

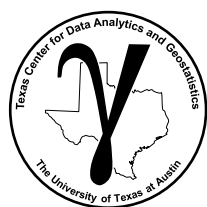
# **PGE 383 Machine Learning**

## **Machine Learning**

**Lecture outline . . .**

- **Machine Learning Overview**
- **Examples of Machine Learning**
- **Energy Machine Learning**

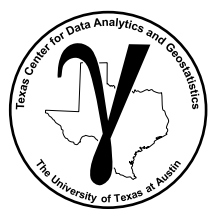
**Michael Pyrcz, The University of Texas at Austin**



# Motivation

Learn the concepts common to a variety of machine learning approaches:

- Inference and prediction
- Training and testing
- Parameters and hyperparameters



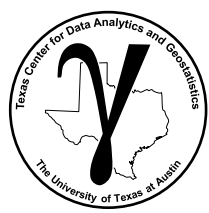
# **PGE 383 Machine Learning**

## **Machine Learning**

**Lecture outline . . .**

- **Machine Learning Overview**

**Michael Pyrcz, The University of Texas at Austin**



# Big Data

Big Data, you have big data if your data has a combination of these:

**Volume:** many data samples, difficult to handle and visualize

**Velocity:** high rate collection, continuous relative to decision making cycles

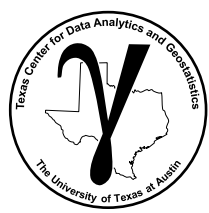
**Variety:** data form various sources, with various types and scales

**Variability:** data acquisition changes during the project

**Veracity:** data has various levels of accuracy

*“Energy has been big data long before tech learned about big data.”*

– Michael Pyrcz

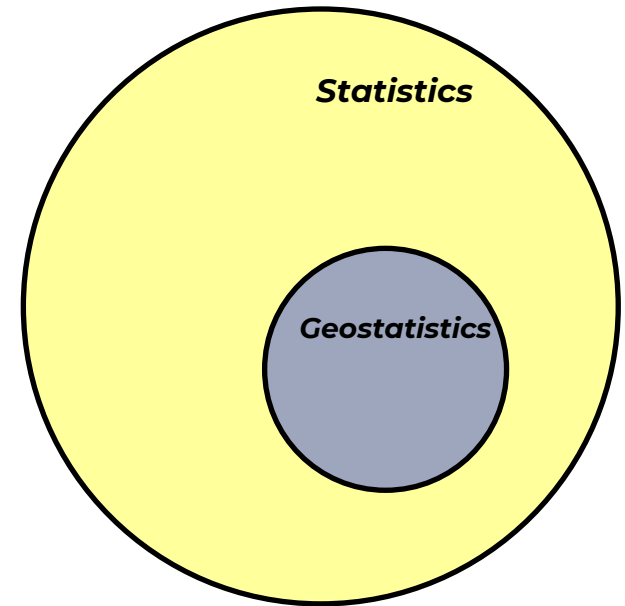


# Big Data Analytics

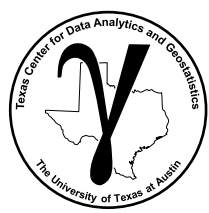
**Statistics** is collecting, organizing, and interpreting data, as well as drawing conclusions and making decisions.

**Geostatistics** is a branch of applied statistics:

1. the spatial (geological) context
2. the spatial relationships
3. volumetric support
4. uncertainty



Proposed Venn diagram for statistics and geostatistics.



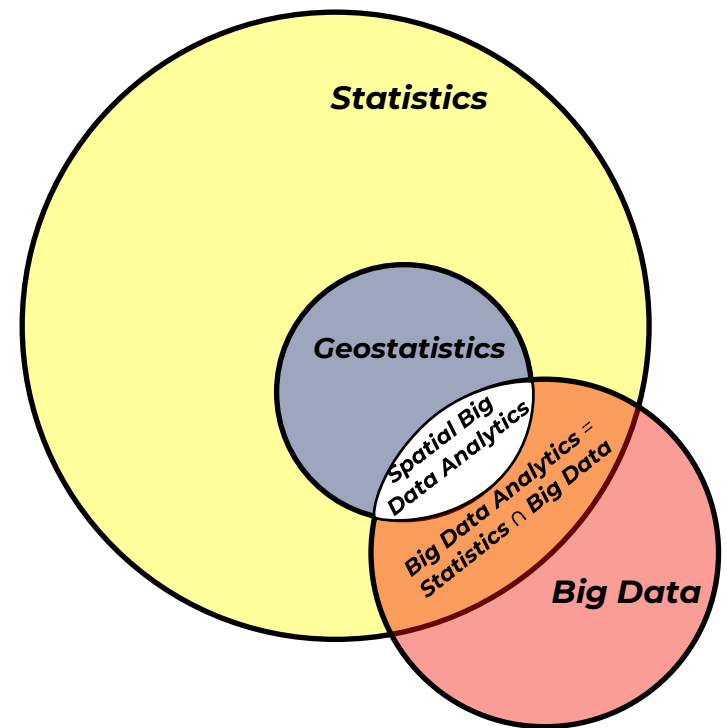
# Big Data Analytics

**Data Analytics** is the analysis of data to support decision making.

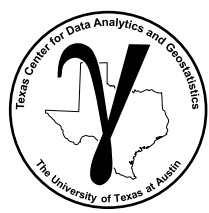
**Big Data Analytics** is the process of examining large and varied data sets to discover patterns and make decisions.

**Spatial Big Data Analytics** is expert use of spatial statistics / geostatistics on big data to support decision making.

*‘Data analytics is the use of statistics and visualization’*

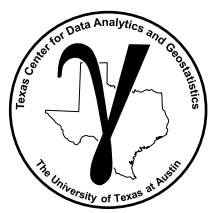


Proposed Venn diagram for spatial big data analytics.



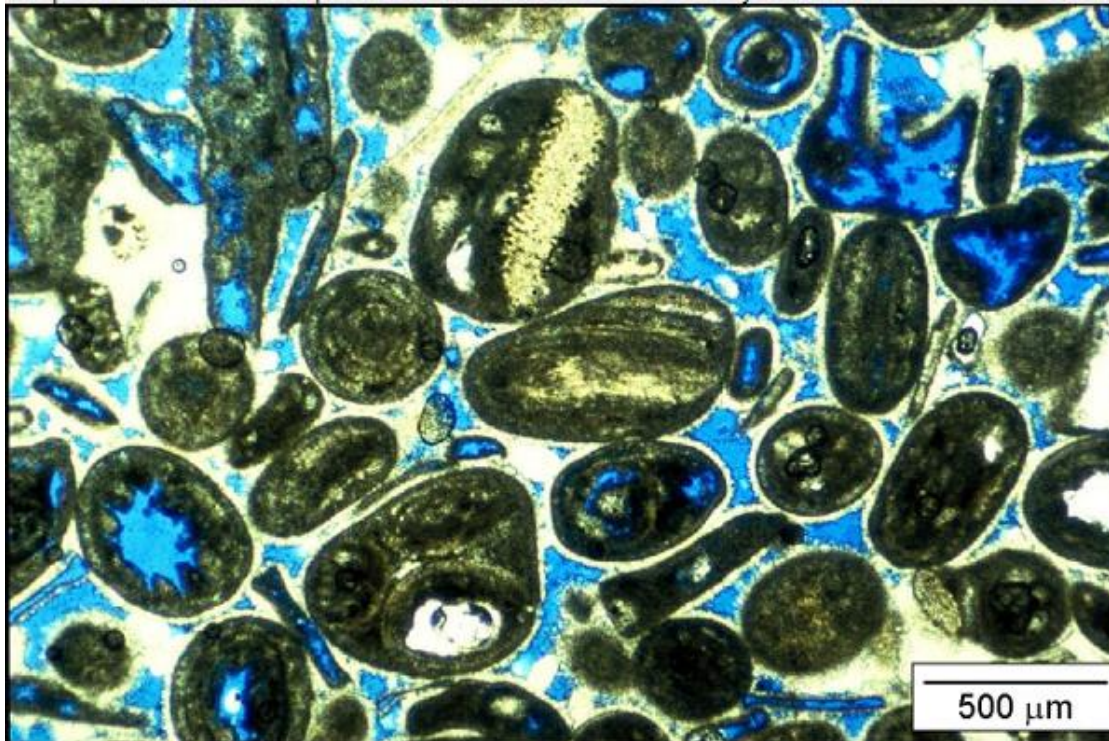
# Machine Learning

**learning** → “... is the study of algorithms and mathematical models **toolkit** →  
that computer systems use to  
progressively improve their performance on a specific task.  
Machine learning algorithms build a mathematical model  
of sample data, known as "training data", **training with data** →  
**general** → in order to make predictions or decisions  
without being explicitly programmed to perform the task.”  
“... where it is **not a panacea** →  
infeasible to develop an algorithm of specific instructions for performing the task.”



# Variables / Features

- **Variable / Feature:** any property measured / observed in a study
  - e.g. porosity, permeability, mineral concentrations, saturations, contaminant concentration
  - in data mining / machine learning this is known as a **feature**
  - measure often requires **significant analysis, interpretation** etc.

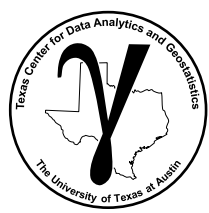


Total Porosity  
all blue area

Effective Porosity  
all connected blue  
area

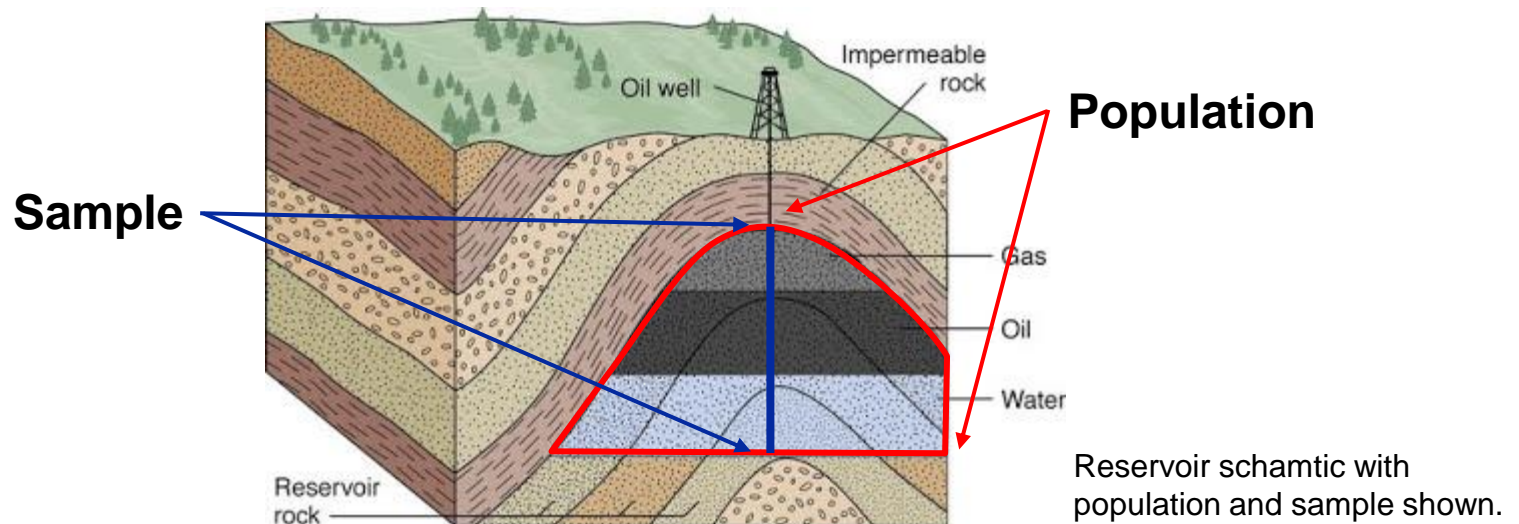
Carbonate thin section from  
BEG, UT Austin from course  
by F. Jerry Lucia.  
[http://www.beg.utexas.edu/lmo/d/\\_IOL-CM07/old-4.29.03/cm07-step05.htm](http://www.beg.utexas.edu/lmo/d/_IOL-CM07/old-4.29.03/cm07-step05.htm)

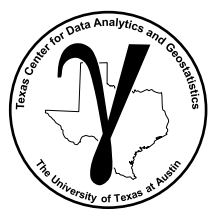




# Population and Sample

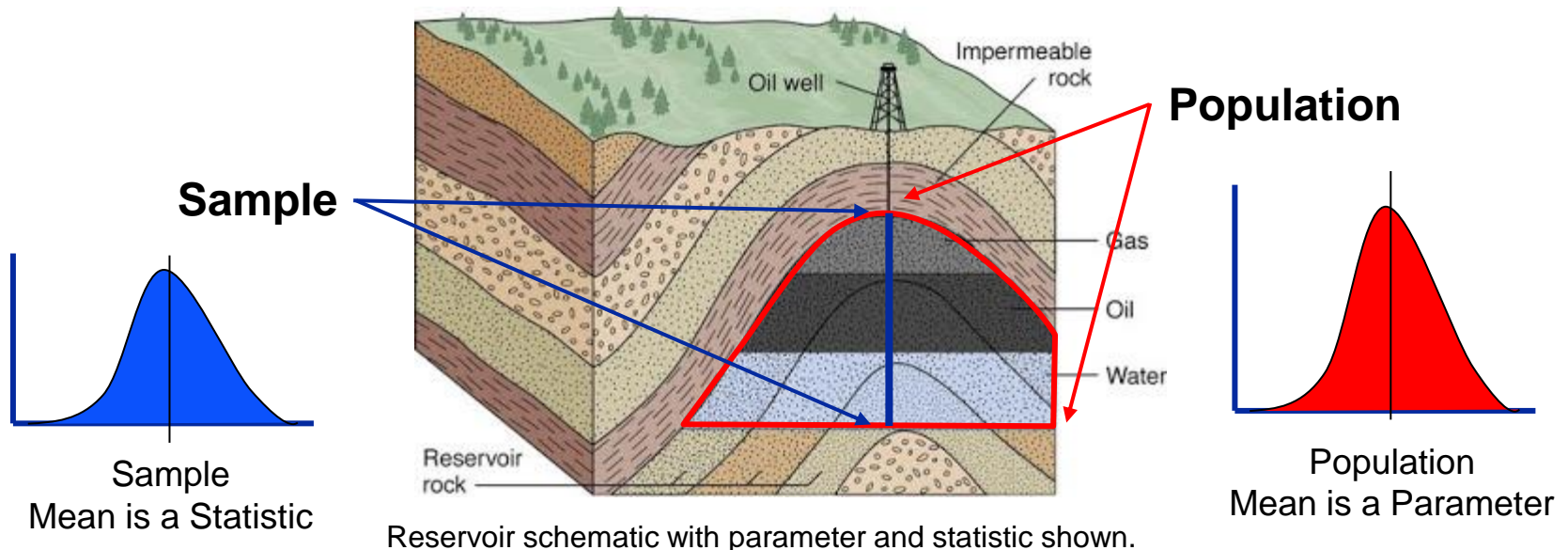
- **Population:** Exhaustive, finite list of property of interest over area of interest. Generally the entire population is not accessible.
  - e.g. exhaustive set of porosity at each location within a reservoir
- **Sample:** The set of values, locations that have been measured
  - e.g. porosity data from well-logs within a reservoir

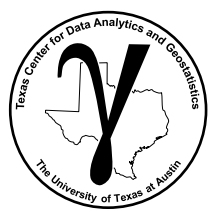




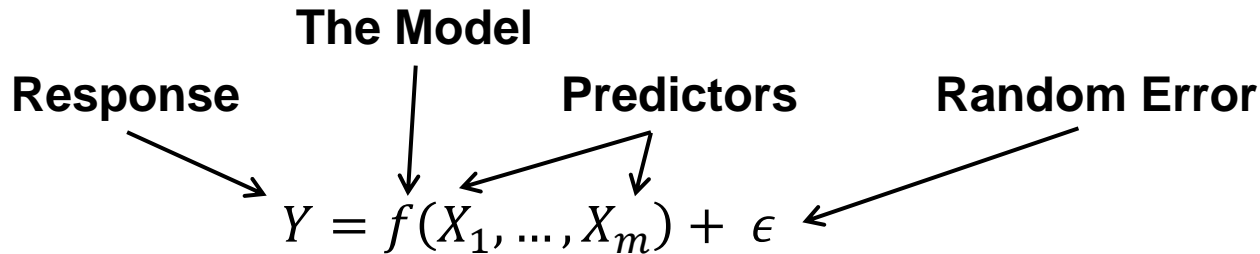
# Parameter and Statistic

- **Parameters:** summary measure of a population
  - e.g. population mean, population standard deviation, we rarely have access to this
  - **model parameters** is different in machine learning, and we will cover later.
- **Statistics:** summary measure of a sample
  - e.g. sample mean, sample standard deviation, we use statistics as estimates of the parameters





# Machine Learning Nuts and Bolts



Predictors (or Independent) Features (or Variables)

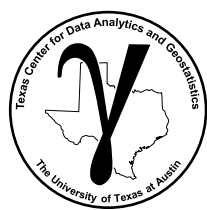
- the model inputs

Response (or Dependent) Features (or Variables)

- the model outputs

Machine Learning is All About Estimating the model,  $f$ , for two purposes:

- Inference or Prediction



# Inference

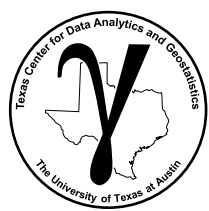
**What is the relationship between each predictor feature?**

$$f(X_1, \dots, X_m)$$

- sense of the relationship (positive or negative)?
- shape of relationship (sweet spots)?
- relationships may depend on values of other predictors!

## **Recall, Inferential Statistics**

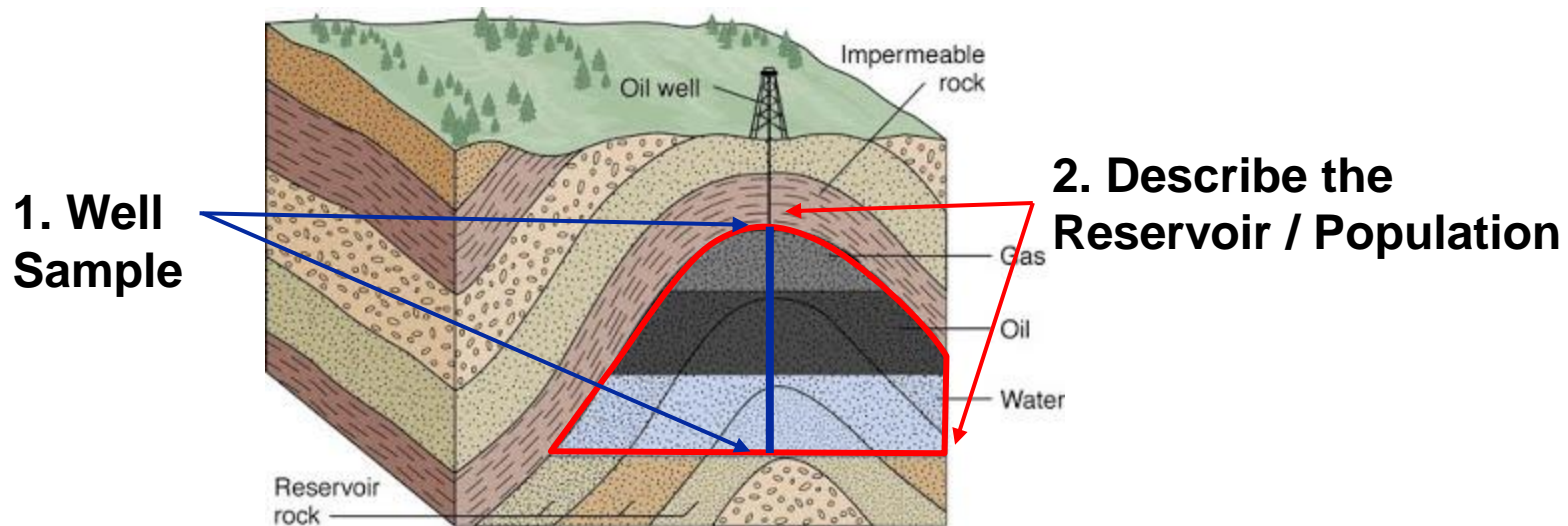
- given a sample, describe the population
- e.g. given 3 heads and 7 tails, what is the probability the coin is fair?



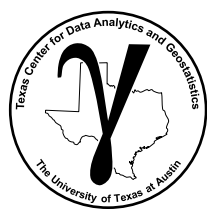
# Inference

- **Inferential Statistics**

- Given a random sample from a population, describe the population
- Given the well(s) samples, describe the reservoir



Reservoir schematic with inference problem, given well sample, describe the reservoir, population.



# Prediction

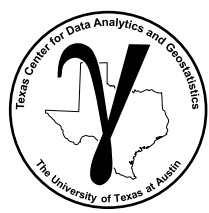
## The best prediction of the response feature

$$\hat{Y} = \hat{f}(X_1, \dots, X_m) + \epsilon$$

- Estimate the function,  $\hat{f}$ , for the purpose of predicting  $\hat{Y}$
- We are focused on getting the most accurate estimates,  $\hat{Y}$ , where  $\hat{Y}$  is an estimate of  $Y$

## Recall, Predictive Statistics

- given an assumption about the population, predict the outcome in the next sample
- e.g. given a fair coin what is the probability of 3 heads and 7 tails?

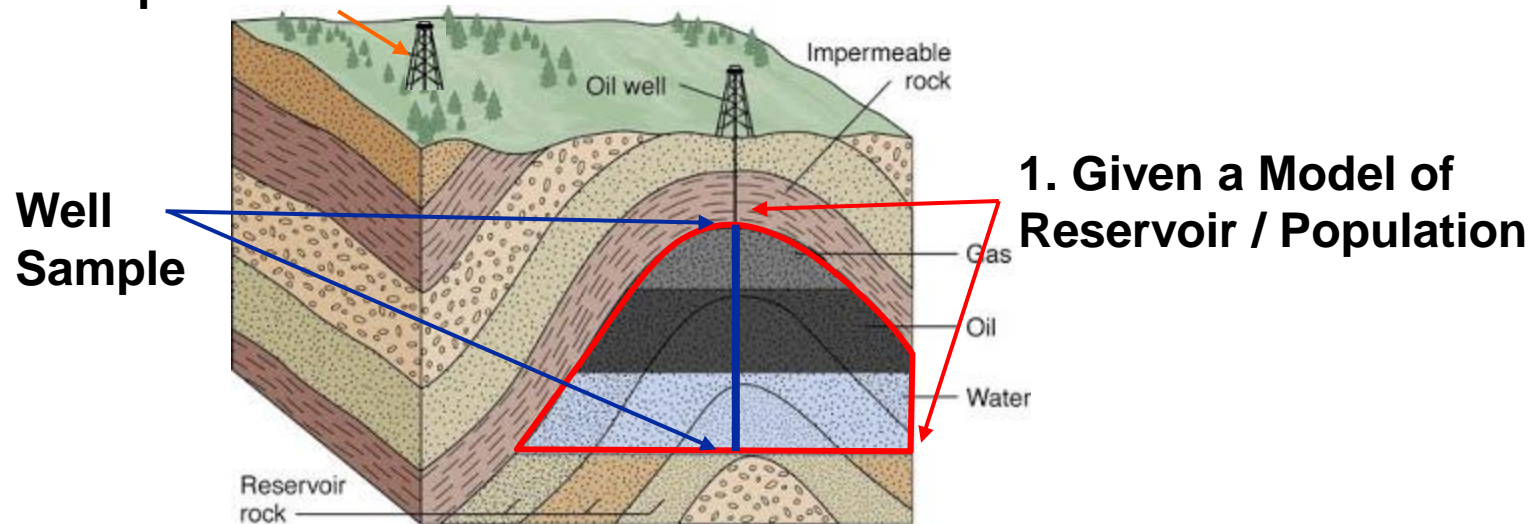


# Prediction

- **Predictive Statistics**

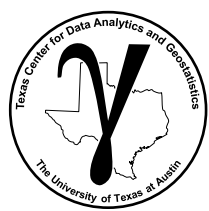
- Predict the samples given assumptions about the population
- Given our model of the reservoir, predict the next well (pre-drill assessment) sample, e.g. porosity, permeability, production rate etc.

## 2. Pre-Drill Prediction for Proposed Well



Reservoir schematic with inference problem, given well sample, describe the reservoir, population.





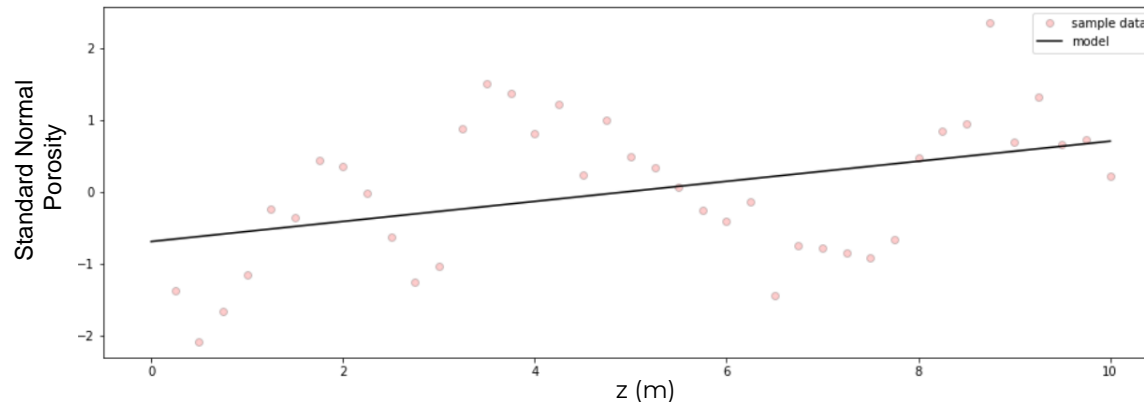
# Parametric Models

## Working with Parametric Models

- Makes an assumption about the functional form, shape
- We gain simplicity and advantage of only a few parameters
- For example, here is a linear model:

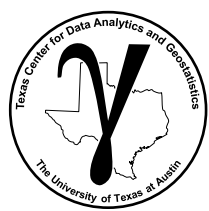
$$Y = f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_m X_m$$

There is a risk that  $\hat{f}$  is quite different than  $f$ , then we get a poor model!



Linear regression model to predict porosity from the z coordinate.



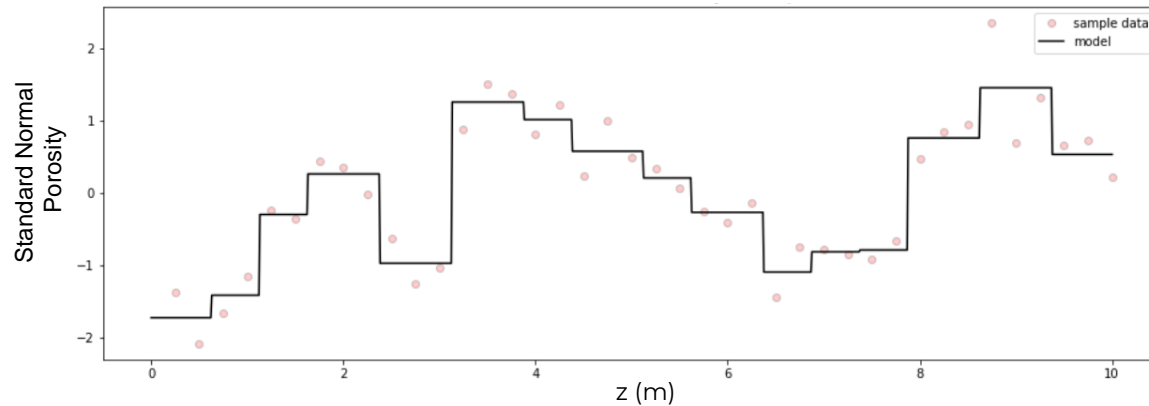


# Nonparametric Models

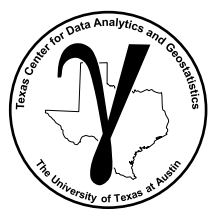
## Working with Nonparametric Models

- Makes no assumption about the functional form, shape
- More flexibility to fit a variety of shapes for  $f$
- Less risk that  $\hat{f}$  is a poor fit for  $f$
- Typically need a lot more data for an accurate estimate of  $f$

*‘Nonparametric is actually parametric rich!’*



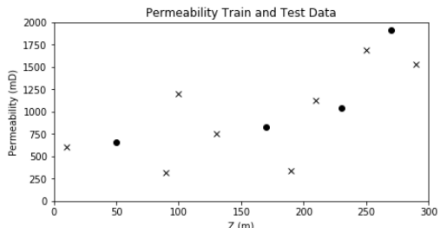
Decision tree regression model to predict porosity from the z coordinate.



# Model Workflow

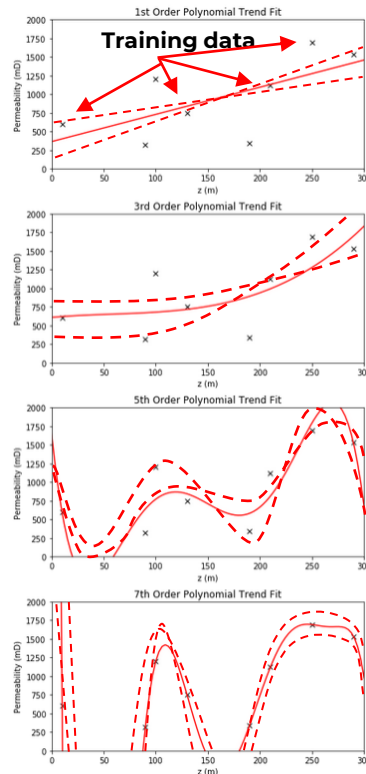
## Building Machine Learning Models

### 1. Split the Data into Train and Test Subsets



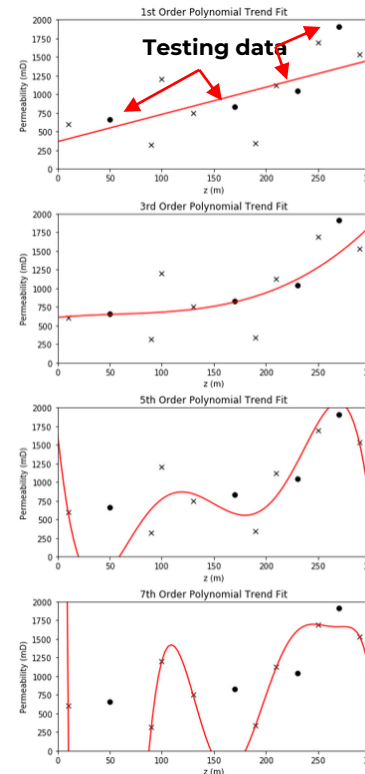
2. Build Models Over a Range of Hyperparameters  
Increasing Model Complexity

### 3. Train the Model Parameters with Training Data



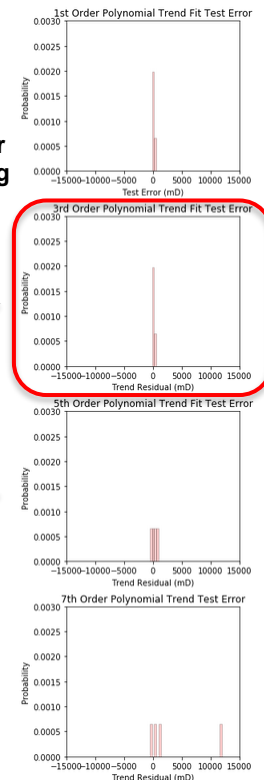
Best Fit Model to Training Data

### 4. Check Each Best Fit Model with Withheld Testing Data

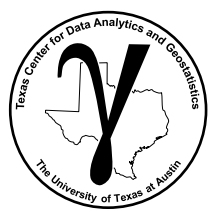


Error Over the Testing Data

### 5. Tune the Model Hyperparameters with Testing Data



Machine learning model building workflow to avoid overfit.



# Model Parameters

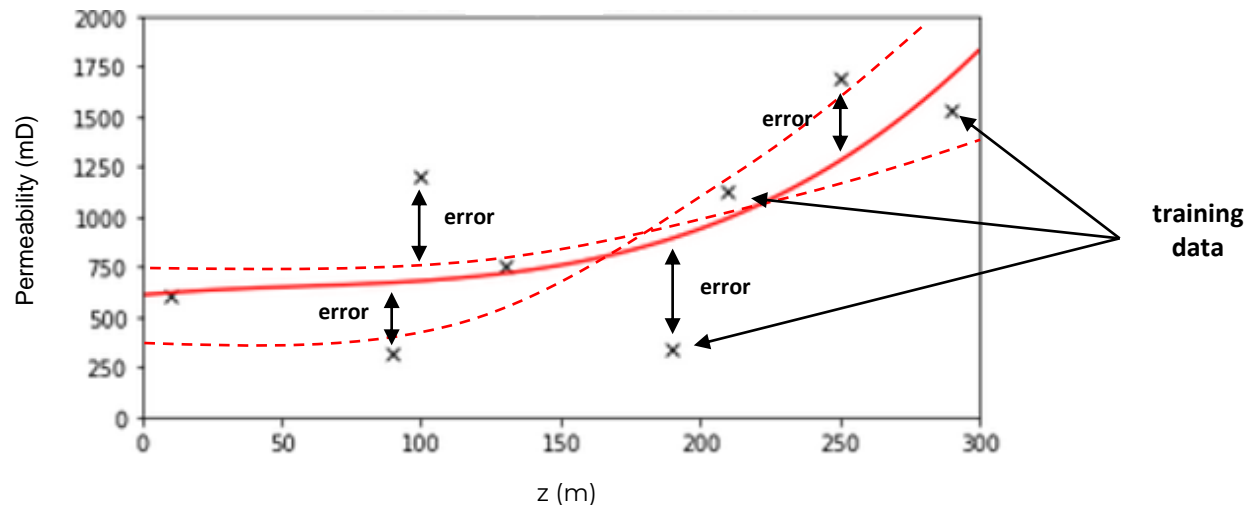
## Machine Learning Model Parameters

### Model Parameters

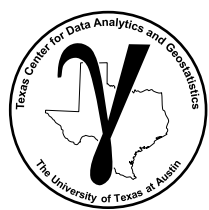
- Fit during training phase to minimize error at the training data
- For this 3<sup>rd</sup> order polynomial:

$$k = b_3 z^3 + b_2 z^2 + b_1 z + c$$

**Parameters:**  
 $b_3, b_2, b_1$  and  $c$



Setting model parameters to minimize the error relative to training data.



# Model Hyperparameters

## Machine Learning Model Hyperparameters

### Model Hyperparameters

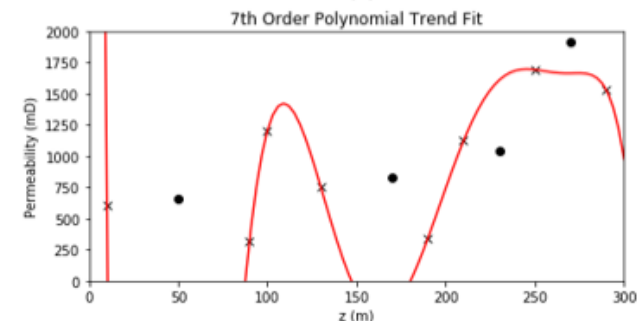
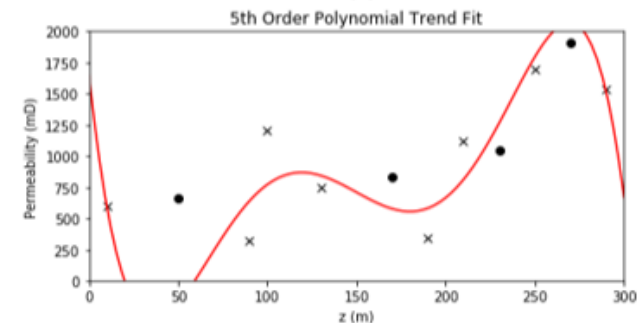
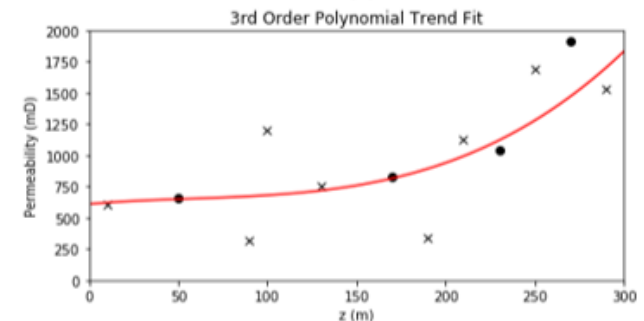
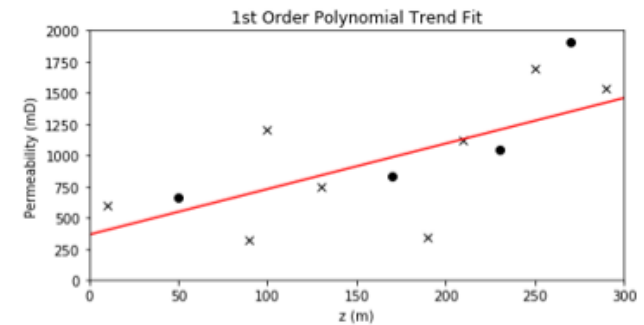
- Constrain the model complexity.
- Select hyperparameters that maximize accuracy with the testing data.
- For a polynomial model:

Increasing  
Complexity

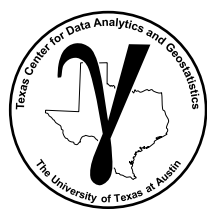
**3<sup>rd</sup> Order:**  $k = b_3z^3 + b_2z^2 + b_1z + c$

**2<sup>nd</sup> Order:**  $k = b_2z^2 + b_1z + c$

**1<sup>st</sup> Order:**  $k = b_1z + c$



Varying the model complexity, model hyperparameter, to maximize fit with testing data.



# Assessing Model Accuracy

## Method Selection is Important

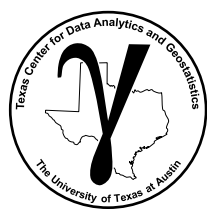
- No one method performs well on all datasets.
- Based on experience, understanding the data and limitations of the methods

## Measuring Quality of Fit in Training

- for regression, the most common measure is the mean square error

$$MSE = \frac{1}{n} \sum_{i=1}^n \left[ (y_i - \hat{f}(x_i^1, \dots, x_i^m))^2 \right] \quad \begin{array}{l} \text{for } i = 1, \dots, n \text{ training data and} \\ \text{for } 1, \dots, m \text{ features.} \end{array}$$

where we have  $n$  observations of training data for response  $y_i$ , and predictor  $x_i^1, \dots, x_i^m$  features.



# Assessing Model Accuracy

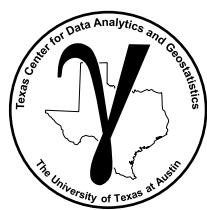
## Measuring Quality of Fit in Testing / Real-world Use

- The challenge is that that real question we have is how well can we predict outside the training data, testing data.

$$MSE = E \left[ (y_0 - \hat{f}(x_0^1, \dots, x_0^m))^2 \right] \quad \begin{array}{l} \text{for } i = 1, \dots, n \text{ training data and} \\ \text{for } 1, \dots, m \text{ features.} \end{array}$$

where we have observations of the response,  $y_0$ , and predictor features not used to train the model,  $x_0^1, \dots, x_0^m$ .

- Recall,  $E$  is the expectation. A probability weighted average, given equal probability the same as the average.
- We want to know how our model performs when we move away from the training data!



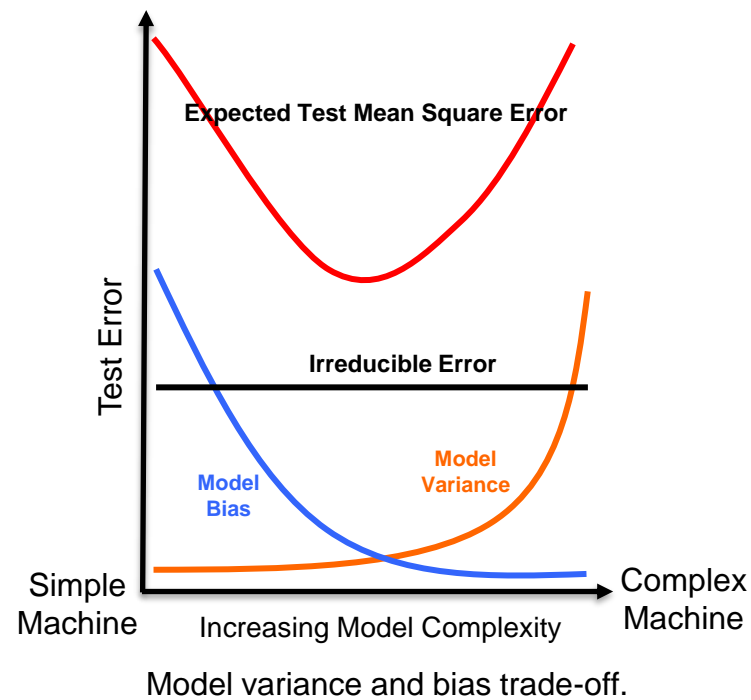
# Model Bias Variance Trade-Off

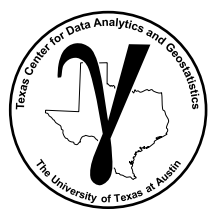
## The Components of Error in Testing / Real-world Use

The Expected Test Square Error components:

$$E \left[ (y_0 - \hat{f}(x_1^0, \dots, x_m^0))^2 \right] = \underbrace{\text{Var}(\hat{f}(x_1^0, \dots, x_m^0))}_{\text{Model Variance}} + \underbrace{[\text{Bias}(\hat{f}(x_1^0, \dots, x_m^0))]^2}_{\text{Model Bias}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}}$$

- **Model Variance** is error due to sensitivity to the dataset
- **Model Bias** is error due to using an approximate model
- **Irreducible Error** is due to missing variables and limited samples



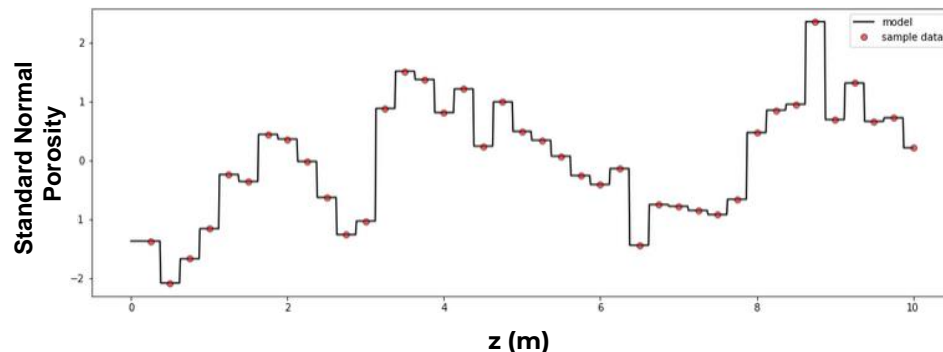
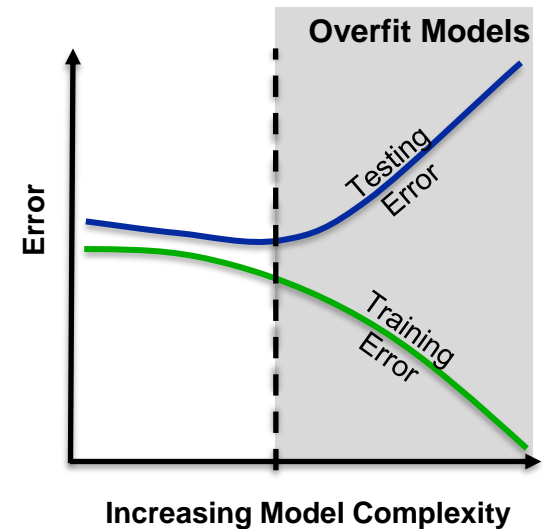


# Model Overfit

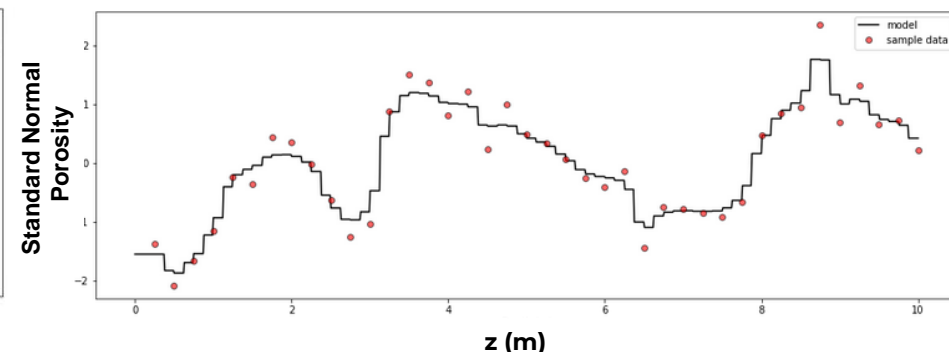
## Machine Learning Model Overfit

### Model Overfit

- Fitting data noise / data idiosyncrasies
- Increased complexity will generally decrease error with respect to the training dataset
- but, may result in increase error with testing data → at this complexity/flexibility we are overfit!

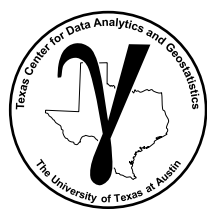


Overfit model to training data.



A more balanced fit model to training data.





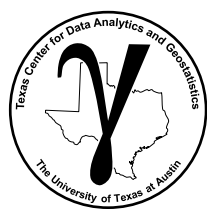
# Training and Testing

## The Training and Testing Split

- the most common approach is random selection
- this may not be fair testing

## Fair Testing

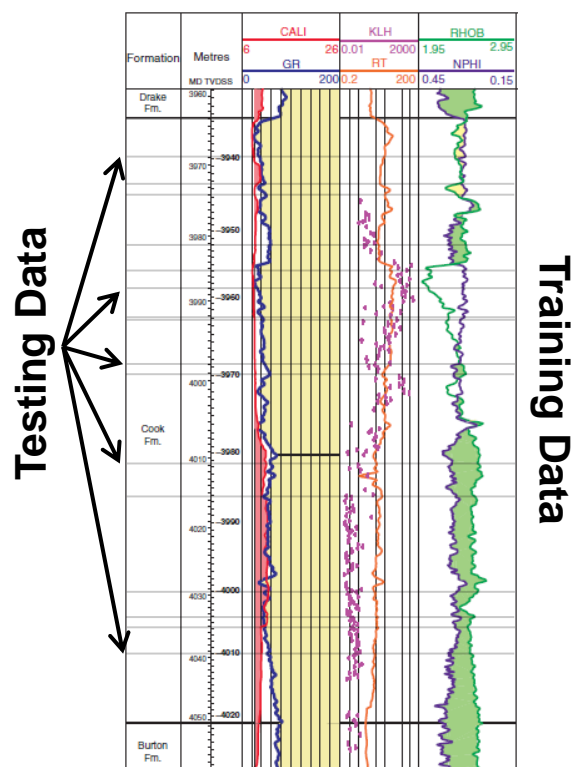
- the range of testing difficulty is similar to the real-world use of the model
- too easy – testing cases are the same or almost the same as training cases, random sampling is often too easy!
- too hard – testing cases are very different from the training cases, the model is expected to severely extrapolate



# Training and Testing

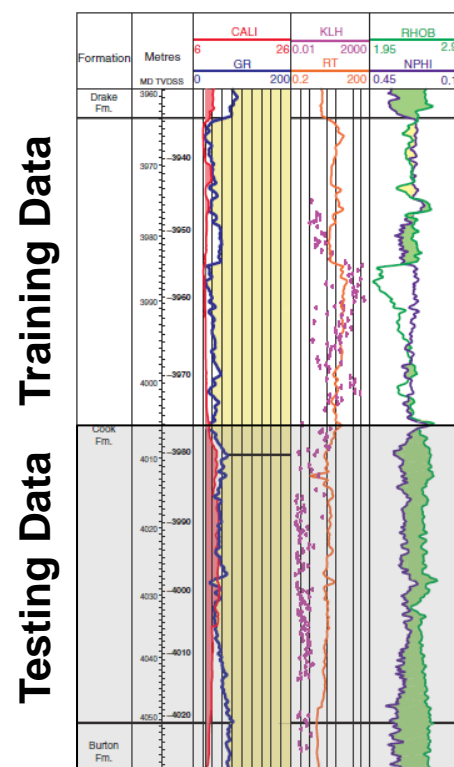
## Fair Testing in Spatial / Temporal Settings

Too Easy

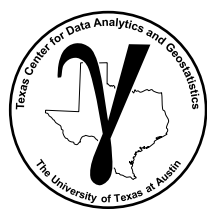


Predictions only at 1/2 ft offsets

Too Hard

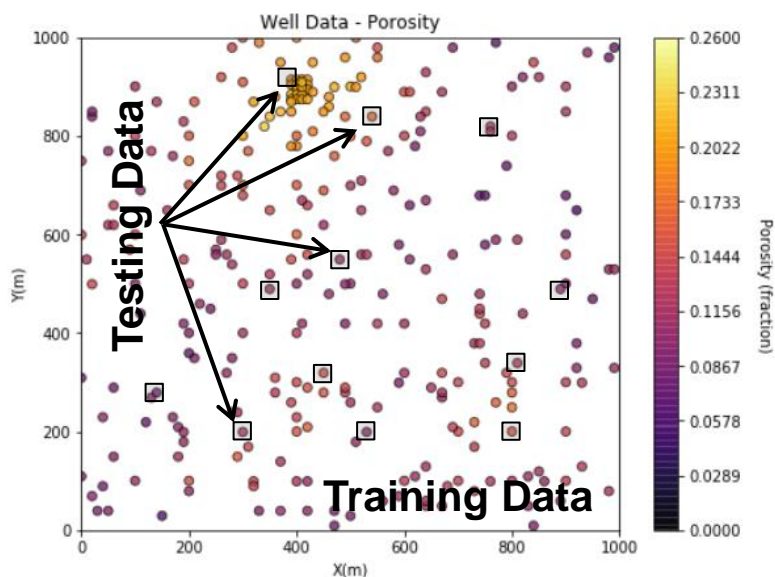


Predictions in a different rock.

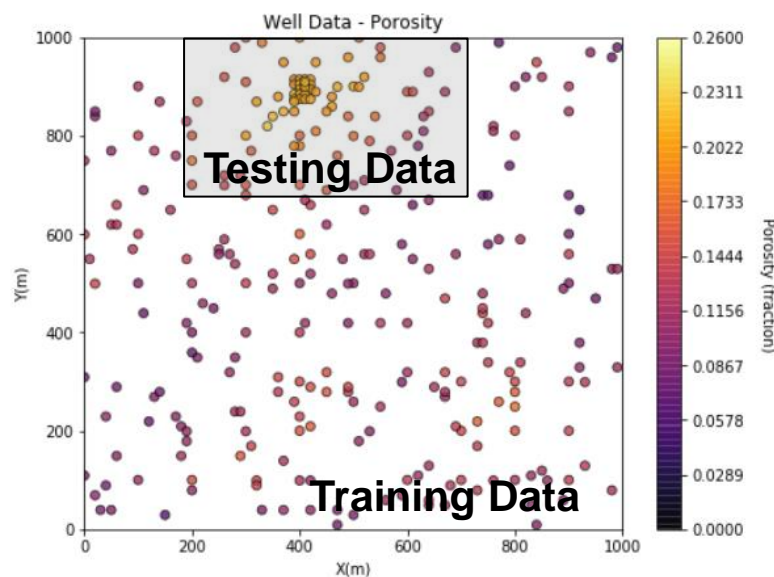


# Training and Testing

## Fair Testing in Spatial / Temporal Settings

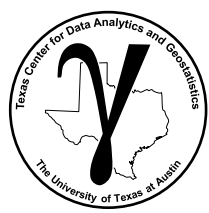


Predictions at only short offsets from training data.



Predictions in a distinctly different range of reservoir values.

We will use random sampling and visualize the training and testing data in Euclidean or feature space. More could be done.

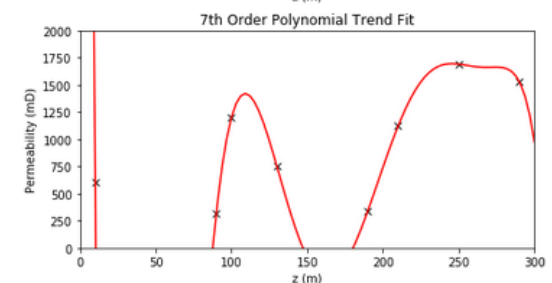
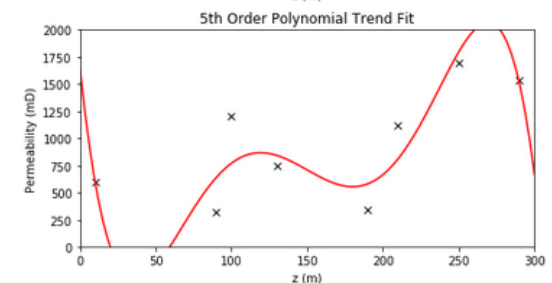
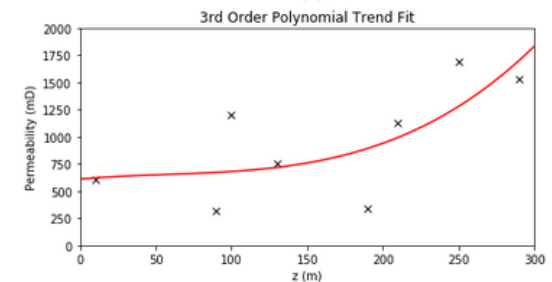
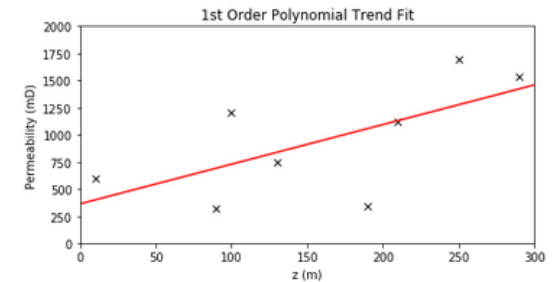


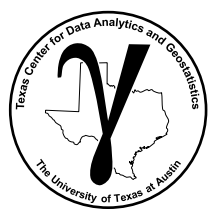
# Model Complexity / Flexibility Definition

## Model Complexity / Flexibility

A variety of concepts may be used to describe model complexity:

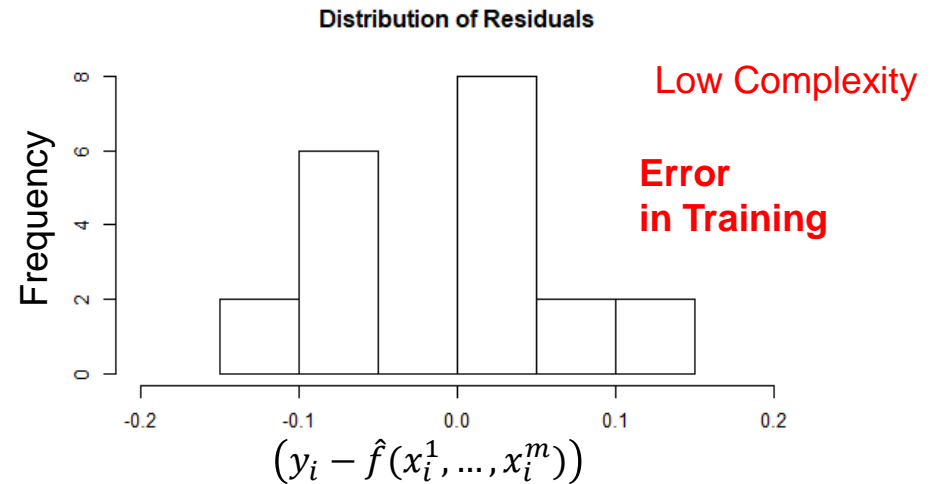
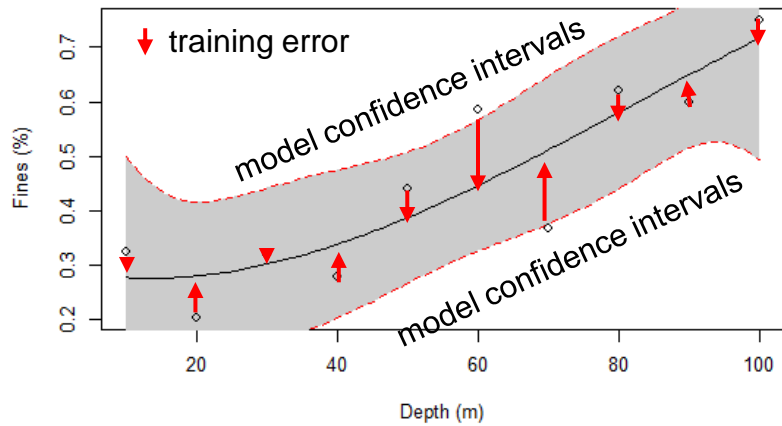
- The number of features:
  - predictor variables are in the model, dimensionality of the model
- The number of terms / parameters
  - the order applied for each term, e.g. linear, quadrature, thresholds
- Expression of the model:
  - Can the model be expressed as:
    - » a compact equation – polynomial regression
    - » nested conditional statements – decision tree
- For example, more complexity with a high order polynomial, larger decision trees etc. →



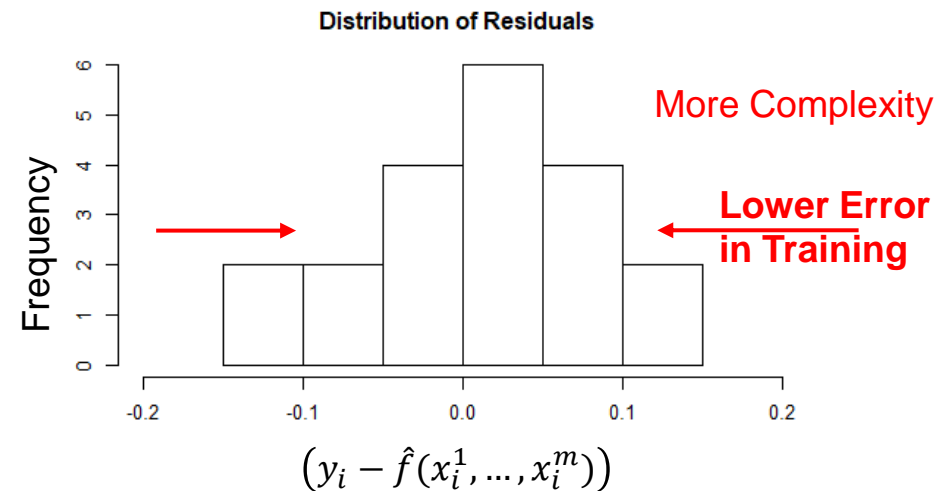
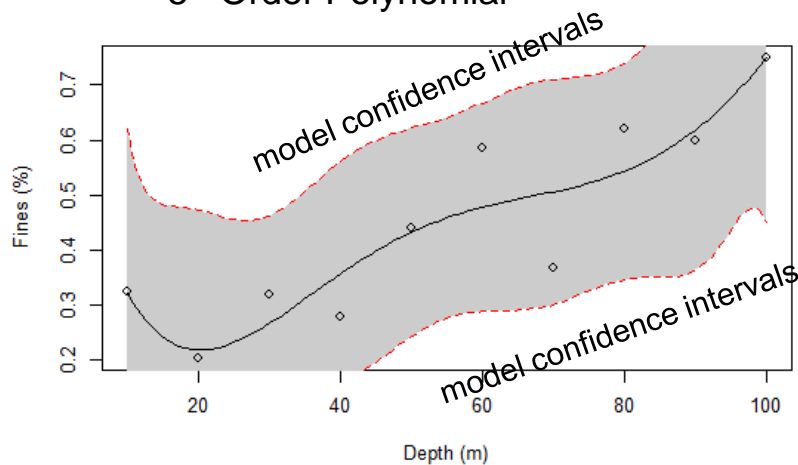


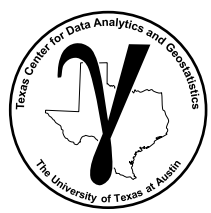
# Simple Statistical Demonstration Overfitting

- Example of trend fits:
  - 3<sup>rd</sup> Ordered Polynomial



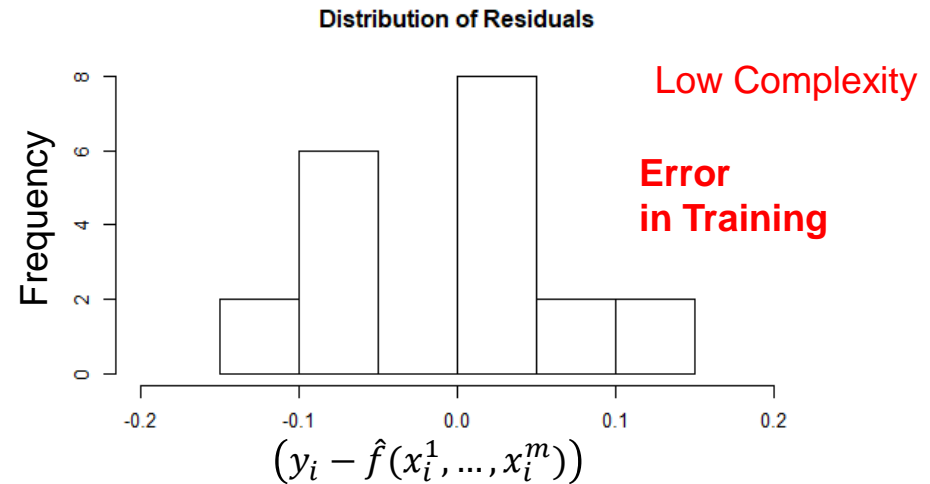
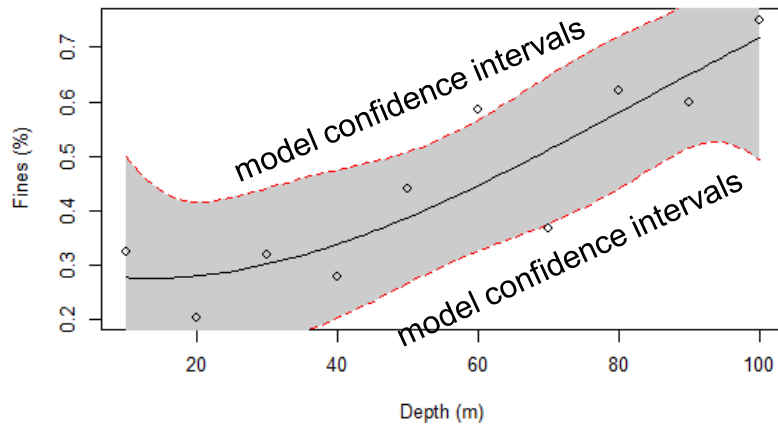
- 5<sup>th</sup> Order Polynomial



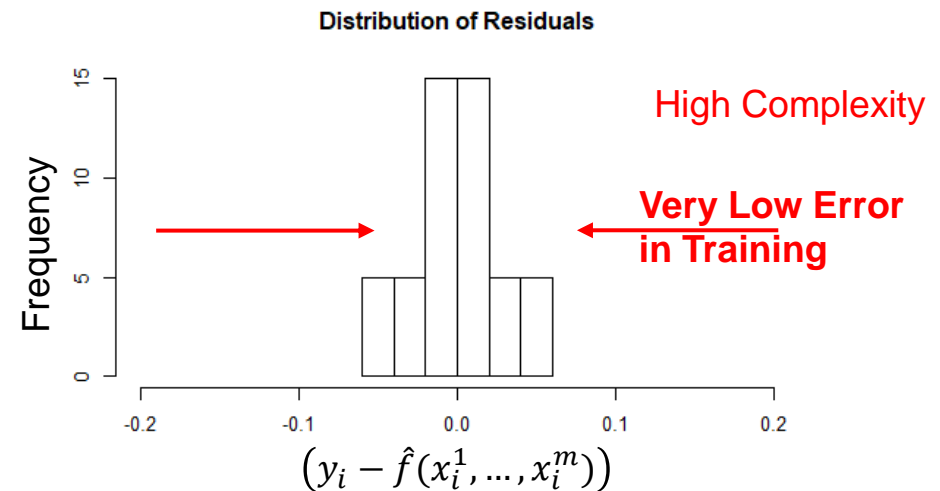


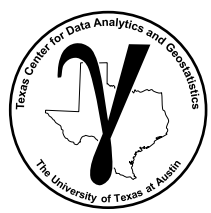
# Simple Statistical Demonstration Overfitting

- Example of trend fits:
  - 3<sup>rd</sup> Ordered Polynomial



- 8<sup>th</sup> Order Polynomial



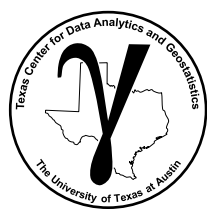


# **PGE 383 Machine Learning**

## **Machine Learning**

**Lecture outline . . .**

- **Model Fitting,  
Overfitting and Model  
Generalization**

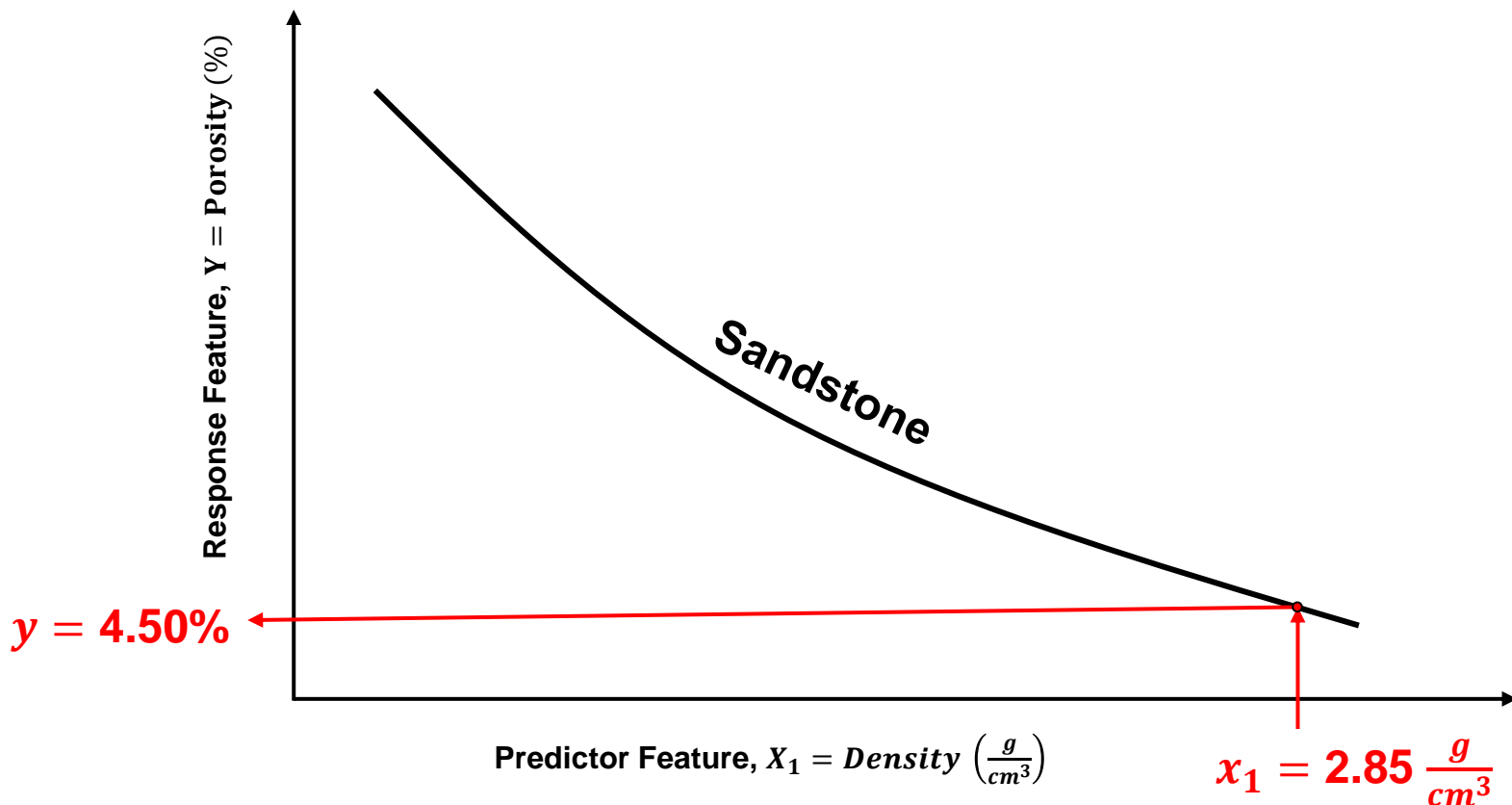


# Fitting, Overfitting and Model Generalization Example

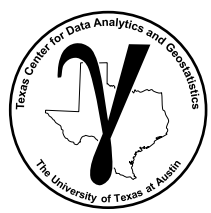
Let's take a simple example from petrophysics to explain fitting, overfitting and generalization

- We need to learn this model, we cannot observe/measure rock porosity in a well bore directly.

*rock porosity from the well log density measure for your sandstone*



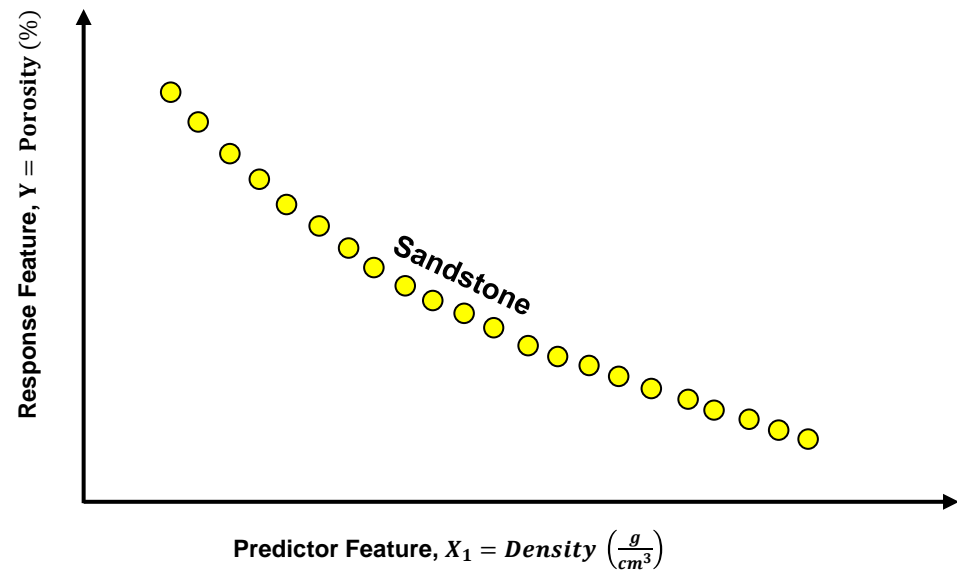
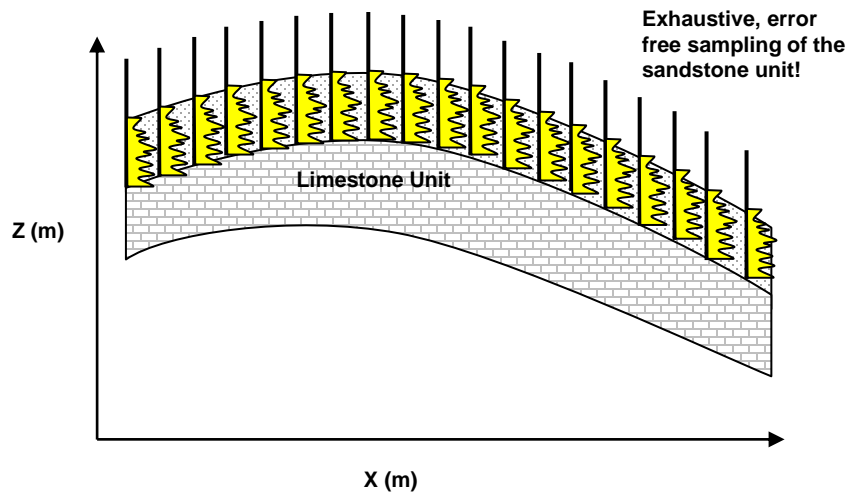
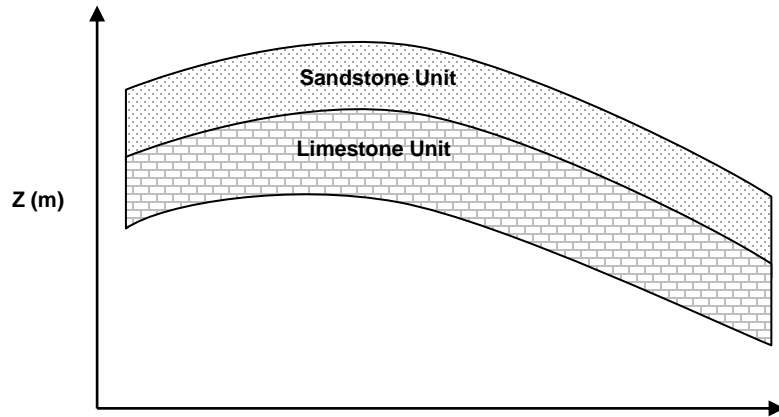


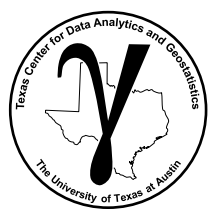


# Overfitting and Model Generalization Example

Assume you are omniscient, and you see the entire natural setting/population!

- If we could see the natural setting at the resolution needed to solve our problem and with complete coverage, we would have the population and know this model between our predictor feature,  $X_1$ , and response feature,  $Y$ , perfectly.

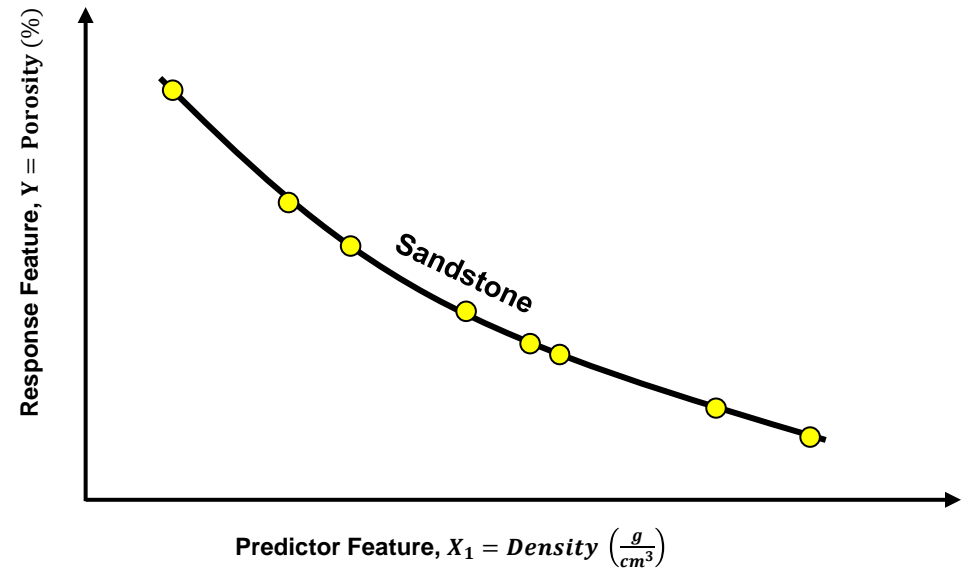
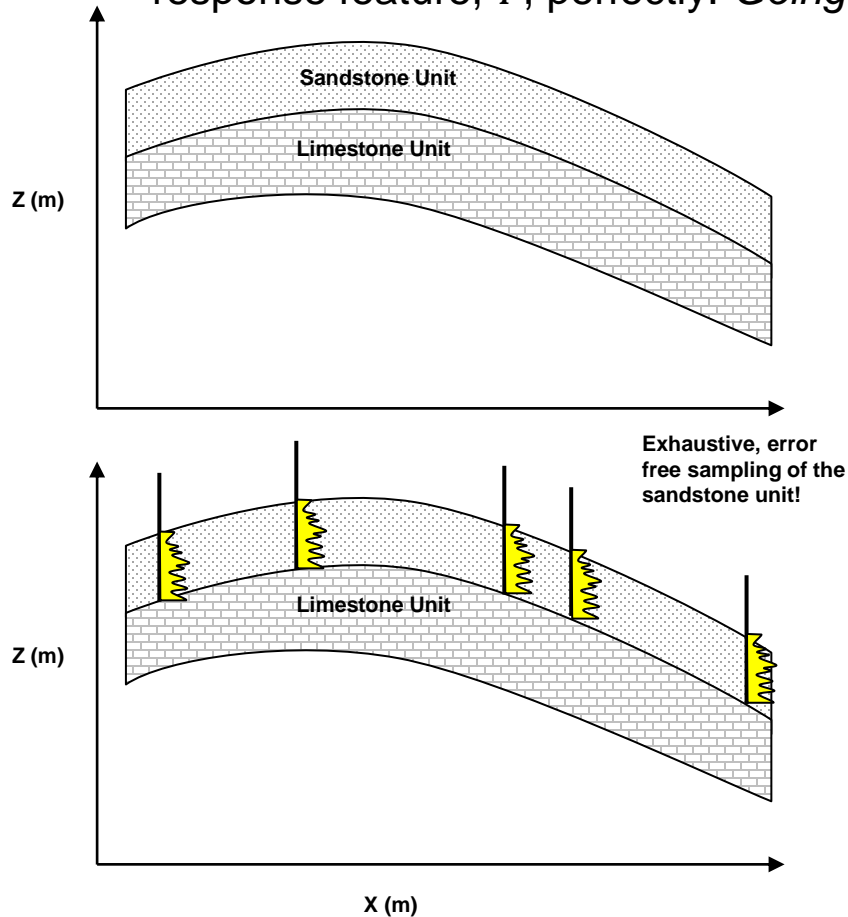


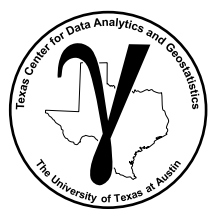


# Overfitting and Model Generalization Example

Assume you integrate physics and limited samples from the population.

- We could build a model with physics (domain information), hinged on limited sample coverage.
- A good (best) model for the relationship between the predictor feature,  $X_1$ , and response feature,  $Y$ , perfectly. *Going forward we will assume data-driven only.*

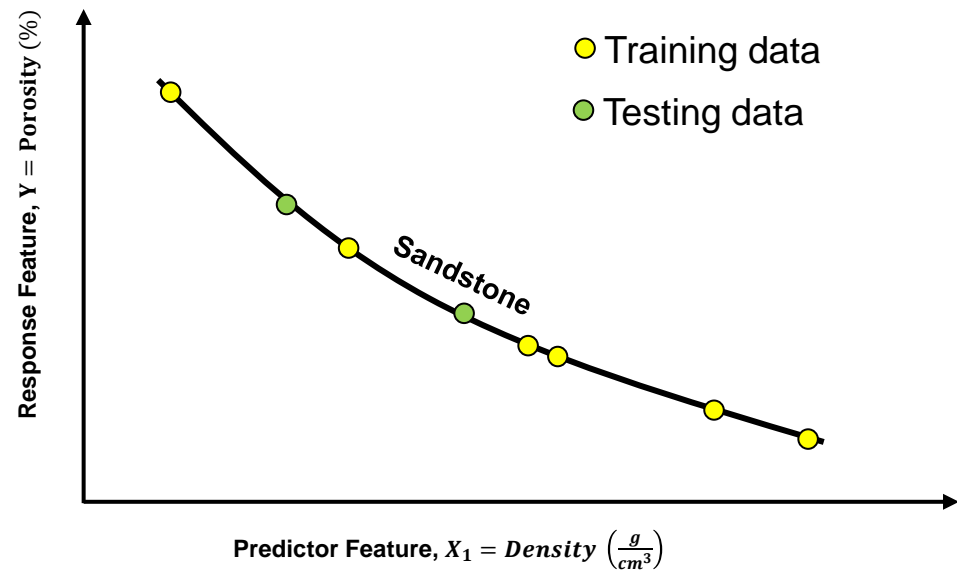
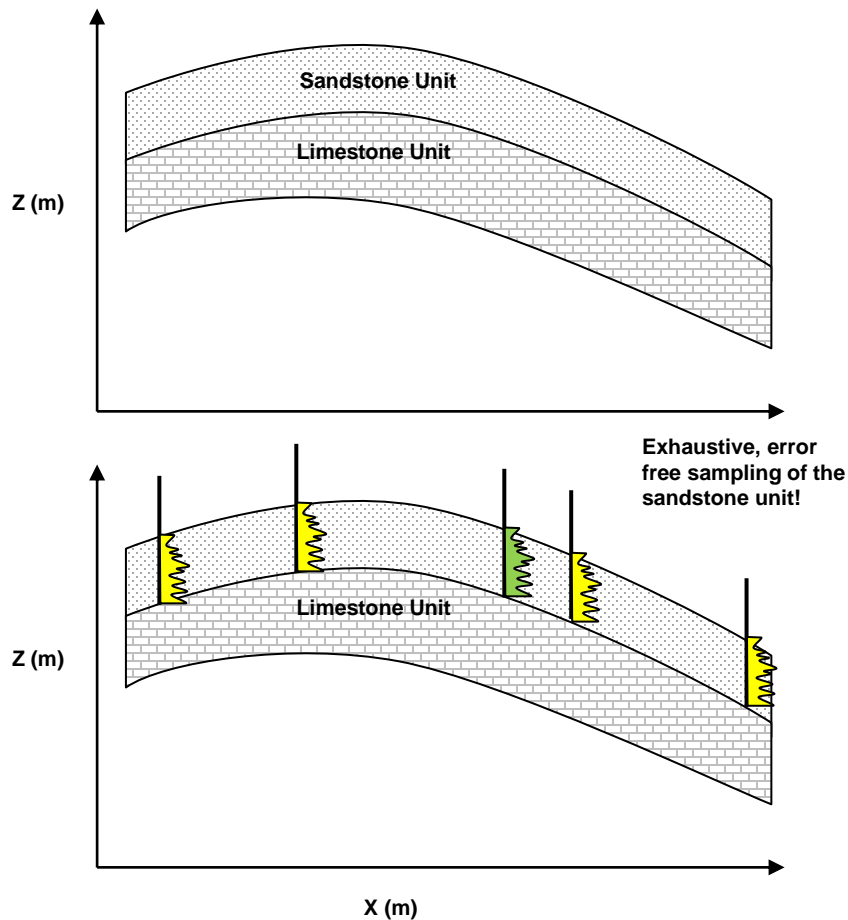


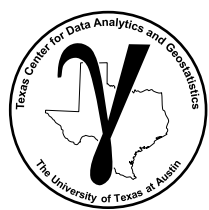


# Overfitting and Model Generalization Example

Assume the data-driven approach, training/tuning a model,  $Y = f(X_1)$ .

- We will separate the data into:
  - Training data to train the model parameters - fit
  - Testing data, withheld from training, to tune the model hyperparameters - complexity

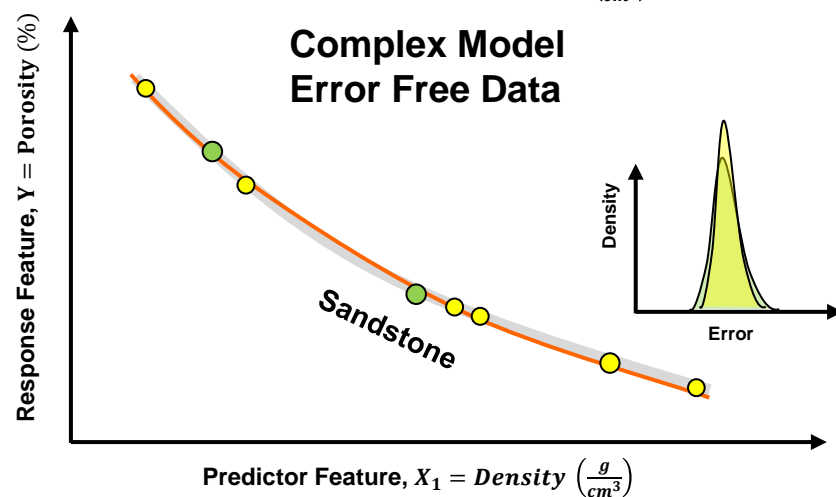
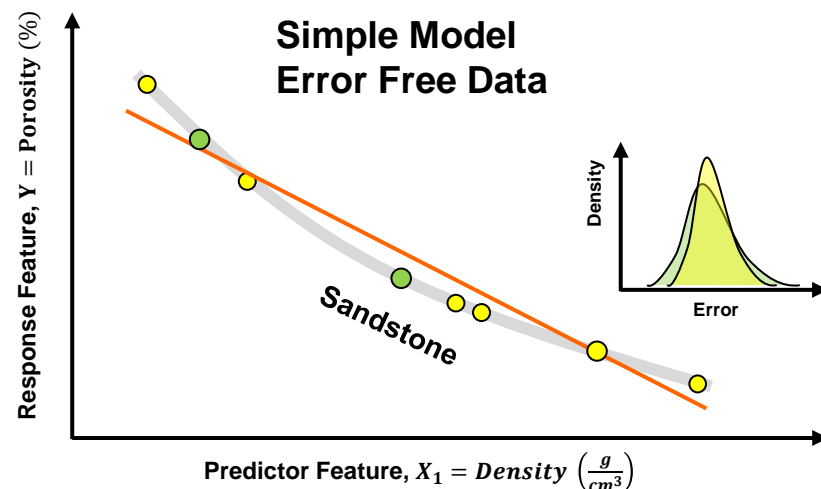
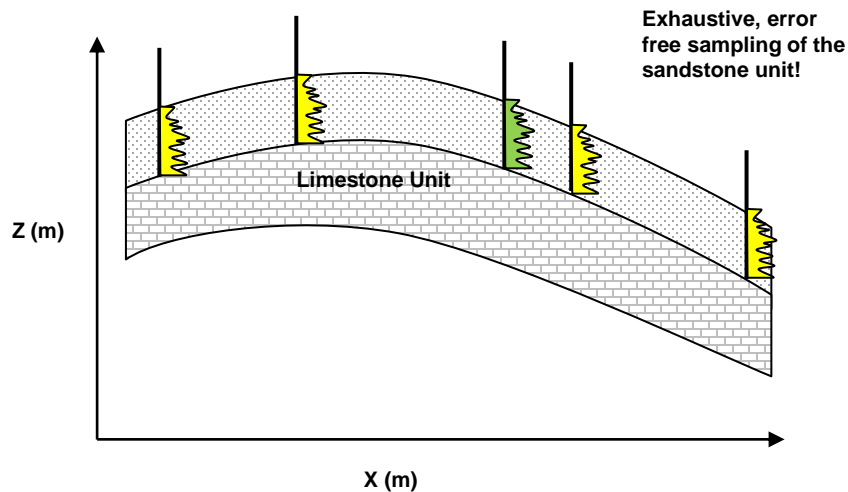
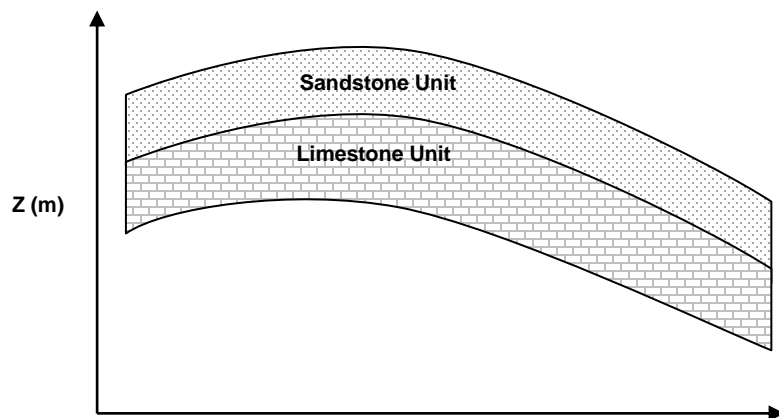


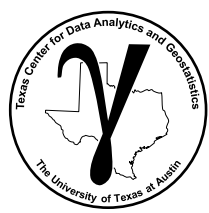


# Overfitting and Model Generalization Example

Assume the data-driven approach, training/tuning a model,  $Y = f(X_1)$ .

- We need to fit an exhaustive model,  $\hat{Y} = \hat{f}(X_1), \forall x_1 \in [x_{min}, x_{max}]$
- As expected, the more complicated model is a better fit. So far it generalizes ok away from training!





# Overfitting and Model Generalization Example

But we don't have error-free measures, we have samples with error

- Error in the measurement of the predictor feature, well log measurement error,  $\epsilon_{X_1}$ .
- Error in the collocated core-based porosity measure,  $\epsilon_Y$ .

Simple Models:

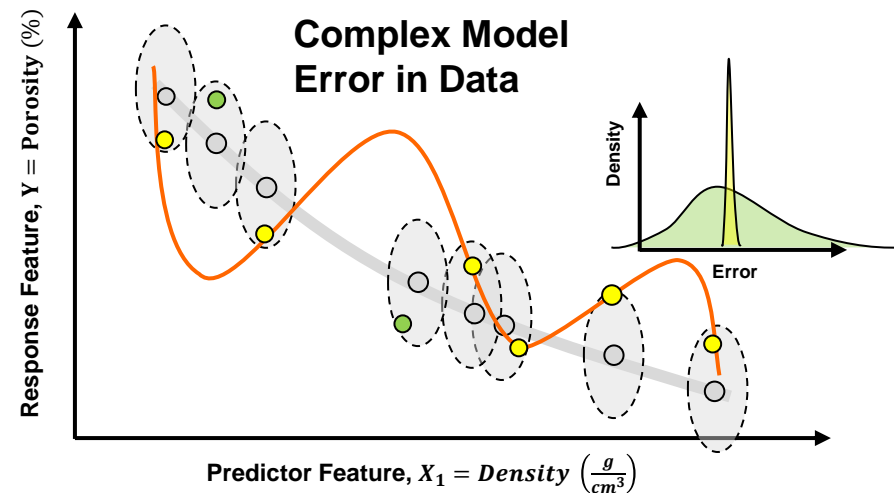
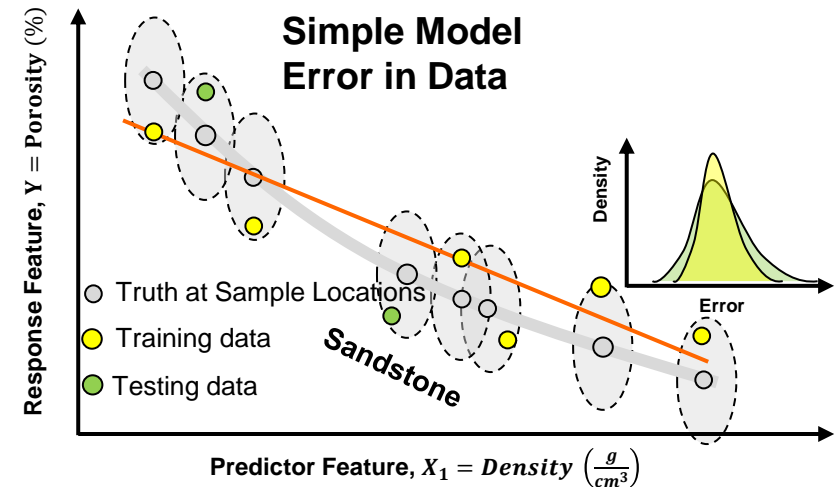
- Less sensitive to error/noise in the data

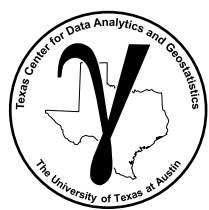
Complexity:

- The ability to flexibly learn the natural system

Complexity + Data Error = Overfit

- Model that fits noise
- Model that poorly generalizes, poor predictions away from training data

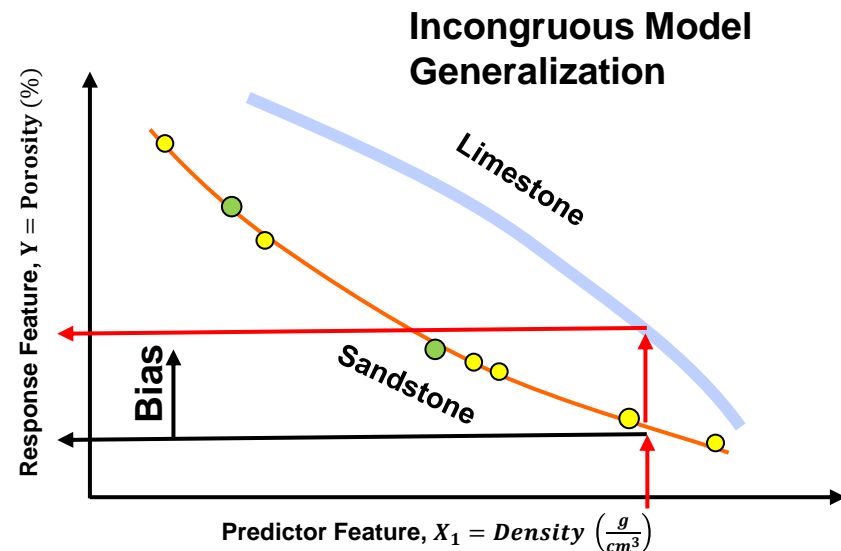
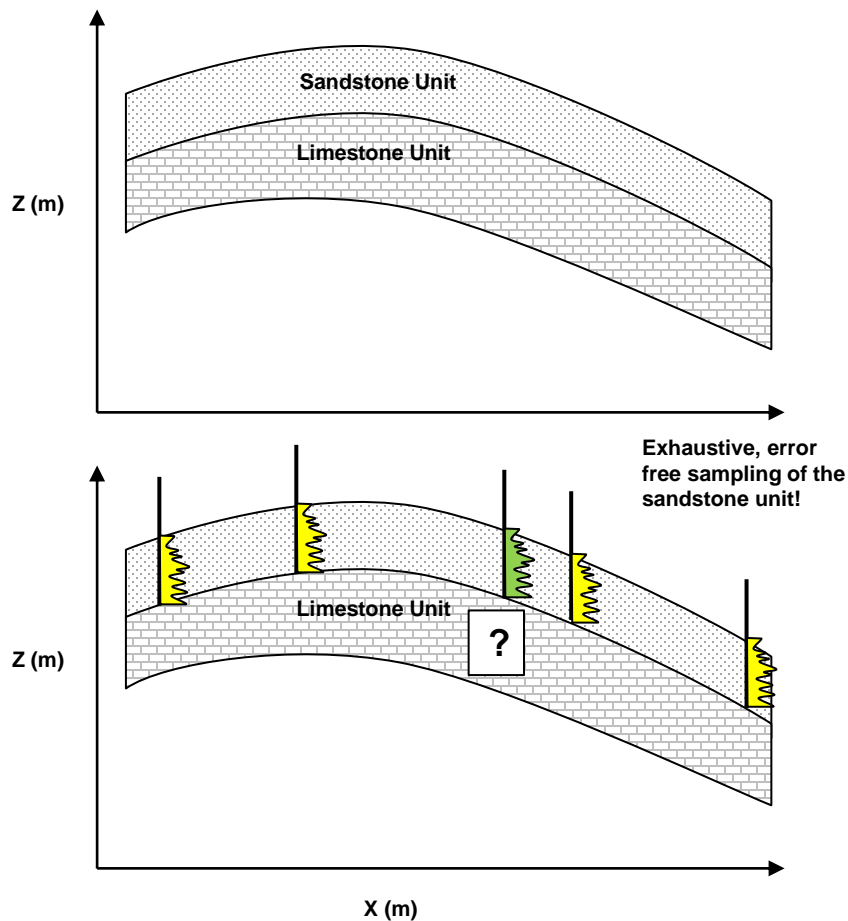




# Overfitting and Model Generalization Example

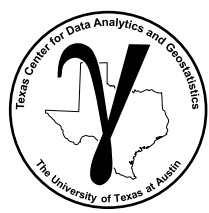
How far can we go with model generalization?

- What if we train and test with sandstone and apply the model to limestone?



There are limits for the congruous application of our machines.

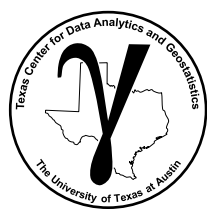
- Training / testing data must be consistent with real-world use.
- As with geostatistics we should be explicit about our decision of stationarity.



# Definition of Overfit

## Overfit:

- More model complexity/flexibility than can be justified with the available data, data accuracy, frequency and coverage
- Model explains “idiosyncrasies” of the data, capturing data noise/error in the model
- High accuracy in training, but low accuracy in testing / real-world use away from training data cases – **poor ability of the model to generalize**



# Building Our Machine, One More Time

Now that we have all the concepts, let's walk through the workflow again.

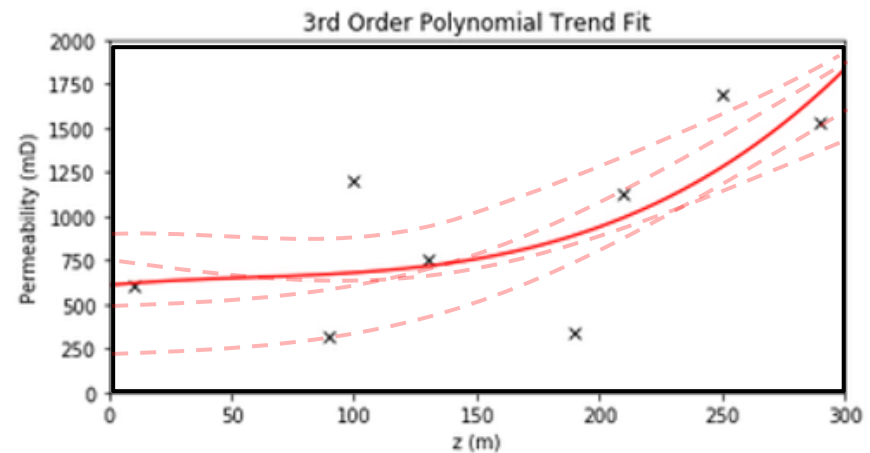
## Apply Training Data to Train the Model Parameters.

- Repeat for all levels of complexity as specified by a range of hyperparameters.
- For example, the parameters of this 3<sup>rd</sup> order polynomial model.

$$b_3, b_2, b_1 \text{ and } c$$

$$k = b_3 z^3 + b_2 z^2 + b_1 z + c$$

- But not appropriate to determine level of complexity (hyperparameter)



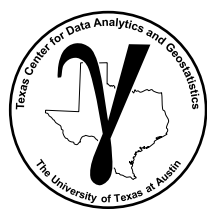
Train model parameters for each level of complexity to maximize fit with training data.

$$MSE = \frac{1}{n} \sum_{i=1}^n \left[ (y_i - \hat{f}(x_1^j, \dots, x_m^j))^2 \right], \text{ for } i = 1, \dots, n_{train}$$

Minimize the summary measure of error over the training data.

**Hyperparameter of our model:**  
1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup> 4<sup>th</sup> ... order polynomial.



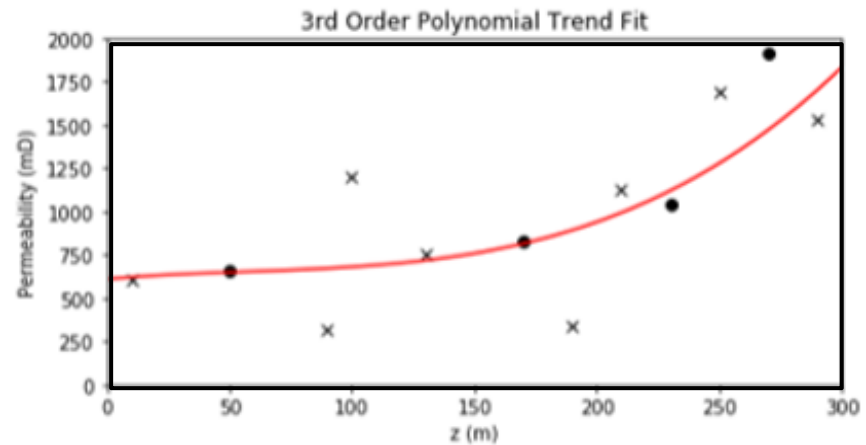


# Building Our Machine, One More Time

Now that we have all the concepts, let's walk through the workflow again.

## Apply Withheld Data to test our Machine.

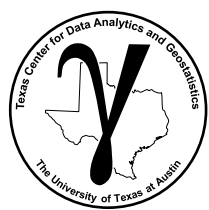
- Calculate the error, over the withheld from training, testing data for all levels of complexity as specified by a range of hyperparameters.
- Select the hyperparameters that minimize error over the withheld testing data.



Test model for each level of complexity against testing data, select the best performing hyperparameters in the testing.

$$MSE = \frac{1}{n} \sum_{i=1}^n \left[ (y_i - \hat{f}(x_1^j, \dots, x_m^j))^2 \right], \text{ for } i = 1, \dots, n_{test}$$

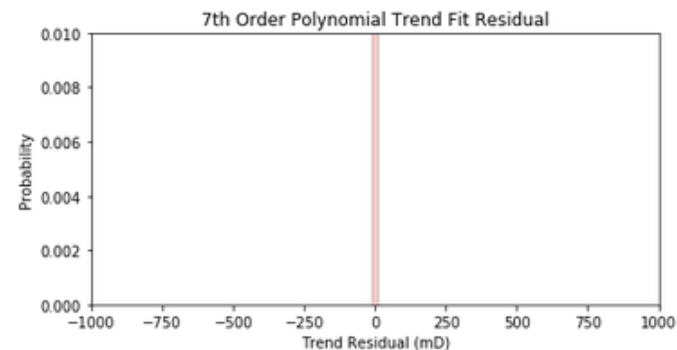
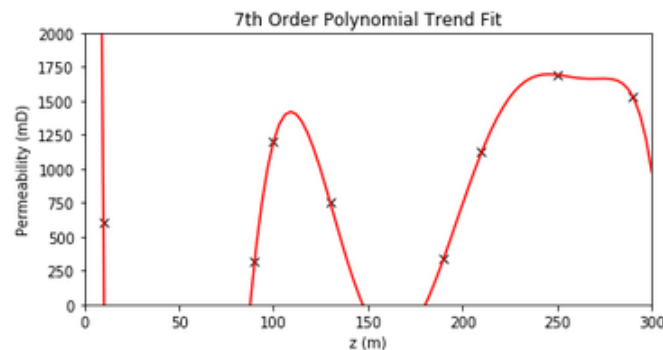
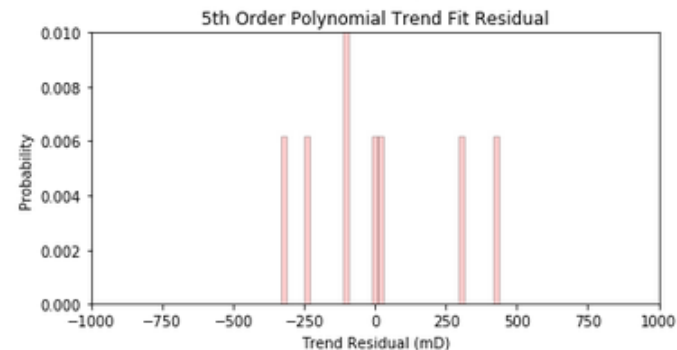
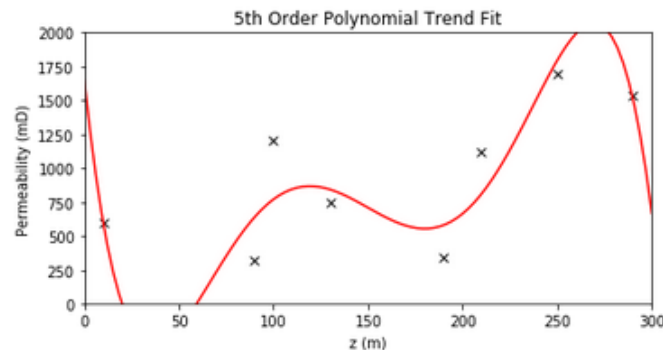
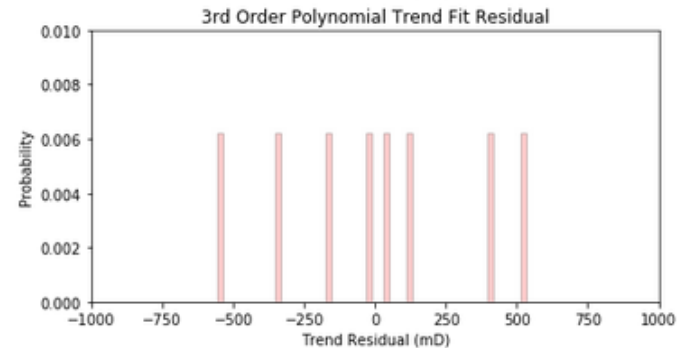
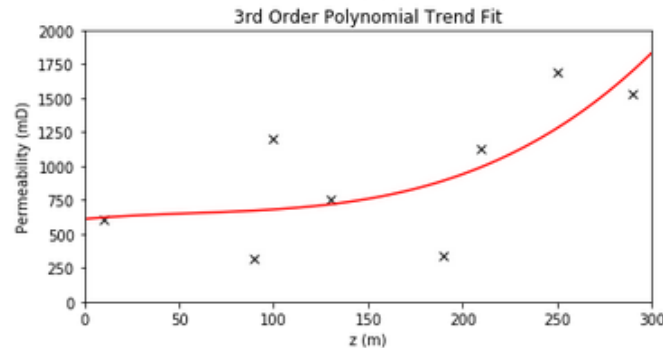
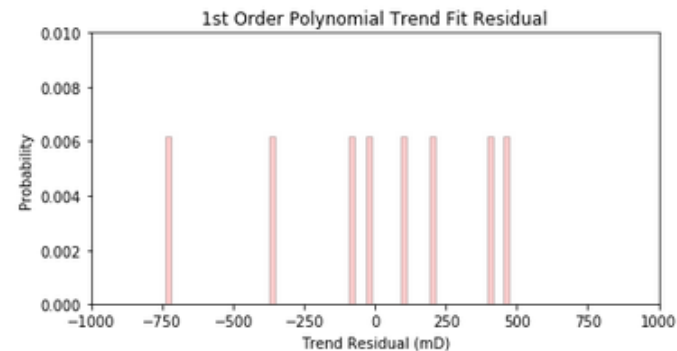
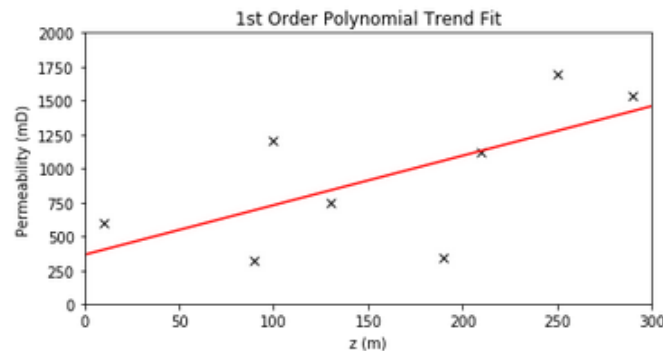
Minimize the summary measure of error over the testing data.

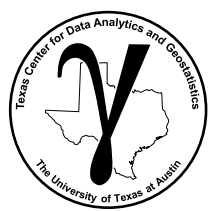


# Building Our Machine, One More Time

What would have happened if we just maximized fit to the data?

- Very complicated model would be best.
- Perfectly [over]fit the data.
- Tuning protects us from overfit, by simulating real-world use of the model.

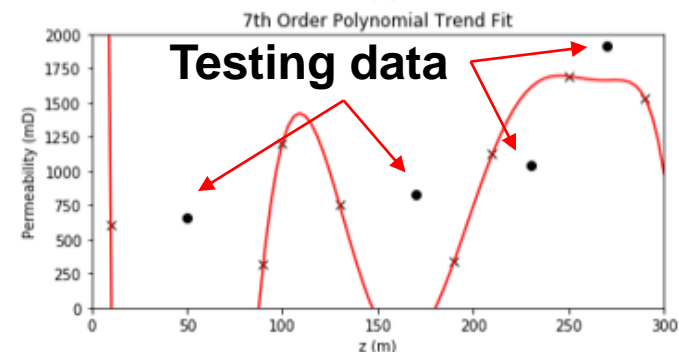
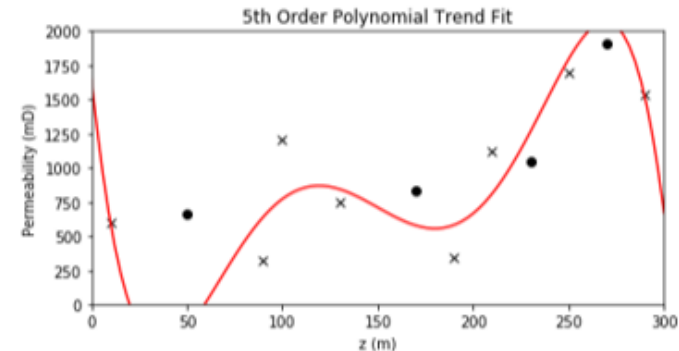
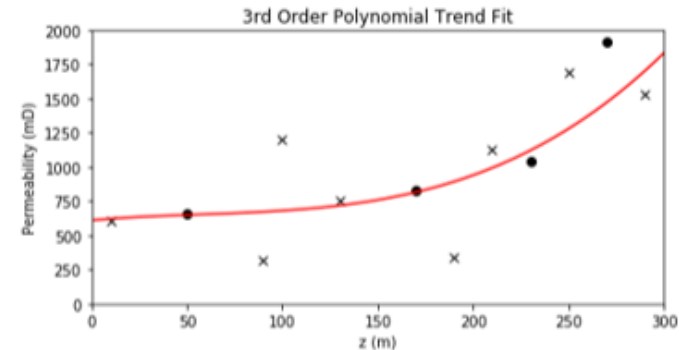
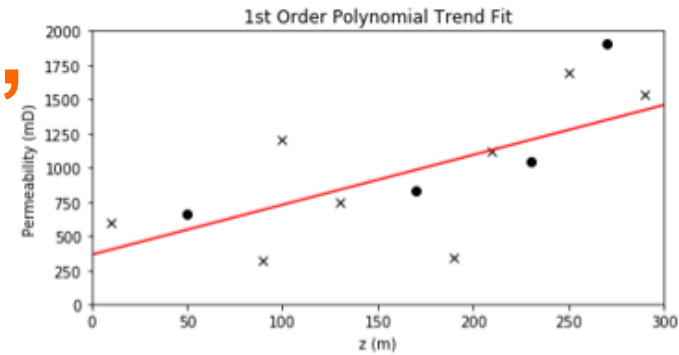


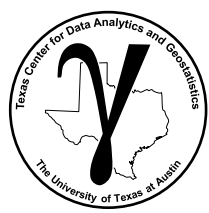


# Building Our Machine, One More Time

## The More Complicated Model Would be Overfit

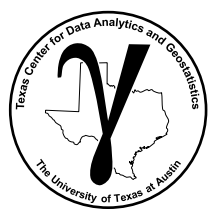
1. Have high accuracy at training data
2. Poor testing accuracy with new observations!
3. Very dangerous with extrapolation.
4. Low model bias, but **high model variance**.





# Now We Begin Machine Learning

- With these concepts established, let's start to get into machine learning / statistical learning methods
  - These methods will allow you to perform inference and prediction
  - Work with complicated data sets / big data analytics
  - Detect patterns in data
- Remember in our business to win:
  - Have the best data
  - Use the data best
- We are at the beginning of the 4<sup>th</sup> paradigm for scientific discovery
  - Data-driven discovery
- Smart fields, 4D seismic surveys, increased computational resources
  - Expanding opportunities for machine learning
- We'll start inferential:
  - Clustering, Principal Component Analysis, Multidimensional Scaling

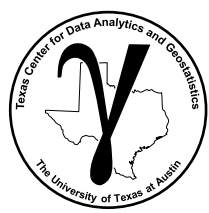


# **PGE 383 Machine Learning**

## **Machine Learning**

**Lecture outline . . .**

- **Examples of Machine Learning**

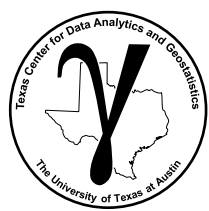


# Examples of Machine Learning

Provides a set of examples with machine learning to address subsurface challenges.

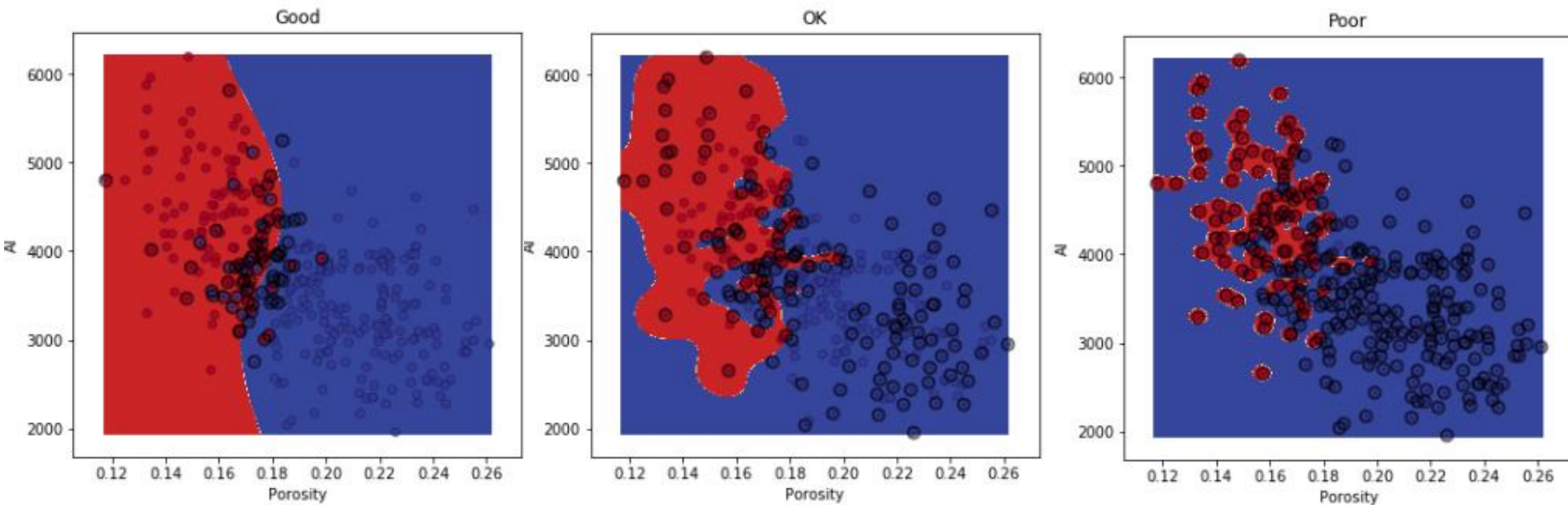
We will cover a wide range of machine learning methods in this class.

This is to motivate and inspire.

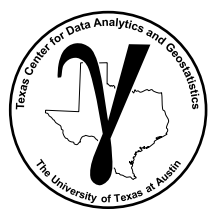


# Support Vector Machines

Support vector machines for interpolating, extrapolating facies from data.



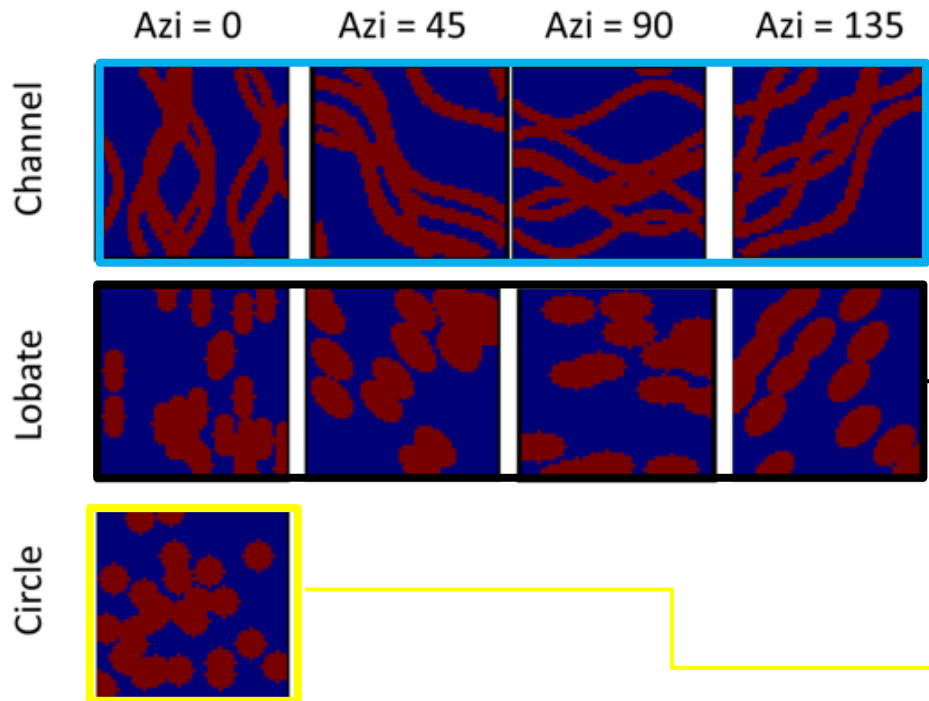
- A range of spatial models with radial basis function  
Kernels with a variety of penalties and Kernel parameters



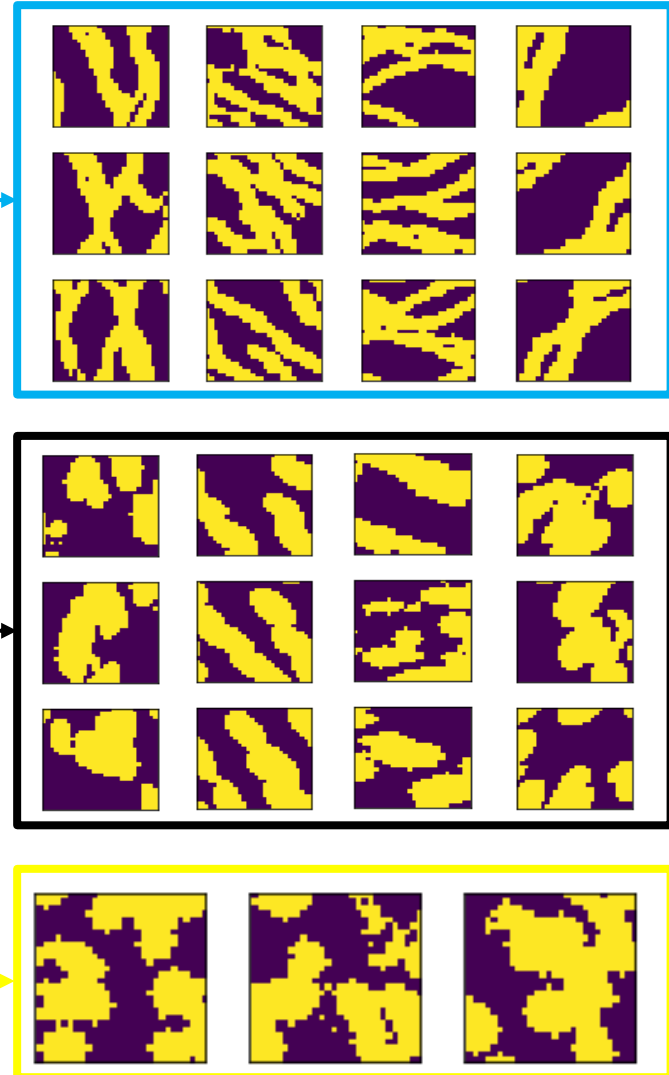
# Neural Networks

- Simple geometric training images for channels, ellipses and circles and realizations.

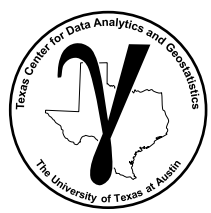
## Training Images



## Realizations

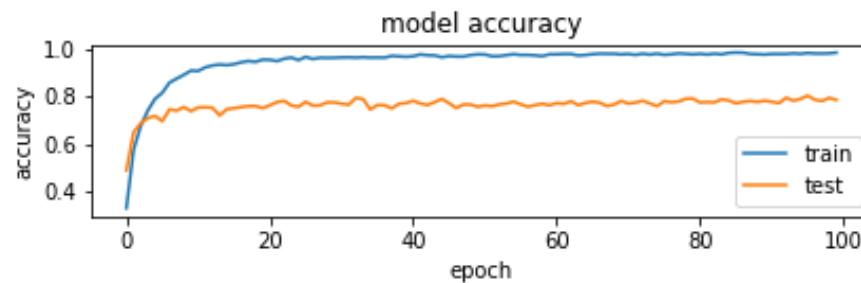




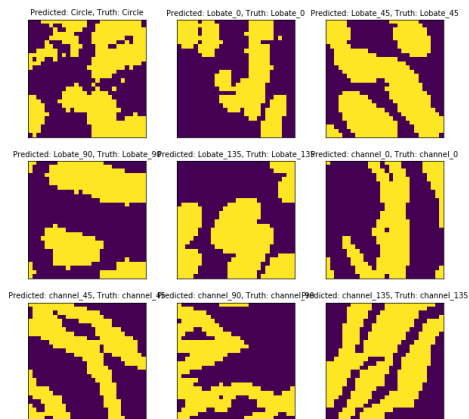


# Neural Networks

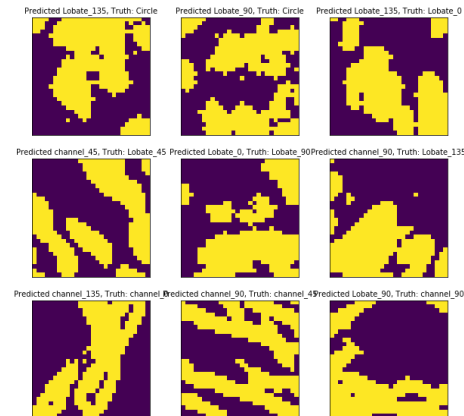
- The training and testing accuracy vs. number of training cycles, **Epoch**
- Levels off at about 78% accuracy

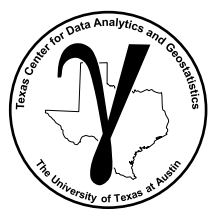


## Correct Identification



## Incorrect Identification

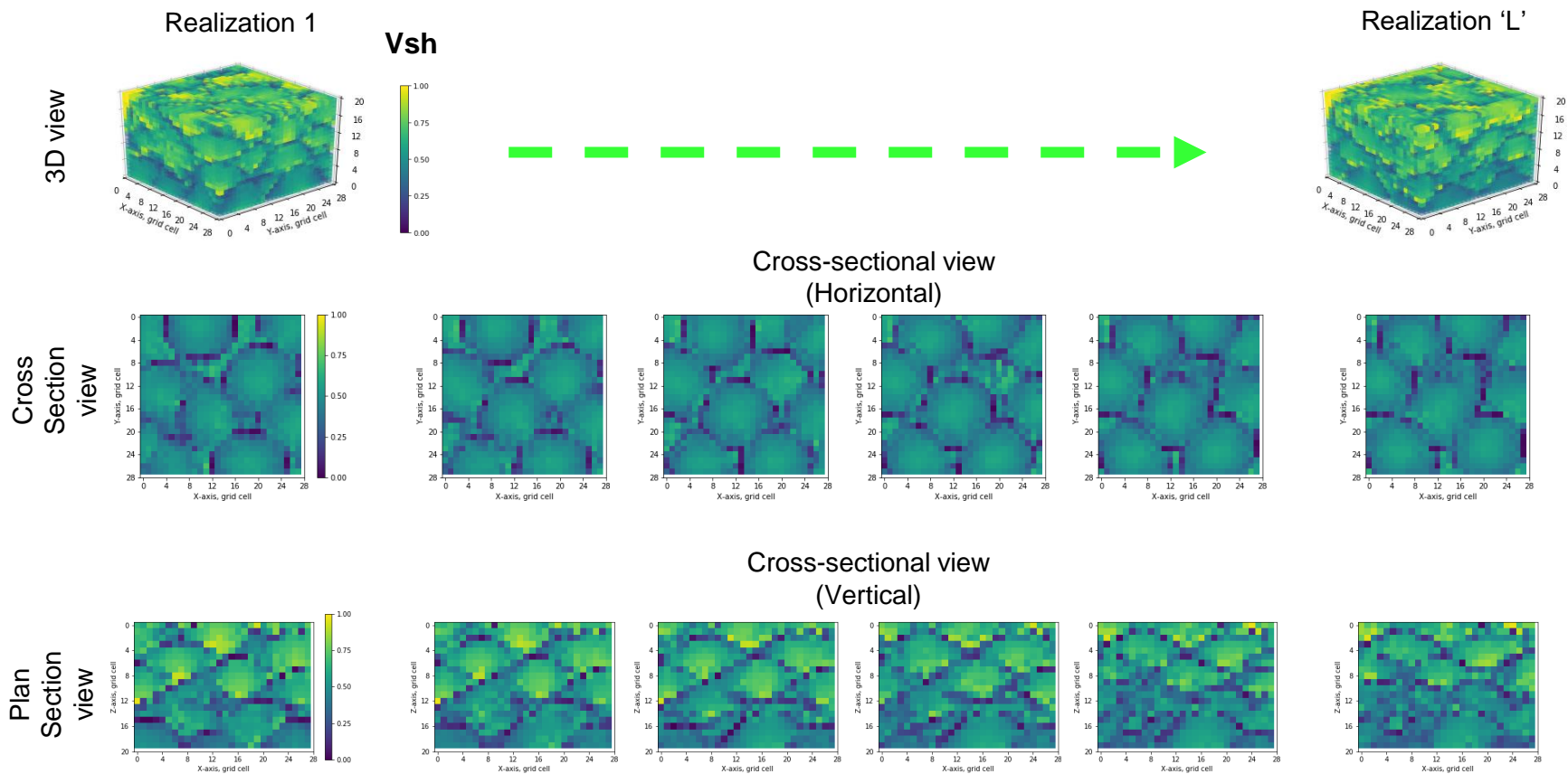




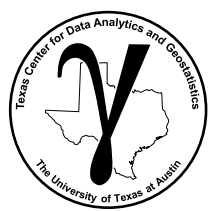
# Convolutional Neural Nets

Can explore the space of uncertainty along a continuous manifold.

- A latent reservoir manifold based on a single parameter



Workflow developed by Honggeun Jo and Javier Santos, PhD student at The University of Texas at Austin.



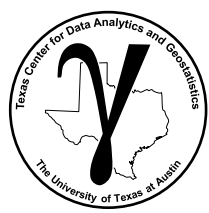
# Convolutional Neural Nets

## Filling In Missing Spatial Information

- Semantic inpainting algorithm (Yeh et al., 2015).
- Using conceptual and perceptual information



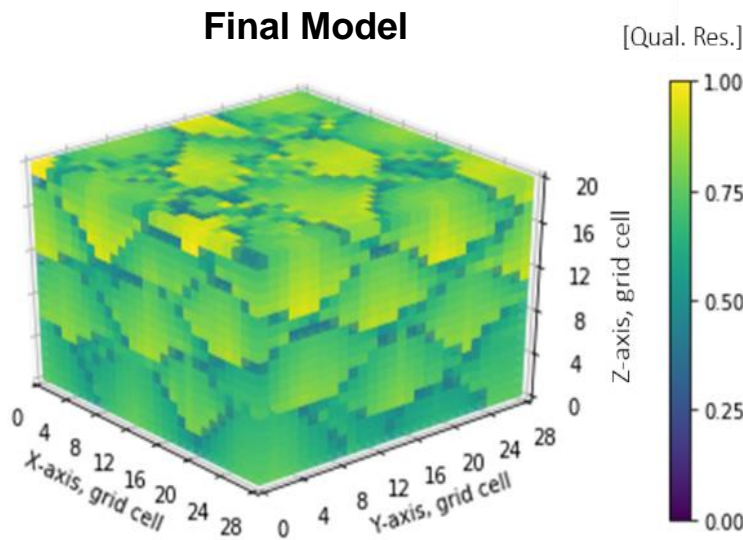
Examples of semantic image inpainting with DCGAN (Yeh et al., 2016)



# Convolutional Neural Nets

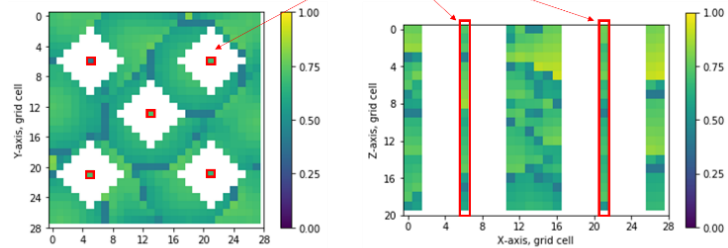
## Conditioning to Well Data?

- Remove model around data
- Use conceptual (model around mask) and perceptual (model elsewhere to fill in missing model consistent with data

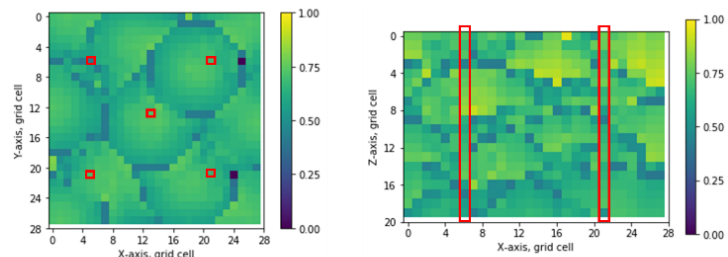


## Original Model

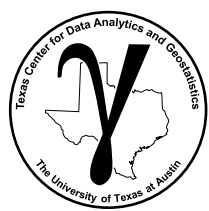
Well data



## Inpainting Around Data



(b)

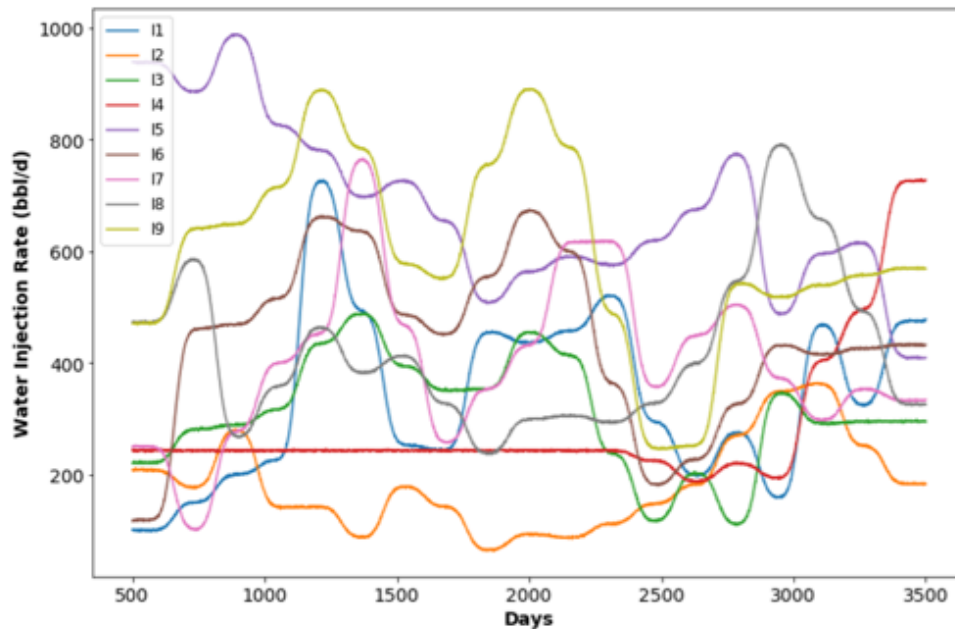


# Long Short-Term Memory Networks

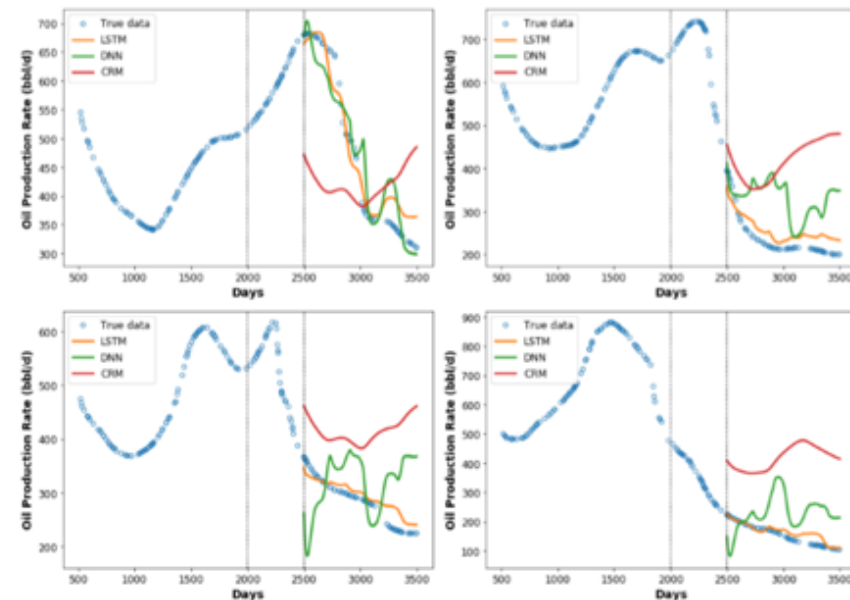
Prediction of producer flow rates based on complicated interactions of injectors.

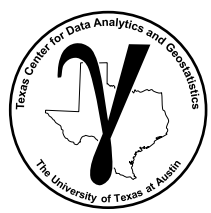
- Train with 2500 days and predict future 100 days.

Injection Rates Over Train and Test Intervals



Production Over Train and Modeled Over Test





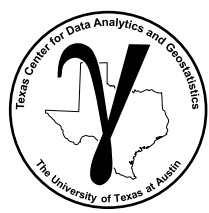
# **PGE 383 Machine Learning**

## **Machine Learning**

**Lecture outline . . .**

- **Energy Machine Learning**

**Michael Pyrcz, The University of Texas at Austin**



# The 4<sup>th</sup> Paradigm

## Welcome to the 4<sup>th</sup> Paradigm of Scientific Discover!

### 1<sup>st</sup> Paradigm Empirical Science

- Experiments

### 2<sup>nd</sup> Paradigm Theoretical Science

- Laws of classical mechanics, electrodynamics

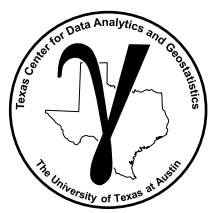
### 3<sup>rd</sup> Paradigm Computational Science Simulation

- Continuum mechanics for heterogeneous systems
- Computational fluid dynamics

### 4<sup>th</sup> Paradigm Data-driven Science

- Detection of patterns and anomalies in big data
- Artificial intelligence





# Energy Digitization

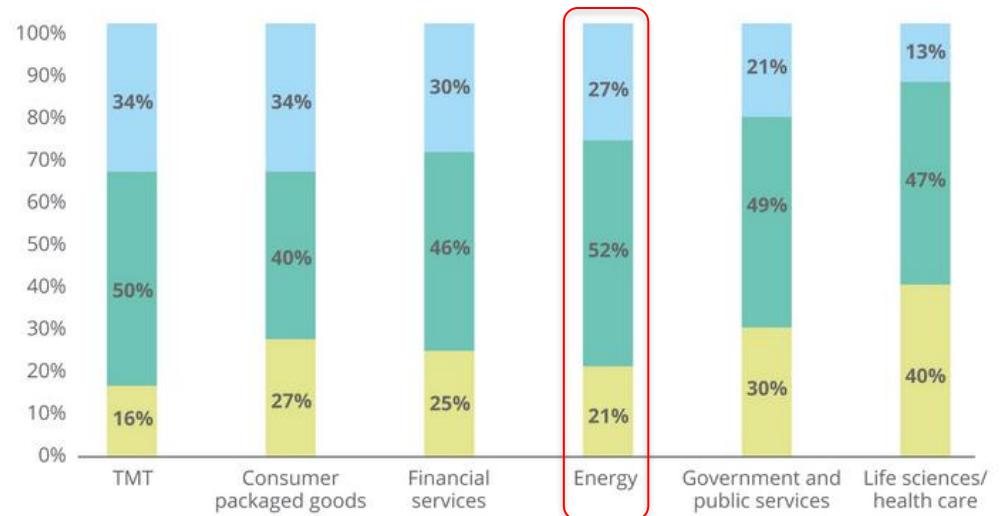
## We Are Not Alone

- Digital transformations are underway in all sectors of our economy
- Every energy company that I visit is working on this right now

FIGURE 14

**TMT companies had the greatest percentage of median- and higher-maturity organizations**

■ Lower maturity ■ Median maturity ■ Higher maturity



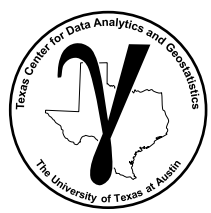
Note: Percentages may not total 100% due to rounding.

Source: Deloitte Digital Transformation Executive Survey 2018.

Deloitte Insights | [deloitte.com/insights](https://deloitte.com/insights)

Digital transformation study by Deloitte, 2019.





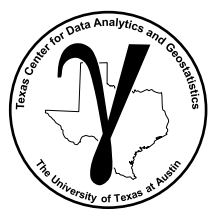
# Energy Digitization

## My Biases:

- Opportunities to do more with our data
- Opportunities to teach data analytics and statistical / machine learning methods to engineers and geoscientists for improved capability
- Geoscience and engineering knowledge & expertise remains core to our business



Digital transformation PricewaterhouseCoopers (PwC) panel  
April 9th, 2019.

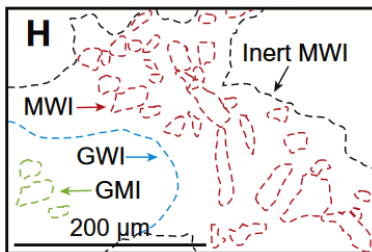
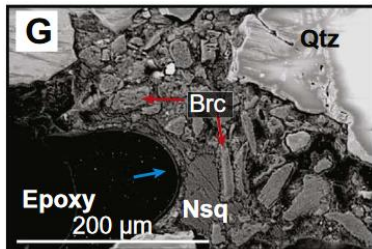


# Working in the 4<sup>th</sup> Paradigm

**We integrate all paradigms, new tools to add value:**

- We augment with new scientific paradigms
- We don't replace older paradigms!

## 1<sup>st</sup> Paradigm Empirical Science



Microfluidics experiment brucite  
carbonation experiment  
(Harrison et al., 2017).

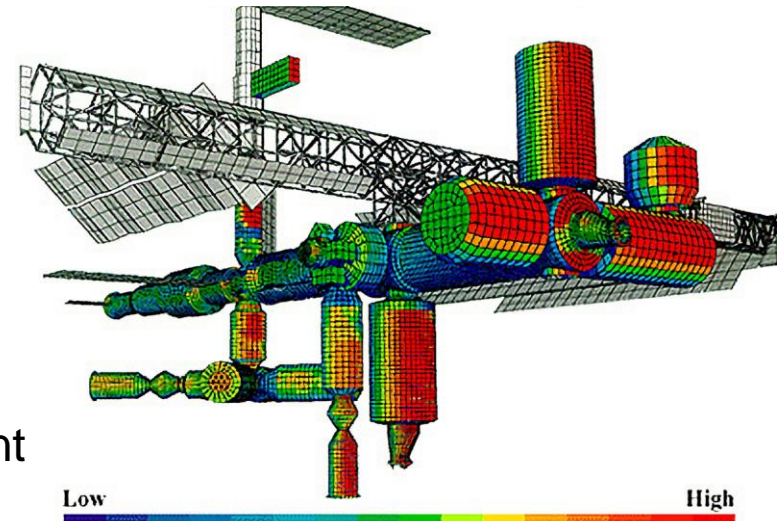
## 2<sup>nd</sup> Paradigm Theoretical Science

$$q = -\frac{k}{\mu} \nabla p$$

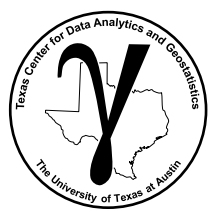
$q$  – flux  
 $k$  – permeability  
 $\mu$  – dynamic viscosity  
 $\nabla p$  – pressure gradient

Darcy's law.

## 3<sup>rd</sup> Paradigm Computational Science Simulation



International space station impact risk from  
computer simulation. Image from  
[https://en.wikipedia.org/wiki/Risk\\_management](https://en.wikipedia.org/wiki/Risk_management).



# Data

## Data-driven Science Needs Data, Data Preparation Remains Essential

>80% of any subsurface study is data preparation and interpretation

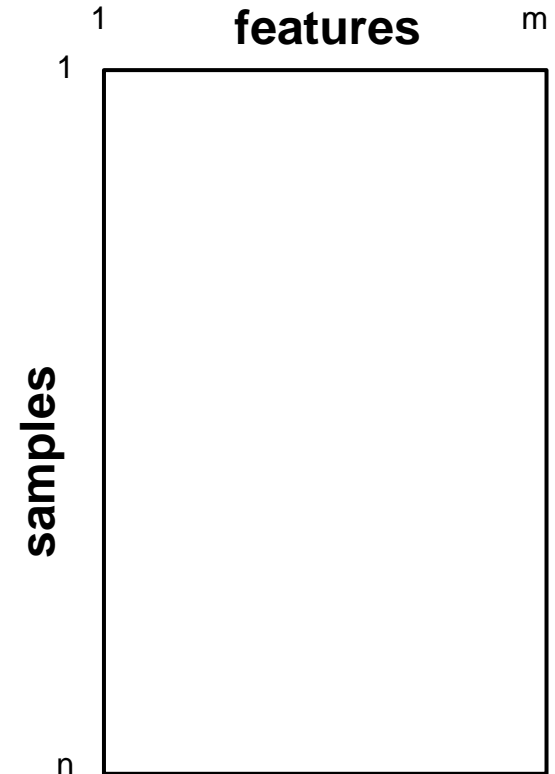
We continue to face a challenge with data:

1. Data curation
2. Large volume
3. Large volumes of metadata
4. Variety of data, scale, collection, interpretation
5. Transmission, controls and security

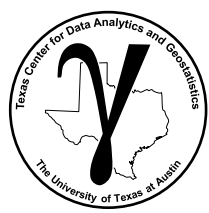
Clean databases are prerequisite to all data analytics and machine learning

Must start with this foundation

Garbage in, garbage out



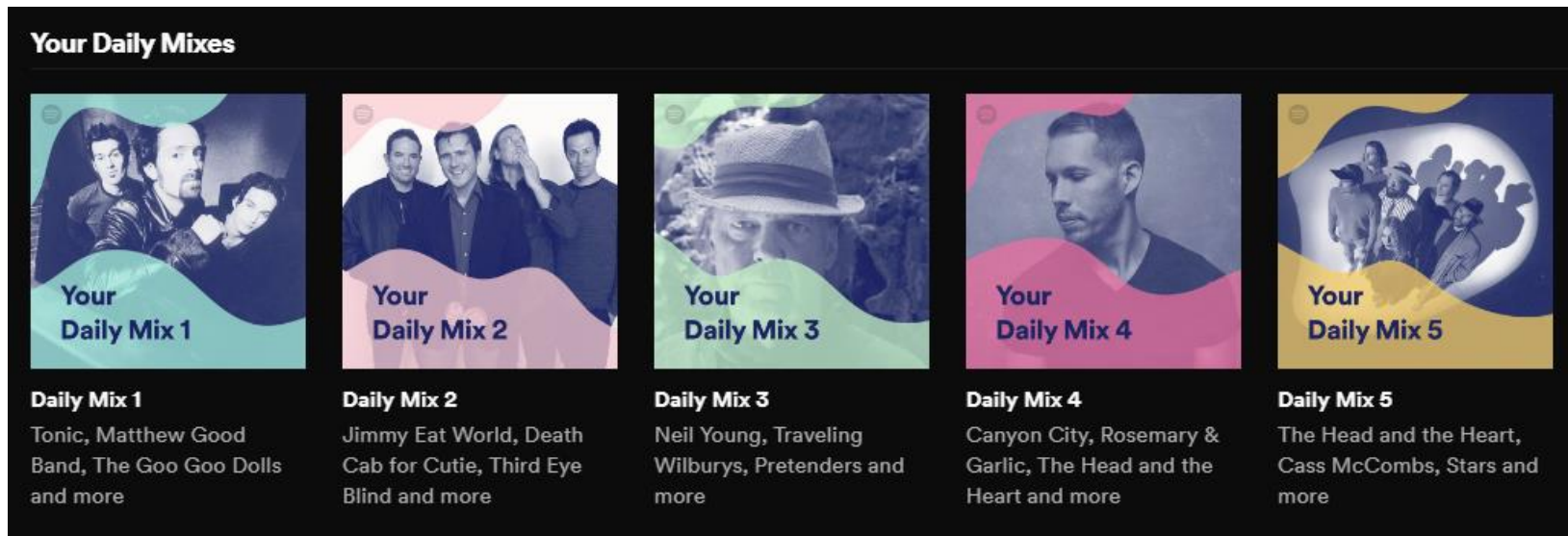
The common data table, samples and features.



# Energy is Unique

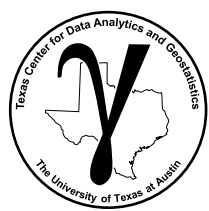
## Energy is Different and May Need New Solutions:

- Sparse, uncertain data, complicated and heterogeneous, open earth systems
- high degree of necessary geoscience and engineering interpretation and physics
- expensive, high value decisions that must be supported



Spotify recommender system from my account summer, 2019.

- Be a critical user / consumer / developer of this technology



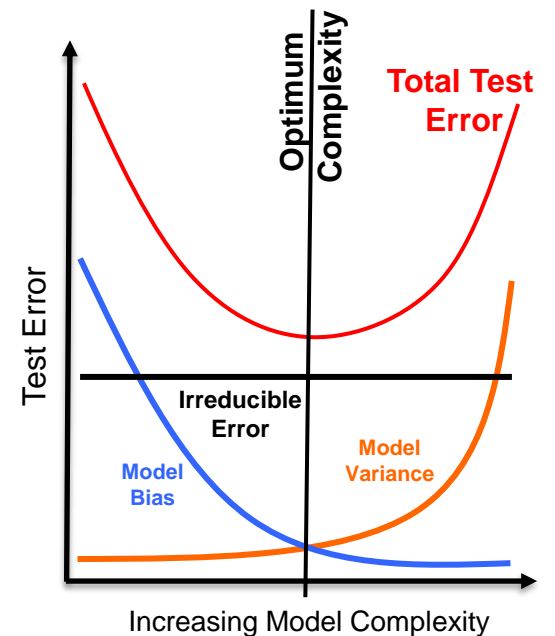
# Don't Jump to Complexity

The Expected Test Mean Square Error may be calculated as:

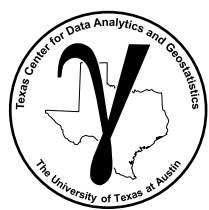
$$E[(y_0 - \hat{f}(x_1^0, \dots, x_m^0))^2] = \underbrace{\text{Var}(\hat{f}(x_1^0, \dots, x_m^0))}_{\text{Model Variance}} + \underbrace{[\text{Bias}(\hat{f}(x_1^0, \dots, x_m^0))]^2}_{\text{Model Bias}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}}$$

Remember:

- **Model Variance, Model Bias and Irreducible Error**
- Often simpler model outperform more complicated models, controlling model variance is critical!
- While providing a more interpretable model to support high value decisions



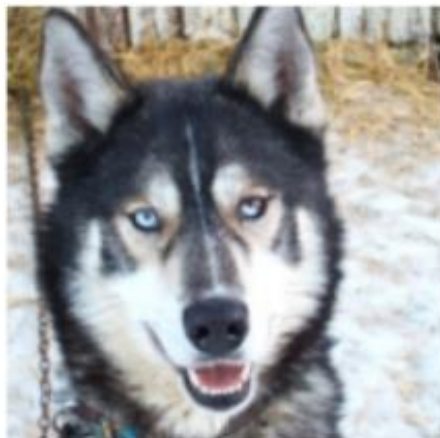
Model variance and bias trade-off.



# Interpretability is Critical

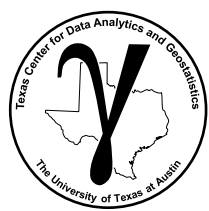
## Develop Methods and Workflows that Provide Useful Diagnostics

- Interpretability may be low
- Application may become routine and trusted
- The machine is trusted, becomes an 'unquestioned authority'



Machine learning-based logistic classifier to identify wolf or dog.

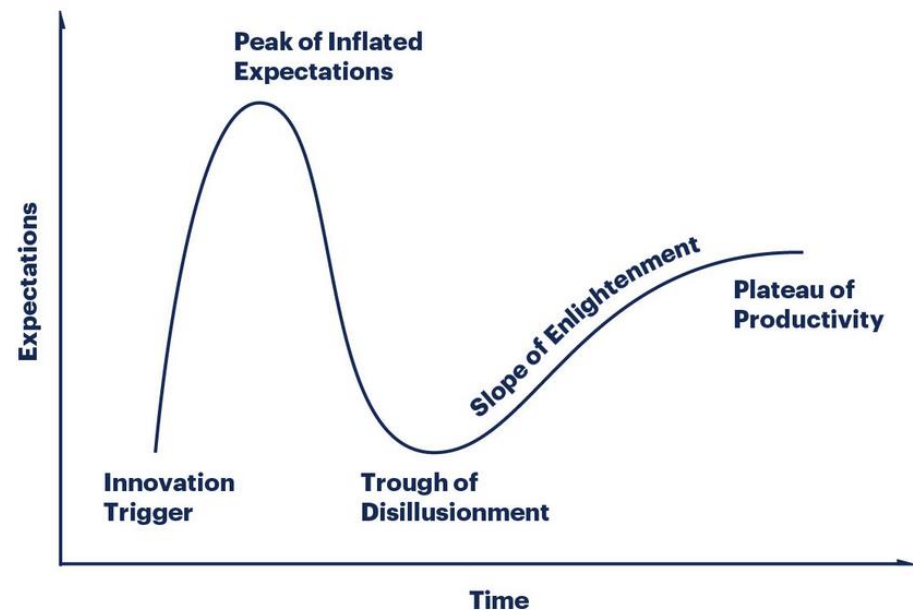
*'Even the developers that work on this stuff have no idea what it is doing' 'These systems do not fail gracefully!' – Peter Haas TED Talk.*



# Meeting Technical Expectations

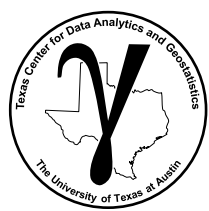
## The Technology Hype Cycle (from Gartner)

- Where are we currently for data analytics and machine learning?
- Varies by company and by group within company.
- Globally, expectations are high!



Technology hype cycle from time of discovery.



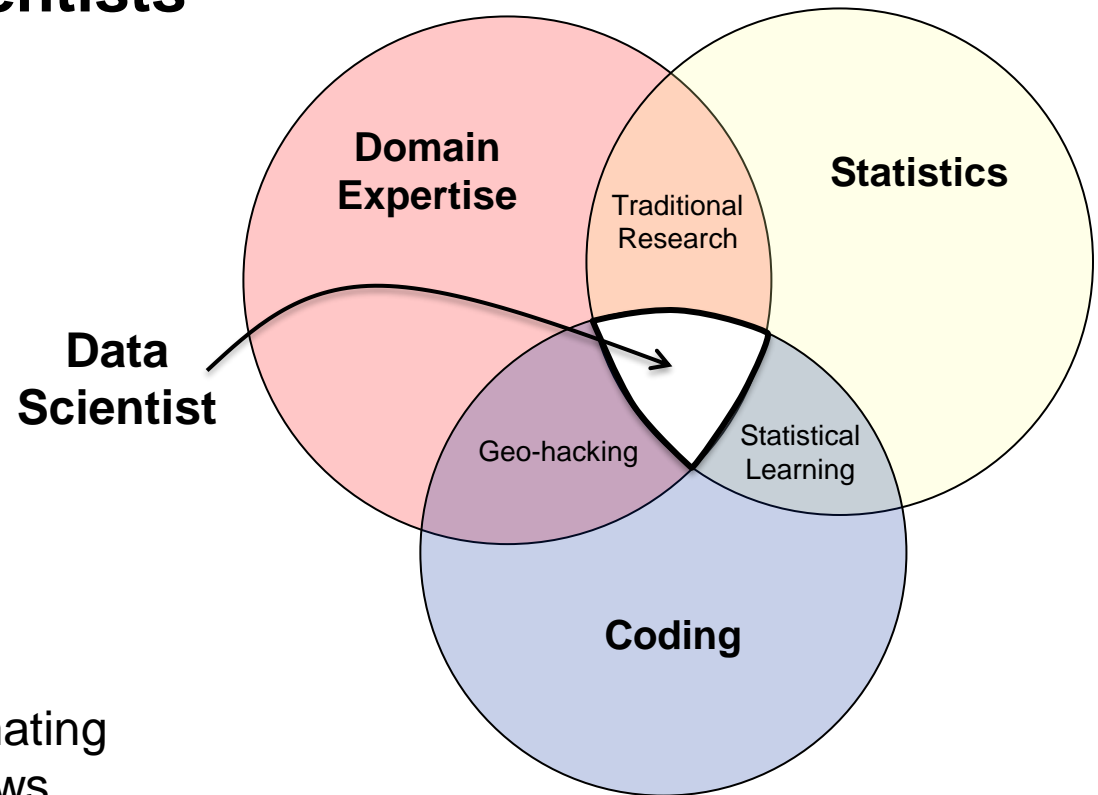


# Developing Operational Capability

## We need Data Scientists

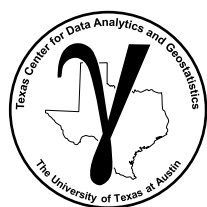
### Intersection of:

- Domain Expertise
  - Geoscience
  - Engineering
- Statistics
  - Probability
  - Data Analytics
- Coding
  - Scripting and Automating
  - Prototyping Workflows



Venn diagram for the data scientist.





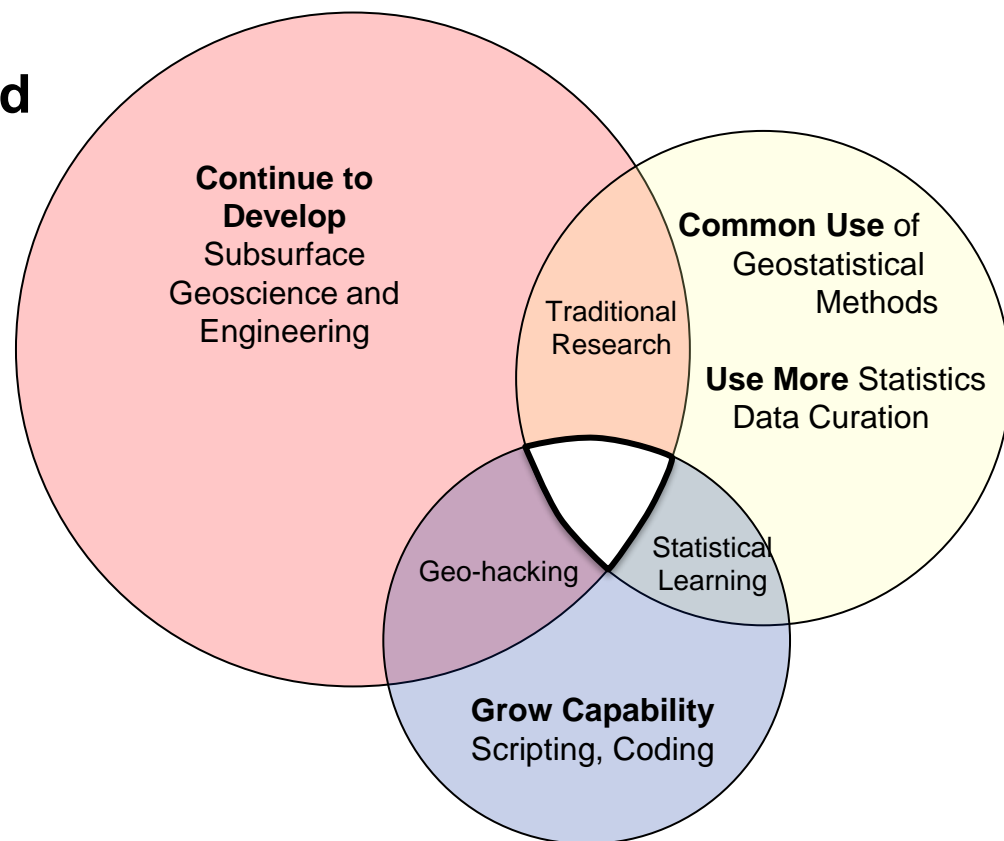
# Developing Operational Capability

## Graduate Geoscientists and Engineers with Data Analytics Capabilities

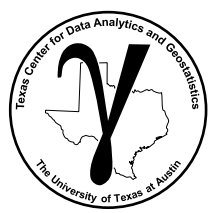
Well-prepared with data-driven knowledge to contribute in our industry

## Build Capability in the Existing Geoscience and Engineering Workforce

Geoscience and engineering capability remains core to our work



Proposed diagram for a path forward for growing data science capabilities among geoscientists and engineers, regions scaled by importance.



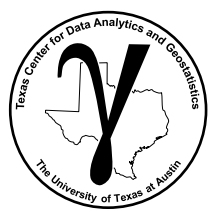
# The Power of Data Analytics

## Statistics to Mitigate Cognitive Biases

- Anchoring Bias: what we see is impacted by anything we have seen recently
- Recency Bias: we weight observations by how recently we saw them
- Confirmation Bias: we tend to see what confirms our current theory

‘I would not have seen it, if I hadn’t believed it!’

- Ashleigh Brilliant



# **PGE 383 Machine Learning**

## **Machine Learning**

**Lecture outline . . .**

- **Machine Learning Overview**
- **Examples of Machine Learning**
- **Energy Machine Learning**

**Michael Pyrcz, The University of Texas at Austin**