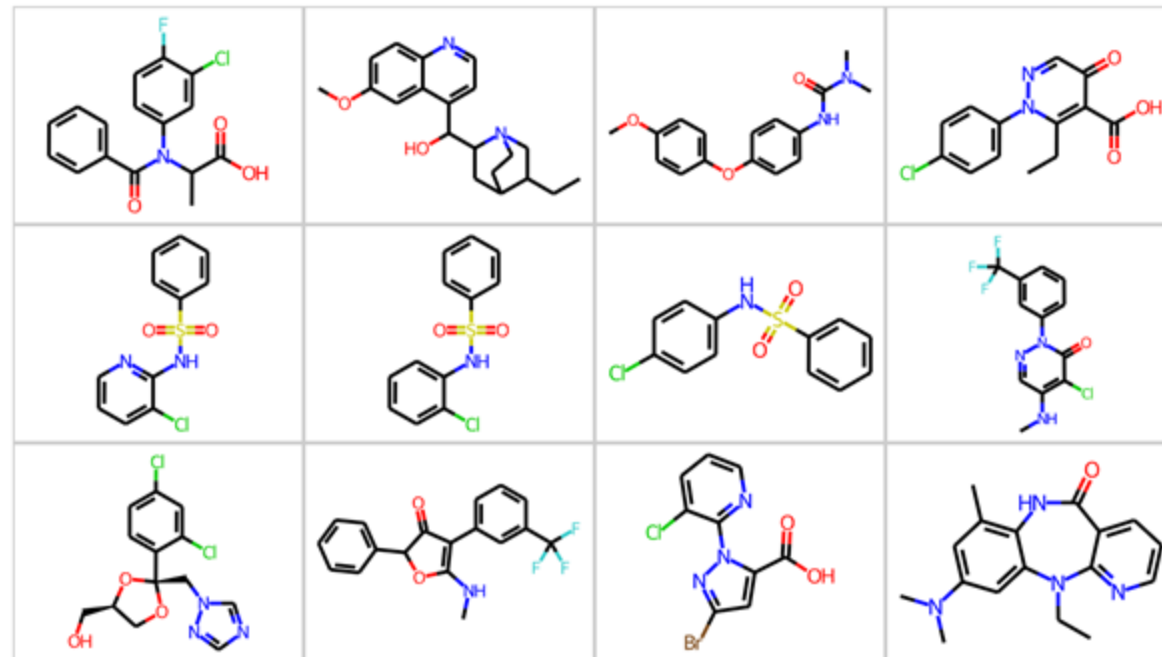
# Labeling Molecules Based on Examples
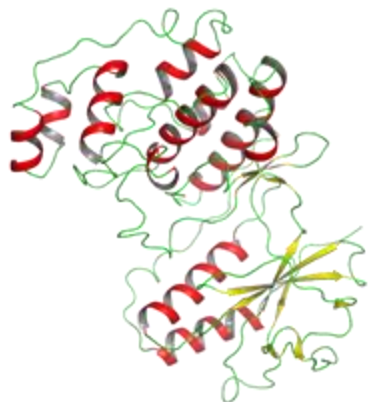
**Molecules with measured data**



-4.54    -2.63    -1.17    -4.05

-3.39    -3.44    -4.17    -4.72

-4.33    -1.90    -4.73    -3.70

**Molecules to be predicted**

**Log10(Molar Aqueous Solubility)**

# Using Predictive Models to Drive Drug Discovery



**On-target Activity**



**Off-target Activity**



Aqueous solubility
Biochemical activity
PK
Potency
PPB
in vivo activity
Vss
CYP inhibition
Selectivity
Cellular activity
Ligand efficiency
t1/2
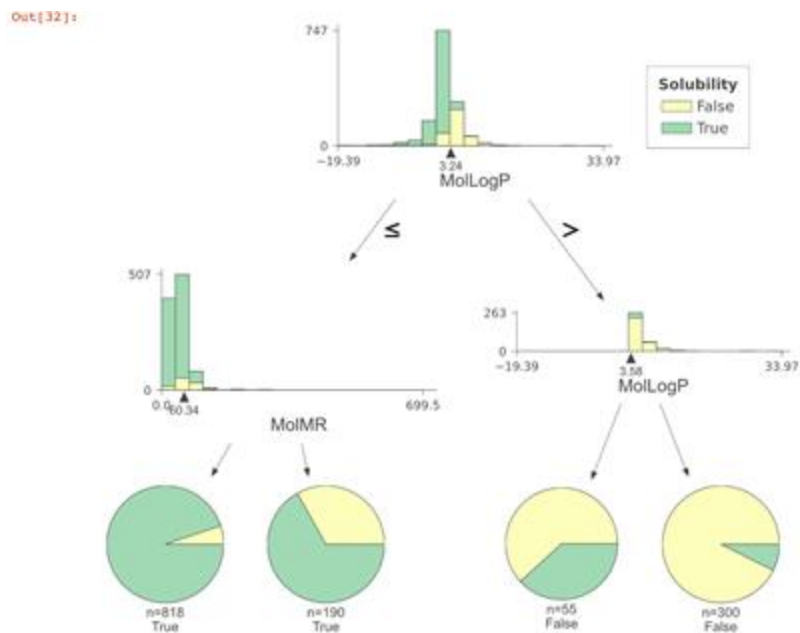FBS Binding
Hepatocyte stability
Dose
FaSSIF solubility
Efflux
hERG
Permeability



**Pharmacokinetics**



**Physical Properties**

# Two Types of ML Models – Classification and Regression

**Classification**



**Regression**



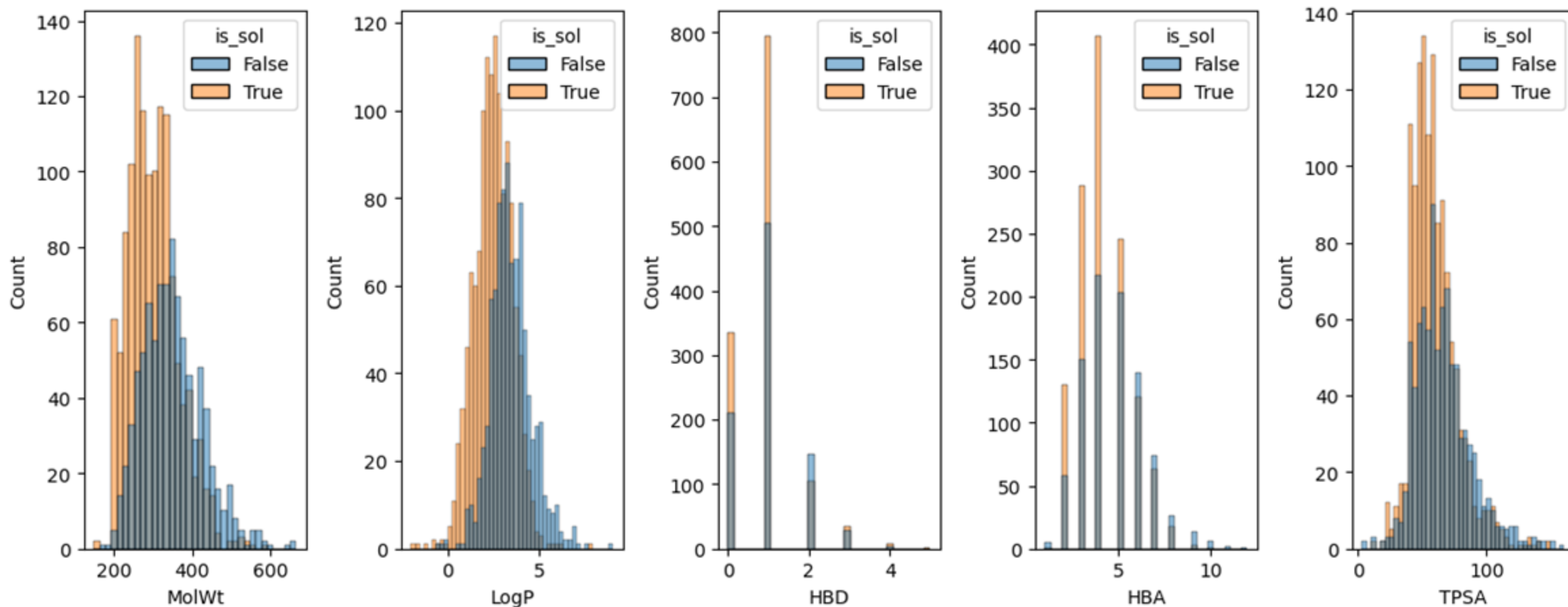**Predict a category, e.g. soluble/insoluble**

**Predict a value, e.g. 4.2**

# Let's Start With a Dataset



Aqueous solubility by CLND
2173 compounds
Min = 0.23µM
Max = 648µM

Fang, Cheng, et al. "Prospective validation of machine learning algorithms for absorption, distribution, metabolism, and excretion prediction: An industrial perspective." *Journal of Chemical Information and Modeling* 63.11 (2023): 3263-3274.
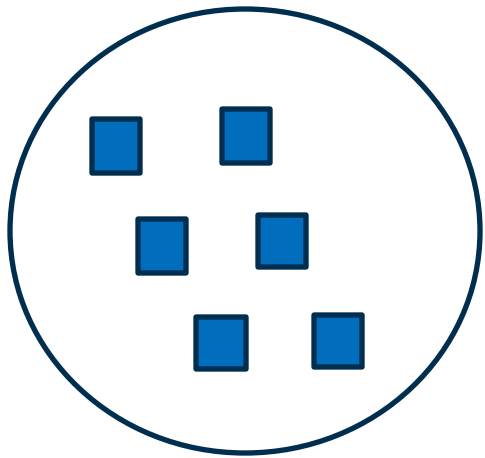
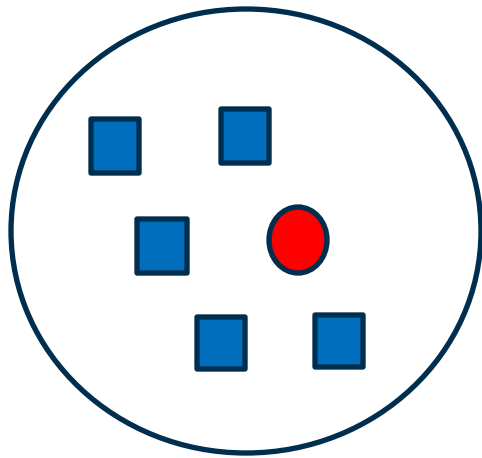# Which Property Best Separates Soluble vs Insoluble Molecules?
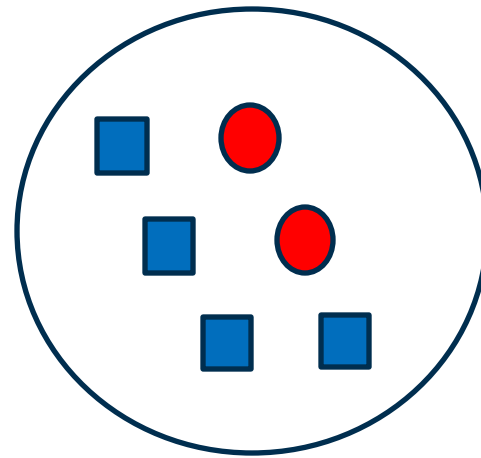
# The Gini Index Quantifies the "Purity" of a Split
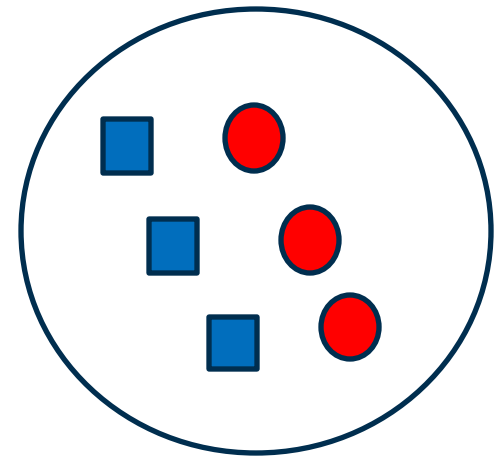
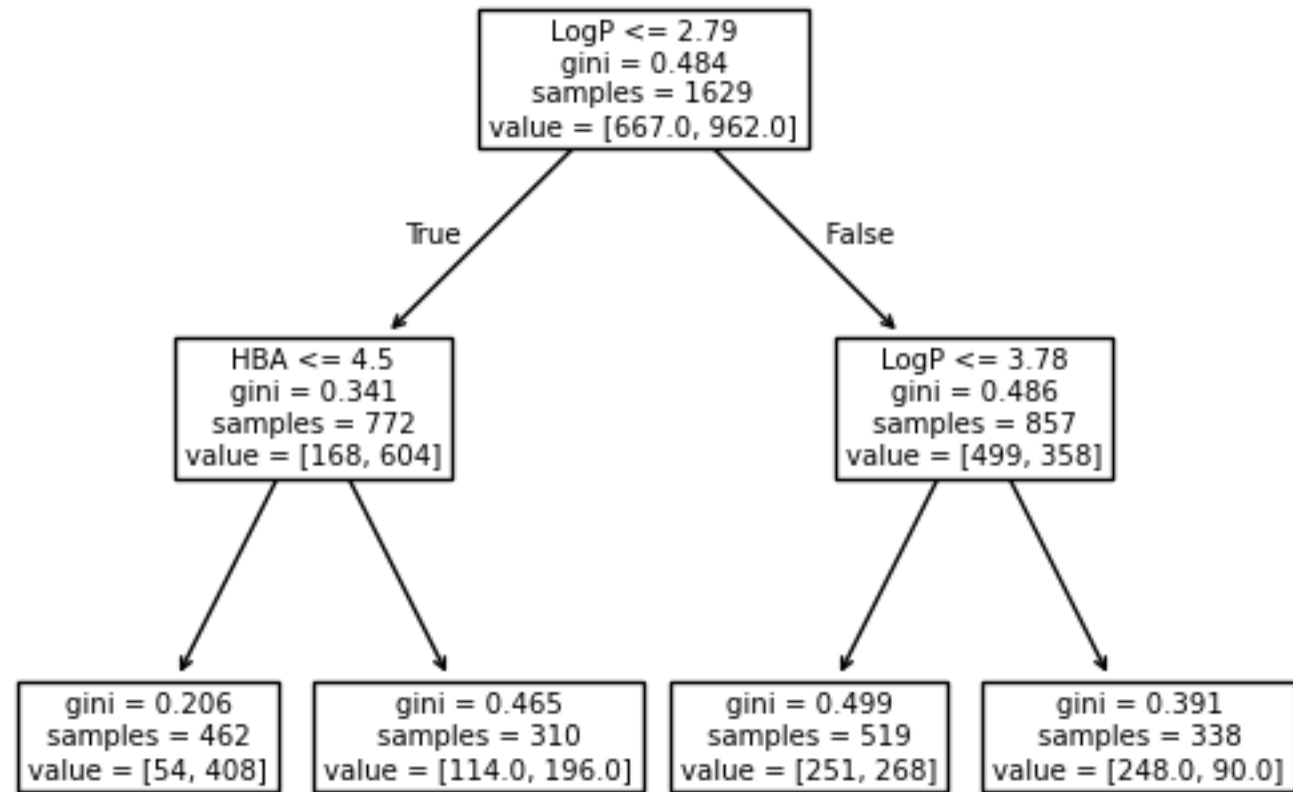$$\text{Gini Index} = 1 - \sum_{i=1}^{n} (P_i)^2$$



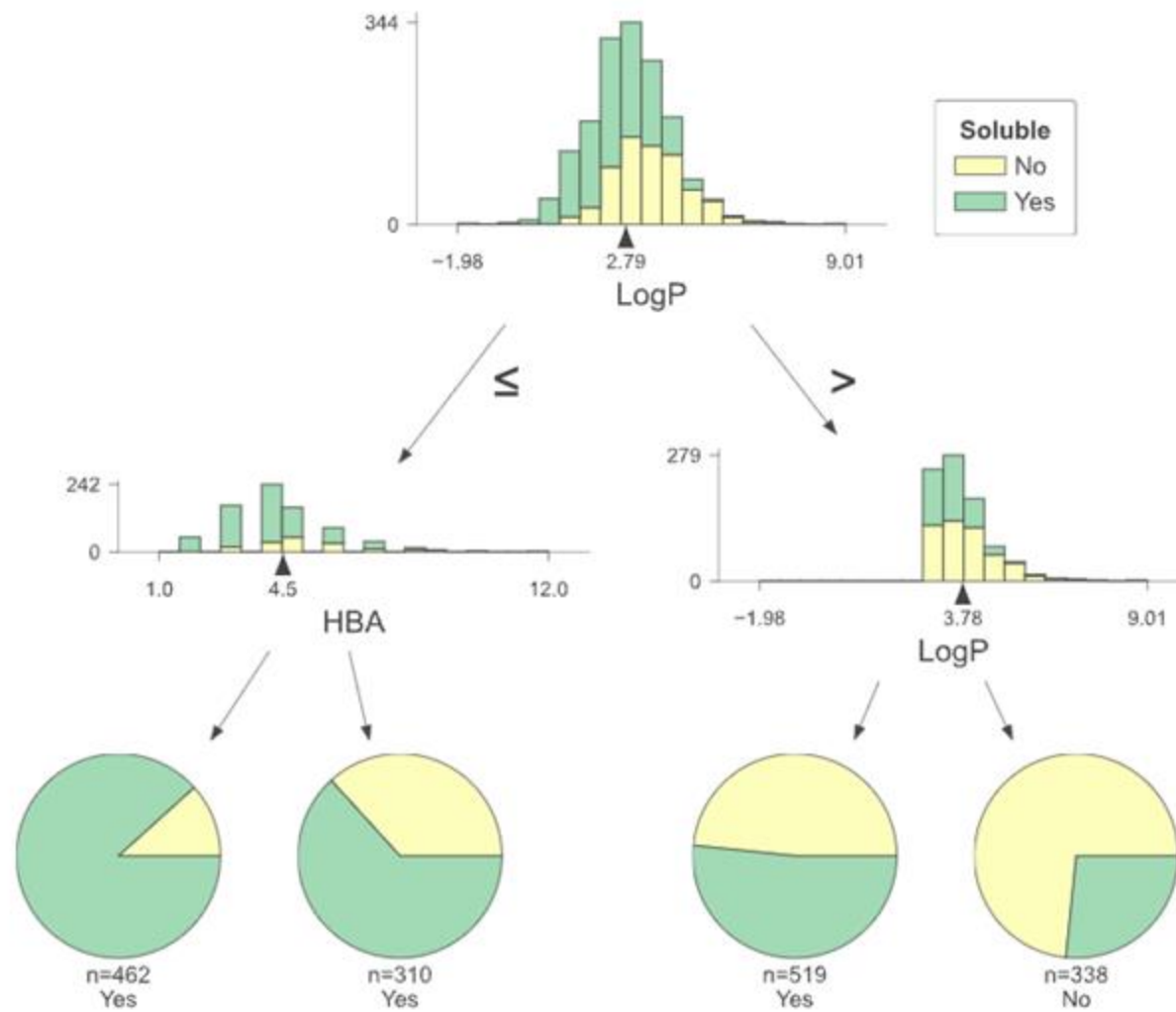| 1.00 | 0.83 | 0.66 | 0.50 |

# Build a Decision Tree
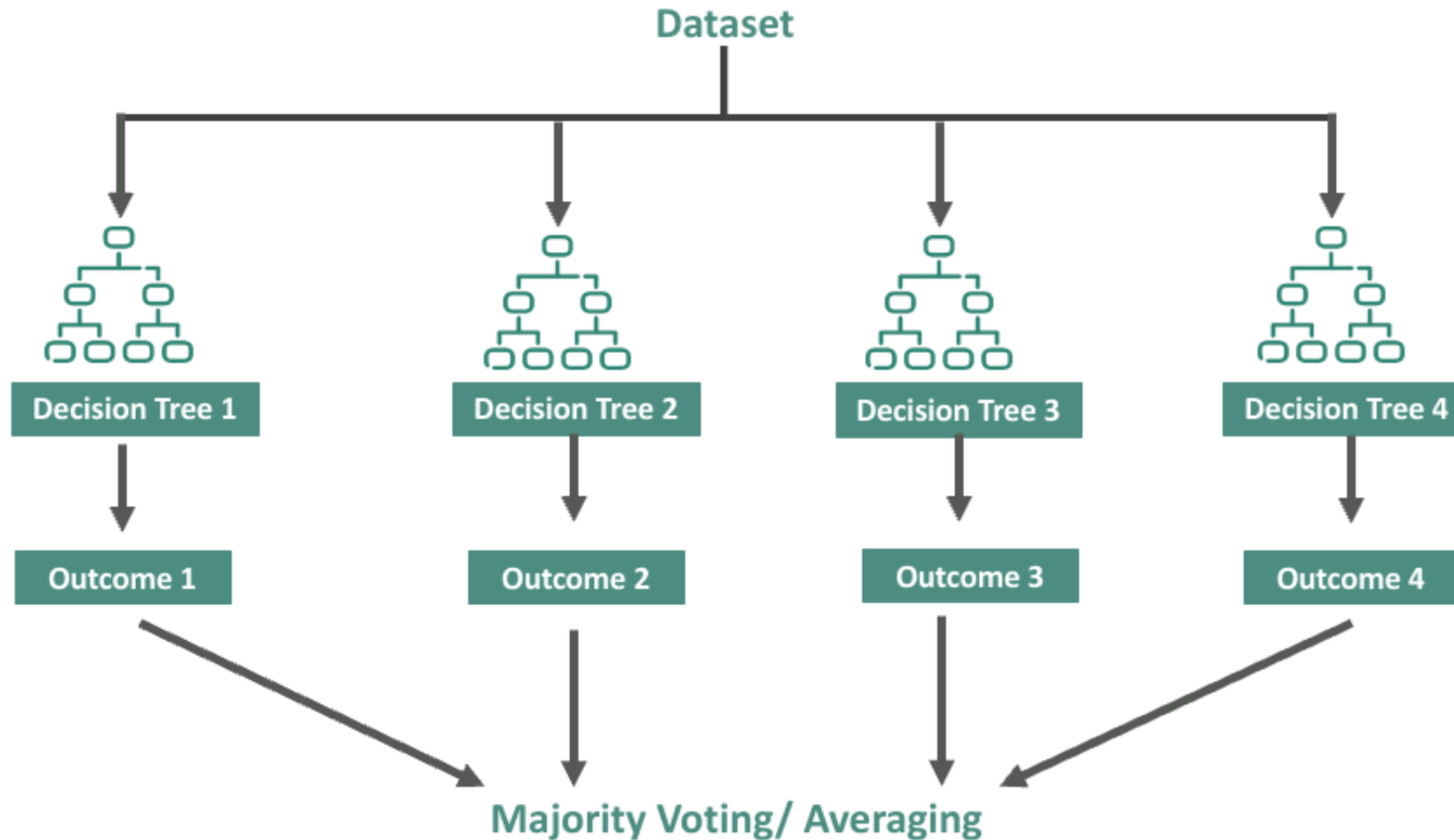
# A Better (IMO) Decision Tree Visualization

# Random Forest Uses an Ensemble of Decision Trees

# There Are Many Tree Ensemble Methods

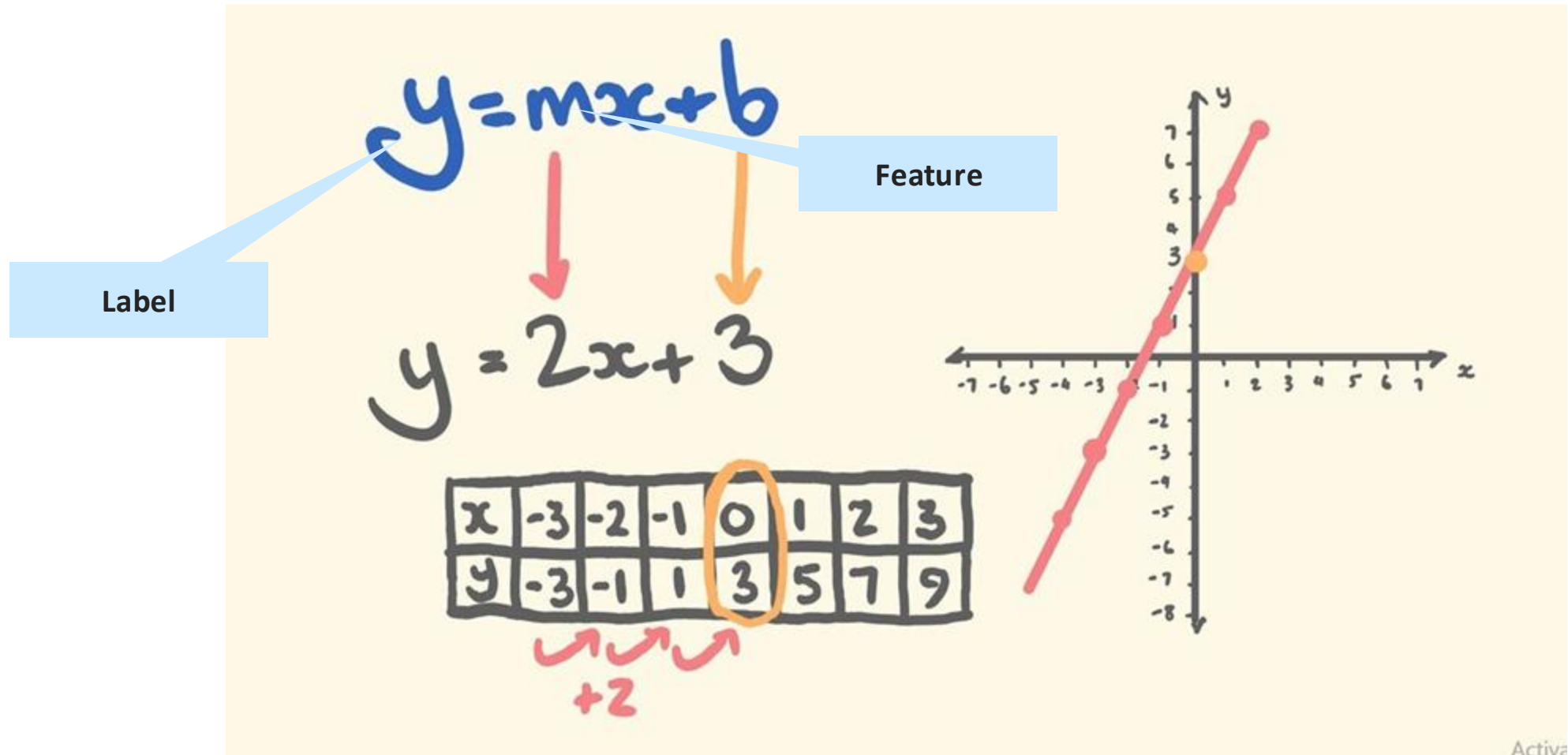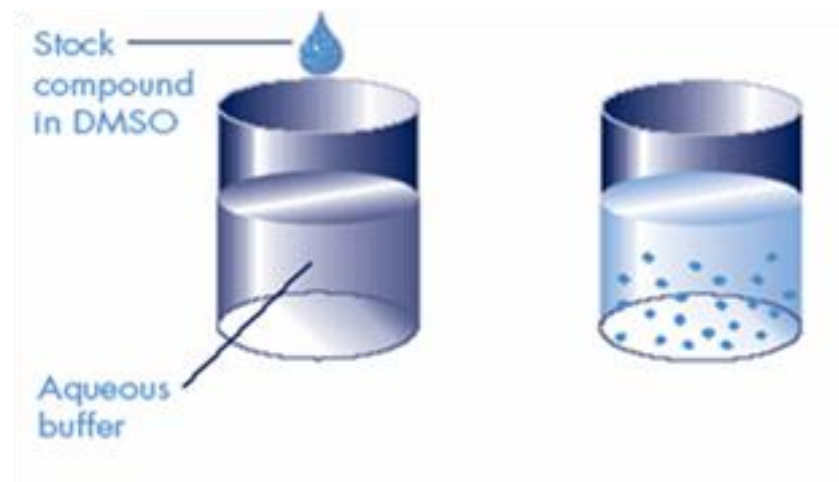# ML Predicts a Set of Labels (y) Based on Features (X)

# Defining Labels (y) and Features (X) - An Example



**y = Log Aqueous Solubility**

**Label**

0101100100010001110010001100101010101011

**X = A Vector Representing a Molecule**

**Features**

# Define Features Based on Books People Have Read



**Fred**

| 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 |

**Sally**

| 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 |

**Jane**

| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 |

**Create vectors representing books purchased by individuals**
**1 = bought book**
**0 = did not buy book**

# Chemical Fingerprints as Molecular Descriptors



| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **1** | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 |
| **2** | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |

# There are A LOT of ways to do this



*2 Volumes*
*6,000 references from 450 journals*

**DRAGON 7 has 5,270 descriptors z**

# Training a Machine Learning Model



Training Set

Assayed Ligands

Feature Vector

| | |
|---|---|
| MOL001 | 1.7 |
| MOL002 | 2.3 |
| MOL003 | 4.1 |
| MOL004 | 5.8 |
| MOL005 | 3.0 |
| MOL006 | 2.1 |
| MOL007 | 1.2 |
| MOL008 | 4.3 |
| MOL009 | 3.1 |
| MOL010 | 2.1 |

0 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 1 1 0...
0 0 0 0 0 0 1 1 0 0 1 0 0 0 0 0 0 0 1 0...
0 0 0 1 0 0 1 1 0 0 1 0 0 0 1 0 0 0 0 1 0...
0 1 0 1 0 0 0 1 0 1 1 0 0 0 1 0 0 0 0 1 0...
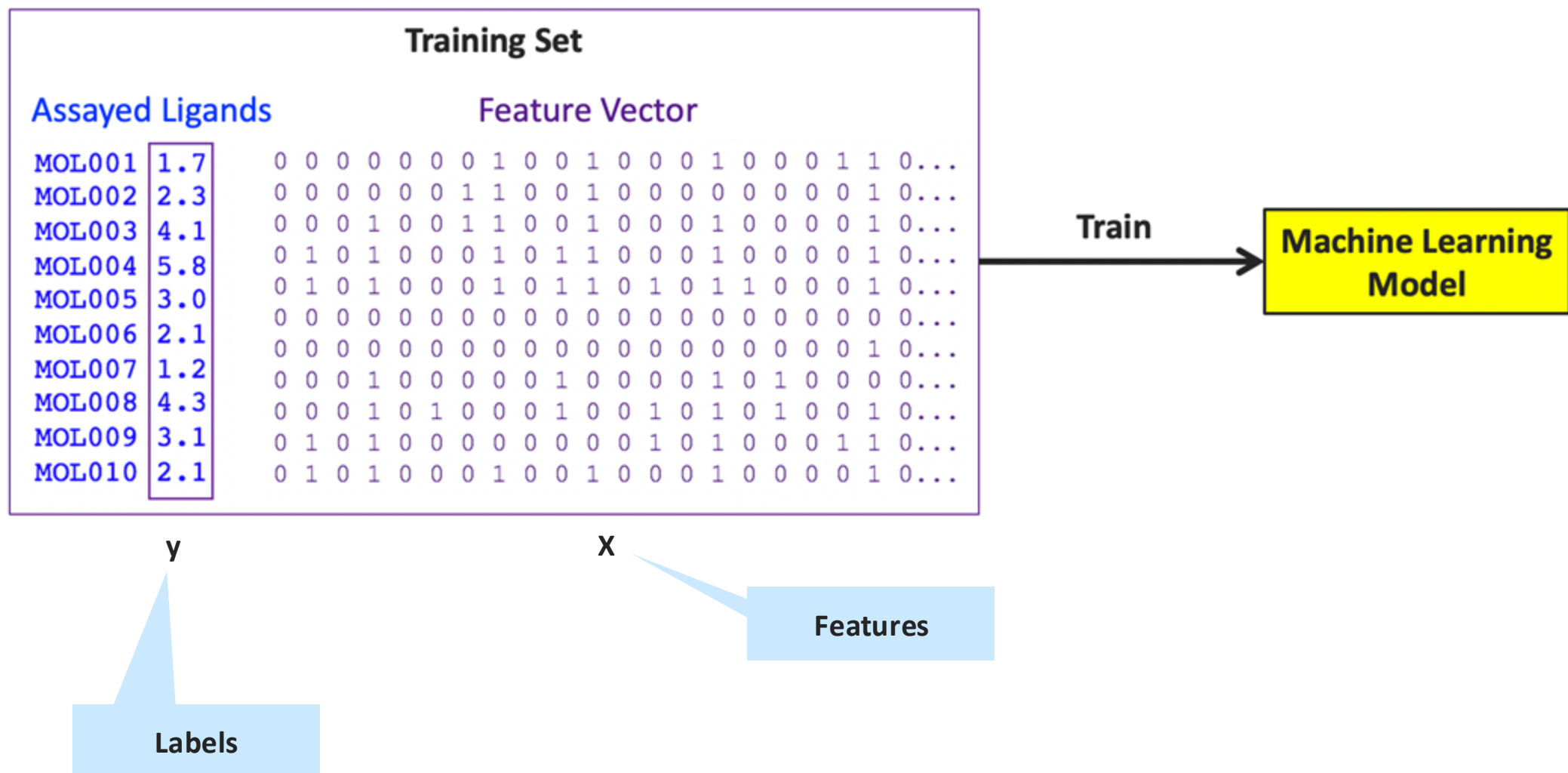0 1 0 1 0 0 0 1 0 1 1 0 1 0 1 1 0 0 0 1 0...
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0...
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0...
0 0 0 1 0 0 0 0 0 1 0 0 0 0 1 0 1 0 0 0 0...
0 0 0 1 0 1 0 0 0 1 0 0 1 0 1 0 1 0 0 1 0...
0 1 0 1 0 0 0 0 0 0 0 1 0 1 0 0 0 1 1 0...
0 1 0 1 0 0 0 1 0 0 1 0 0 0 1 0 0 0 0 1 0...

**Train** → **Machine Learning Model**

y

X

Features

Labels

# Making Predictions With a Machine Learning Model
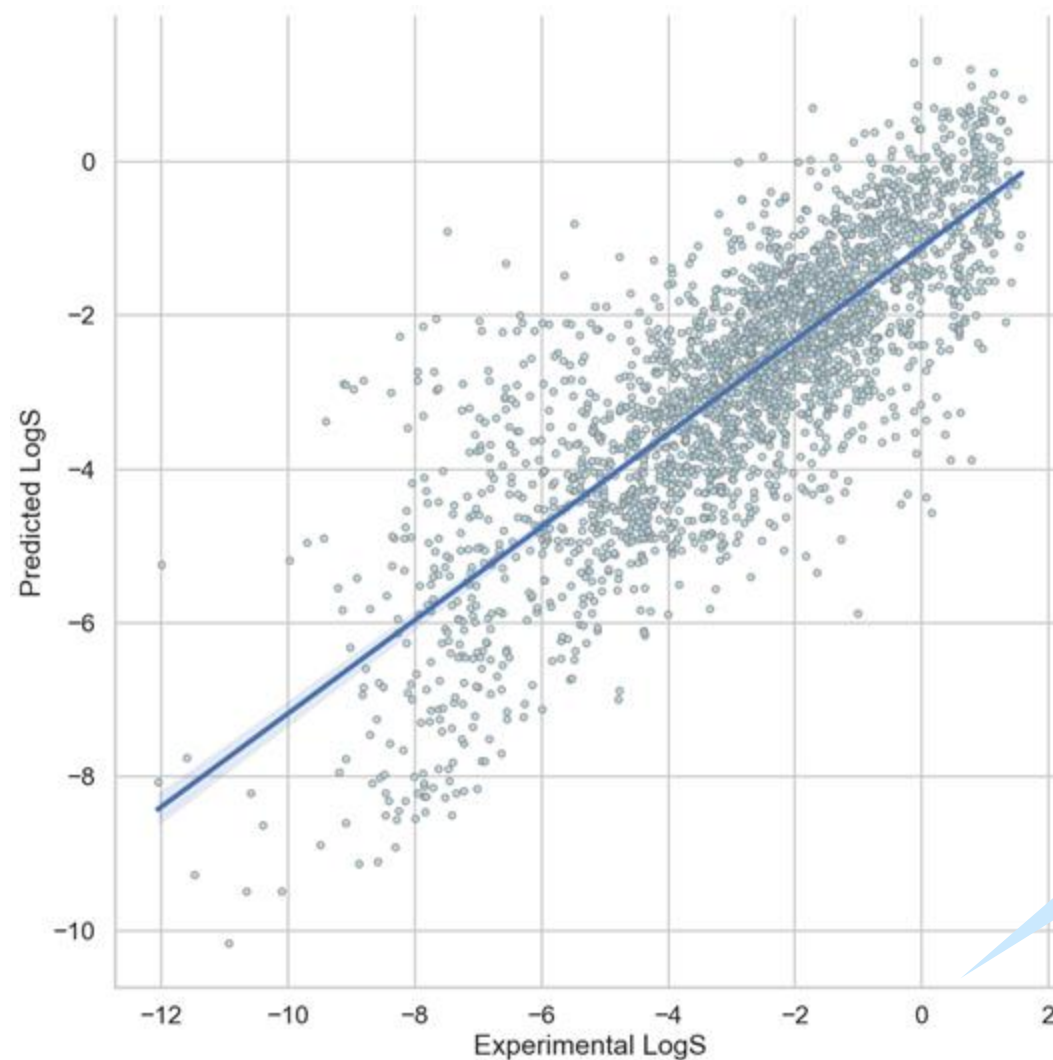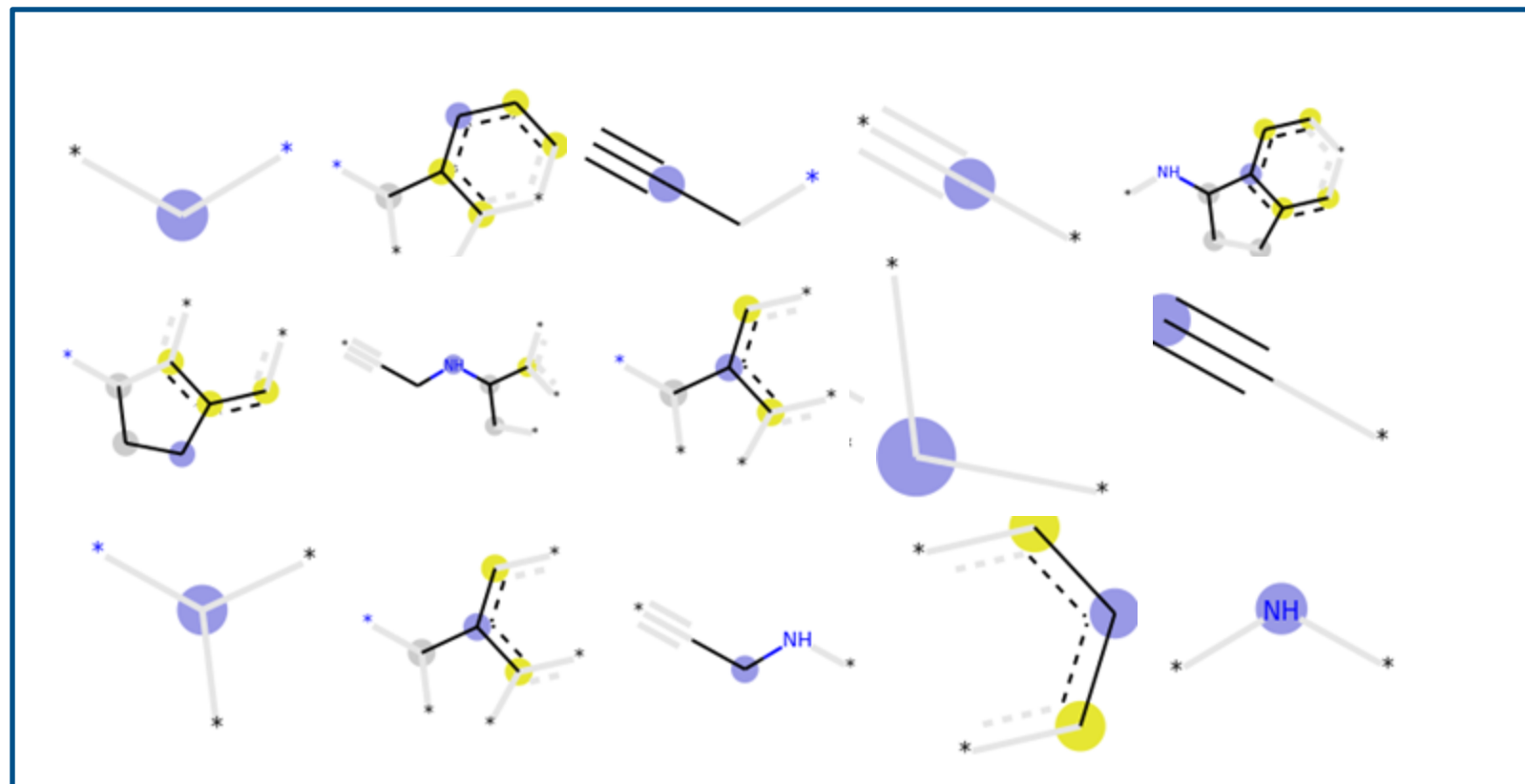
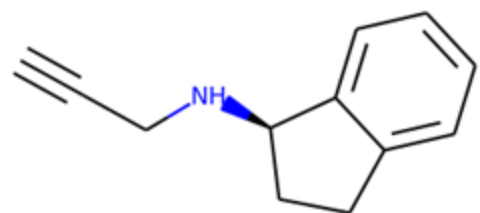# Prediction Performance of an Aqueous Solubility Model



Dynamic range is very large, does not represent a typical use case

Nodes are functions that accept values from connected nodes and output a value between 0 and 1

Predicted value = 3.2
True value = 2.8
Error = 0.4

Backpropagate to minimize loss (e.g. RMSE)

Input Layer ∈ R⁹⁶   Hidden Layer ∈ R¹²   Hidden Layer ∈ R¹⁰   Output Layer ∈ R¹

**Graph Convolutions**

*J. Comput. Aided Mol. Des*. 2016, 595–608

**Message Passing Neural Network**

*J. Chem. Inf. Model.* 2019, 59, 3370–3388

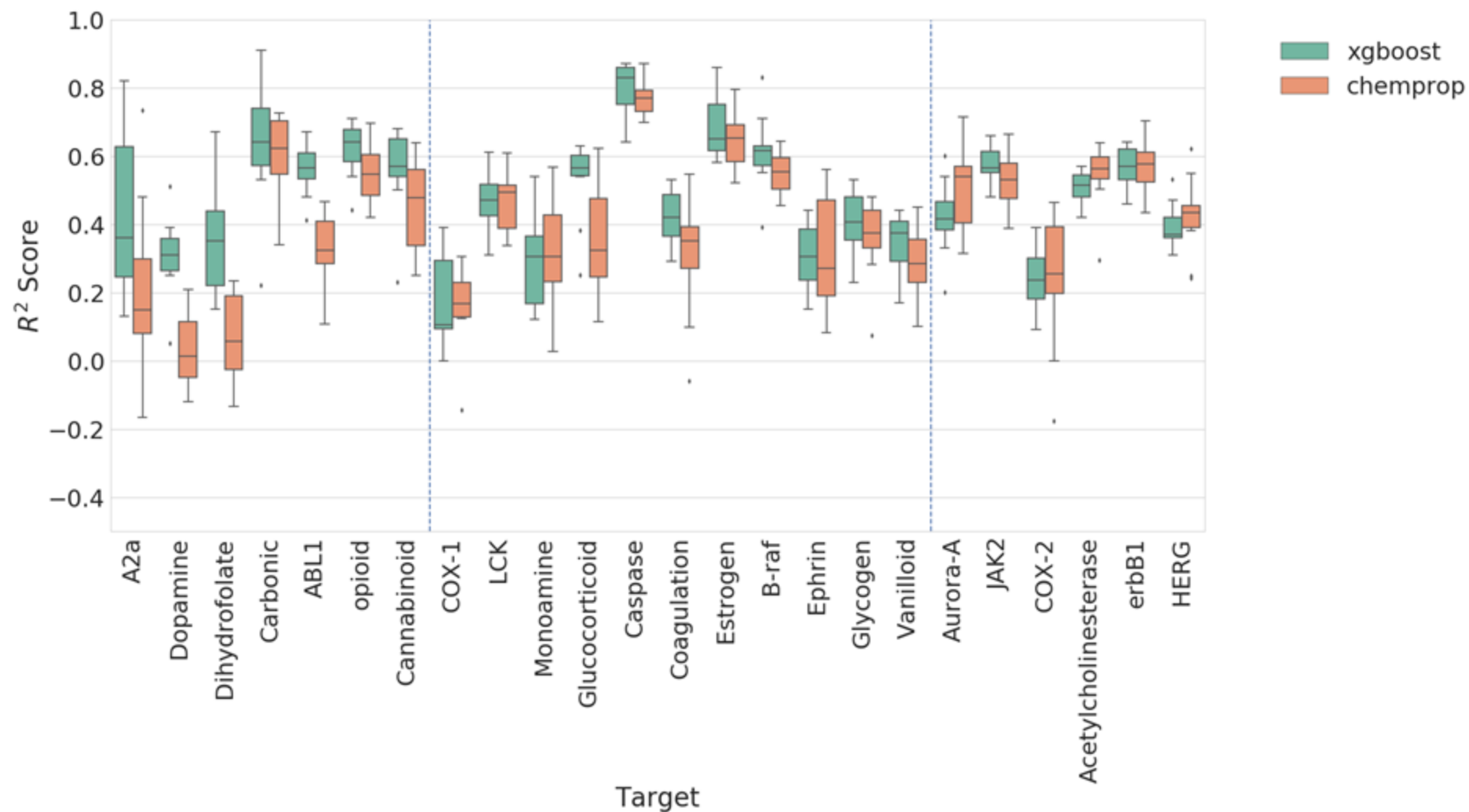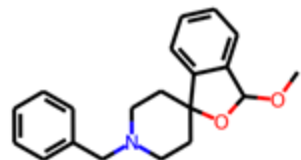# Are Neural Network Representations Better?

# Incorporating 3D into Molecular Machine Learning is an Unsolved Problem



**2D single-instance**

0 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 1 1 0... → 5.21 ✓

**3D single-instance**

0 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 1 1 0... → 5.21 ✗

**3D multiple-instance**

0 0 0 0 0 0 0 1 0 0 1 0 0 0 1 0 0 0 1 1 0...
0 0 0 0 0 0 1 1 0 0 1 0 0 0 0 0 0 0 1 0...
0 0 0 1 0 0 1 1 0 0 1 0 0 0 1 0 0 0 0 1 0...
0 1 0 1 0 0 0 1 0 1 1 0 0 0 1 0 0 0 0 1 0...
0 1 0 1 0 0 0 1 0 1 1 0 1 0 1 1 0 0 0 1 0...
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0...
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0...
0 0 0 1 0 0 0 0 0 1 0 0 0 0 1 0 1 0 0 0 0...
0 0 0 1 0 1 0 0 0 1 0 0 1 0 1 0 1 0 0 1 0...
0 1 0 1 0 0 0 0 0 0 0 0 1 0 1 0 0 0 1 1 0...
0 1 0 1 0 0 0 1 0 0 1 0 0 0 1 0 0 0 0 1 0...

→ 5.21 ?

# How to Find Me

https://pwalters.github.io



**My Boring Website**    Publications    Tutorials    Blog    Videos

**Pat Walters**

Cheminformatics, ML

📍 Cambridge, MA

✉ Email

𝒮 Google Scholar

◯ Github

in LinkedIn

✗ X (formerly Twitter)

Pat Walters is Chief Data Officer at Relay Therapeutics in Cambridge, MA. Prior to joining Relay, he spent more than 20 years at Vertex Pharmaceuticals where he was Global Head of Modeling & Informatics. Pat is the 2023 recipient of the Herman Skolnik Award for Chemical Information Science from the American Chemical Society. He is a member of the editorial advisory boards for the Journal of Chemical Information and Modeling and Artificial Intelligence in the Life Sciences, and previously held a similar role with the Journal of Medicinal Chemistry. Pat is co-author of the book "Deep Learning for the Life Sciences", published in 2019 by O'Reilly and Associates. He received his Ph.D. in Organic Chemistry from the University of Arizona where he studied the application of artificial intelligence in conformational analysis. Prior to obtaining his Ph.D., Pat worked at Varian Instruments as both a chemist and a software developer. He received his B.S. in Chemistry from the University of California, Santa Barbara.