

Data Analytics and Geostatistics: Machine Learning



Lecture outline . . .

- General Comments
- Prediction and Inference
- Decision Tree

Introduction

Modeling Prerequisites

Spatial Estimation

Spatial Uncertainty

Multivariate, Spatial

Multivariate Analysis

Machine Learning

Novel Workflows

Conclusions

Instructor: Michael Pyrcz, the University of Texas at Austin

Data Analytics and Geostatistics: Machine Learning



Other Resources:

- Statistical Learning, Dimensional Reduction and Decision Tree

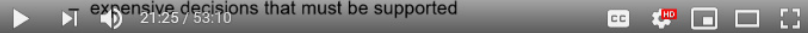


Machine Learning / Statistical Learning



To better utilize data to improve decision-making with consistency and speed.

- Applications in Energy
 1. Feature detection / Guided interpretation in dense data sets like seismic, smart fields / Big data analytics
 2. Optimization of field development decisions
 3. Exploration prioritization
 4. Fast proxies for forecasting
 - Why is Energy different?
 - sparse and uncertain data
 - complicated and heterogeneous systems
 - high degree of irreversible interpretation, engineering physics
- expensive decisions that must be supported



Introduction

Modeling Prerequisites

Spatial Estimation

Spatial Uncertainty

Multivariate, Spatial

Multivariate Analysis

Machine Learning

Novel Workflows

Conclusions

Instructor: Michael Pyrcz, the University of Texas at Austin

Data Analytics and Geostatistics: Machine Learning



Lecture outline . . .

- General Comments

Introduction

Modeling Prerequisites

Spatial Estimation

Spatial Uncertainty

Multivariate, Spatial

Multivariate Analysis

Machine Learning

Novel Workflows

Conclusions

Instructor: Michael Pyrcz, the University of Texas at Austin

Machine Learning / Statistical Learning

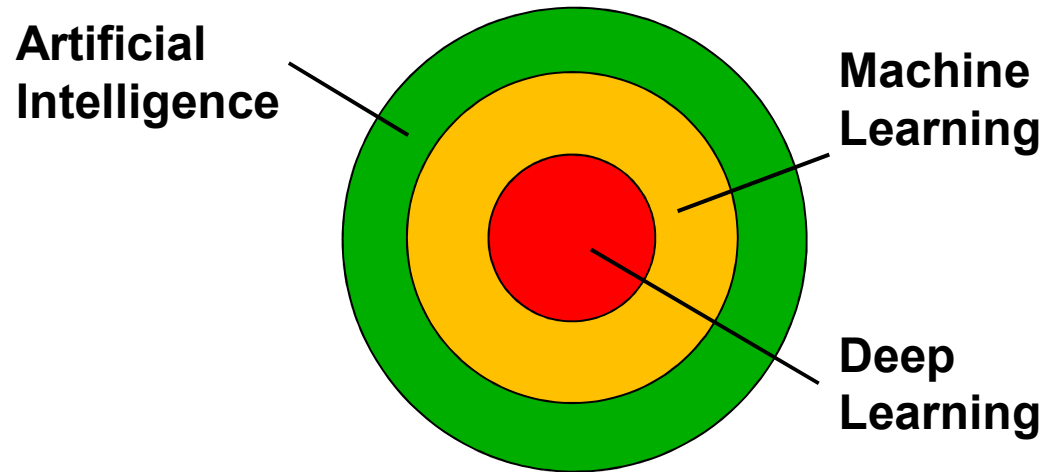


- I'm using this book for this section of the class: An Introduction to Statistical Learning with Applications in R, 2013, James et al., Springer.
(<http://www-bcf.usc.edu/~gareth/ISL/ISLR%20First%20Printing.pdf>)

- Statistical Learning
 - vast set of tools for learning from data
 - based on initial assumptions and hypothesis
- Machine Learning vs. Statistical Learning
 - vast set of tools for learning patterns
 - very little if any prior assumptions
- Supervised Learning
 - building a predictive model for estimating an output given one or more inputs
- Unsupervised Learning
 - all inputs, no output
 - learn from the structures of the data alone

Note: Some consider statistical learning and machine learning to be the same
I'll use them interchangeably

Machine Learning / Statistical Learning



Artificial Intelligence: the theory and development of computer systems able to perform tasks that normally require human intelligence, such as visual perception, speech recognition, decision-making, and translation between languages (Google Dictionary)

Machine Learning: is an application of artificial intelligence (AI) that provides systems the ability to automatically learn and improve from experience without being explicitly programmed (Google Dictionary). Access data and learn for themselves.

Deep Learning: subset of machine learning for unsupervised learning from unstructured, unlabeled data.

Machine Learning / Statistical Learning



Machine Learning:

toolkit

**training
with data**

“is the study of algorithms and mathematical models 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", in order to make predictions or decisions without being explicitly programmed to perform the task.”

learning

general

“where it is infeasible to develop an algorithm of specific instructions for performing the task.”

not a panacea

Machine Learning - Wikipedia

Machine Learning / Statistical Learning



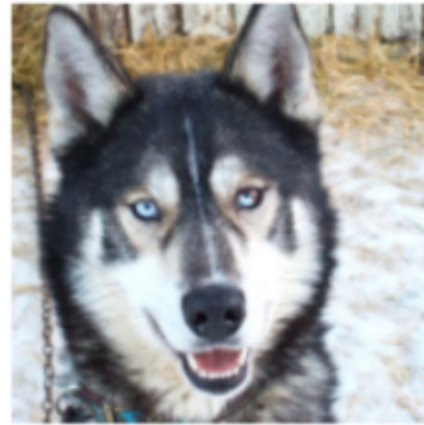
Concerns:

Biased training data

Ribeiro et al. (2016) trained a logistic regression classifier with 20 wolves and dogs images to detect the difference between wolves and dogs.

The problem is:

- interpretability may be low
- application may become routine and trusted
- the machine is trusted, becomes an authority



(a) Husky classified as wolf



(b) Explanation

Figure 11: Raw data and explanation of a bad model's prediction in the "Husky vs Wolf" task.

Image and example from Ribeiro et al., (2016)
<https://arxiv.org/pdf/1602.04938.pdf>

Big Data



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

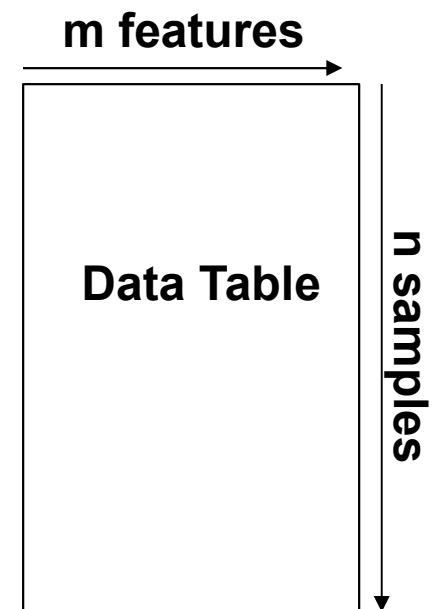
Volume: large number of data samples, large memory requirements and difficult to visualize

Velocity: data is gathered at a high rate, continuously 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 before tech learned about big data.”

– Michael Pyrcz

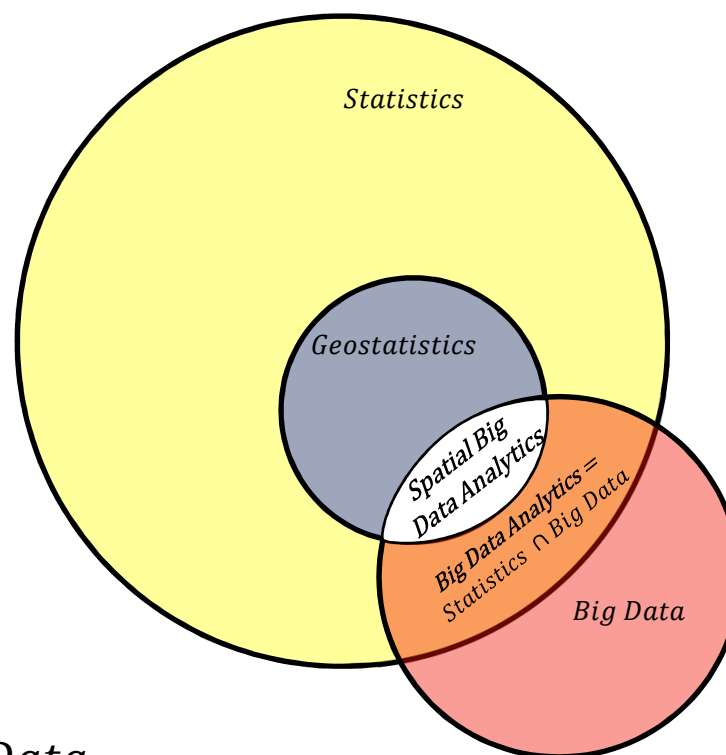
Big Data Analytics – methods to explore and detect patterns, trends and other useful information from big data to improve decision making.

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, and (4) uncertainty.

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



Proposed Venn diagram for spatial big data analytics.

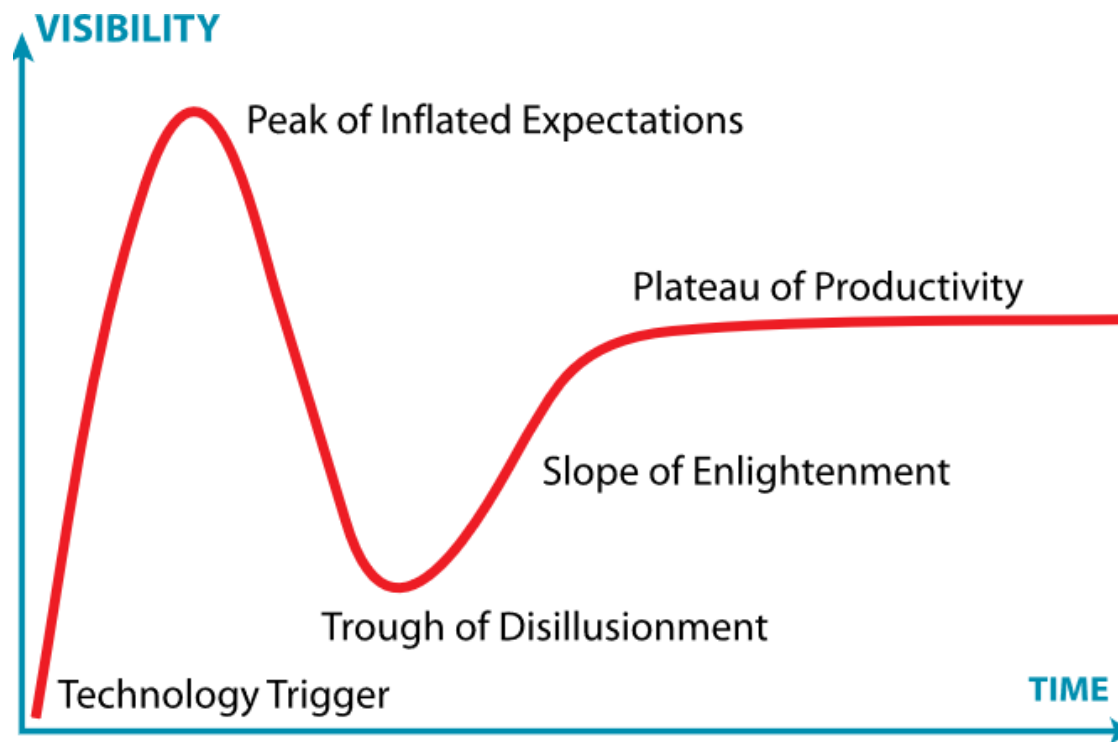
$$\textit{Spatial Big Data Analytics} = \textit{Geostatistics} \cap \textit{Big Data}$$

Big data analytics is expert use of (geo)statistics on big data.

Machine Learning / Statistical Learning



- Hype Cycle – from information technology firm Gartner (https://en.wikipedia.org/wiki/Hype_cycle)



Where are we currently for machine learning?

Machine Learning / Statistical Learning



- Applications Around You / Societal Impacts
 1. Driving directions that crowd source and update improve traffic flow
 2. Air traffic routing
 3. Spam filters
 4. Plagiarism checkers
 5. Translation / computer reading
 6. Credit card fraud detection
 7. Face recognition (Facebook, Snapchat etc.)
 8. Recommendations (Amazon, Netflix, YouTube)
 9. Smart personal assistants

Machine Learning / Statistical Learning



To better utilize data to improve decision-making with consistency and speed.

- Applications in Energy
 1. Feature detection / Guided interpretation in dense data sets like seismic, smart fields / Big data analytics
 2. Optimization of field development decisions
 3. Exploration prioritization
 4. Fast proxies for forecasting
- Why is Energy different?
 - sparse and uncertain data
 - complicated and heterogeneous systems
 - high degree of irreversible interpretation, engineering physics
 - expensive decisions that must be supported

Machine Learning / Statistical Learning



- Just like spatial statistics / geostatistics, statistical learning is a set of tools to add to your tool box as an engineer
- Each is very dangerous to use as a black box. You will need to understand what's under the hood
 - methods, workflows, assumptions and limitations.
 - scope and trade offs between alternative methods
- Imagine you are a carpenter (all geostatistics workflows) (Pyrz and Deutsch, 2014).
 - You would have a tool box
 - You would know each tool perfectly well
 - Understand performance over a variety of applications
 - You would understand the range of applications, weaknesses, strengths, limits.
 - Choice between tools would be based on expert judgement of circumstances and goals of a project
 - You would choose specific tools to have ready for use and other for more rare circumstances
 - Too few tools and a box overwhelmed with obscure tools are both issues.

Data Analytics and Geostatistics: Machine Learning



Lecture outline . . .

- Prediction and Inference

Introduction

Modeling Prerequisites

Spatial Estimation

Spatial Uncertainty

Multivariate, Spatial

Multivariate Analysis

Machine Learning

Novel Workflows

Conclusions

Instructor: Michael Pyrcz, the University of Texas at Austin



The Model

- Predictors, Independent Variables, Features
 - input variables
 - for a model $Y = f(X_1, \dots, X_m) + \epsilon$, these are the X_1, \dots, X_m
 - note ϵ is a random error term
- Response, Dependent Variables
 - output variable
 - for a model $Y = f(X_1, \dots, X_m)$, this is Y
- Statistical / Machine Learning is All About
 - Estimating f for two purposes
 1. Prediction
 2. Inference



Prediction

- Estimating, \hat{f} , for the purpose of predicting \hat{Y}
 - We are focused on getting the most accurate estimates, \hat{Y}
 - We may not even understand what is happening between the X 's!
 - We are concerned about the relationships between X and Y
- Accuracy of \hat{Y} depends on reducible and irreducible error
 - \hat{f} is not a perfect model. Error due to the estimate of f is reducible error
 - but even if we had f , $\hat{Y} = f(X)$, prediction would still have error
 - This is because Y is a function of ϵ , $Y = f(X) + \epsilon$
 - and ϵ is irreducible

Inference

- There is value in understanding the relationships
 - for $Y = f(X_1, \dots, X_m) + \epsilon$ we can understand the influence / interactions of each X_α on Y *and* each other.
- 1. Which predictors are associated with the response?
 - a) What data to collect? Value of information.
 - b) What data to focus on? Simplification of the model. Communication. Big hitters.
- 2. What is the relationship between each response and each predictor?
 - a) sense of the relationship (positive or negative)?
 - b) shape of relationship (sweet spot)?
 - c) relationships may depend on values of other predictors!
- 3. Can the relationship be modeled linearly?
 - a) much simplified
 - b) very low parametric representation
 - c) use multiGaussian?



Estimating f

- Parametric Methods

- make an assumption about the functional form, shape
- we gain simplicity and advantage of only a few parameters
- use training data to fit or train the model
- test the model with withheld test data
- 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!

- Model fitting

- Apply training data
- Solve for least squares solution for coefficients $\beta_0, \beta_1, \dots, \beta_m$



Estimating f

- Nonparametric Methods
 - make no assumption about the functional form, shape
 - estimate f that approaches the data without being too rough
 - more flexibility to fit a variety of shapes for f
 - less risk that \hat{f} is a poor fit for f
 - does not reduce the problem to estimating a small set of parameters
 - » Typically need a lot more data for an accurate estimate of f
 - » Risk of overfitting is greater
 - » May also be parameter rich (e.g. a decision tree as a set of thresholds and averages)
 - » Lacks a compact expression for the model



Training and Testing

Training Phase

- The training subset of the data is applied to select the model parameters (fit the model) usually optimized to minimize the mean square error.

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

Testing Phase

- Apply the model to the testing data (data withheld from training)
- Optimize the model hyperparameters (e.g. complexity) to minimize mean square error with the testing data

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

Do not use all data to train or you will likely overfit to the data and not predict well with new data. Various methods, **k-fold cross validation** is common.

Prediction Accuracy vs. Model Interpretability / Explainability



- **Interpretability / Explainability**

- is the ability to understand the model
- how each predictor is associated with the response
- for example, with a linear model is very easy to observe the influence of each predictor on the response
- but for an artificial neural net it is very difficult

Complexity / Flexibility



Complexity / Flexibility

- Consider these potential polynomials \hat{f} to predict \hat{Y}

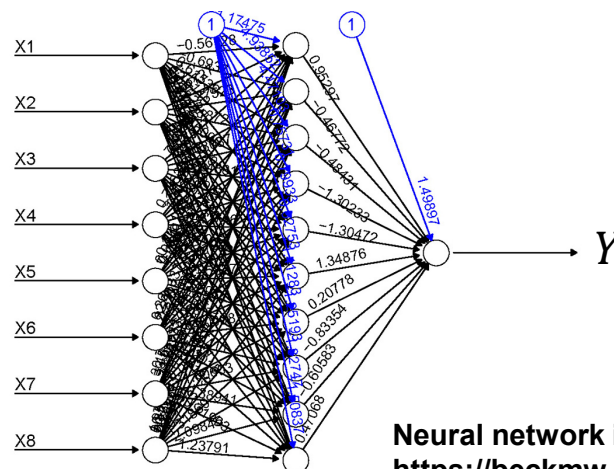
$$Y = \beta_0 + \beta_1 X$$

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_4 X^4 + \beta_5 X^5 + \beta_6 X^6$$

- The 6th order polynomial is more complicated and more flexible to fit the relationship between feature, X , and response, Y
- Now, what if we use 8 bins on X and 10 nodes in a hidden layer of a neural net?:

Indicator Code X into Bins

$$I(x; x_k) = \begin{cases} 1, & \text{if } x \in X_k \\ 0, & \text{otherwise} \end{cases}$$



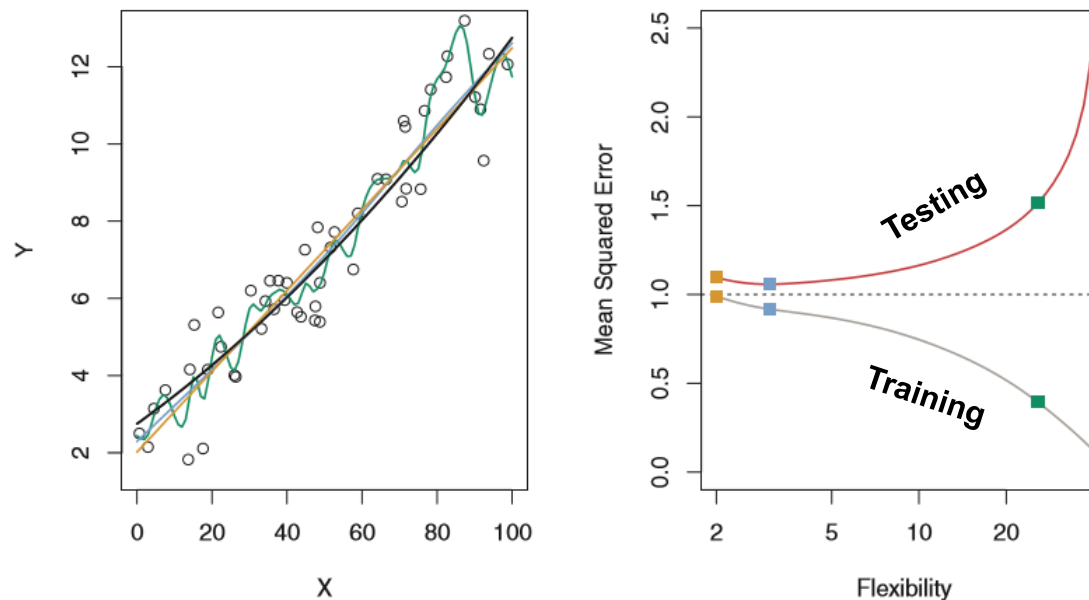
Neural network in R image from:
https://beckmw.files.wordpress.com/2013/11/neuralnet_plot.jpg

Assessing Model Accuracy



- **Flexibility vs. Accuracy**

- Increased flexibility will generally decrease MSE on the **training dataset**
- May result in increase MSE with **testing data**
- Not generally a good idea to select method only to minimize training MSE



Data and model fits (left) and MSE for training and testing (right) from James et al. (2013).

- High flexibility + minimize MSE = likely overfit.

Bias and Variance Trade-off



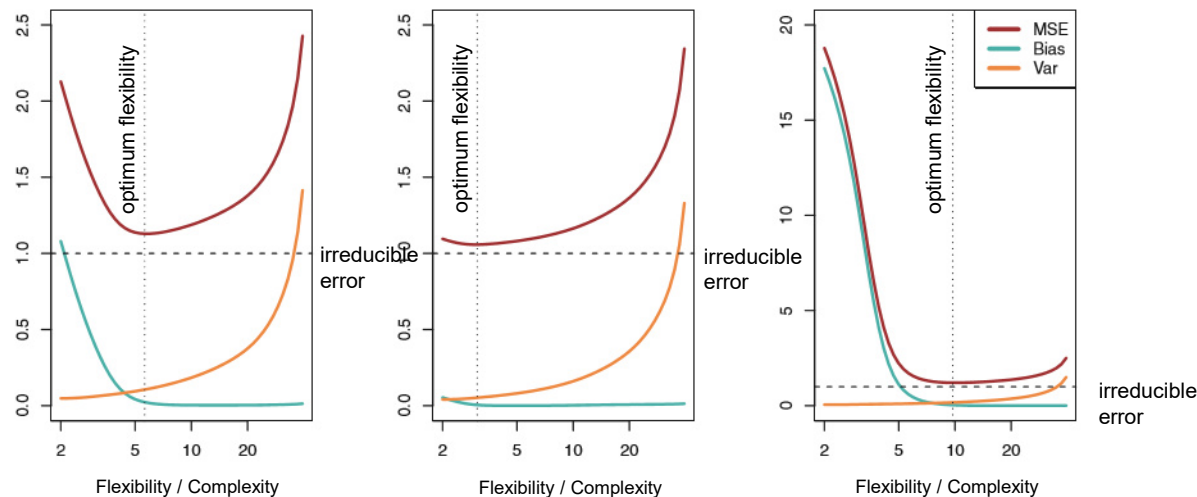
- The **Expected Test Mean Square Error** may be calculated as:

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

Model Variance is the variance if we had estimated the model with a different training set (simpler models \searrow lower variance)

Model Bias is error due to using an approximate model (simpler models \nearrow higher bias)

Irreducible error is due to missing variables and limited samples \Rightarrow can't be fixed with modeling



Model variance, model bias and test MSE for 3 datasets with variable flexibility (Fig 2.12, James et al., 2013), labels added for clarification.

James, G, Witten, D., Hastie, T. and Tibshirani, R., 2013, An Introduction to Statistical Learning with Applications in R, Springer, New York

Data Analytics and Geostatistics: Machine Learning



Lecture outline . . .

- Decision Tree

Introduction

Modeling Prerequisites

Spatial Estimation

Spatial Uncertainty

Multivariate, Spatial

Multivariate Analysis

Machine Learning

Novel Workflows

Conclusions

Instructor: Michael Pyrcz, the University of Texas at Austin



Decision Trees

- Decision trees are used for supervised learning.

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

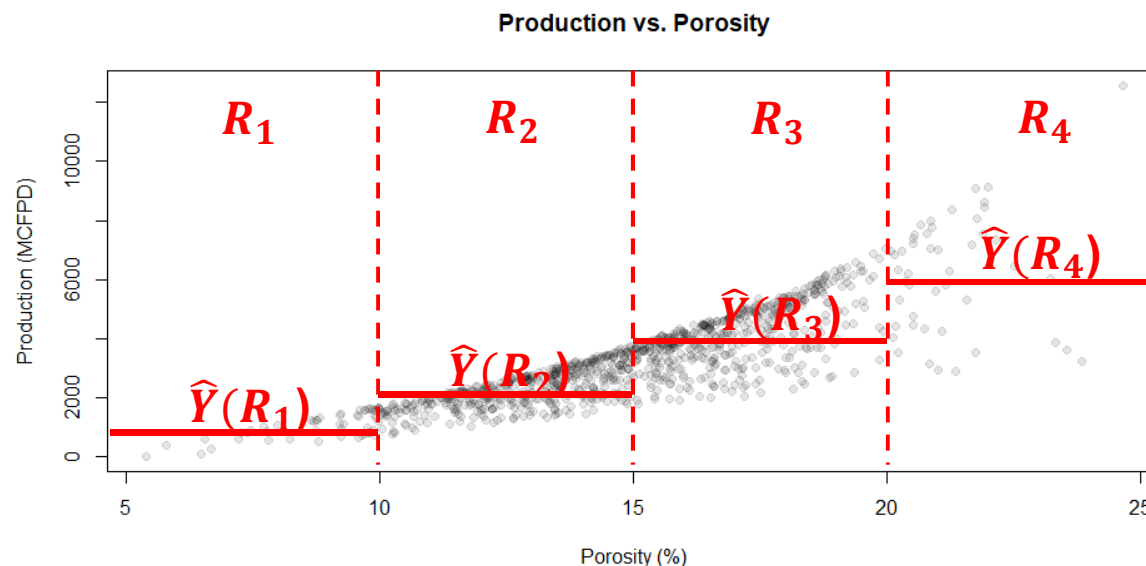
we are predicting a response, Y , from a set of features, X_1, \dots, X_m

- May work with continuous Y for regression tree or categorical Y for classification tree.
- Why cover decision trees?
 - They are not the most powerful, cutting edge method in machine learning
 - But they are likely the most understandable, interpretable
 - Decision trees are expanded with random forests, bagging and boosting to be cutting edge.

Let's learn first about a single tree and then we can comprehend the forest.

Decision Trees

- The fundamental idea is to divide the predictor space, X_1, \dots, X_m , into J mutually exclusive, exhaustive regions
 - mutually exclusive – any combination of predictors only belongs to a single region, R_j
 - exhaustive – all combinations of predictors belong a region, R_j , regions cover entire feature space (range of the variables being considered)
- For every observation in a region, R_j , we use the same prediction, $\hat{Y}(R_j)$
- For example predict production, \hat{Y} , from porosity, X_1

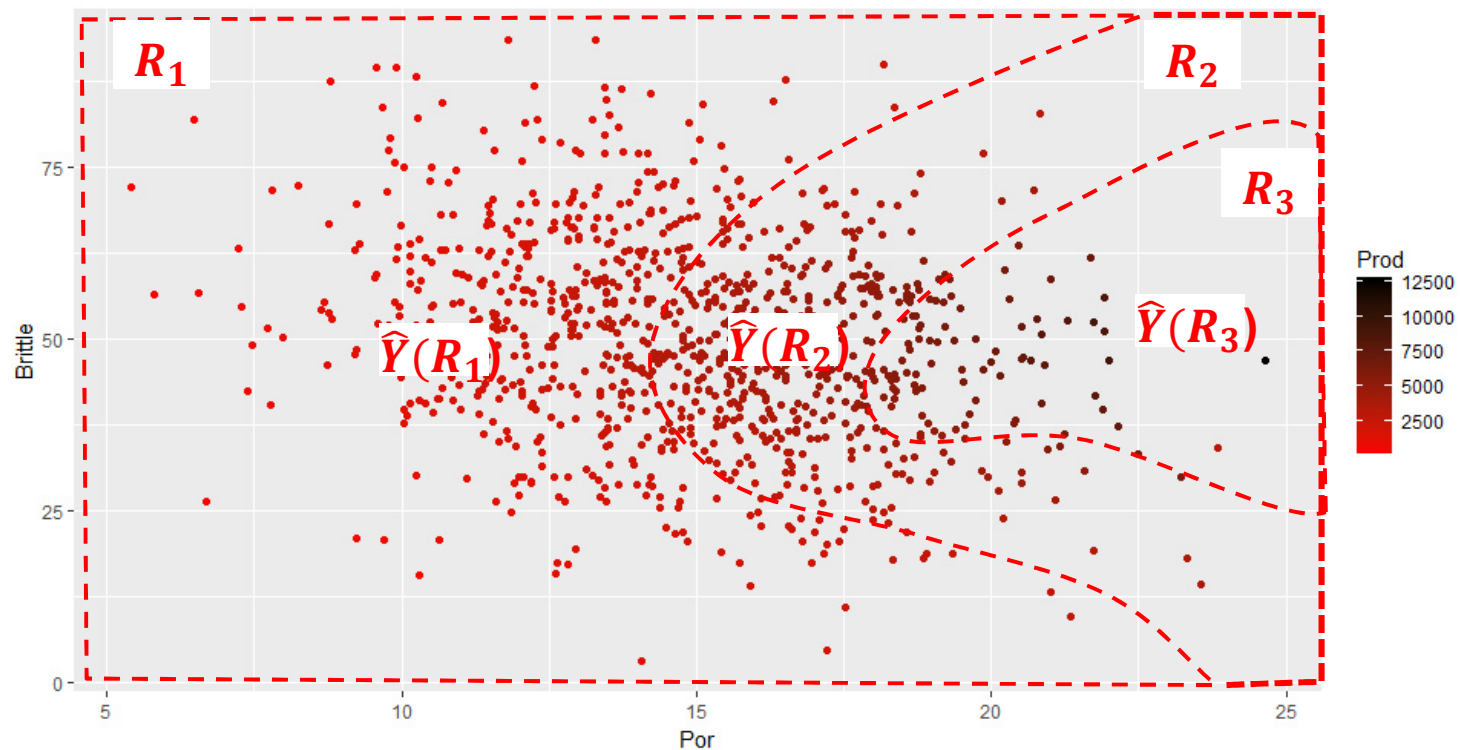


Decision Trees

The Regions



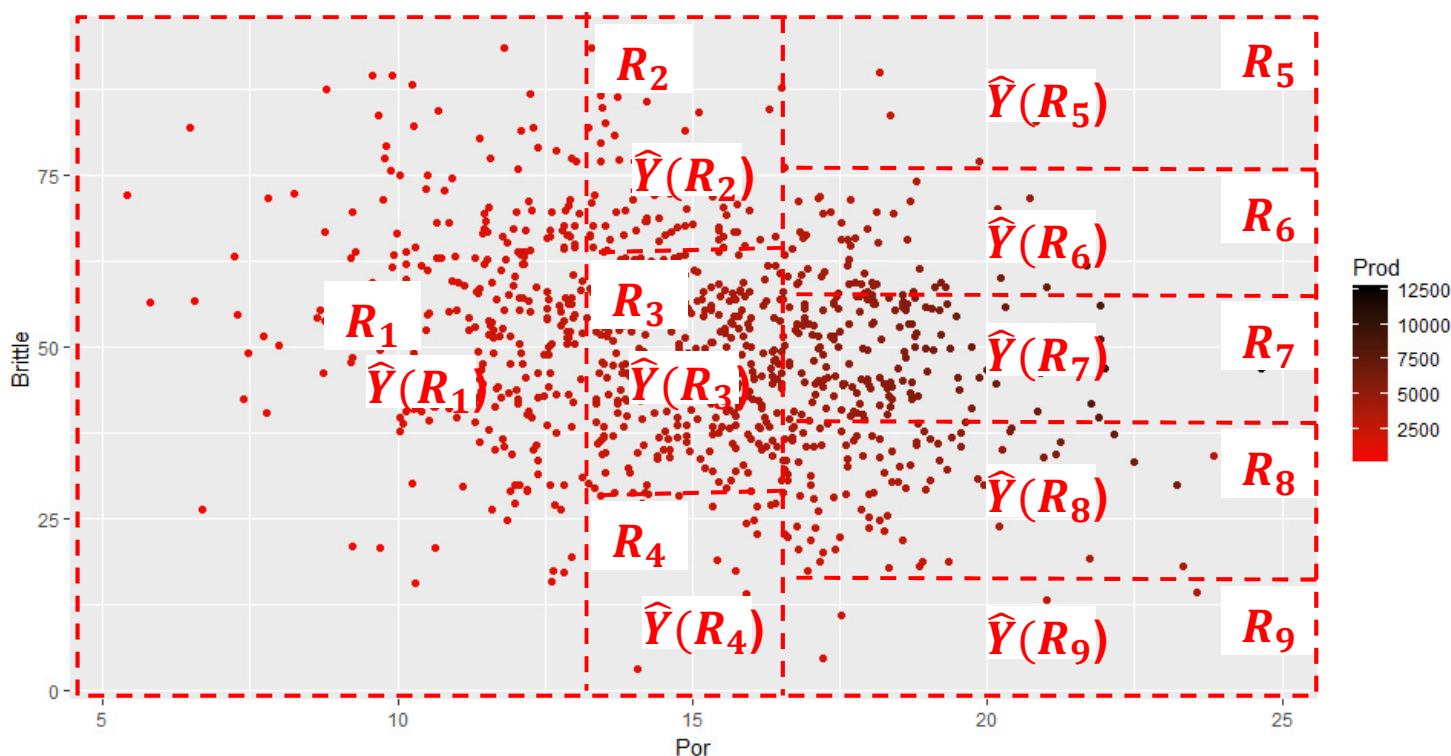
- How do we construct the Regions, R_1, R_2, \dots, R_J ?
 - They could be any shape!
 - Consider the 3 variable problem below.
- Prediction of unconventional well production (MCFPD) from porosity (%) and brittleness (%)



Decision Trees

The Regions

- How do we construct the Regions, R_1, R_2, \dots, R_J ?
 - They could be any shape!
 - Consider the trivariate (3 variable) problem below.
 - We decide to use high-dimensional rectangles or boxes \Rightarrow simple interpretation / rules
 - » Hierarchical segmentation over the features – **very flexible, compact model!**



Prediction of unconventional well production (MCFPD) from porosity (%) and brittleness (%)

Decision Trees

The Regions



How do we construct the Regions, R_1, R_2, \dots, R_J ?

- We want to minimize the Residual Sum of Squares:

$$RSS = \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

looping over J regions and data in each region, $i \in R_j$

- This is the sum of squares of all the data vs. the estimate in their region (the mean of the training data in the region)
- Hint: somehow we need to account for the cost of complexity
 - » We do this through cross validation and pruning

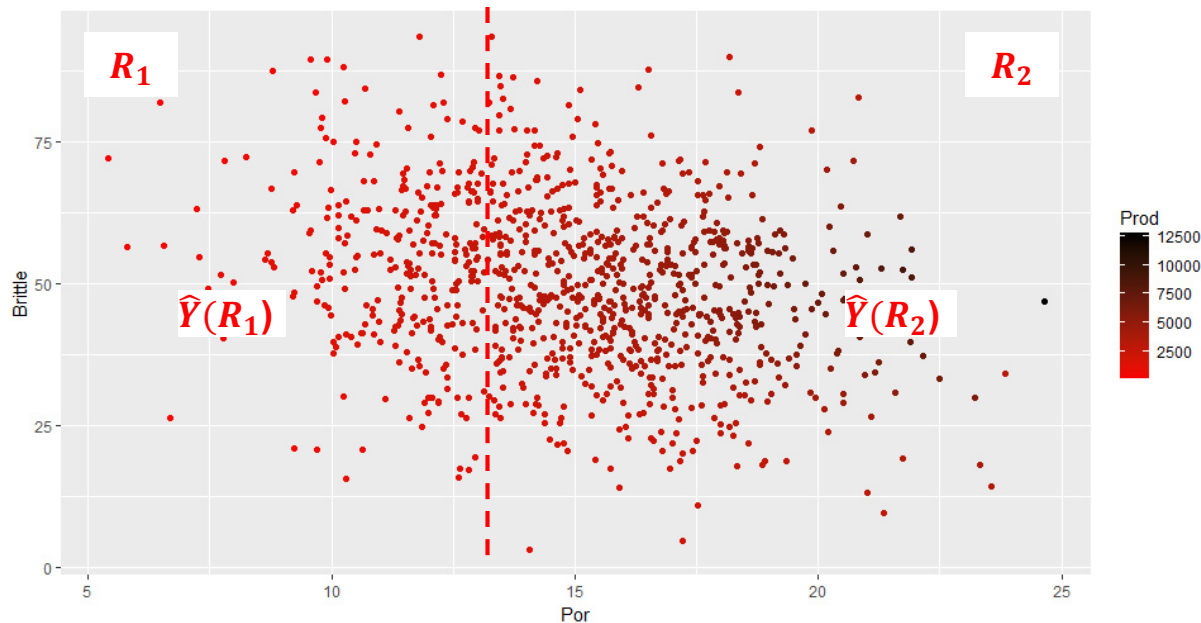
Decision Trees

The Regions



How do we construct the Regions, R_1, R_2, \dots, R_J ?

- Recursive, binary splitting
 - Greedy - at each step the method selects the choice that minimizes RSS. There is no attempt to look ahead, jointly optimize over multiple choices
 - Top-down - at the beginning all data belong to a single region, top of the tree, greedy selection of the single best split over any feature that best reduces the RSS



Prediction of unconventional well production (MCFPD) from porosity (%) and brittleness (%)

Decision Trees

The Regions



How do we construct the Regions, R_1, R_2, \dots, R_J ?

- Let's start with one region, R_1 , with all the training data in it
 - We will place the region boundary based on a threshold, s , inside a this region, j , such that it minimize the RSS.
 - This requires search over all possible thresholds over all features within that region.
 - This is not computationally impossible (not a big space to search)

$$R_{1(m,s)} = \{X | X_m < s\} \text{ and } R_{2(m,s)} = \{X | X_m \geq s\}$$

- X_m are the features and s is the threshold for the segmentation into R_1 and R_2
- We segment such that we minimize the Residual Sum of Squares:

$$RSS = \sum_{i: x_i \in R_1(m,s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: x_i \in R_2(m,s)} (y_i - \hat{y}_{R_2})^2$$

- Then we just repeat for over the two region to find the next best segmentation.

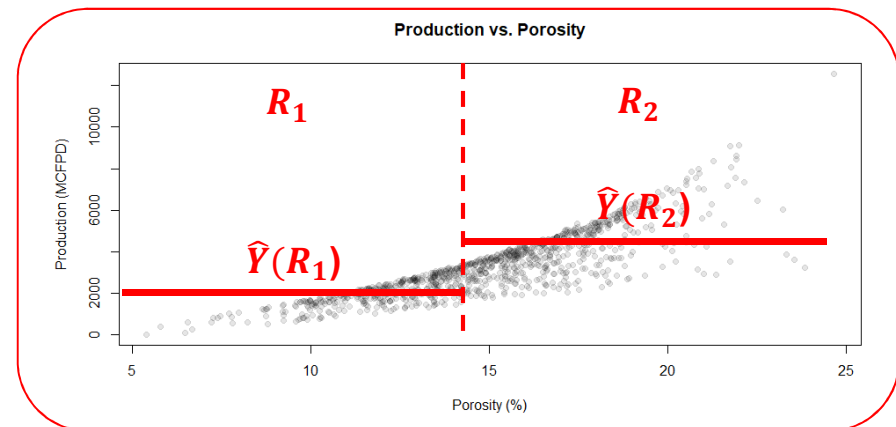
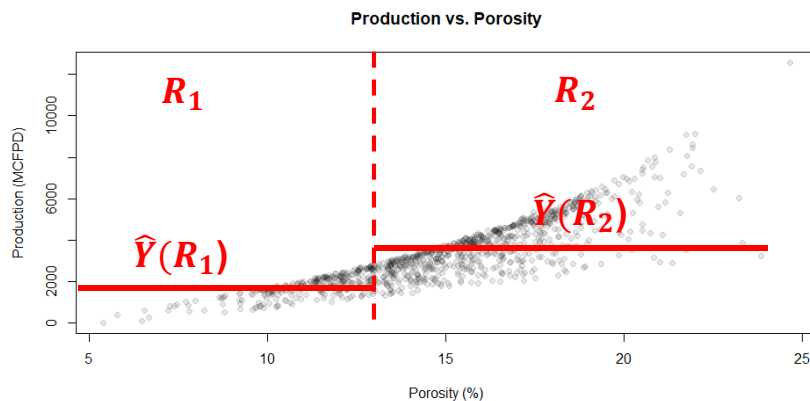
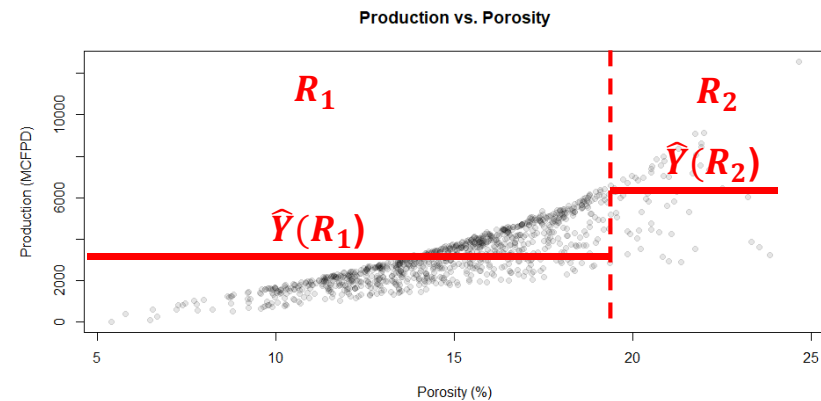
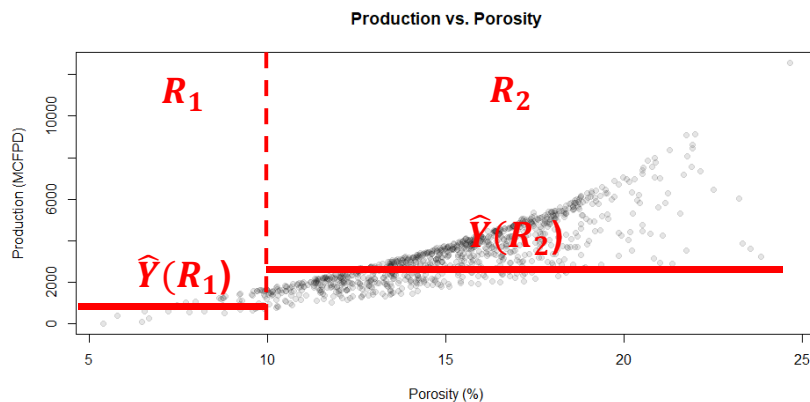
Decision Trees

The Regions



Let's pause and go back to our initial univariate problem and make a tree by hand!

- Where should we split to minimize the error in a tree-based estimate (minimize the residual sum of square)?

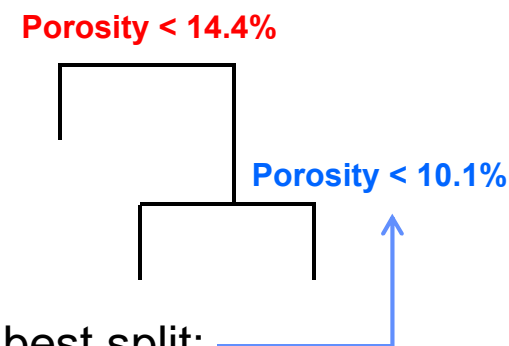
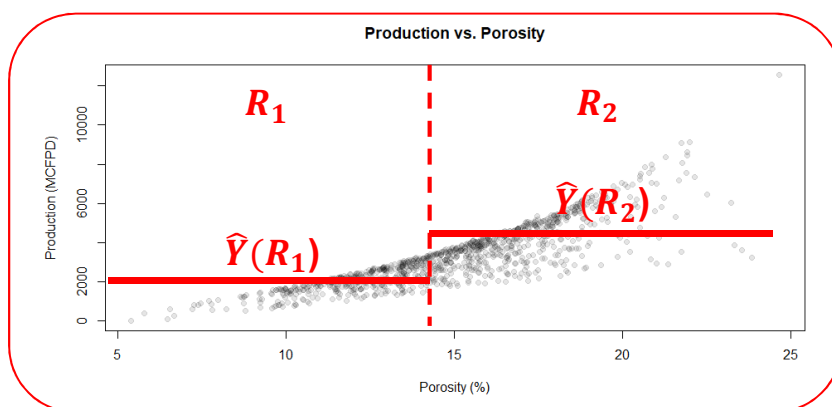


Decision Trees

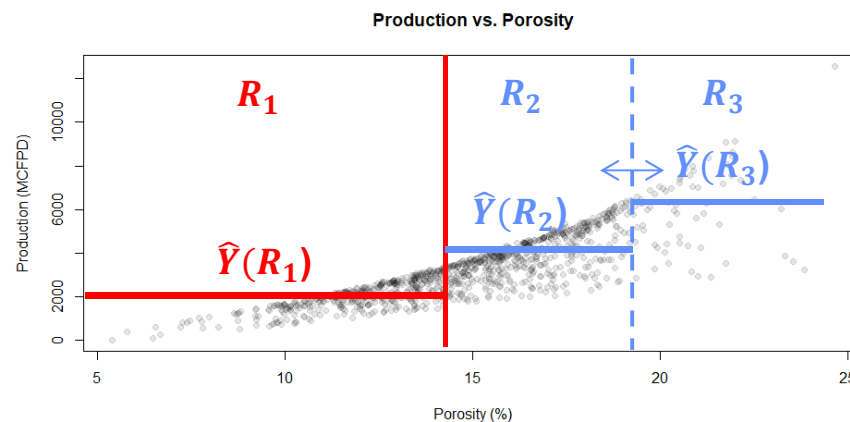
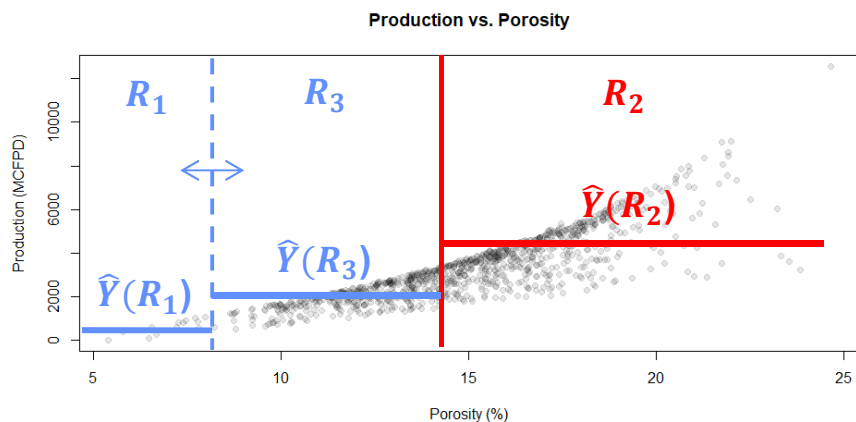
The Regions

Let's pause and go back to our initial bivariate problem and make a tree by hand!

- Found first split, now check for next split the maximizes accuracy



- Search over all regions and variables, to find the next best split:



Decision Trees

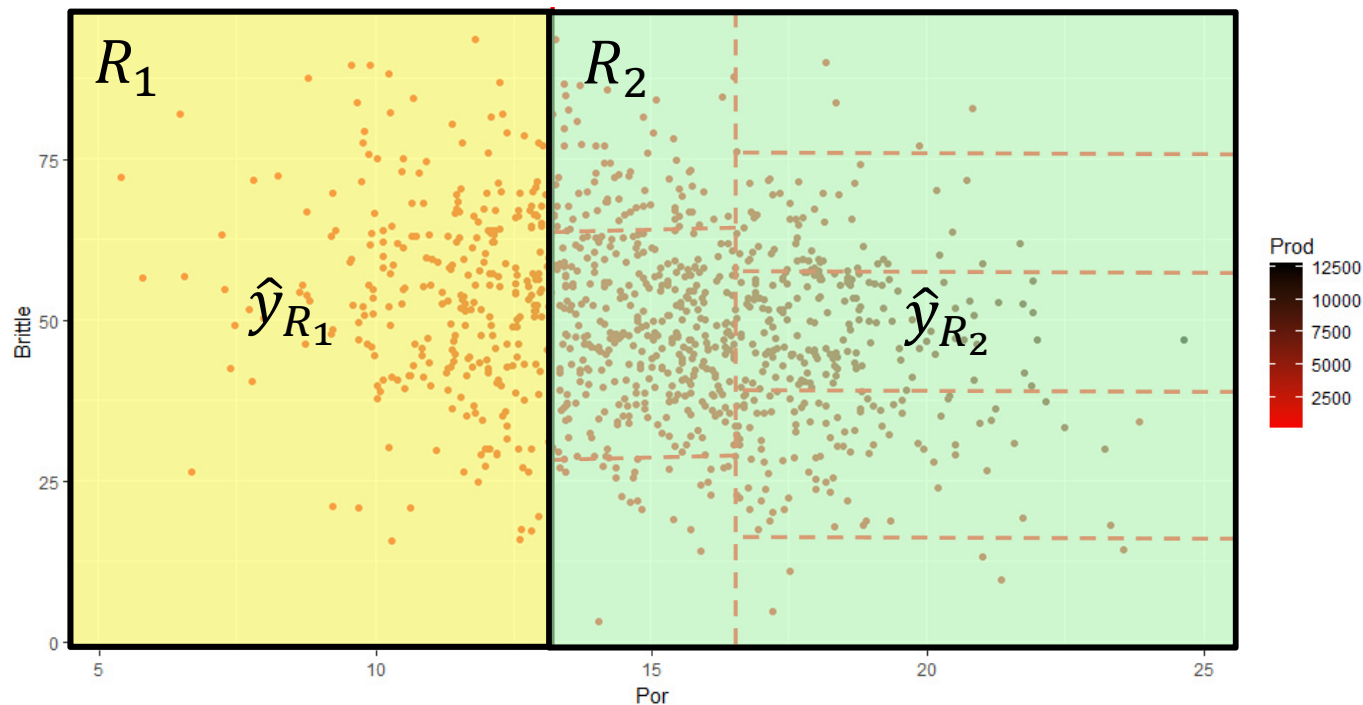
The Regions



How do we construct the Regions, R_1, R_2, \dots, R_J ?

- The we continue sequentially segmenting region with threshold.
 - We will place the region boundaries based on a threshold, s , inside a previous

$$RSS = \sum_{i: x_i \in R_1} (y_i - \hat{y}_{R_1})^2 + \sum_{i: x_i \in R_2} (y_i - \hat{y}_{R_2})^2 + \dots + \sum_{i: x_i \in R_J} (y_i - \hat{y}_{R_J})^2$$



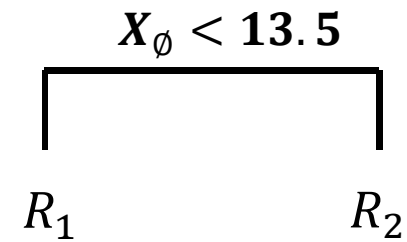
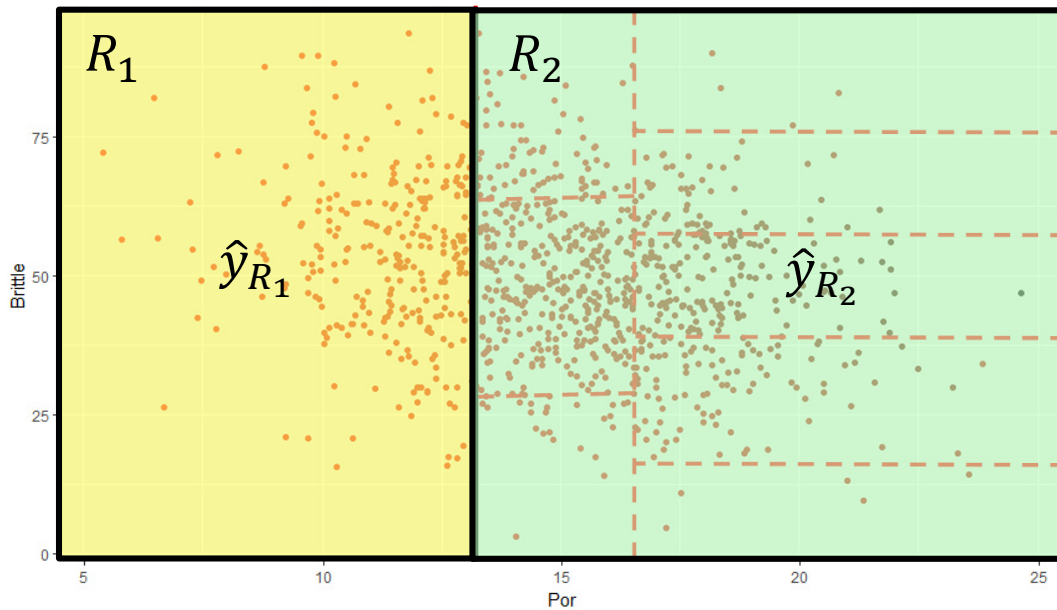
Prediction of unconventional well production (MCFPD) from porosity (%) and brittleness (%)

Decision Trees

The Regions



How do we construct the Regions, R_1, R_2, \dots, R_J ?

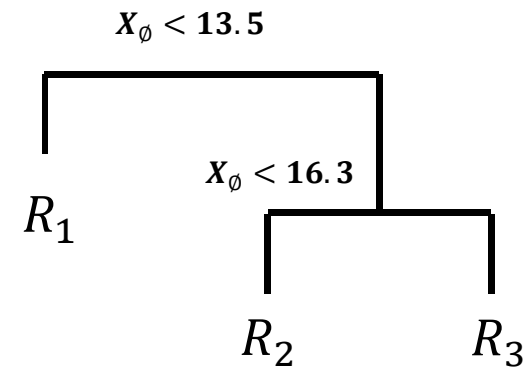
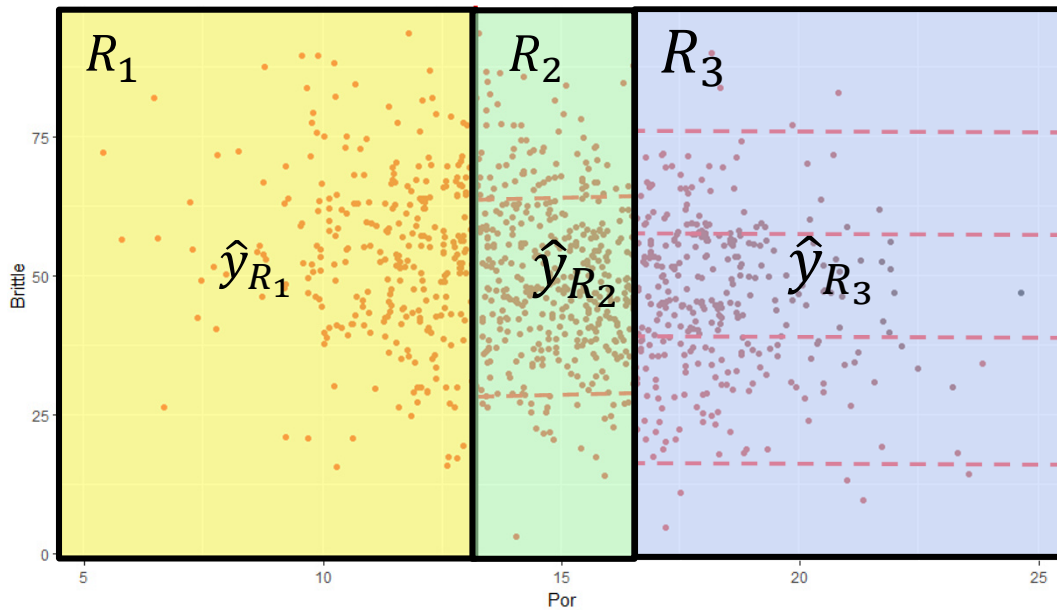


Decision Trees

The Regions



How do we construct the Regions, R_1, R_2, \dots, R_J ?

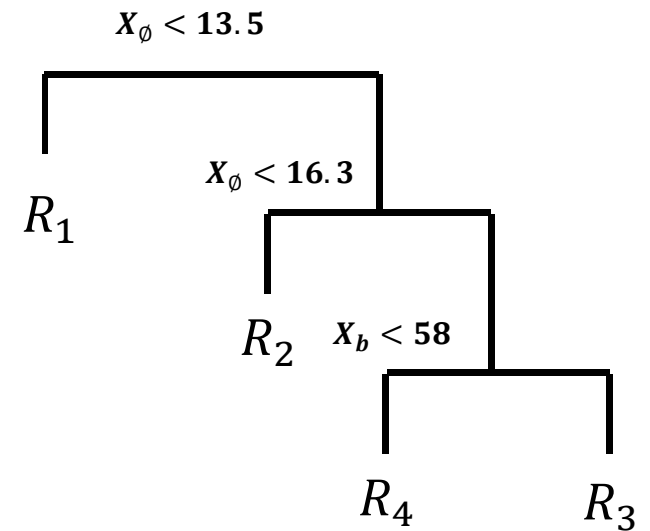
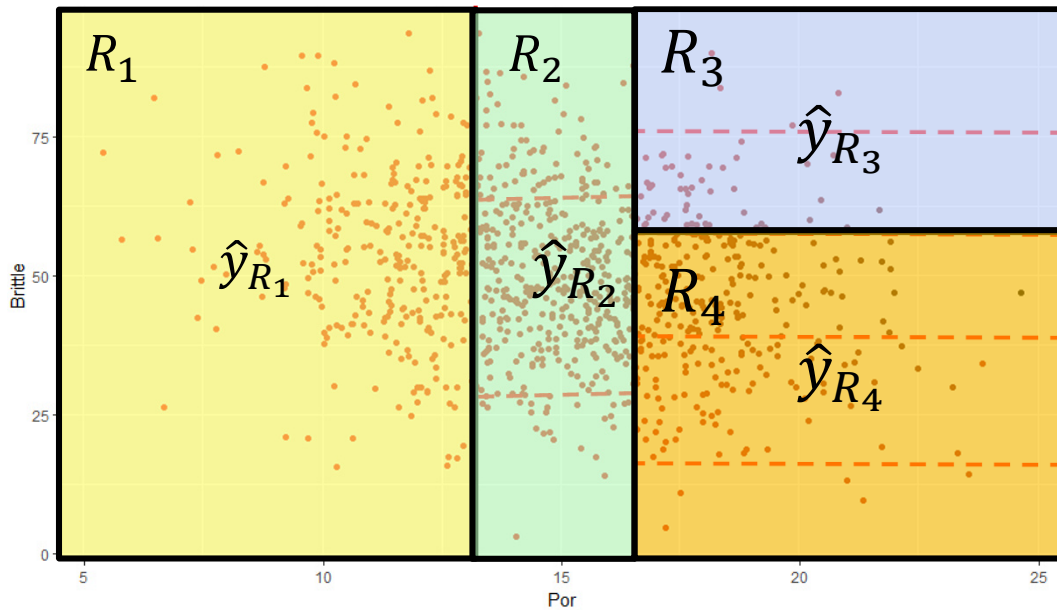


Decision Trees

The Regions



How do we construct the Regions, R_1, R_2, \dots, R_J ?



Decision Trees

The Regions

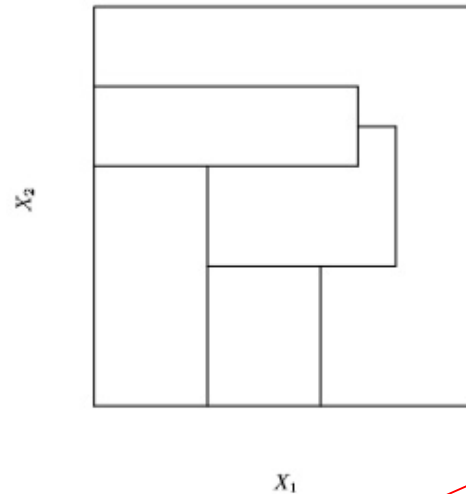


Not from recursive binary splitting

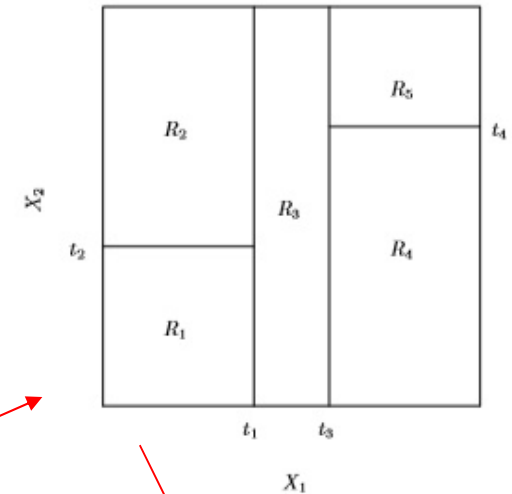
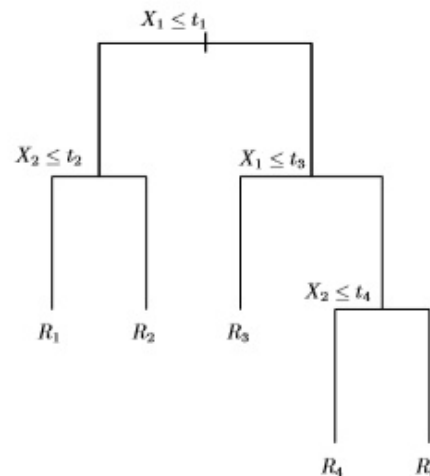
Segmented Feature Space

Example from James et al. (2017)

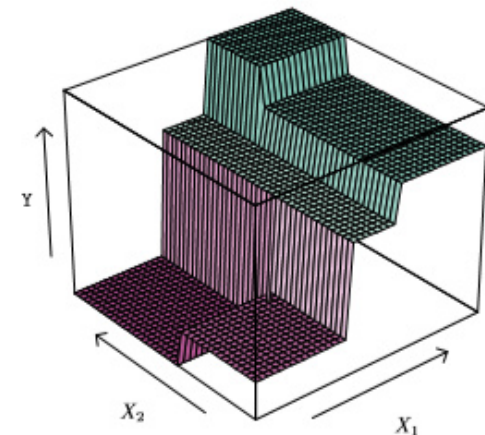
- Top-left 2D feature space partitioning that could not result from recursive binary splitting
- Top-right feature space partitioning, decision tree and estimation surface for feature space.



Decision Tree



Prediction Surface



Decision Trees Termination



When do we stop recursive binary splitting?

- We could continue until every training data value is in its own region!
 - This would be overfit!
- The typical approach is to apply a minimum training data in each region criteria
 - The algorithm stops when all boxes have reached the minimum
- We could continue until we cannot not significantly reduce RSS
 - But the current split could lead to an even better split \Rightarrow short sighted

Decision Trees Pruning



Why do we want a less complicated tree?

- Decision trees, if allowed to grow very complicated are generally overfit.
- It is better to simplify the tree to a smaller tree with fewer splits
 - » lower model variance
 - » better interpretation
 - » with little added model bias
- Limiting tree growth with a high decrease in RSS hurdle is short sighted
- Best strategy is to build a large, complicated tree and then to prune the tree.
 - » We then select the sub tree to provides the lowest test error rate
 - » We cannot consider all possible sub trees (too vast of a solution space)

Decision Trees – Steps



Building a Regression Tree

1. Apply recursive binary splitting to train / grow a large tree with training data, stop when each terminal node has fewer than a minimum number of data or insufficient RSS decrease.
2. Obtain the sequence of best subtrees as a function of complexity (number of terminal nodes) and RSS with training.
3. Use k-fold cross validation to choose the best complexity value. Divide the training observations into K folds. For each fold, $k = 1, \dots, K$:
 - a) Repeats steps from 1-2 on all training excluding those in k fold.
 - b) Evaluate the RSS on the left out data in the k fold.
4. Average the error for each α (K results over each fold) and select complexity (number of terminal nodes) that provides low enough RSS.

K-fold Cross Validation



Cross Validation

- Withhold subset of the data during model training
- Then testing the trained model with withheld subset dataset
- Must make sure cross validation is fair
- Training data set (used for training), Testing data set (withheld for testing)

K-fold Approach

- Select K, for example
- Break data set into K subsets
- Loop over K subsets:
 - use data outside the K part to predict inside the K subset
- Average to summarize the result

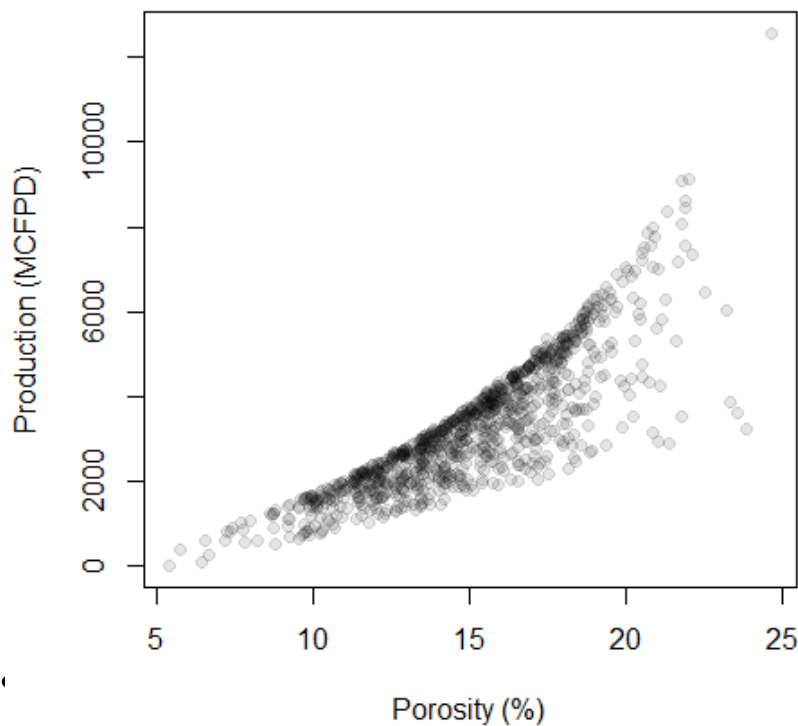
Decision Trees Example



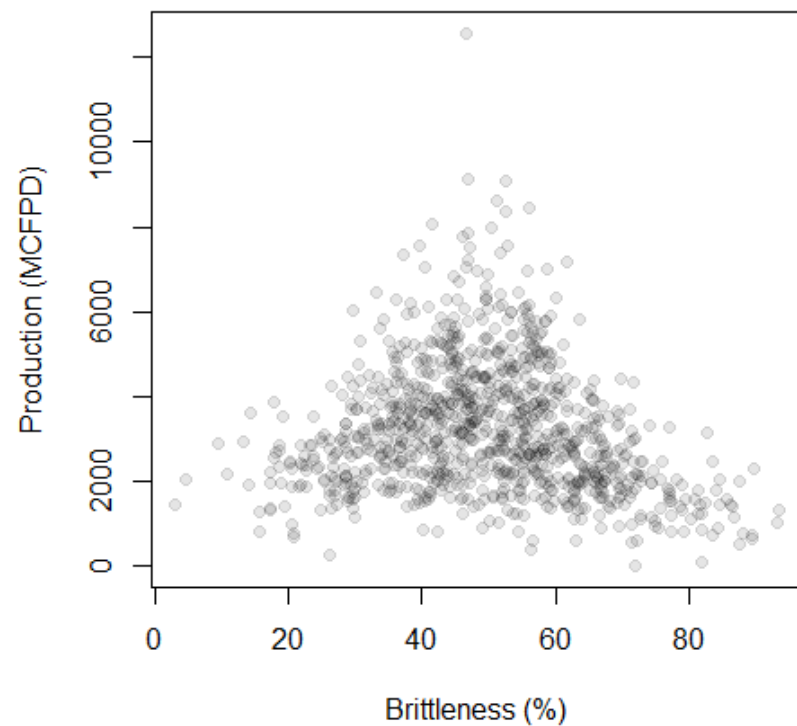
Let's use our Unconventional Multivariate Data

- We added in a production variable for prediction
- Both porosity and brittleness have interesting relationships with production

Production vs. Porosity



Production vs. Brittleness

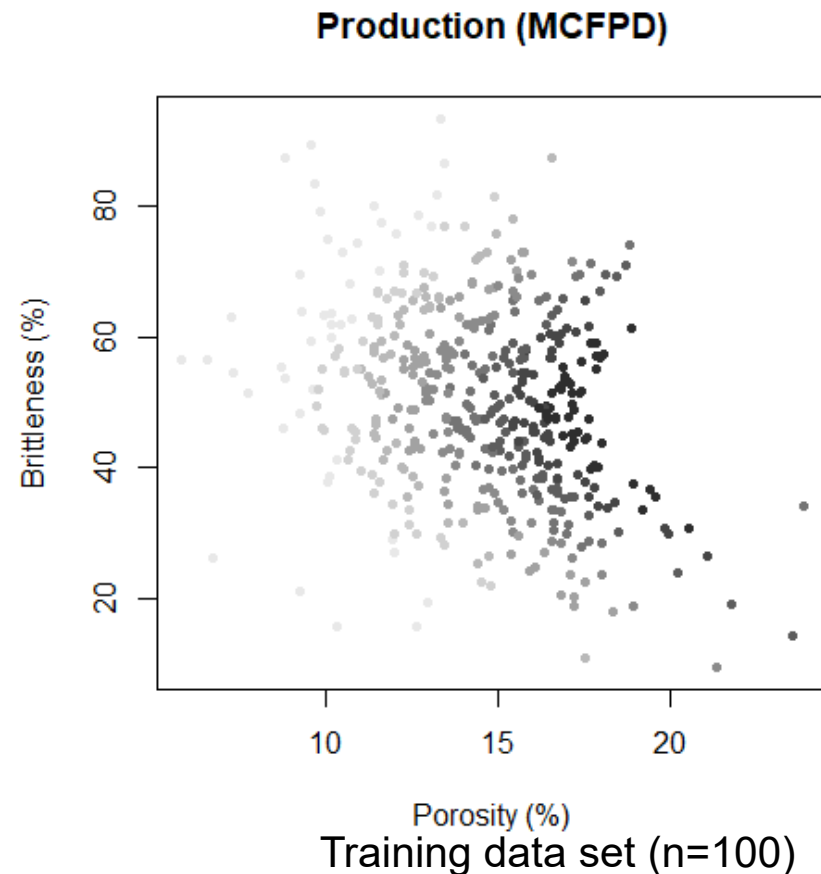
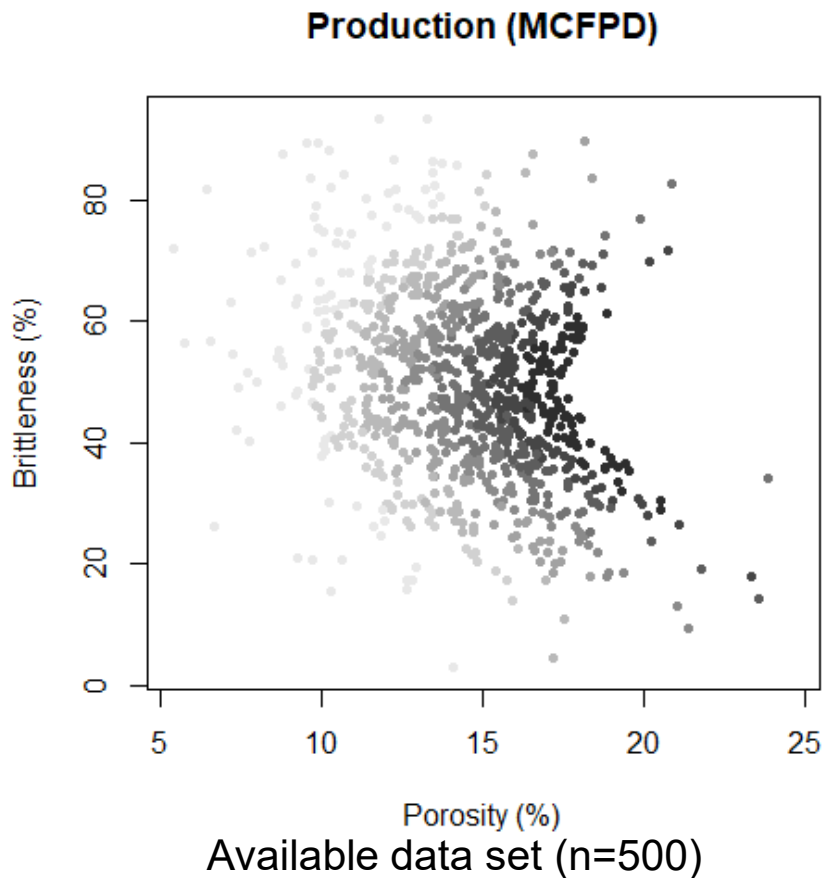


Decision Trees Example



Let's use our Unconventional Multivariate Data

- There is a complicated relationships between porosity, brittleness and production.



Decision Trees Example



Build the initial reasonably complicated tree

- By using the default tree controls we get an 10 terminal node tree.
- We can use the summary command to:

```
Regression tree:
tree(formula = Prod ~ Por + Brittle, data = train, control = tree.control)
Number of terminal nodes: 10
Residual mean deviance: 302900 = 148400000 / 490
Distribution of residuals:
      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-2298.00 -303.50   57.16    0.00  327.50  3668.00
```

- check the complexity of the resulting tree (number of terminal nodes)
- check the summary statistics of the residuals and ensure that the model is not biased (mean = 0.0)
- residual mean deviance is the total residual deviance divided by (the number of observations – number of terminal nodes)
- for a regression trees the total residual deviance is the RSS, reminder:

$$RSS = \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

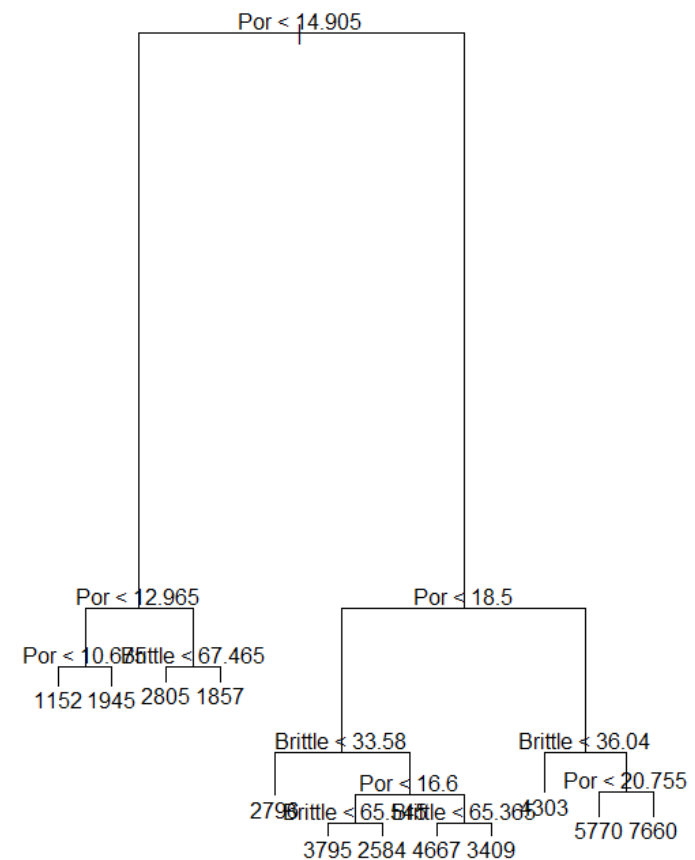
Decision Trees Example



Build the initial reasonably complicated tree

Here's the tree:

- first choice is porosity $<$ or $>$ 14.9%
- we get to the 3rd decision before brittleness is considered
- length of the branches is proportional to decrease in impurity
 - decrease in RSS of the model for regression tree
 - a measure of node heterogeneity for classification trees

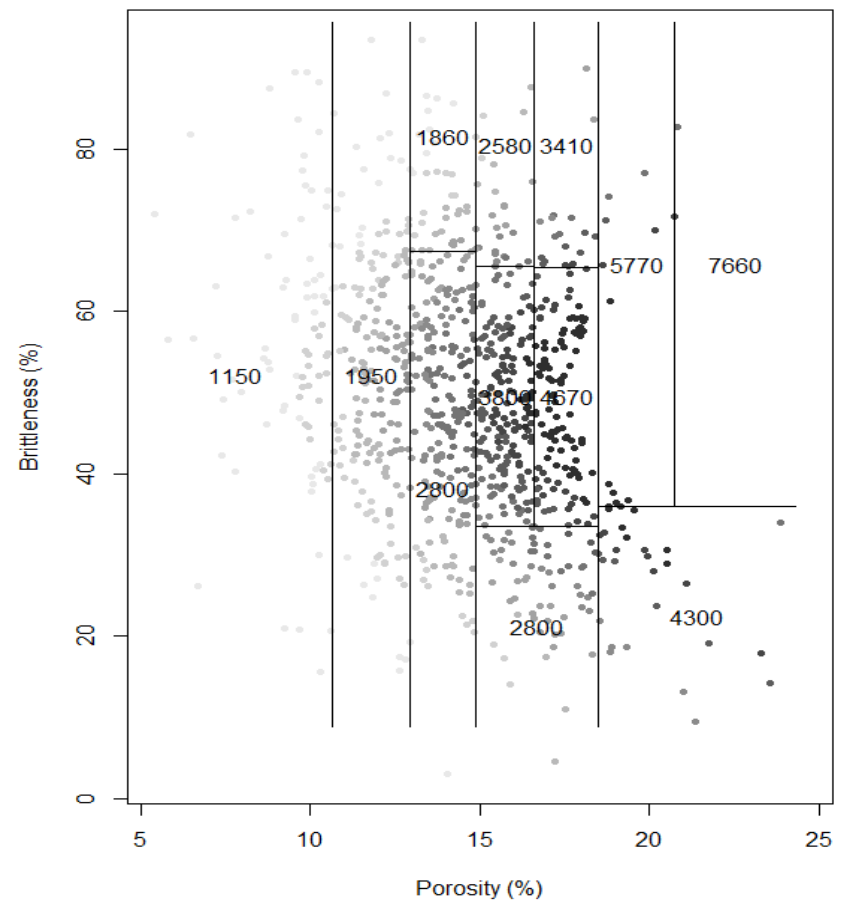


Decision Trees Example



Build the initial reasonably complicated tree

- We can plot the original data and the binary recursive boundaries outlining the various regions and the mean values in each region used as the estimate.



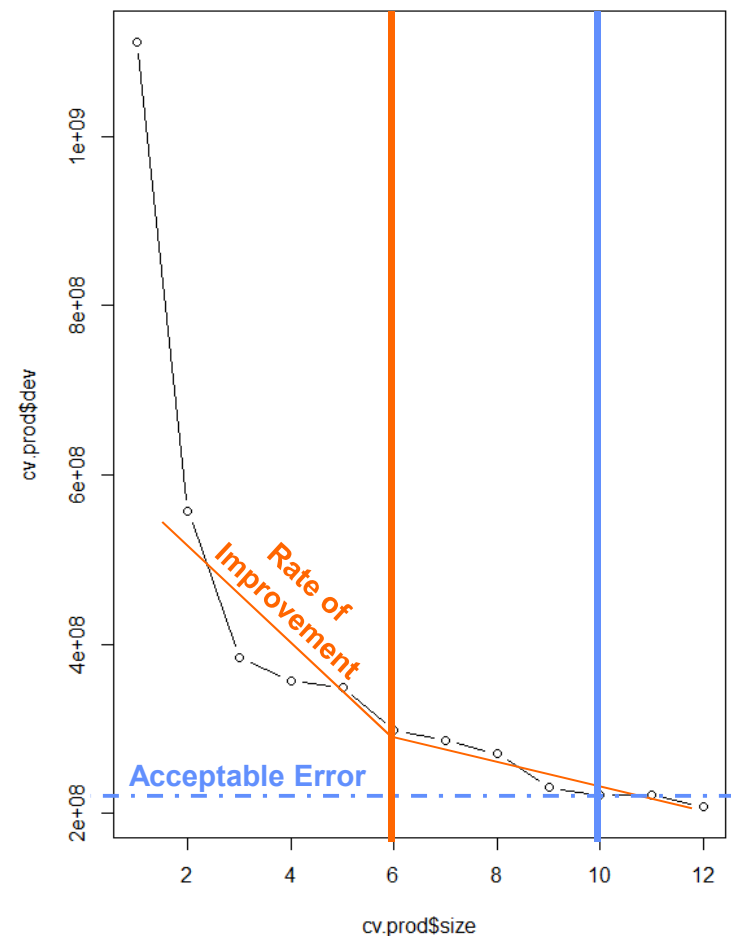
Decision Trees Example



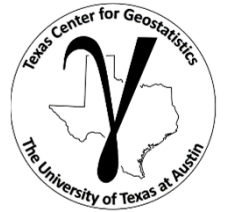
Build the initial reasonably complicated tree

Then we perform k fold cross validation.

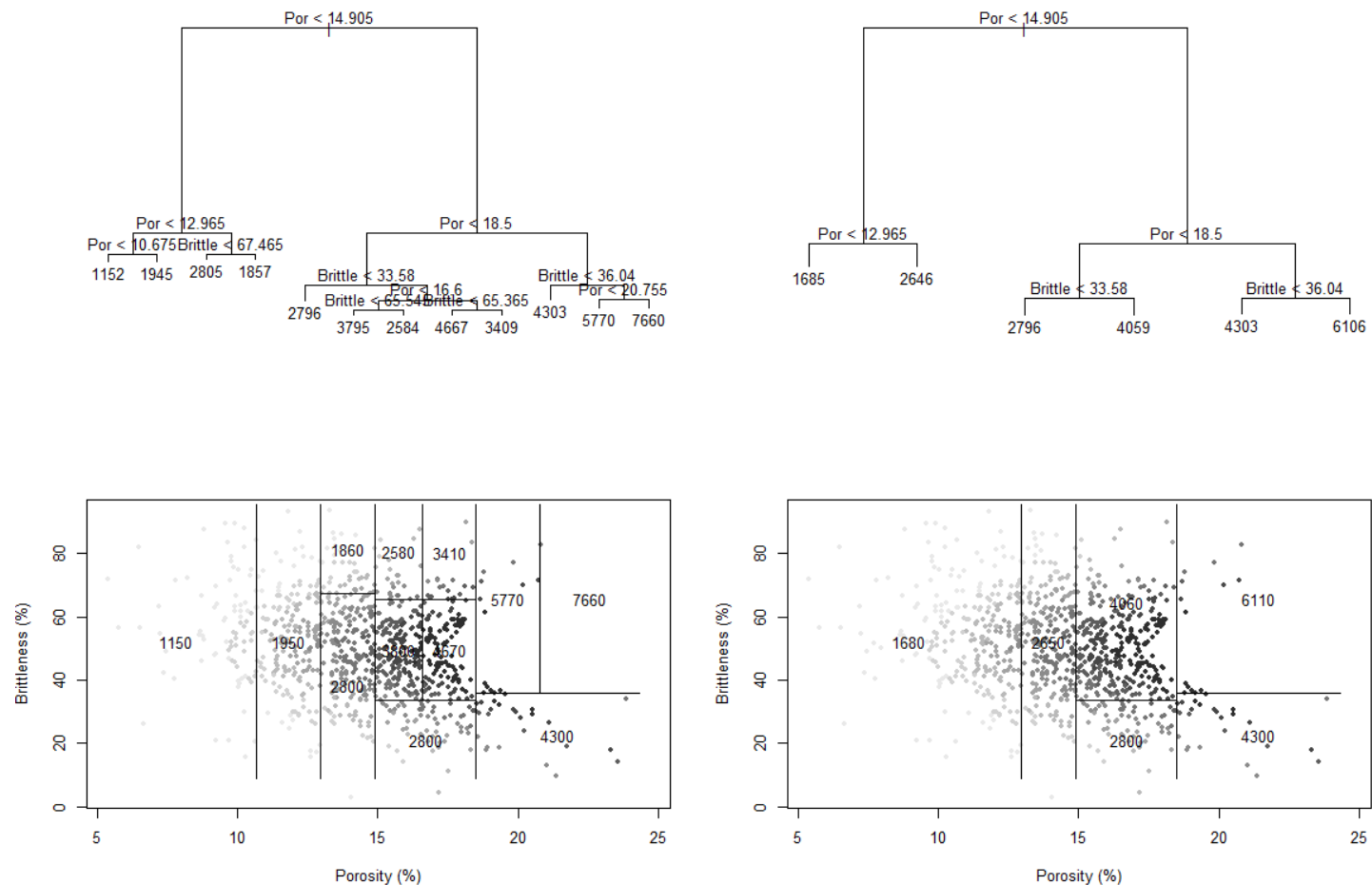
- Decrease tree complexity from 12 nodes (current model) to 1 node (uniform model)
- Calculate the RSS by averaging over k folds of the training data
- We can observe that each additional node improves the model
- Prune complexity based on:
 - Diminishing returns
 - Acceptable level of error



Decision Trees Example



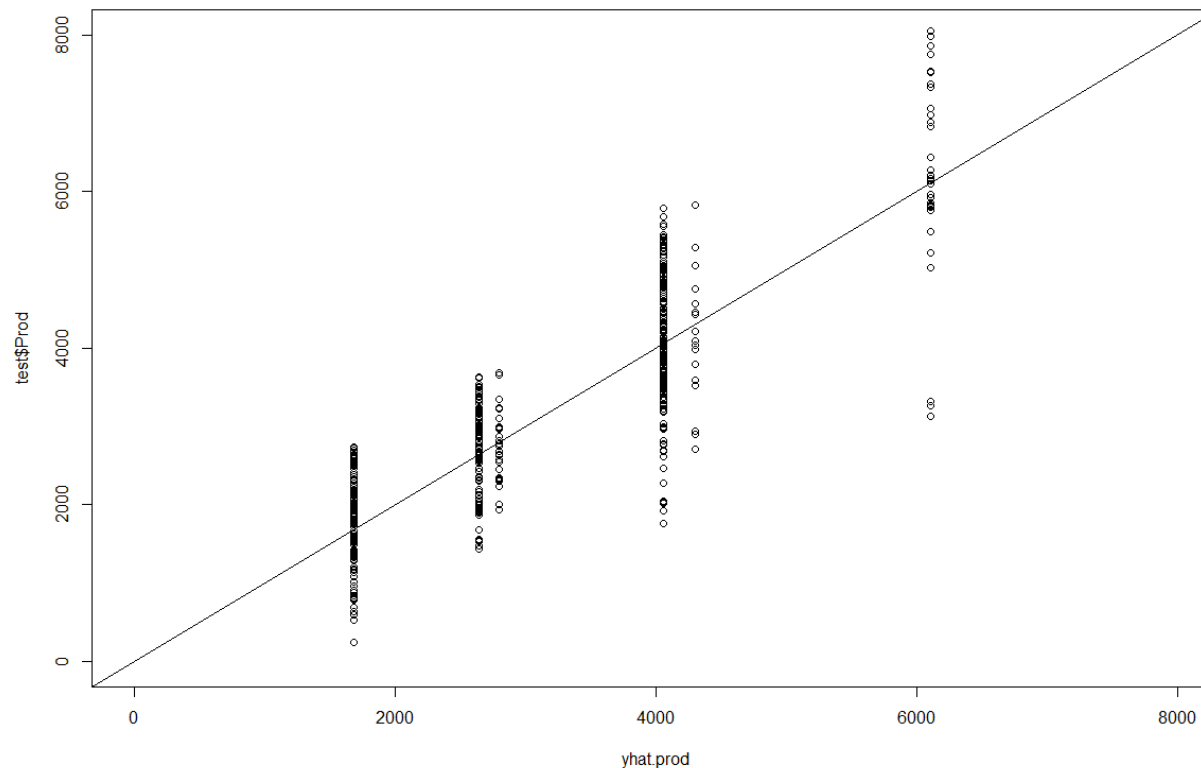
Original and pruned tree:



Decision Trees Example



Cross validation with the testing data set



- Note: the binning due to estimation with the mean of only 6 regions
- We can calculate MSE to assess model accuracy

Decision Trees Demonstration in Python



Decision Tree Tutorial with Subsurface Demonstration in **Python** for Geoscientists and Geo-engineers

Michael Pyrcz, University of Texas at Austin (@GeostatsGuy)



Decision Tree is one of the most simple, explainable and interpretable predictive models in machine learning; therefore, it is a great introduction to regression and classification with machine learning. In addition, the recursive binary segmentation is analogous to human decision making. Finally, understanding a decision tree is a prerequisite for more powerful bagging, random forest and boosting. This tutorial is in Jupyter with Markdown and a realistic unconventional dataset. There is enough documentation that any geoscientists or engineer should be able to try out machine learning.

Decision Tree in Python for Engineers and Geoscientists

Michael Pyrcz, Associate Professor, University of Texas at Austin

Contacts: [Twitter@GeostatsGuy](https://twitter.com/GeostatsGuy) | [GitHub@GeostatsGuy](https://github.com/GeostatsGuy) | www.michaelpyrcz.com | [Google Scholar](https://www.michaelpyrcz.com) | [Book](https://www.michaelpyrcz.com)

This is a tutorial for demonstration of building decision trees in Python with `scikit-learn`. Decision trees are one of the easiest machine learning, prediction methods to explain, apply and integrate. In addition, understanding decision tree-based prediction is a prerequisite to more complicated and powerful methods such as random forest and tree-based bagging and boosting. For this demonstration we use a 1,000 well 7 variable unconventional dataset (file: `uncconv_MV.csv`) that is available on GitHub at <https://github.com/GeostatsGuy/GeostatsData>. The dataset includes 6 predictors (features) and 1 response. We take this multivariate dataset and only retain the three variables (2 predictors and 1 response) for a simple demonstration of the decision tree method. We break the data set into 500 training data and 500 testing data. I used the tutorial in my introduction to Geostatistics undergraduate class (POE337 at UT Austin) as part of a first introduction to geostatistics and Python for the engineering undergraduate students. It is assumed that students have no previous Python, geostatistics nor machine learning experience; therefore, all steps of the code and workflow are explained and described. This tutorial is augmented with course notes in my class. The Python code and markdown was developed and tested in Jupyter.

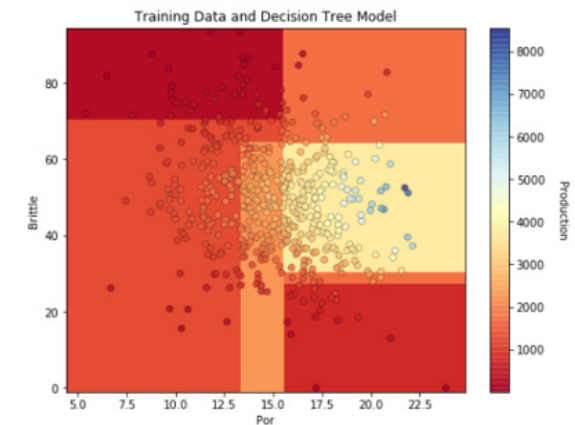
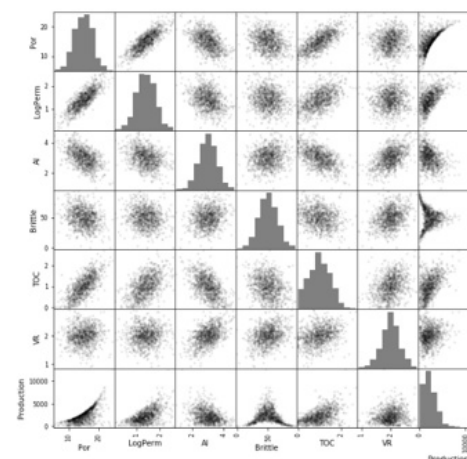
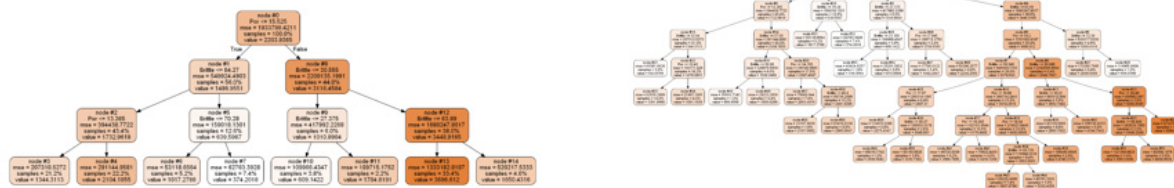
What is a decision tree?

It's make a couple of points about decision trees. For greater detail there are a lot of online resources on decision trees along with the book "An Introduction to Statistical Learning" by James et al. (my favorite).

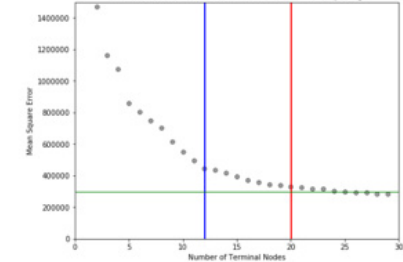
1. method for supervised learning
2. categorical prediction with a classification tree and continuous prediction with a regression tree
3. fundamental idea is to divide feature space into exhaustive, mutually exclusive regions (terminal or leaf nodes in the tree)
4. estimate with the average of data in each region for continuous prediction or the majority category for the data in each region for categorical prediction
5. segment the feature space with hierarchical, binary splitting that may be represented as a decision tree
6. apply a greedy method to find the sequential splits for any feature that minimizes the residual sum of squares

Let's build some decision trees together. You'll get a chance to see the trees and the divided feature space graphically.

| | count | mean | std | min | 25% | 50% | 75% | max |
|------------|--------|-------------|-------------|------------|------------|------------|-------------|-------------|
| Por | 1000.0 | 14.950480 | 3.029634 | 5.400000 | 12.85750 | 14.98500 | 17.080000 | 24.95000 |
| LogPerm | 1000.0 | 1.369880 | 0.405098 | 0.120000 | 1.13000 | 1.39000 | 1.680000 | 2.58000 |
| AI | 1000.0 | 2.982810 | 0.577820 | 0.980000 | 2.57750 | 3.01000 | 3.380000 | 4.70000 |
| Brittle | 1000.0 | 49.719480 | 15.077008 | -10.500000 | 39.72250 | 49.88000 | 59.170000 | 93.47000 |
| TOC | 1000.0 | 1.003810 | 0.504078 | -0.280000 | 0.84000 | 0.90500 | 1.360000 | 2.71000 |
| VR | 1000.0 | 1.991170 | 0.308194 | 0.900000 | 1.81000 | 2.00000 | 2.172500 | 2.90000 |
| Production | 1000.0 | 2247.295809 | 1464.250312 | 2.713535 | 1191.36956 | 1978.48782 | 3023.594214 | 12568.84413 |



Decision Tree Cross Validation Error (MSE) vs. Complexity



Decision Trees Comments



General Comments on Decision Trees

- Easy to explain
- Analog to human decision making
- Graphically displayed
- Continuous or categorical variables
- Lower predictive accuracy than other machine learning methods
- Model bias may be high

Bagging, Random Forest and Boosting



These are all methods to improve the prediction accuracy of trees

- Bagging (used with many types of models)
 - the use of bootstrap on the training dataset to get B training sets
 - train a tree on each data set
 - then use all models and average the result to get the prediction

$$\hat{f}_{bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)$$

- the trees are allowed to grow large
- 100s to 1,000s of trees (forest of mediocre estimates!)
- classification by majority vote
- out-of-bag data (about 1/3 for each tree) are used as a test data set!

Bagging, Random Forest and Boosting



These are all methods to improve the prediction accuracy of trees

- Random Forest
 - same as bagging, but we randomize selection of on about \sqrt{m} of the features!
 - prevents a single strong predictor from dominating the entire set of trees – forces diversity among the trees
 - decorrelating the trees

Bagging, Random Forest and Boosting



These are all methods to improve the prediction accuracy of trees

- Boosting (used with many types of models)
 - sequential modeling of a simple tree
 - build a tree, calculate residual
 - build a tree to model residual from 1st tree
 - build a tree to model the residual from 2nd tree
 - etc.

Statistical Learning New Tools



| Topic | Application to Subsurface Modeling |
|-------------------------------------|--|
| Consider inference and prediction | <p>Value of working with inference and prediction.</p> <p><i>Permeability was deemed to be redundant with porosity and was removed from the prediction model.</i></p> |
| Consider model training vs. testing | <p>Maximize model prediction accuracy with testing not training.</p> <p><i>A reduced complexity model was adopted for predicting porosity from acoustic impedance due improved testing accuracy.</i></p> |

Data Analytics and Geostatistics: Machine Learning



Lecture outline . . .

- General Comments
- Prediction and Inference
- Decision Tree

Introduction

Modeling Prerequisites

Spatial Estimation

Spatial Uncertainty

Multivariate, Spatial

Multivariate Analysis

Machine Learning

Novel Workflows

Conclusions

Instructor: Michael Pyrcz, the University of Texas at Austin