

A decorative banner at the top of the slide features a dark blue background with various molecular models. On the left, there are several colorful, semi-transparent carbon cages (C60, C70, etc.) in shades of yellow, green, and purple. In the center, there are blue wireframe structures. On the right, there are grey ball-and-stick models of organic molecules, including one with a pink atom. A large, detailed grey ball-and-stick model of a graphene sheet is on the far right.

Cornell Future Of Care

Karl Leswing

karl.leswing@gmail.com

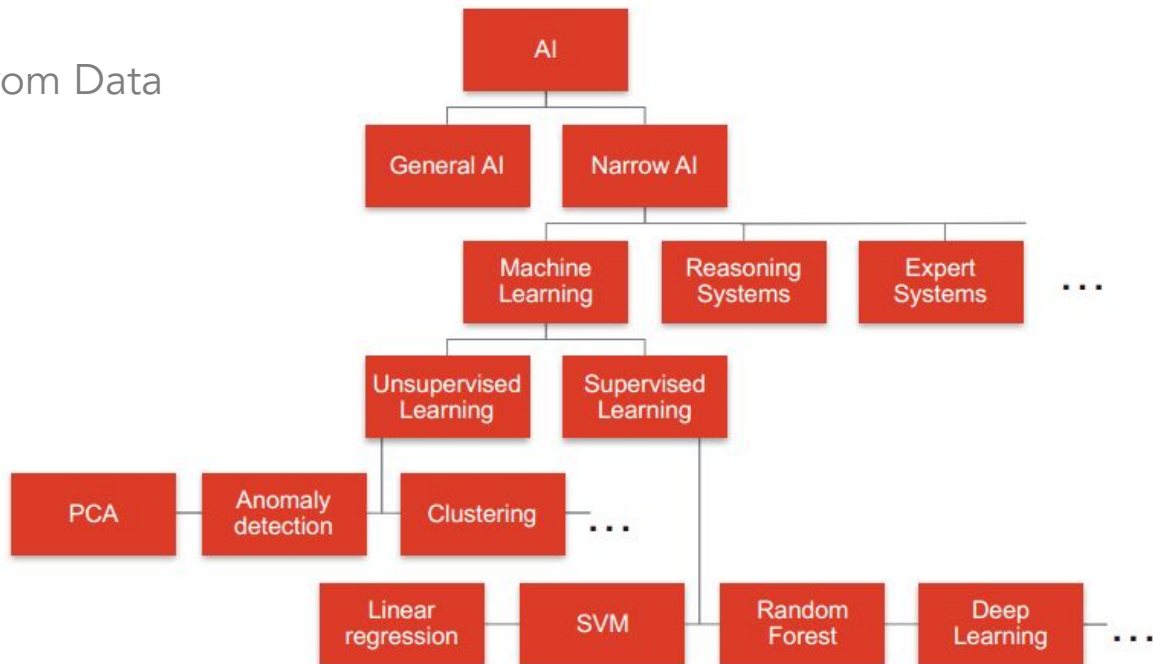
FB/Insta/Twitter/Github: lilleswing

October 2018

https://github.com/lilleswing/future_of_care

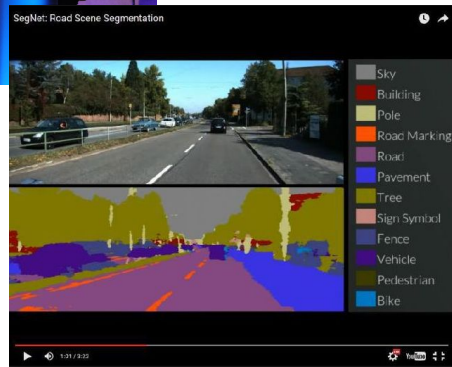
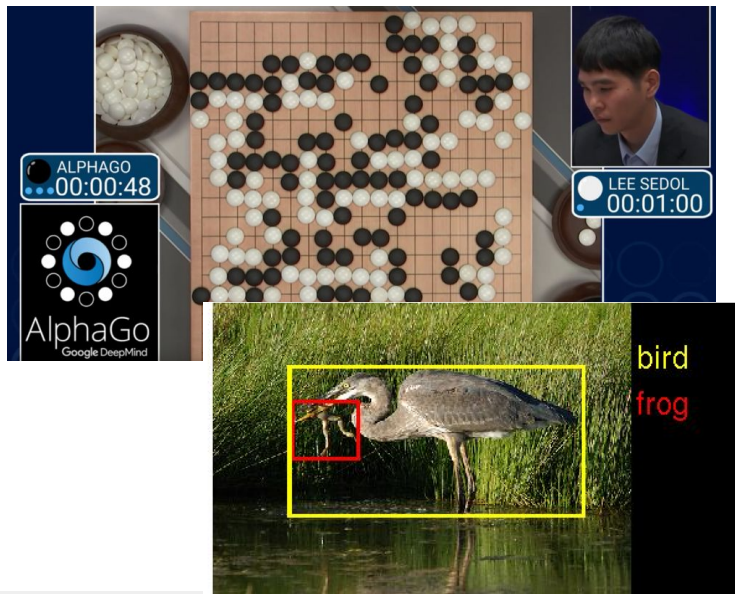
What Is Machine Learning

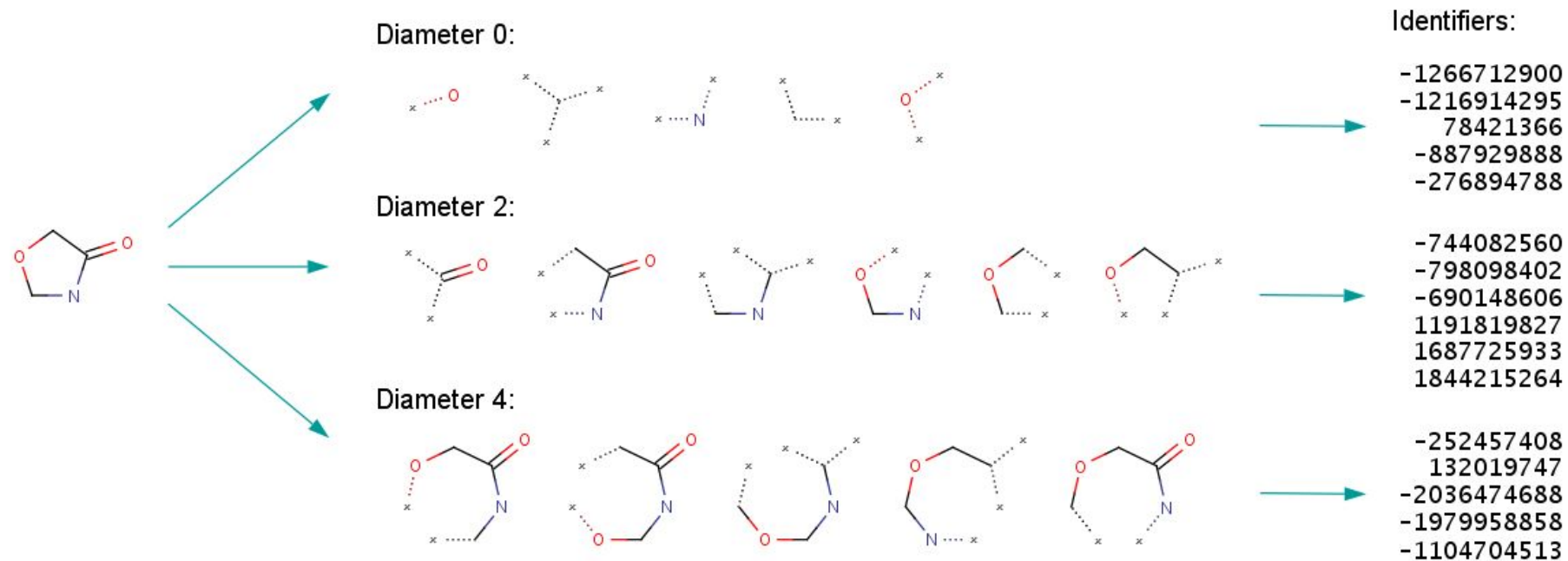
- Supervised Learning
 - Generate Functions From Data



Deep Learning

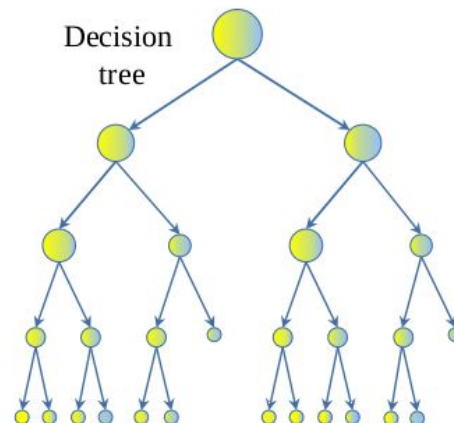
- Deep learning methods are becoming very popular in image recognition, game playing, and question and answer systems.





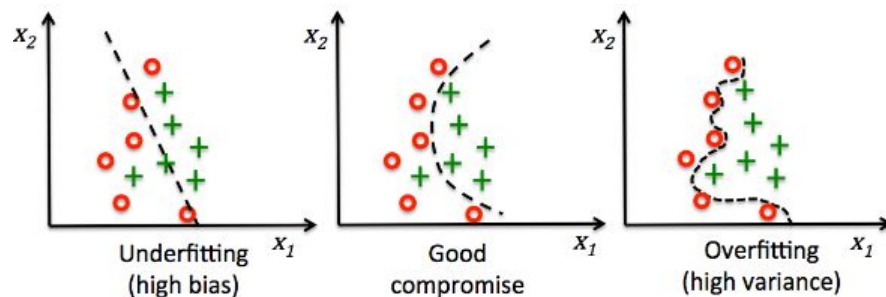
Decision Trees

- Based on the data “split” on a cut off of a single feature
- Can use the most informative feature of all the samples
- Leaf nodes hold predicted values



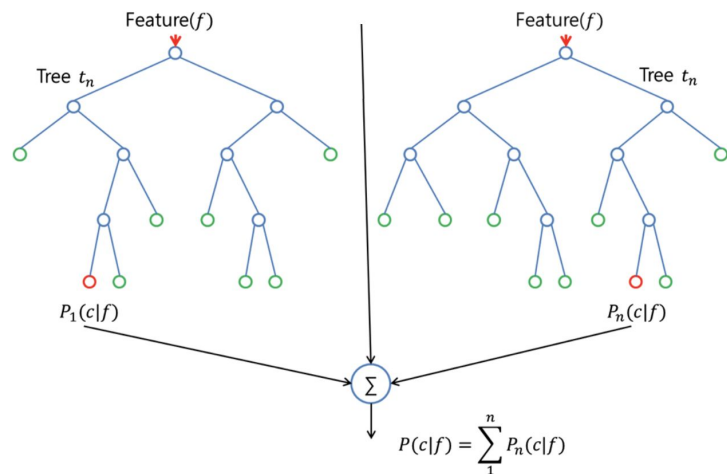
High Variance Low Bias

- Standard decision trees are prone to overfitting
- When there is “noise” in the response with respect to the input we need a more general model



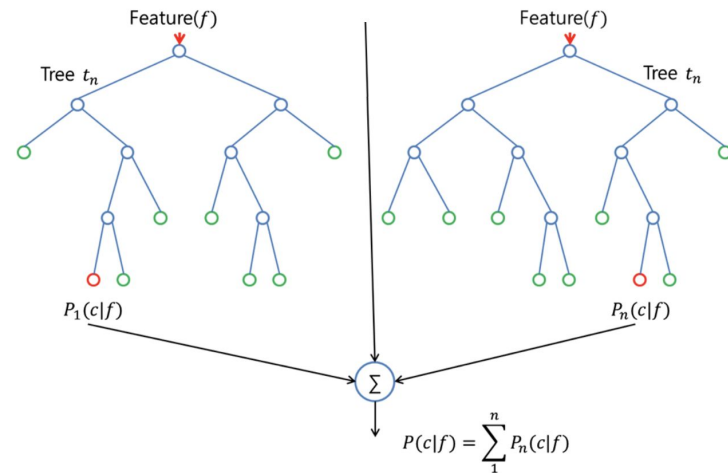
Bootstrap Aggregating (Bagging)

- Make many decision trees!
 - Each one on a subset of the training data, selected uniformly random WITH replacement
- Average results from all decision trees
- More robust to outliers

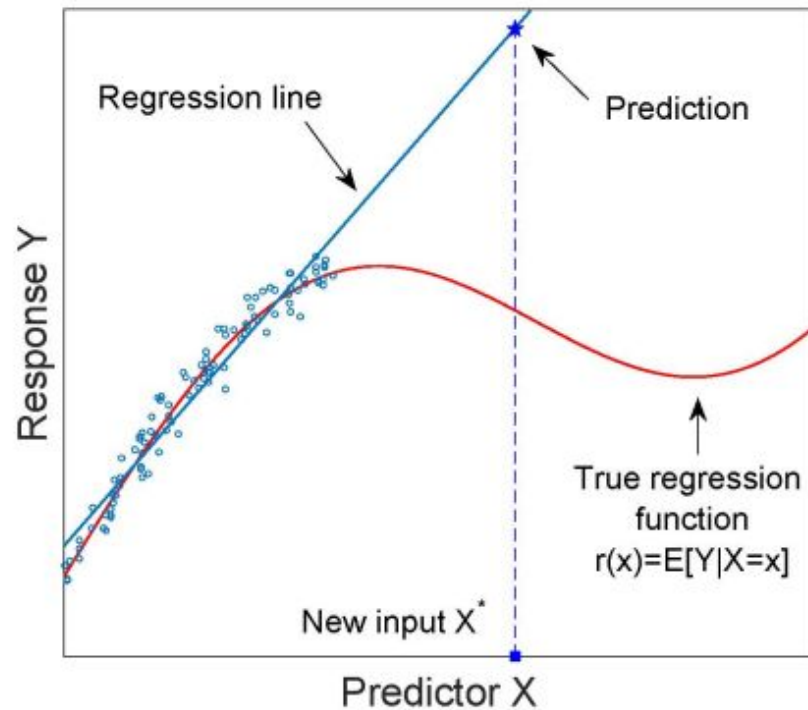


Feature Bagging (Random Forests)

- Many of the trees can be very similar
- E.X if one feature is very predictive all trees use this feature
 - No longer have good ensembling to lower variance
- Solution: At each split only select from a random subset of features



Interpolation

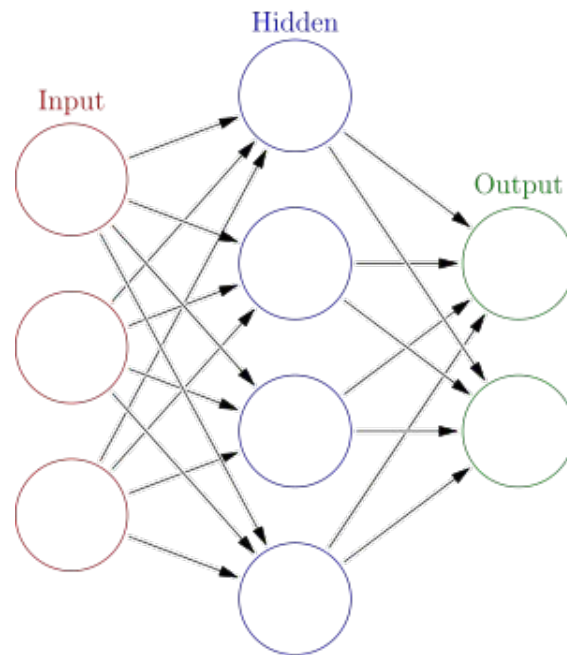


Deep Learning

- Lots of excitement to try to use these methods in other contexts
- Should deep learning be used in materials?
- Where does it provide the greatest benefit?

Artificial Neural Network Overview

- Collection of units called neurons (Circles Here)
- Each neuron computes a function over its inputs (real numbers)
- Each neuron can be connected to multiple outputs
- Trained using back propagation

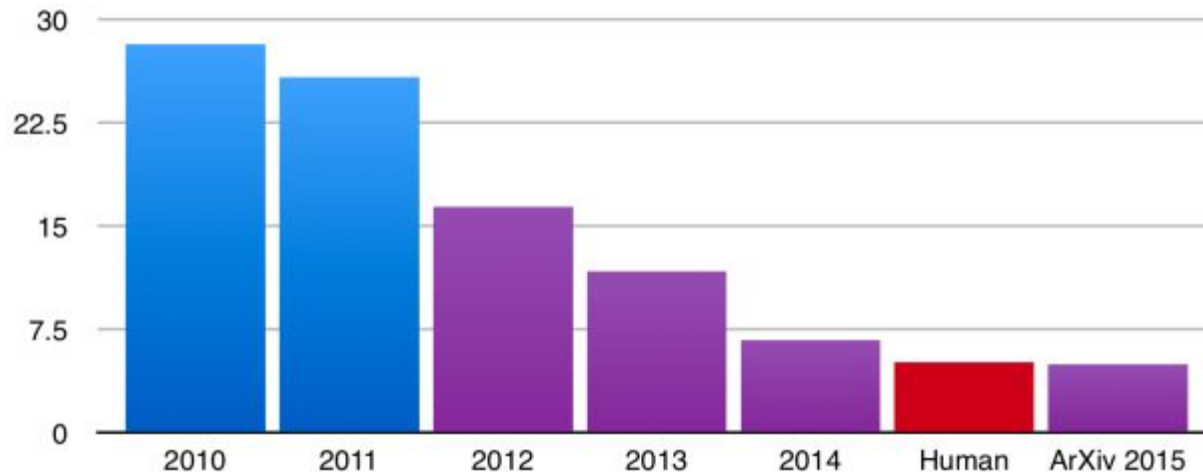


Universal Function Approximation Theorem

- Artificial Neural Networks can represent ANY function
- This does not pan out in practice
 - Limited data and compute power
- Requires us to create data and compute efficient models.

Deep Neural Network Image Classification

ImageNet Large Scale Visual Recognition Challenge Model Accuracy



As of 2015, a 27 layer DNN was more accurate than a human (Stanford student) at sorting 100,000 images into 1,000 different pre-specified categories

Deep Neural Network Image Classification

- The ImageNet classification challenge is very difficult:



Ruler



King crab



Sidewinder



Salt shaker

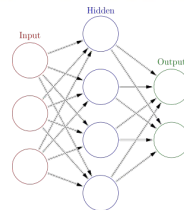
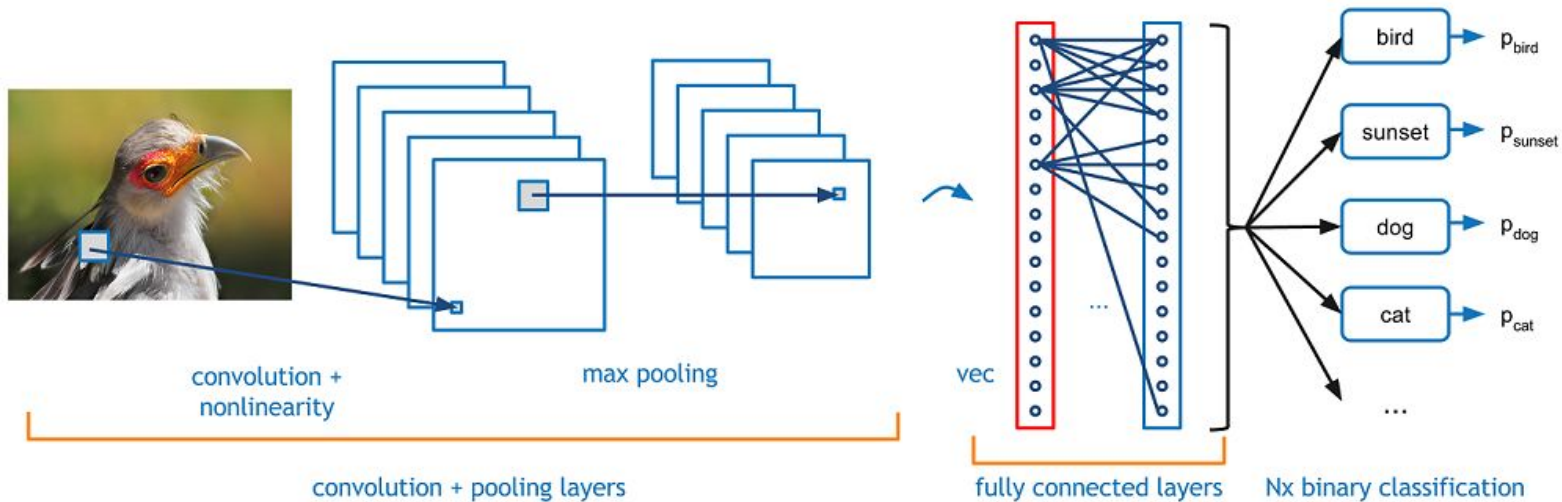


Reel



Hatchet

Convolutional Neural Networks



Convolution Layer

- Slide a learnable mask across the image.

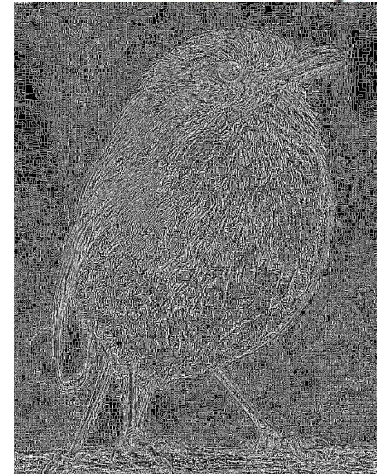
Input image



Convolution
Kernel

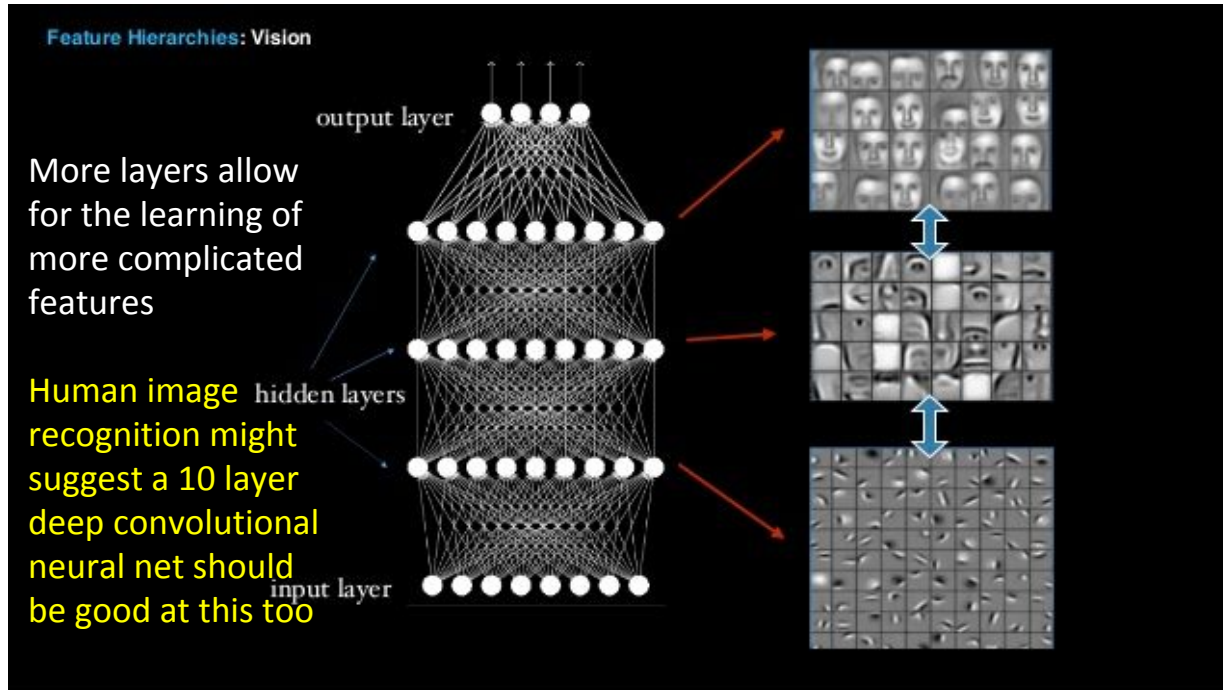
$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

Feature map



Deep Neural Network Image Classification

- A unique aspect of Deep Learning is the ability learn new features as the network is trained:



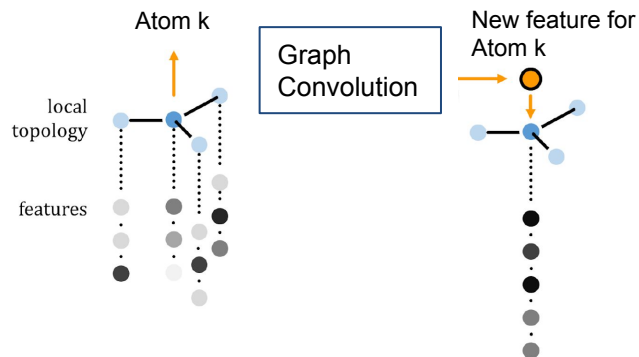
AutoQSAR w/ DeepChem Feature Generation

2D Graphic description of molecules

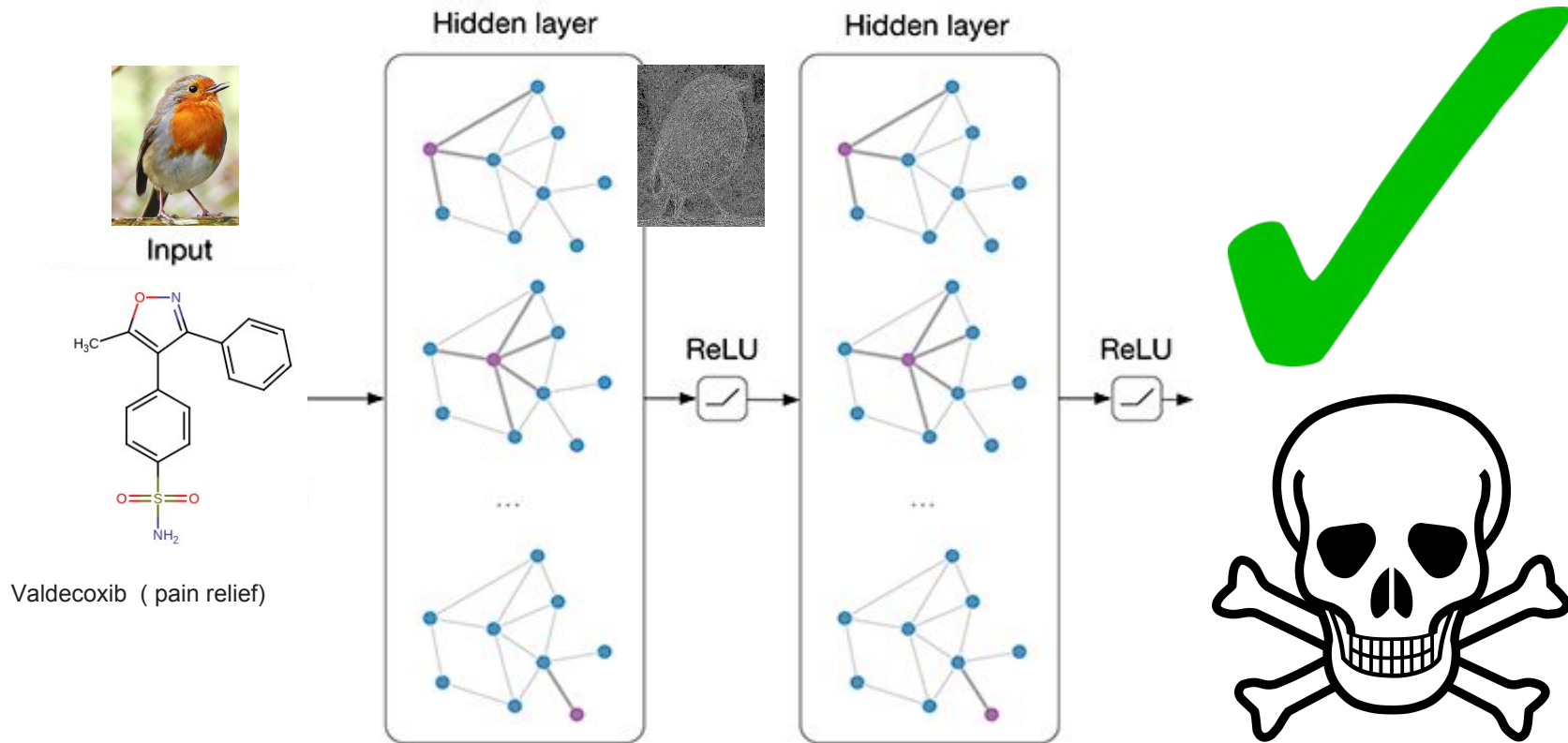
- Each node represents an atom
- Each edge represents a bond
- Atom features include atoms-type, valences, formal charges, and hybridization

Graph Convolution

- **Automatically learn new local features that suit the endpoint**
- These new features are then converted to molecular feature which is feed to dense neural network for model building



Graph Convolutions

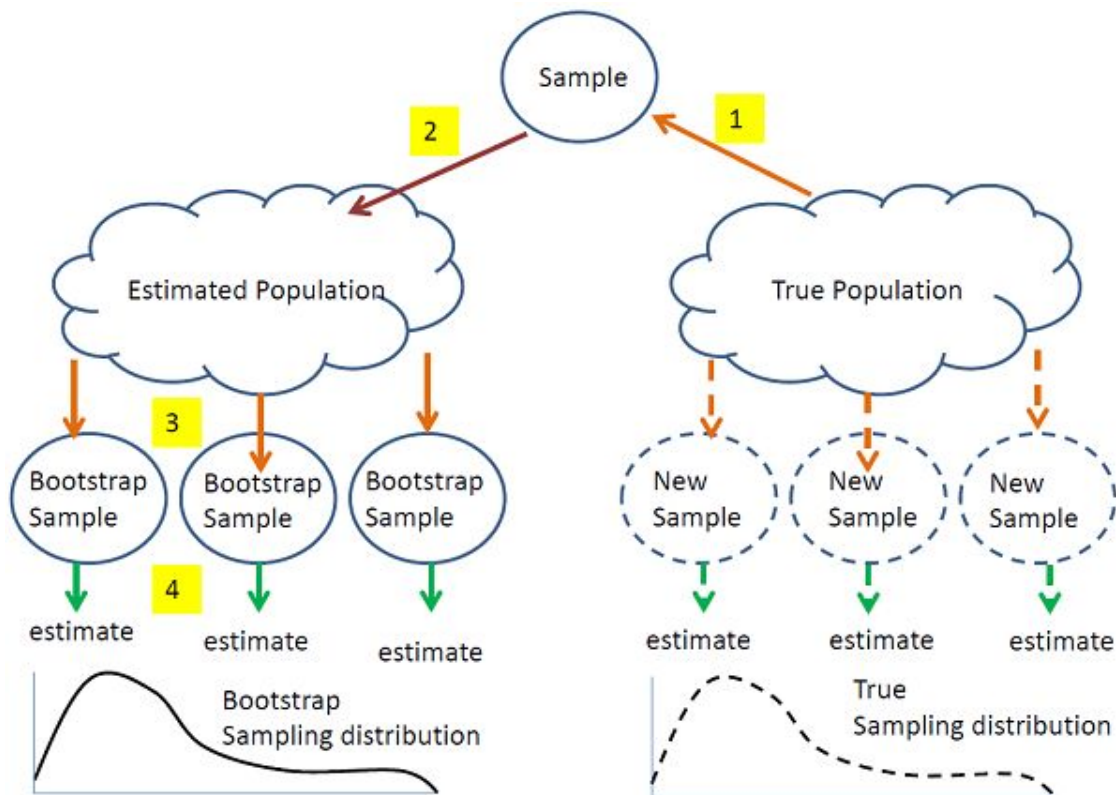


Can We Detect Bad/Interesting Labels In Chemical Data?

08/28/2018

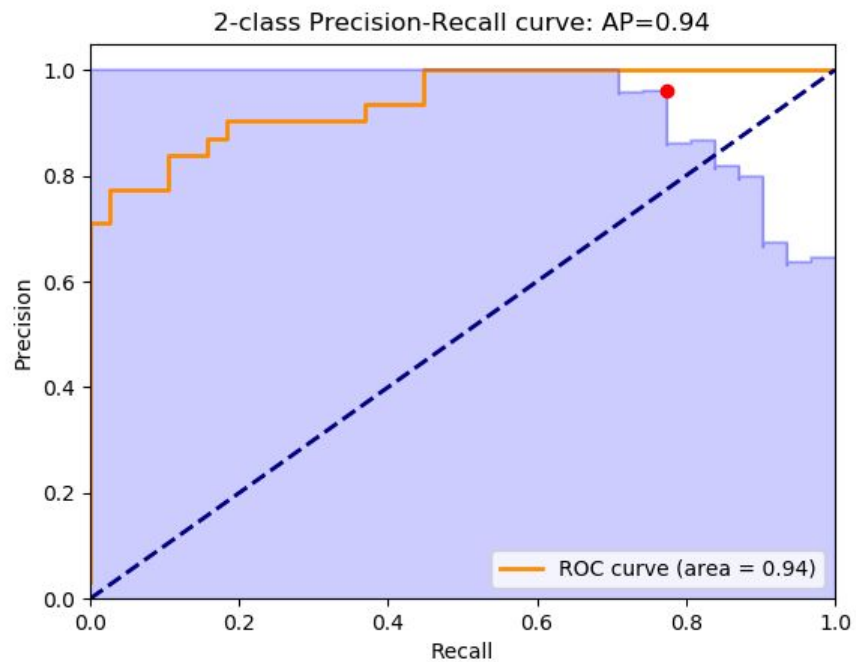
Bootstrapping Interesting Datapoint Identification

- Repeatedly train a model on a random 80% of the data
- Predict on remaining 20%
- Find samples whose predictions are farthest from labels

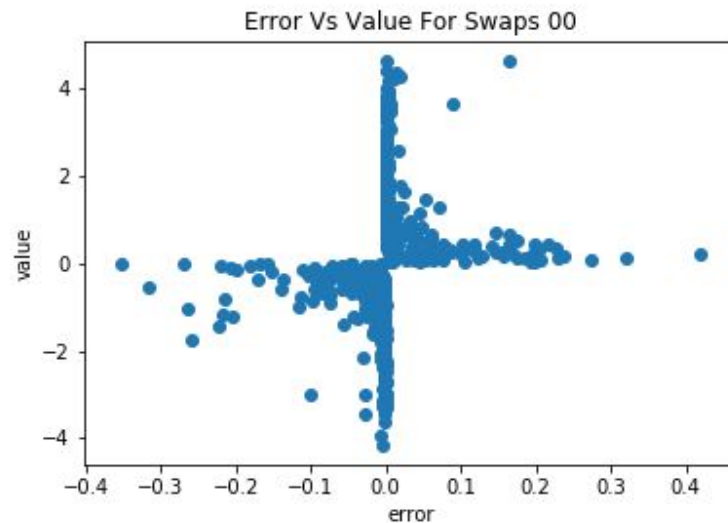
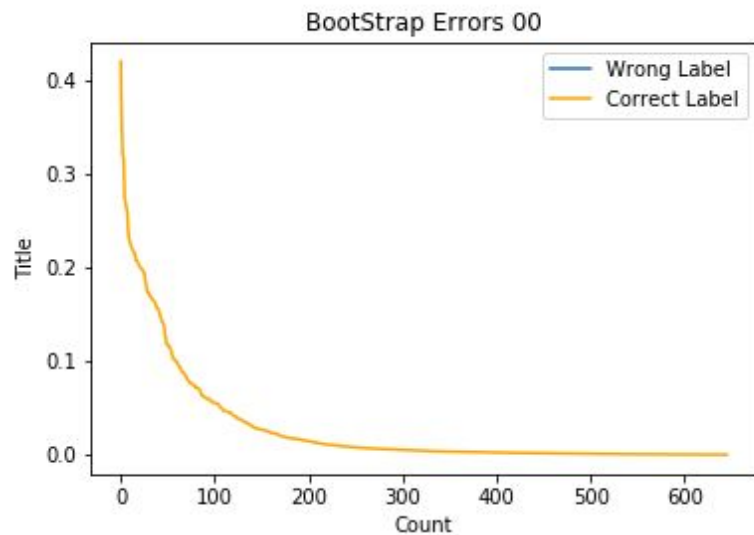


- Classification Model **Descriptors.MolLogP(m) > 0**
- We will randomly incorrectly label x% [0,10,20,50] of compounds and see if we can find the molecules we incorrectly labeled

No Flips



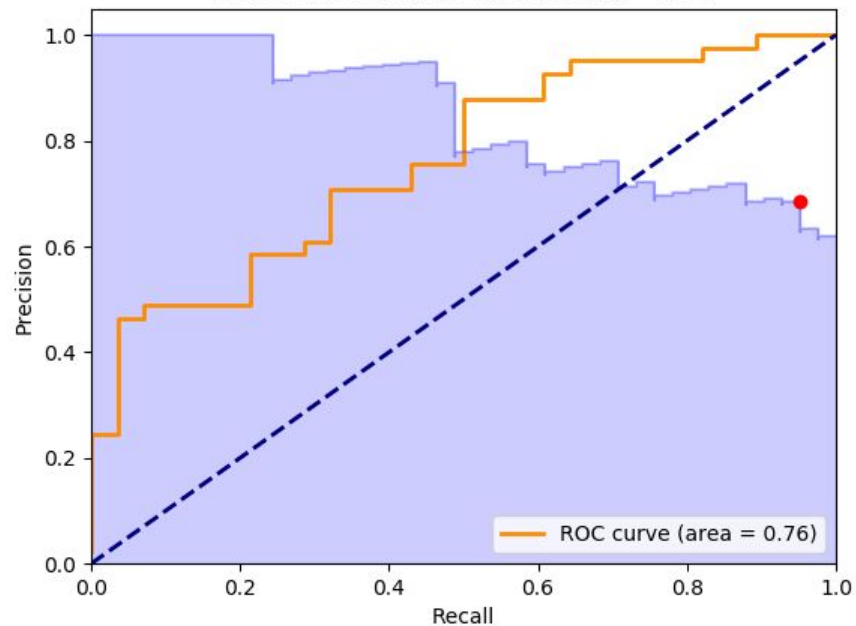
No Flips



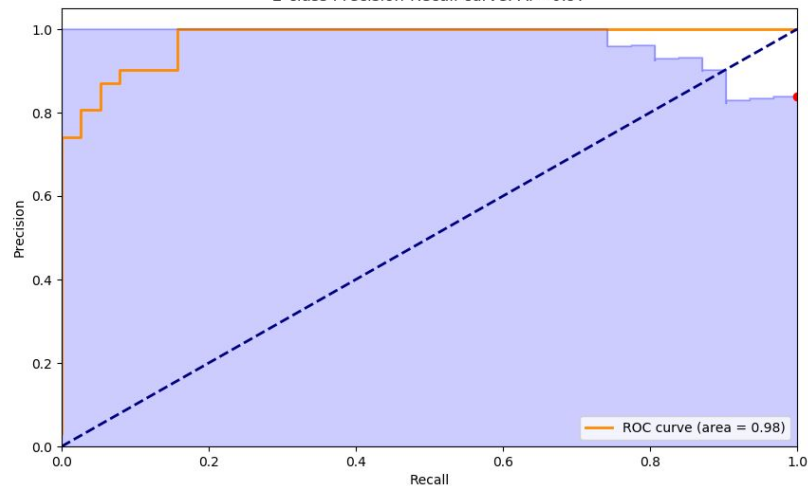
10% Flips

Errored Holdout

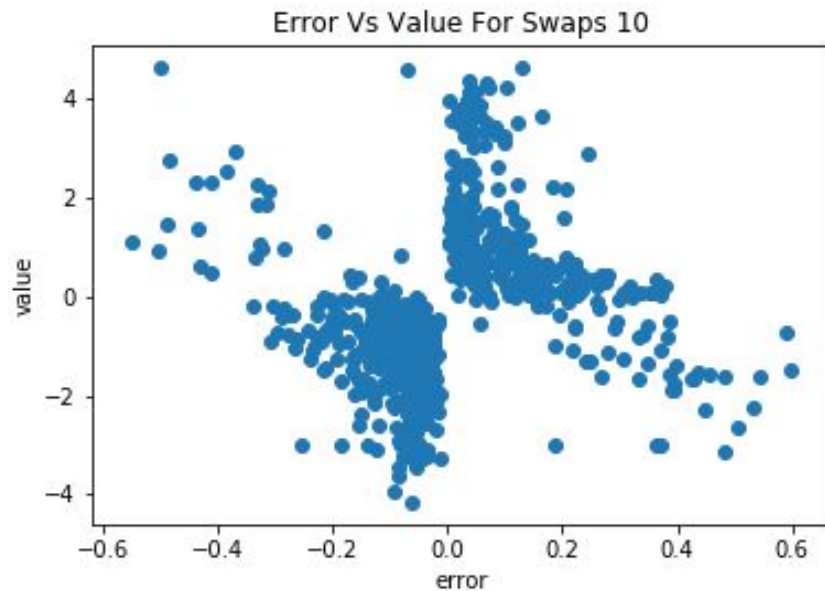
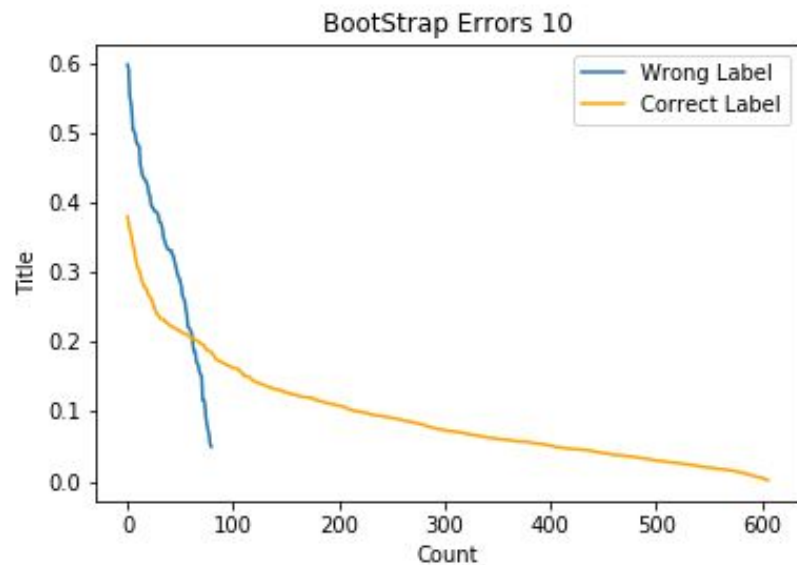
2-class Precision-Recall curve: AP=0.84



2-class Precision-Recall curve: AP=0.97

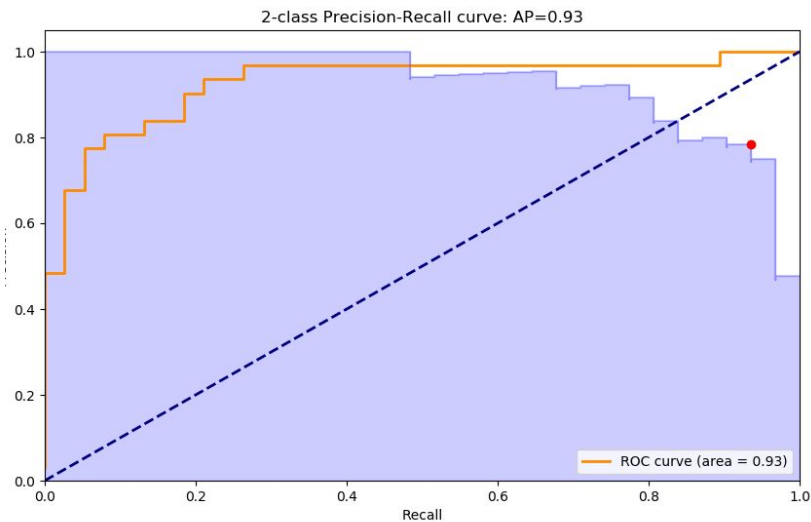
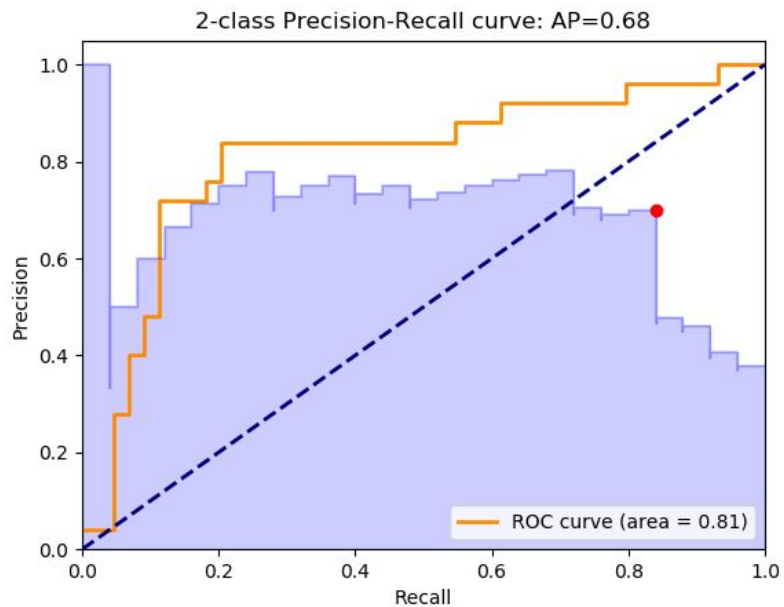


10% Flips

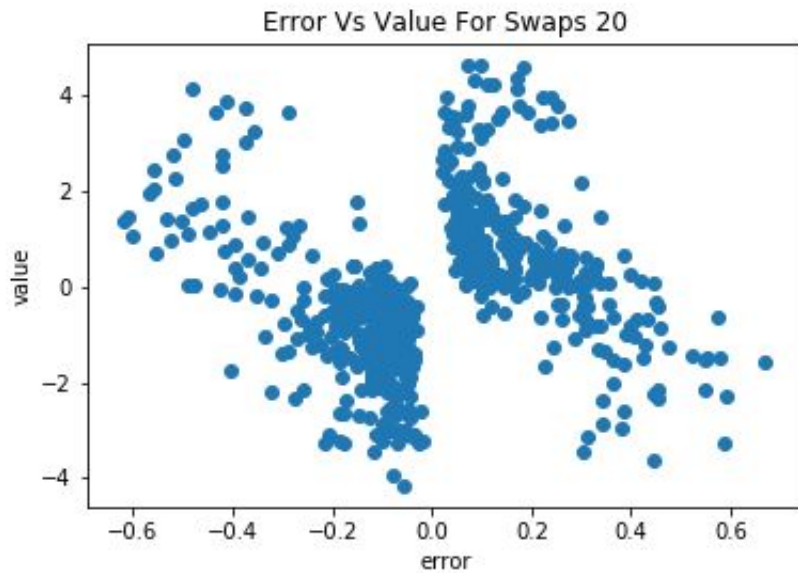
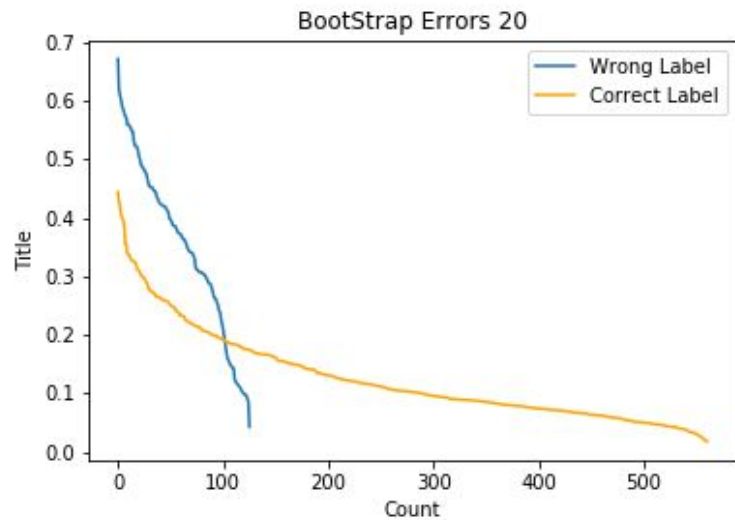


20% Flips

Errored Holdout

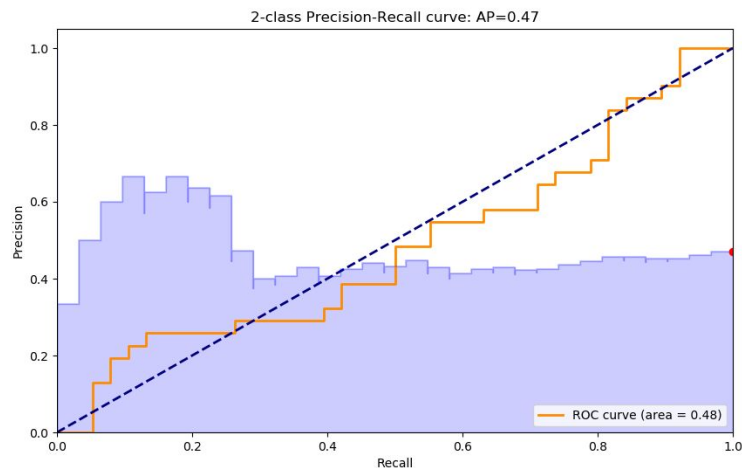
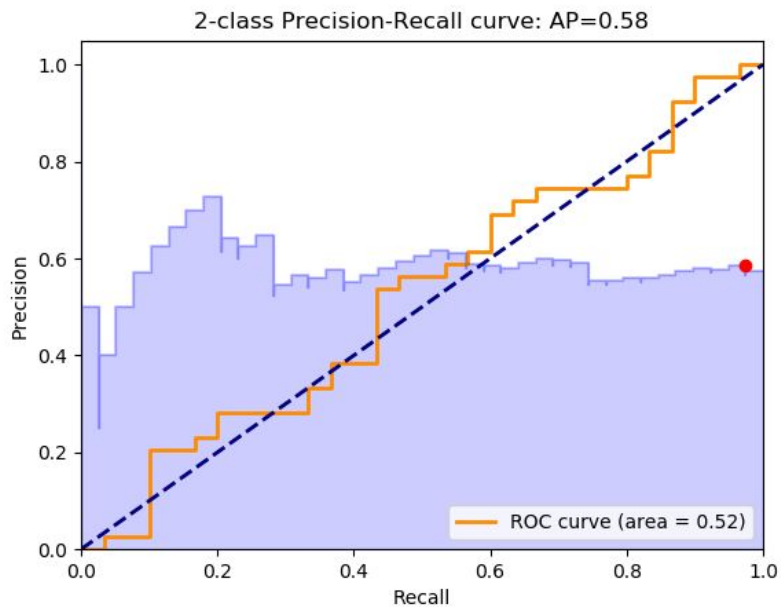


20% Flips

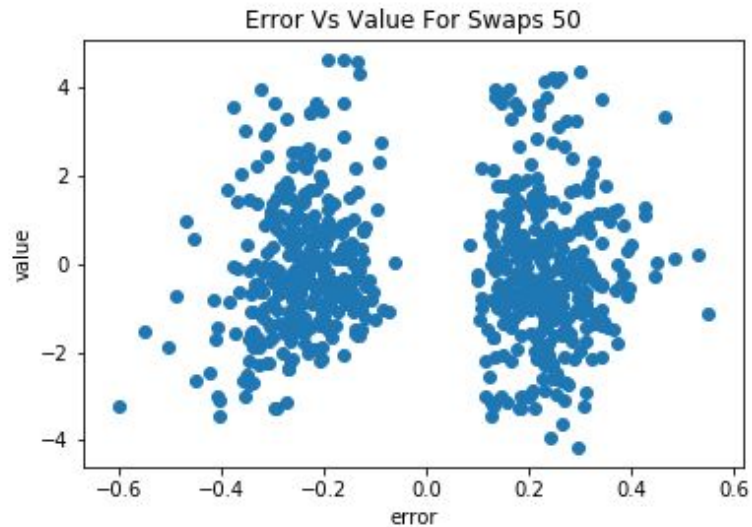
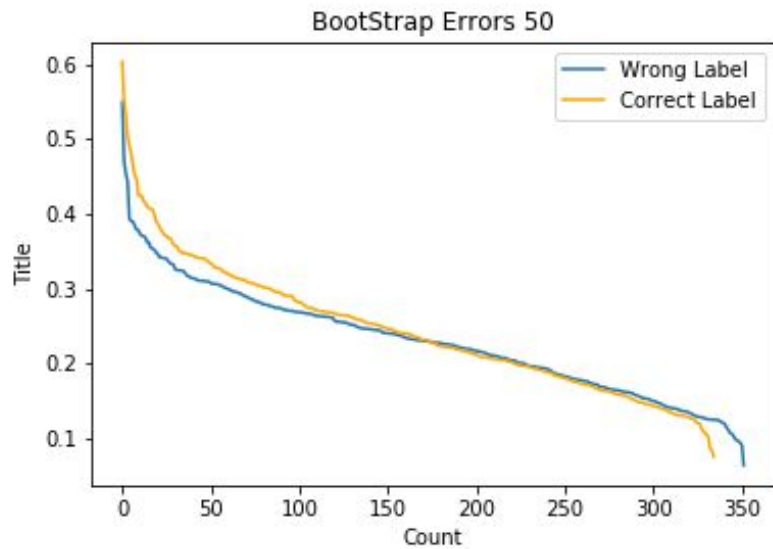


50% Flips

Errored Holdout



50% flips



Interested in Learning More?

- Book is in pre-release looking for feedback!
- <https://www.facebook.com/groups/1362916627160962/>
 - Facebook Group
- <https://gitter.im/deepchem/Lobby>

