

## CHEN 5802: Guest lecture on Active Learning



#### bartel.cems.umn.edu

Chemical Engineering & Materials Science

University of Minnesota

Nathan J. Szymanski



Existing data  $(x_i, y_i)$ 











## What you've learned so far in this course









Domain (without data)



- We have some domain over which new data (*x<sub>i</sub>*, *y<sub>i</sub>*) can be collected
- In practice, data collection tends to be costly and time consuming





**Density Functional Theory** 



Required time: **min to hours** 





## An example: ML trained on DFT calculations



#### **Density Functional Theory**



Required time: **min to hours** 

Machine learning potentials



Required time: **seconds** 





## An example: ML trained on DFT calculations



#### **Density Functional Theory**



#### Required time: min to hours



#### Machine learning potentials



#### Required time: **seconds**

- We can train on existing DFT calculations, but these tend be biased in their sampling of certain chemistries...
- New DFT calculations can be run to supplement the training data



#### An example: autonomous (AI-driven) experiments





• The A-Lab tries to make new materials whose recipes are unknown



#### An example: autonomous (Al-driven) experiments





- The A-Lab tries to make new materials whose recipes are unknown
- But *thousands* of synthesis recipes are often possible, and each recipe requires **hours to days** of experiments → we <u>cannot</u> test them all!





#### Domain (without data)



- We have some domain over which new data (*x<sub>i</sub>*, *y<sub>i</sub>*) can be collected
- In practice, data collection tends to be costly and time consuming
- In which parts of the domain should we collect training data? The goal is to achieve good model performance with limited data



## We can iterate between data collection with model training



Suppose we are collecting data to train a **classification** model in 2D (x-y) space

Initial dataset





## We can iterate between data collection with model training



Suppose we are collecting data to train a **classification** model in 2D (x-y) space



- Based on this initial training data, our model predicts a boundary to separate the two classes
- But there is likely a high amount of *uncertainty* in this prediction owing to limited data



## We can iterate between data collection with model training



Suppose we are collecting data to train a **classification** model in 2D (x-y) space



- With more data, the boundary changes
- This prediction should have **less uncertainty**





Suppose we are collecting data to train a **classification** model in 2D (x-y) space





Class 1
Class 2

## We can iterate between data collection with model training

- Ideally, we could focus our data collection on areas near the decision boundary
- For this, we need the *model uncertainty* as a function of the domain space





#### First, let's distinguish epistemic vs. aleatoric uncertainty





#### **Aleatoric uncertainty:**

Due to inherent randomness (noise) in the data

#### **Epistemic uncertainty:**

Due to a lack of knowledge, usually caused by the absence of training data



#### First, let's distinguish epistemic vs. aleatoric uncertainty

![](_page_17_Picture_1.jpeg)

![](_page_17_Figure_2.jpeg)

![](_page_17_Picture_3.jpeg)

![](_page_17_Picture_4.jpeg)

## Method 1: ensemble methods for <u>classification</u>

![](_page_18_Picture_1.jpeg)

![](_page_18_Figure_2.jpeg)

![](_page_18_Figure_3.jpeg)

These models are each trained **independently**, either with **different initializations** or on **different parts of the training data** 

![](_page_18_Picture_5.jpeg)

![](_page_18_Picture_6.jpeg)

## Method 1: ensemble methods for <u>classification</u>

![](_page_19_Picture_1.jpeg)

![](_page_19_Figure_2.jpeg)

**Class 1** is predicted most often (3/5 times)  $\rightarrow$  60% confidence Higher confidence  $\leftrightarrow$  Lower uncertainty

![](_page_19_Picture_4.jpeg)

## Method 1: ensemble methods for <u>regression</u>

![](_page_20_Picture_1.jpeg)

![](_page_20_Figure_2.jpeg)

We can use the standard deviation

as a measure of uncertainty

$$\boldsymbol{\sigma} = \sqrt{\frac{\sum_{i} (\boldsymbol{y}_{i} - \boldsymbol{\mu})^{2}}{N}}$$

 $y_i$  = individual predictions

 $\mu$  = mean of the predictions

N = number of predictions (this should be as large as possible)

![](_page_20_Picture_10.jpeg)

## Method 1: ensemble methods for <u>regression</u>

![](_page_21_Picture_1.jpeg)

![](_page_21_Figure_2.jpeg)

# We can use the standard deviation

as a measure of uncertainty

$$\boldsymbol{\sigma} = \sqrt{\frac{\sum_{i} (\boldsymbol{y}_{i} - \boldsymbol{\mu})^{2}}{N}}$$

#### In this case:

Mean prediction = **6.43** Std. deviation = **0.56** 

![](_page_21_Picture_9.jpeg)

![](_page_22_Picture_1.jpeg)

![](_page_22_Figure_3.jpeg)

![](_page_22_Picture_5.jpeg)

![](_page_22_Picture_6.jpeg)

![](_page_23_Picture_1.jpeg)

![](_page_23_Figure_3.jpeg)

![](_page_23_Picture_6.jpeg)

![](_page_24_Picture_1.jpeg)

![](_page_24_Figure_3.jpeg)

![](_page_24_Picture_5.jpeg)

![](_page_24_Picture_6.jpeg)

![](_page_25_Picture_1.jpeg)

![](_page_25_Figure_3.jpeg)

![](_page_25_Picture_6.jpeg)

## **Building a Gaussian process requires a prior**

![](_page_26_Picture_1.jpeg)

A *prior* describes an initial set of assumptions or beliefs we have about our system

To build a Gaussian process, we'll need **two things** to define our prior:

![](_page_26_Picture_5.jpeg)

## **Building a Gaussian process requires a prior**

![](_page_27_Picture_1.jpeg)

A prior describes an initial set of assumptions or beliefs we have about our system

To build a Gaussian process, we'll need **two things** to define our prior:

The **mean** function:  $\mu(x)$ 

This defines the average value of all functions at each point in the domain space

Usually, we start with  $\mu(x) = 0$ and then modify it as more data is collected

![](_page_27_Picture_7.jpeg)

## **Building a Gaussian process requires a prior**

![](_page_28_Picture_1.jpeg)

A prior describes an initial set of assumptions or beliefs we have about our system

To build a Gaussian process, we'll need **two things** to define our prior:

The **mean** function:  $\mu(x)$ 

This defines the average value of all functions at each point in the domain space

Usually, we start with  $\mu(x) = 0$ and then modify it as more data is collected The **covariance** function: K(x, x')

This defines how function values at different points in the domain are correlated with one another

In simpler terms, it defines how "smooth" our functions are

![](_page_28_Picture_10.jpeg)

![](_page_28_Picture_11.jpeg)

#### **Commonly used for covariance: the radial basis function**

![](_page_29_Picture_1.jpeg)

$$K(x, x') = \exp\left(-\frac{|x - x'|}{2\sigma^2}\right)$$

• When x and x' are close to one another, the RBF kernel  $K(x, x') \rightarrow 1$ , which means f(x) and f(x') are highly correlated and should exhibit similar values

![](_page_29_Picture_4.jpeg)

![](_page_29_Picture_5.jpeg)

## **Commonly used for covariance: the radial basis function**

![](_page_30_Picture_1.jpeg)

$$K(x, x') = \exp\left(-\frac{|x - x'|}{2\sigma^2}\right)$$

- When x and x' are close to one another, the RBF kernel  $K(x, x') \rightarrow 1$ , which means f(x) and f(x') are highly correlated and should exhibit similar values
- When x and x' are far from one another, the RBF kernel  $K(x, x') \rightarrow 0$ , which means f(x) and f(x') are **not** correlated

![](_page_30_Picture_6.jpeg)

## **Commonly used for covariance: the radial basis function**

![](_page_31_Picture_1.jpeg)

$$K(x,x') = \exp\left(-\frac{|x-x'|}{2\sigma^2}\right)$$

- When x and x' are close to one another, the RBF kernel  $K(x, x') \rightarrow 1$ , which means f(x) and f(x') are highly correlated and should exhibit similar values
- When x and x' are far from one another, the RBF kernel  $K(x, x') \rightarrow 0$ , which means f(x) and f(x') are **not** correlated
- $\sigma$  is a hyperparameter that measures the correlation length. Essentially, this controls how smooth we want our functions to be.

![](_page_31_Picture_6.jpeg)

![](_page_31_Picture_7.jpeg)

#### $\sigma$ is a user-chosen parameter that controls "smoothness"

![](_page_32_Picture_1.jpeg)

![](_page_32_Figure_2.jpeg)

12/20

![](_page_32_Picture_4.jpeg)

#### We can get uncertainty from the prediction variance

![](_page_33_Picture_1.jpeg)

![](_page_33_Figure_2.jpeg)

The variance is the square of the standard deviation:  $Var(x) = \sigma^2(x)$ 

![](_page_33_Picture_4.jpeg)

![](_page_33_Picture_5.jpeg)

## We can get uncertainty from the prediction variance

![](_page_34_Picture_1.jpeg)

![](_page_34_Figure_2.jpeg)

The variance is the square of the standard deviation:  $Var(x) = \sigma^2(x)$ 

**Large variance** exist in regions without much data  $\rightarrow$  high uncertainty!

![](_page_34_Picture_5.jpeg)

![](_page_34_Picture_6.jpeg)

## We can get uncertainty from the prediction variance

![](_page_35_Picture_1.jpeg)

![](_page_35_Figure_2.jpeg)

The variance is the square of the standard deviation:  $Var(x) = \sigma^2(x)$ 

Little to no variance around the known data → low uncertainty!

![](_page_35_Picture_6.jpeg)

## To improve model accuracy: sample high-uncertainty areas

![](_page_36_Picture_1.jpeg)

![](_page_36_Figure_2.jpeg)

If all we care about is model accuracy, then we should sample the part of the domain with the **highest uncertainty** 

![](_page_36_Picture_4.jpeg)

![](_page_36_Picture_5.jpeg)

## To improve model accuracy: sample high-uncertainty areas

![](_page_37_Picture_1.jpeg)

![](_page_37_Figure_2.jpeg)

![](_page_37_Picture_4.jpeg)

## To improve model accuracy: sample high-uncertainty areas

![](_page_38_Picture_1.jpeg)

![](_page_38_Figure_2.jpeg)

![](_page_38_Picture_4.jpeg)

## Now, what if we want to use our model for optimization?

![](_page_39_Picture_1.jpeg)

![](_page_39_Figure_2.jpeg)

## **Objective:** Find the value of x where f(x) is maximal

To accomplish this, we need to balance **exploration** with **exploitation** 

![](_page_39_Picture_5.jpeg)

![](_page_39_Picture_6.jpeg)

## Now, what if we want to use our model for optimization?

![](_page_40_Picture_1.jpeg)

![](_page_40_Figure_2.jpeg)

**Objective:** Find the value of x where f(x) is maximal

To accomplish this, we need to balance **exploration** with **exploitation** 

#### **Exploration:**

Collecting new training data in parts of the domain that are under-sampled and have high model uncertainty

This is basically what we've been prioritizing so far

![](_page_40_Picture_8.jpeg)

![](_page_40_Picture_9.jpeg)

## Now, what if we want to use our model for optimization?

![](_page_41_Picture_1.jpeg)

![](_page_41_Figure_2.jpeg)

#### **Objective:** Find the value of x where f(x) is maximal

To accomplish this, we need to balance *exploration* with *exploitation* 

#### **Exploration:**

Collecting new training data in parts of the domain that are under-sampled and have high model uncertainty

#### **Exploitation:**

Collecting new training data in parts of the domain where f(x) is expected to be optimal

![](_page_41_Picture_9.jpeg)

![](_page_41_Picture_10.jpeg)

## **Acquisition functions balance exploration/exploitation**

![](_page_42_Picture_1.jpeg)

Acquisition functions quantify the **anticipated benefit** of sampling a point (*x*)

![](_page_42_Picture_4.jpeg)

![](_page_43_Picture_1.jpeg)

Acquisition functions quantify the **anticipated benefit** of sampling a point (*x*)

**Upper Confidence Bound (UCB):** 

$$a(x,\beta) = \mu(x) + \beta\sigma(x)$$

Mean of the predicted function

Uncertainty of the predicted function

**β** is a hyperparameter which we can use to control the exploration-exploitation tradeoff

![](_page_43_Picture_8.jpeg)

![](_page_44_Picture_1.jpeg)

![](_page_44_Figure_2.jpeg)

 $\beta = 0$ 

#### **Pure exploitation**

![](_page_44_Picture_5.jpeg)

![](_page_45_Picture_1.jpeg)

![](_page_45_Figure_2.jpeg)

Sample the point corresponding to the maximum in the acquisition function

 $\beta = 0$ 

**Pure exploitation** 

![](_page_45_Picture_6.jpeg)

![](_page_45_Picture_7.jpeg)

![](_page_46_Picture_1.jpeg)

![](_page_46_Figure_2.jpeg)

 $\beta = 0$ 

**Pure exploitation** 

![](_page_46_Figure_5.jpeg)

 $\beta = 1$ 

#### **Some exploration**

![](_page_46_Picture_8.jpeg)

![](_page_47_Picture_1.jpeg)

![](_page_47_Figure_2.jpeg)

 $\beta = 2$ 

#### **More exploration**

![](_page_47_Figure_5.jpeg)

 $\beta = 4$ 

And even more exploration

![](_page_47_Picture_8.jpeg)

![](_page_47_Picture_9.jpeg)

#### **Putting it all together: Bayesian optimization!**

![](_page_48_Picture_1.jpeg)

![](_page_48_Figure_2.jpeg)

Widely used for optimization in the **physical sciences**, where experiments require a lot of time, money, and effort

From "Bayesian Hyperparameter Optimization" by Matti Karppanen

![](_page_48_Picture_6.jpeg)

![](_page_49_Picture_1.jpeg)

#### Surrogate models:

- Gaussian processes are most popular since they provide uncertainty
- But any ML model can be used neural networks are increasingly common

#### **Acquisition functions:**

• Many different ones exist: upper confidence bound, expected improvement, entropy search, and so on...All balance exploration/exploitation differently

#### **Prior:**

- The mean and covariance should be set to reflect your system
- For example, a periodic system should use periodic covariance (trig functions)

![](_page_49_Picture_10.jpeg)

![](_page_50_Picture_1.jpeg)

#### Surrogate models:

- Gaussian processes are most popular since they provide uncertainty
- But any ML model can be used neural networks are increasingly common

#### **Acquisition functions:**

• Many different ones exist: upper confidence bound, expected improvement, entropy search, and so on...All balance exploration/exploitation differently

**Prior:** 

- The mean and covariance should be set to reflect your system
- For example, a periodic system should use periodic covariance (trig functions)

![](_page_50_Picture_10.jpeg)

![](_page_51_Picture_1.jpeg)

#### Surrogate models:

- Gaussian processes are most popular since they provide uncertainty
- But any ML model can be used neural networks are increasingly common

#### **Acquisition functions:**

• Many different ones exist: upper confidence bound, expected improvement, entropy search, and so on...All balance exploration/exploitation differently

#### **Prior:**

- The mean and covariance should be set to reflect your system
- For example, a periodic system should use periodic covariance (trig functions)

![](_page_51_Picture_10.jpeg)

![](_page_51_Picture_11.jpeg)