

# What Does a Geoscientist, Engineer or Data Scientist Need to Know About Geostatistics? and Why It Would Be Helpful?

**Michael Pyrcz, Associate Professor**

Texas Center for Geostatistics

Hildebrand Department of Petroleum and Geosystems Engineering

Bureau of Economic Geology, Jackson School of Geosciences

The University of Texas at Austin



# Who am I?

- New professor in UT PGE, Fall 2017
- Teaching geostatistics, data analytics and machine learning
- over 17 years of experience in consulting, teaching and industrial R&D in statistical modeling, reservoir modeling and uncertainty characterization, statistical / machine learning
- associate editor with Computers and Geosciences, editorial board for Mathematical Geosciences
- author of the textbook “Geostatistical Reservoir Modeling” and > 45 peer reviewed publications, patents etc.
- program chair for SPE Petroleum Data Driven Analytics Technical section (PD<sup>2</sup>A)
- **“Committed to professional development and collaboration to support our industry.”**
- **Expanded role of quantification, data analytics**



Michael Not Working



Michael Working



# Goals

1. **High grade and cover those fundamental concepts** that can impact your daily work.
  - ✓ **New Concepts → New Opportunities**
  - ✓ **Cross Discipline Expertise → Communication / Integration**
  - ✓ **Impact Your Work → New Workflows and Research**
2. **Preview more in-depth discussions and training**, that we will cover in March together.
3. **Shameless promotion** of collaboration opportunities.



# Outline

## 1. Definitions

- geostatistics / reservoir modeling
- big data
- complexity

## 2. Fundamental concepts.

- big data and machine learning
- stationarity
- bias
- uncertainty
- facies
- spatial continuity

## 3. Some Motivational Examples



# (Geo)statistics

## Some Definitions

**Statistics** is concerned with mathematical methods for collecting, organizing, and interpreting data, as well as drawing conclusions and making reasonable decisions on the basis of such analysis.

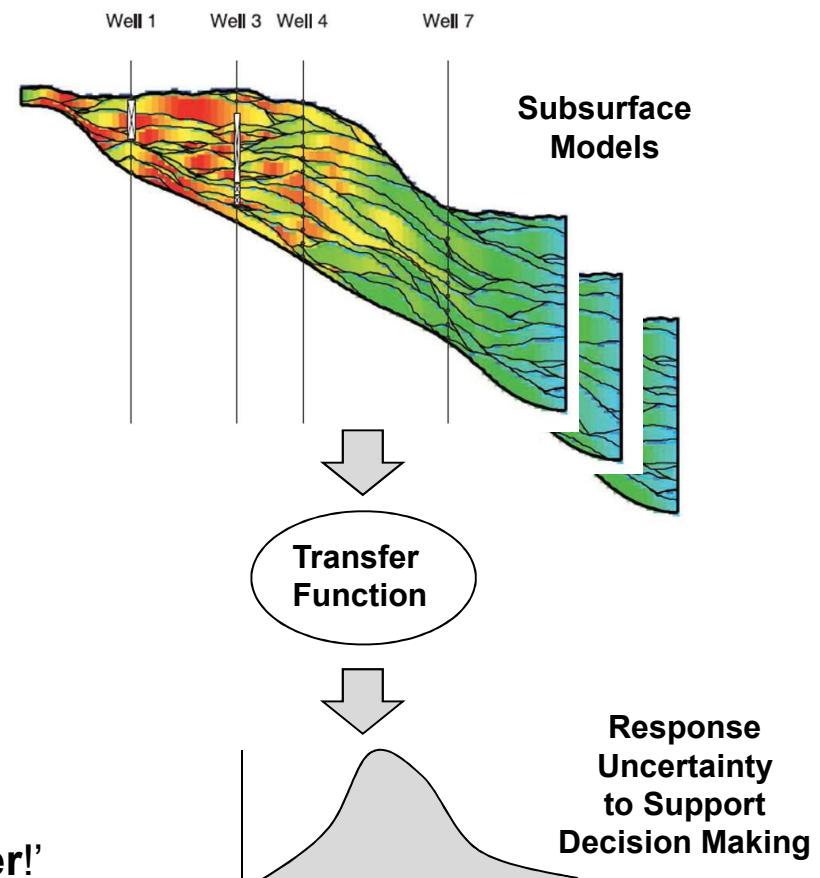
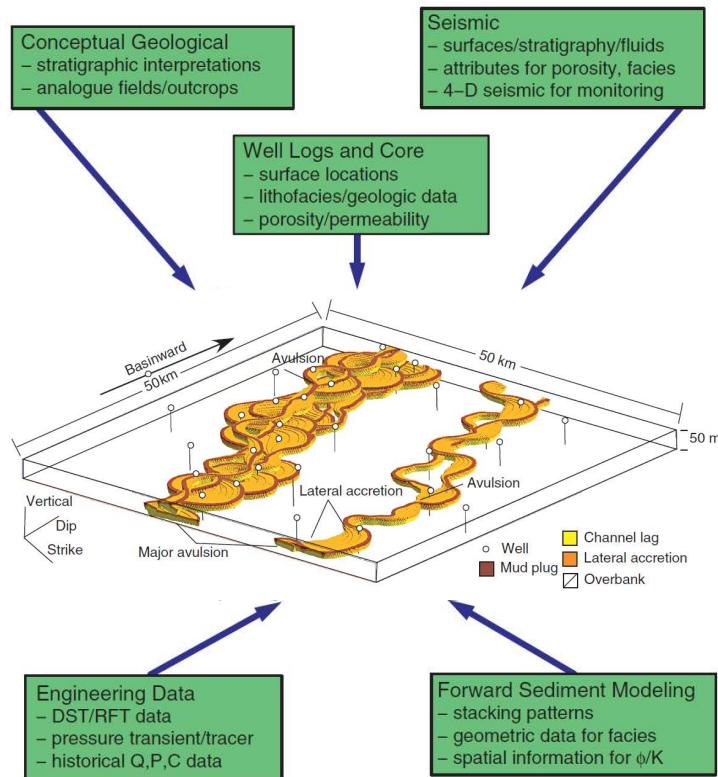
**Geostatistics** is a branch of applied statistics that emphasizes (1) the geological context of the data, (2) the spatial relationship between data, (3) spatial uncertainty and (4) the different volumetric support and precision of the data.

Why do we work with geostatistics in Geosciences?

- ✓ **Geological Context**
- ✓ **Spatial Relationships**
- ✓ **Variable Scale of Data**
- ✓ **Variable Data Precision**
- ✓ **Highly Multivariable**

# (Geo)statistics Some Definitions

**Reservoir / Subsurface Modeling** is the integration of all subsurface information to build a suite of models representing uncertainty to support decision making.



'If it doesn't get in the model, it doesn't matter!'



# (Geo)statistics Some Definitions

Geostatistics developed from **practice of subsurface estimation and modeling** in mining, theory added later.

**TABLE 2.1. RESERVOIR CONCEPTS AND ASSOCIATED GEOLOGICAL AND GEOSTATISTICAL EXPRESSIONS**

Concept	Geological Expression	Geostatistical Expression
Major changes in relationships between reservoir bodies	Architectural complexes and complex sets	Regions—separate units and model with unique methods and input statistics
Changes in reservoir properties within reservoir bodies	Basinward and landward stepping Fining/Coarsening up	Nonstationary mean
Stacking patterns reservoir bodies	Organization, disorganization, compartmentalization, compensation	Attraction, repulsion, minimum and maximum spacing distributions, interaction rules
Major direction of continuity	Paleo-flow direction	Major direction of continuity, locally variable azimuth model
Relationship between vertical and horizontal continuity	Walther's Law	Geometric and zonal anisotropy
Distinct reservoir property groups	Lithofacies, depositional facies, and architectural elements	Reservoir categories, stationary regions
Heterogeneity	Architecture	Spatial continuity model geometric parameters, training image patterns

<sup>a</sup>Most geostatistical constructs can be directly mapped to geological constructs that describe the reservoir.

Common concepts, it all translates!

**Geostatistics is the practical quantification of the subsurface to support decision making.**

# Big Data, Machine Learning

## Big Data Criteria:

- Volume
  - Velocity
  - Variety
  - Veracity
- } We have this!

We've been big data before there was big data.



image from <https://www.humansatsea.com>

## Machine / Statistical Learning

Training a computer detect features, find complicated relationships with complicated, multivariate, large datasets

We've been doing that too.

But we can learn from the best practice in big data and machine learning.

Webinar - Big Data Analytics for Petroleum Engineering: Hype or Panacea?  
December 8, 2017: Big Data Analytics for Petroleum Engineering: Hype or Panacea? Little Data + Simple Model = Big Data?



Webinar at <http://www.cpge.utexas.edu/?q=node/385>

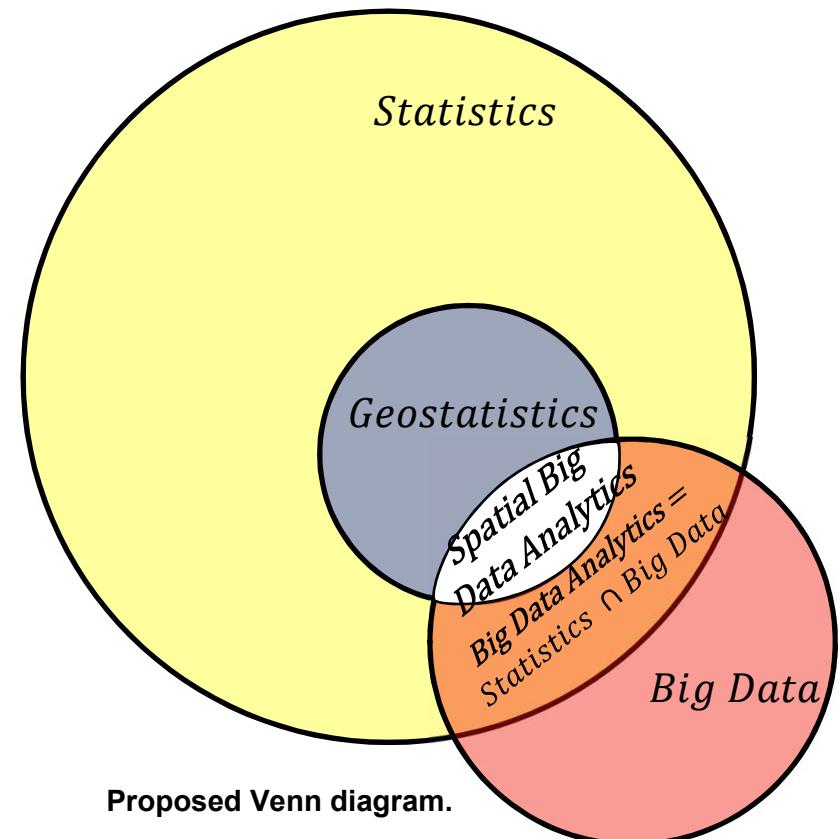
# Big Data, Machine Learning and Geostatistics



**Statistics** is concerned with mathematical methods for collecting, organizing, and interpreting data, as well as drawing conclusions and making reasonable decisions on the basis of such analysis.

**Geostatistics** is a branch of applied statistics that emphasizes: (1) the spatial (geological) context of the data, (2) the spatial relationship between data, (3) the different volumetric support and precision of the data, and (4) spatial and data uncertainty.

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



Given this:

**Spatial big data analytics is the expert use of (geo)statistics to learn from our spatial data set.**

# Model Accuracy and Complexity

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

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

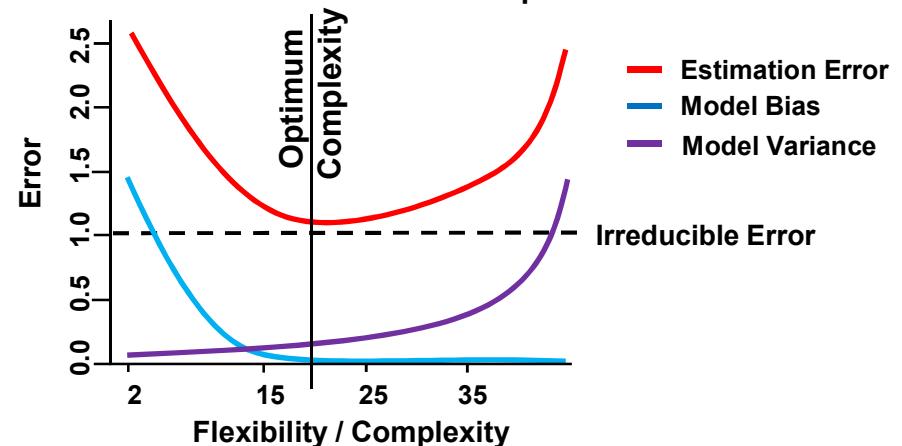
**Model Variance** - variance due to limited data (simpler models  $\downarrow$  lower variance)

**Model Bias** – error due to simple model (simpler models  $\uparrow$  higher bias)

**Irreducible error** - due to missing variables and limited samples

There are **trade-offs**, resulting in an **optimum level of complexity**.

**The most complicated model is not necessarily the best!**

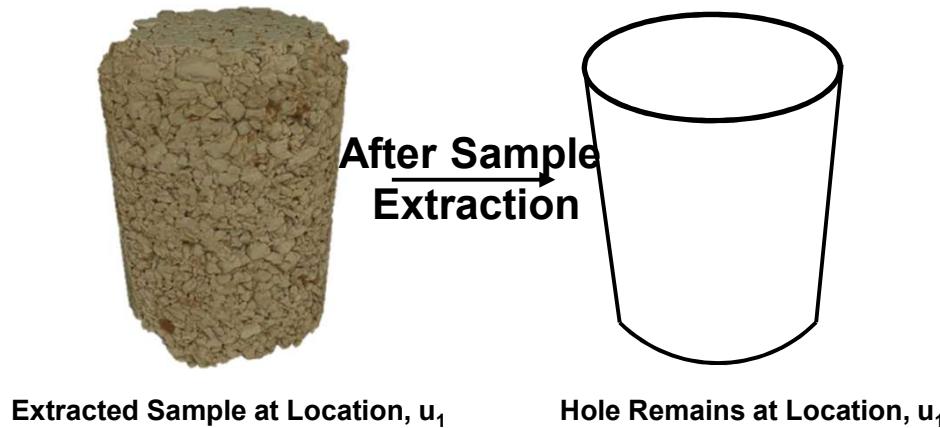


# Stationarity

## Substituting time for space



Any statistic requires replicates, repeated sampling (e.g. air or water samples from a monitoring station). In our geospatial problems repeated samples are not available at a location in the subsurface.



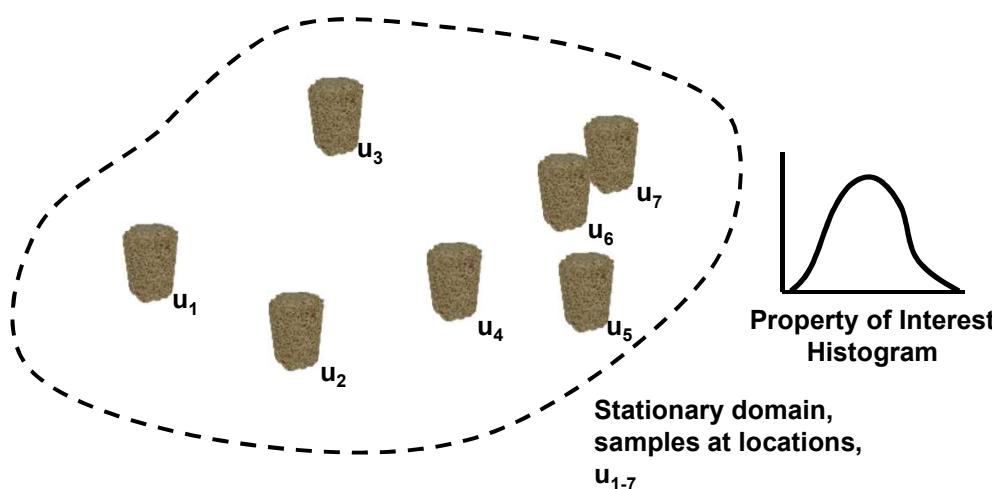
Instead of time, **we must pool samples over space** to calculate our statistics. This decision to pool is the decision of stationarity. It is the decision that the subset of the subsurface is all the “same stuff”.

# Stationarity

## Substituting time for space



The decision of the stationary domain for sampling is an expert choice. Without it we are stuck in the “hole” and **cannot calculate any statistics** nor say anything about the behavior of the subsurface **between the sample data**.



**Import License:** choice to pool specific samples to evaluate a statistic.

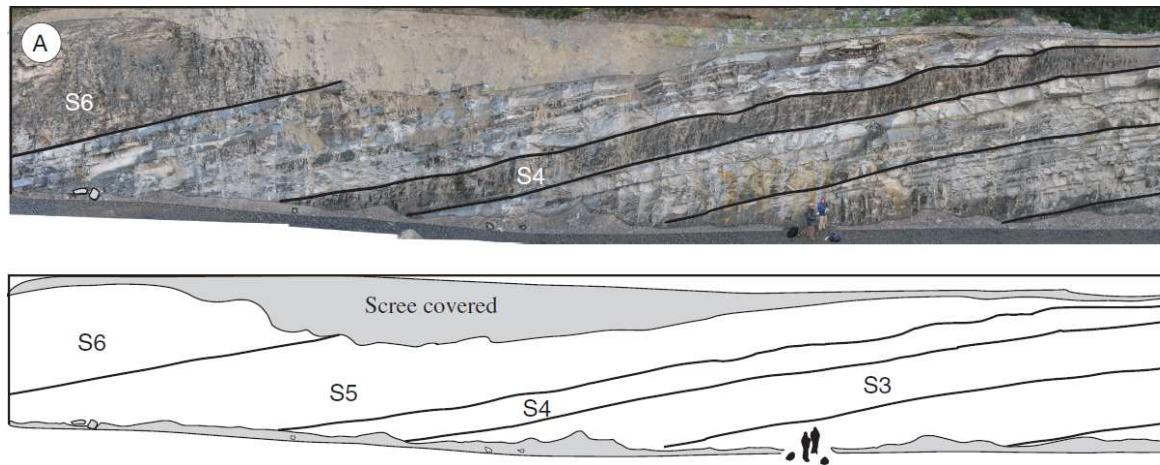
**Export License:** choice of where in the subsurface this statistic is applicable.

# Stationarity

## Definition 1: Geologic



**Geological Definition:** e.g. 'The rock over the stationary domain is sourced, deposited, preserved, and postdepositionally altered in a similar manner, the domain is map-able and may be used for local prediction or as information for analogous locations within the subsurface; therefore, it is useful to pool information over this expert mapped volume of the subsurface.'



Photomosaic, line drawing Punta Barrosa Formation sheet complex (Fildani et al. (2009)).

# Stationarity

## Definition 2: Statistical



**Statistical Definition:** The metrics of interest are invariant under translation over the domain. For example, one point stationarity indicates that histogram and associated statistics do not rely on location,  $\mathbf{u}$ . Statistical stationarity for some common statistics:

**Stationary Mean:**  $E\{Z(\mathbf{u})\} = m, \forall \mathbf{u}$

**Stationary Distribution:**  $F(\mathbf{u}; z) = F(z), \forall \mathbf{u}$

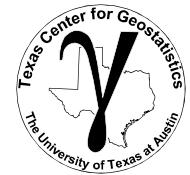
**Stationary Semivariogram:**  $\gamma_z(\mathbf{u}; \mathbf{h}) = \gamma_z(\mathbf{h}), \forall \mathbf{u}$

Stationarity: *What metric / statistic? Over what volume?*

May be extended to any statistic of interest including, facies proportions, bivariate distributions and multiple point statistics.

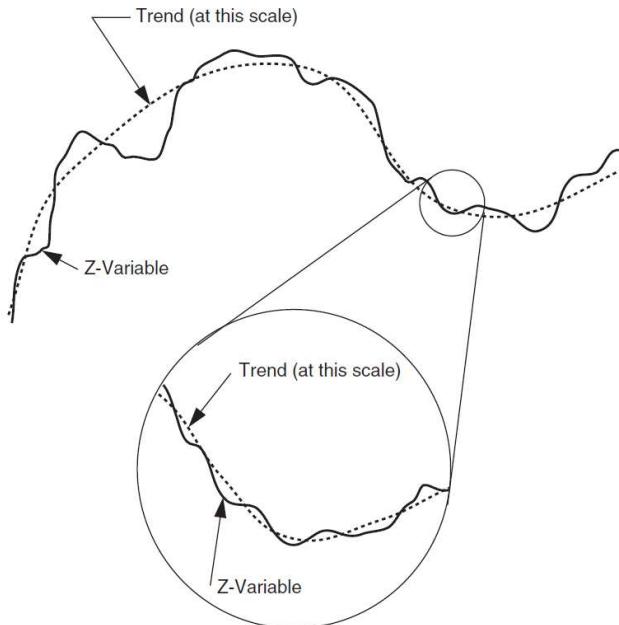
# Stationarity

## Comments on Stationarity



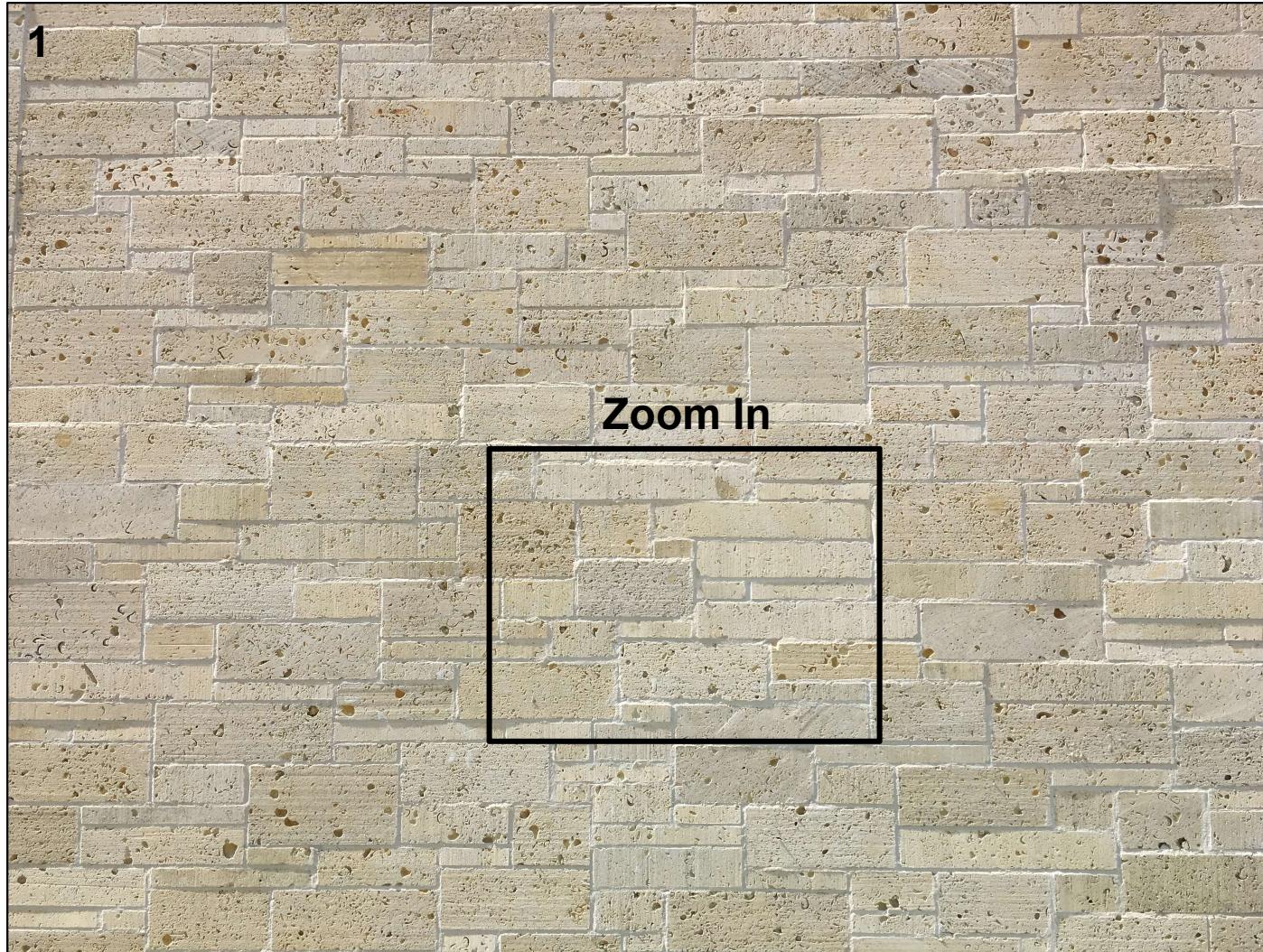
**Stationarity is a decision, not an hypothesis;** therefore it cannot be tested. Data may demonstrate that it is inappropriate.

**The stationarity assessment depends on scale.** This choice of modeling scale(s) should be based on the specific problem and project needs.



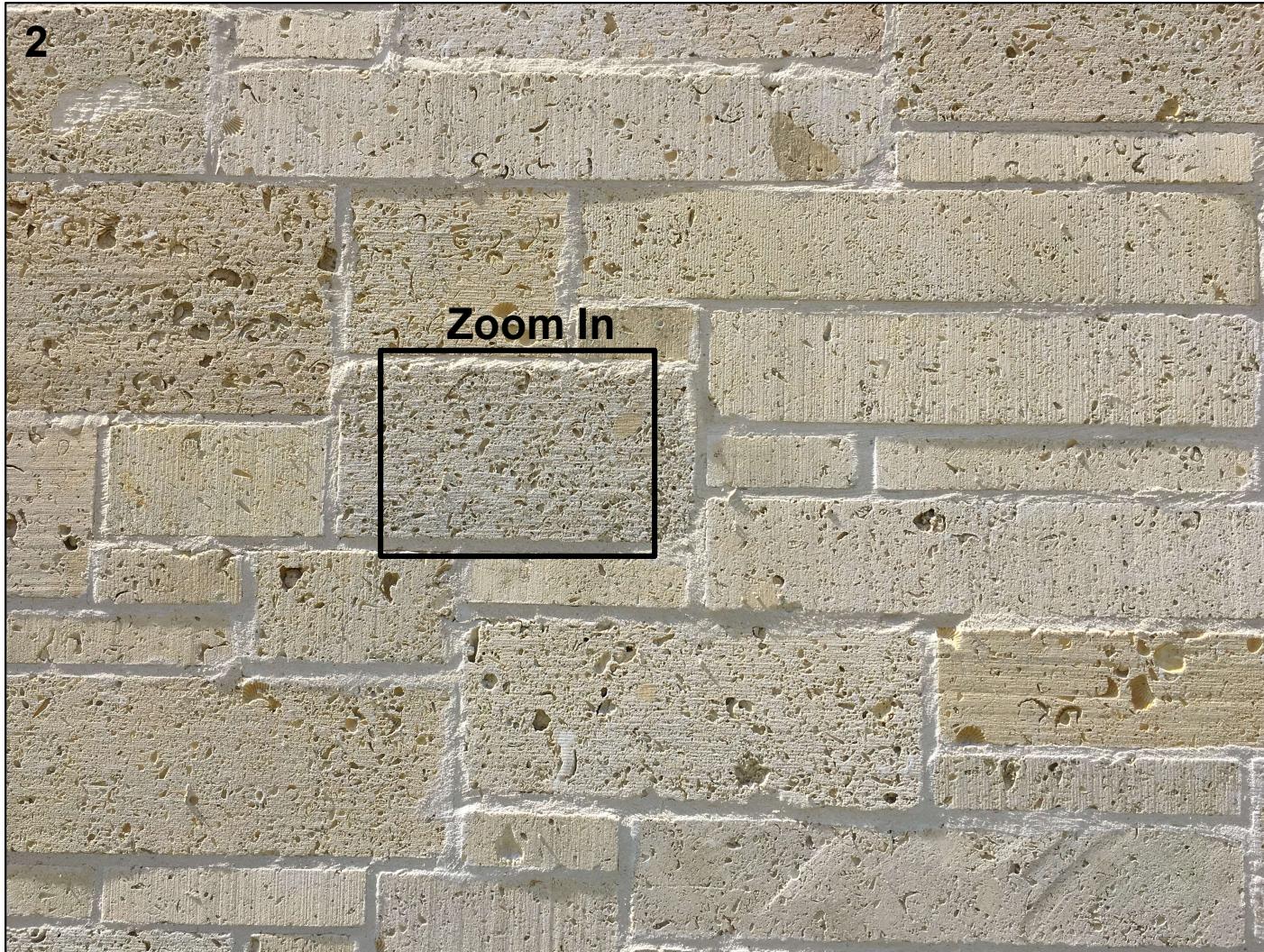
Scales of stationarity from Pyrcz and Deutsch (2014).

# Stationarity and Scale Example



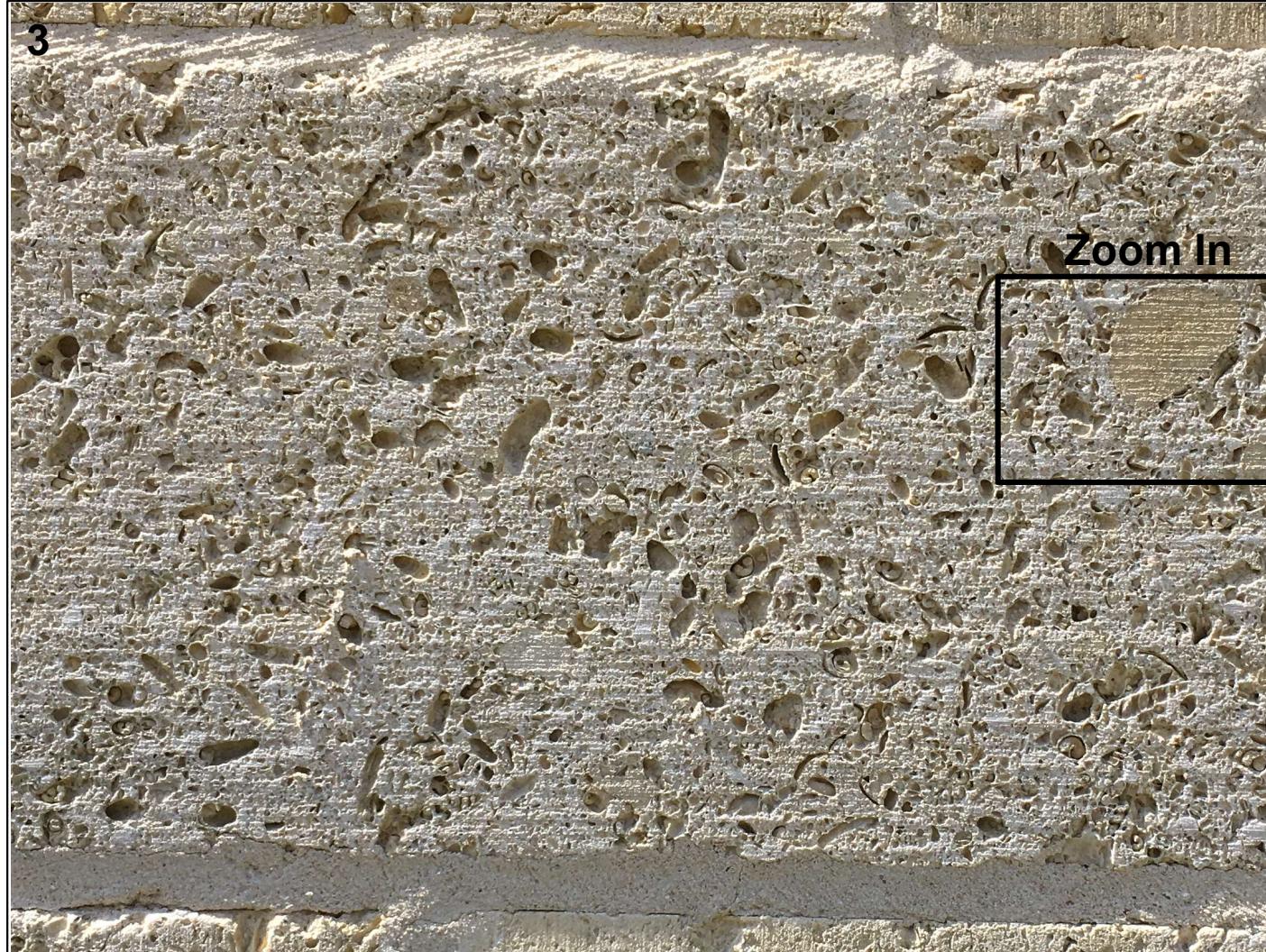
- Is this image stationary? What metric do you consider?

# Stationarity and Scale Example



- A smaller group of bricks?

# Stationarity and Scale Example



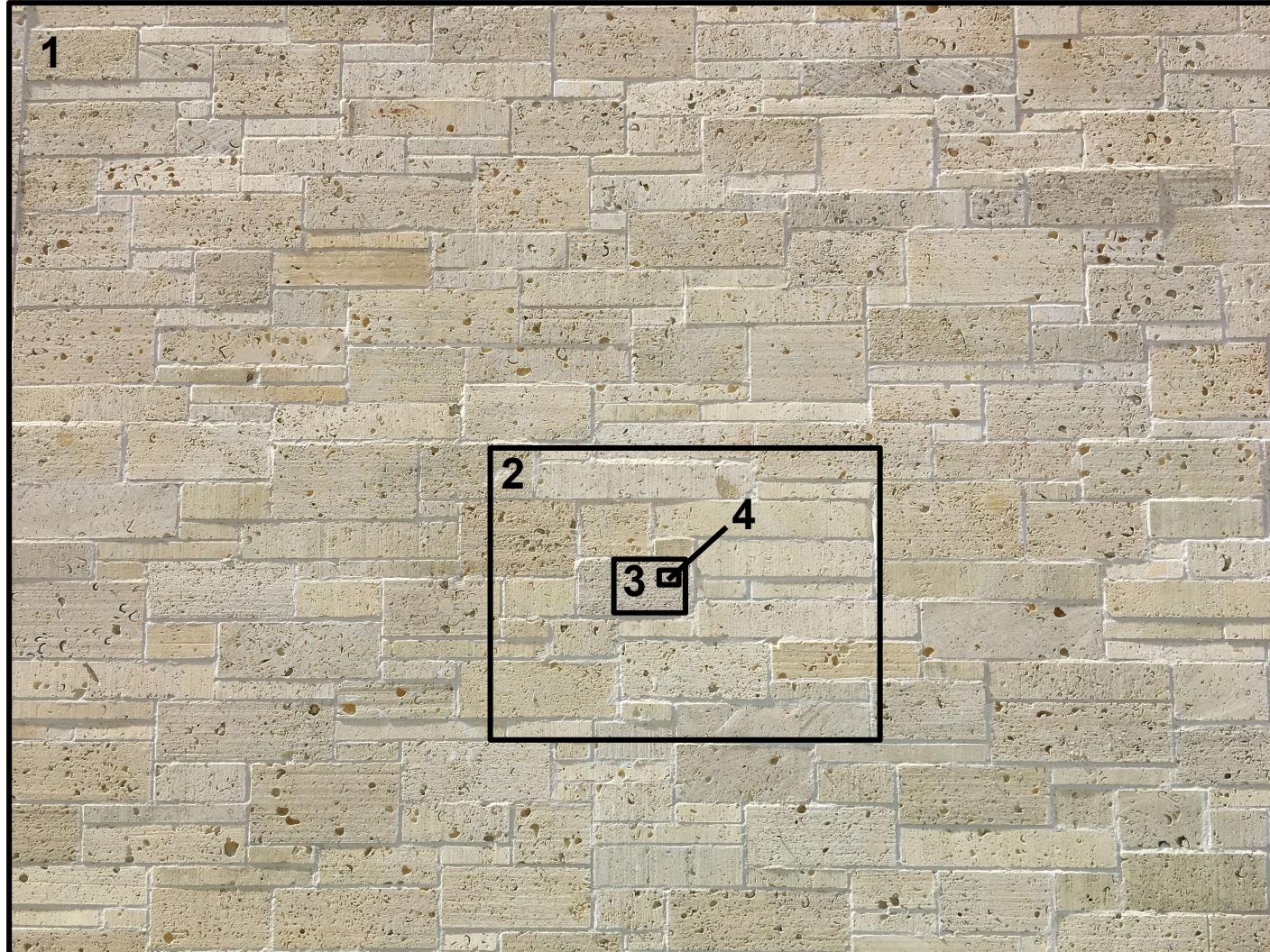
- A single brick?

# Stationarity and Scale Example



- Small part of a brick?

# Stationarity and Scale Example



- Is this image stationary? What metric do you consider?



# Comments on Stationarity

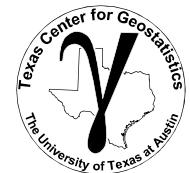
**We cannot avoid a decision of stationarity.** No stationarity decision and we cannot move beyond the data. Conversely, assuming broad stationarity over all the data and over large volumes of the earth is naïve.

**Geomodeling stationarity is the decision:** (1) over what region to pool data (import license) and (2) over what region to use the resulting statistics (export license).

**Nonstationary trends** may be mapped and the remaining stationary residual modelled statistically / stochastically, trends may be treated uncertain.

**Good geological mapping and data integration is essential!**

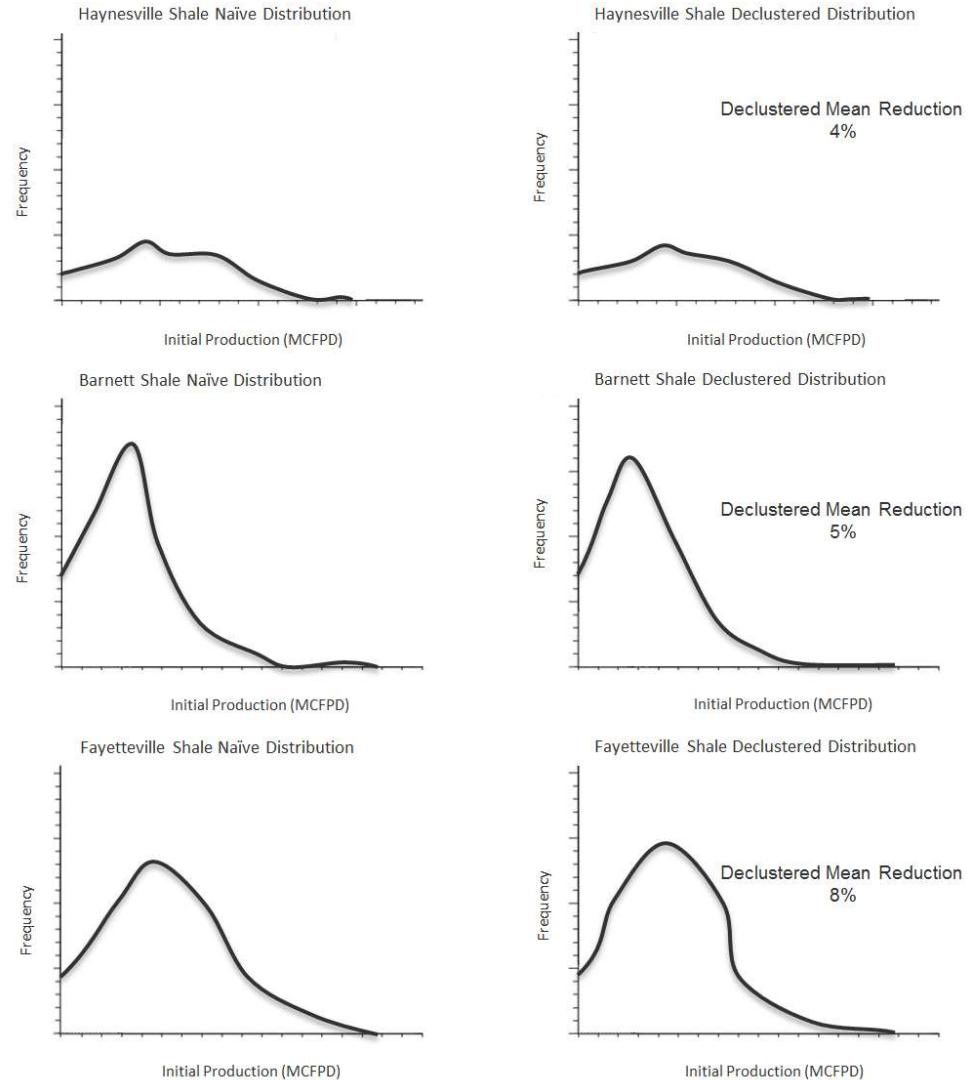
**Stationarity is the framework of any subsurface data and models.**



# Bias is Ubiquitous

## Representative Statistics

- Compiled IP datasets for domestic shale plays
  - Filtered datasets to reduce influence of completions
- Representativity an issue even with large datasets and relatively good coverage
  - Observed changes in naïve to declustered means of 4 – 8%



Naïve and declustered distributions from cell-based declustering (modified from Pyrcz et al, 2017)

# One Source of Bias Data Collection



Data is collected to answer questions:

- how far does the contaminant plume extend? – *sample peripheries*
- where is the fault? – *drill based on seismic interpretation*
- what is the highest mineral grade? – *sample the best part*
- who far does the reservoir extend? – *offset drilling*

and to maximize NPV directly:

- maximize production rates

# Data Collection

There are also limits to our data collection:

- accessibility to the sample – obstruction, reliable drilling, subsalt imaging
- inability to process the sample – may not be able to recover shale core samples
- can't run permeability evaluation on low permeability rock

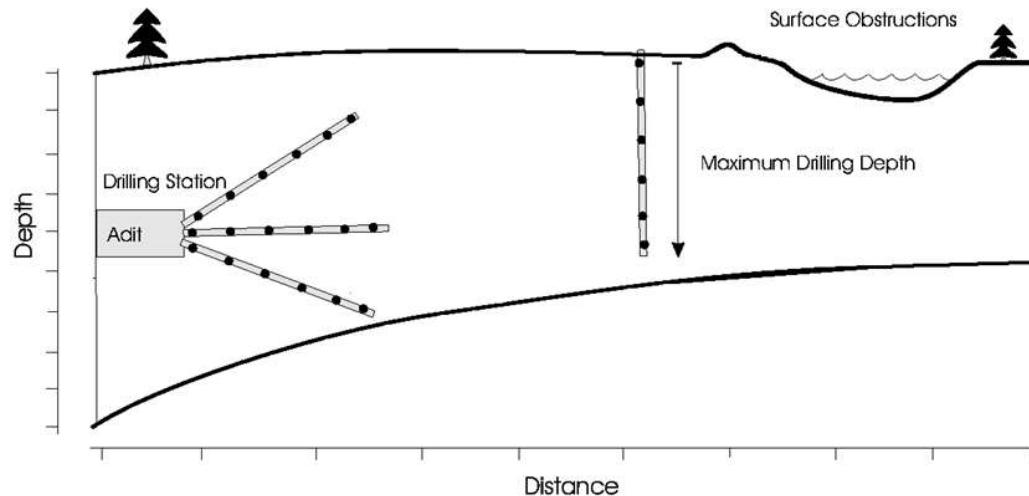


Image from Pyrcz and Deutsch (2003) <http://gaa.org.au/pdf/DeclusterDebias-CCG.pdf>



# Data Collection

If we were sampling for representativity of the sample set and resulting sample statistics, by theory we have 2 options:

1. random sampling
2. regular sampling (as long as we don't align with natural periodicity)

**What would happen if you proposed random sampling in the Gulf of Mexico at \$150M per well?**

We should not change current sampling methods as they result in best economics, we should address sampling bias in the data.

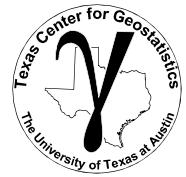
**Never use raw spatial data without access sampling bias / correcting.**



# Solutions to Biased Sampling

- There is a need, however, to adjust the histograms and summary statistics to be representative of the entire volume of interest. We use statistics to make decisions!
1. **Declustering techniques** assign each datum a weight based on closeness to surrounding data
    - $w_i, i = 1, \dots, n$  (weights are greater than 0 and sum to n)
    - Histogram and cumulative histogram use  $w_i, i = 1, \dots, n$  instead of equal weighted,  $w_i = 1.0$ .
  2. **Debiasing techniques** derive an entirely new distribution based on a secondary data source such as geophysical measurements or expert interpretation

# Declustering for Spatial Sampling Bias



- Split up the area of interest with Voronoi partition.
  - Intersected perpendicular bisectors between adjacent data points
  - The result is data weights, all summary statistics can be weighted

$$w(\mathbf{u}_j) = \frac{A_j}{\sum_{j=1}^n A_j} \text{ for } \sum_{j=1}^n w(\mathbf{u}_j) = 1$$

$$w(\mathbf{u}_j) = n \frac{A_j}{\sum_{j=1}^n A_j} \text{ for } \sum_{j=1}^n w(\mathbf{u}_j) = n$$

- This method is sensitive to boundary
- Commonly applied in a variety of scientific fields for weighted averages of spatial phenomenon with irregular sampling.

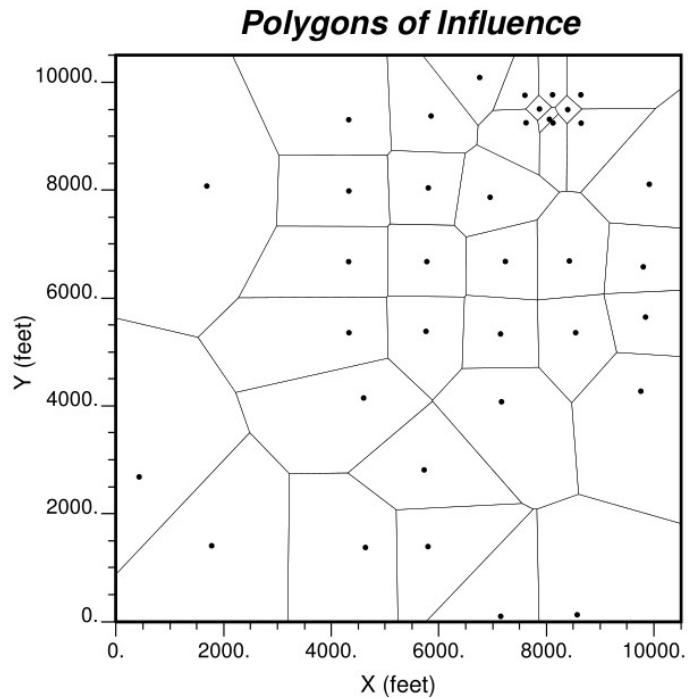
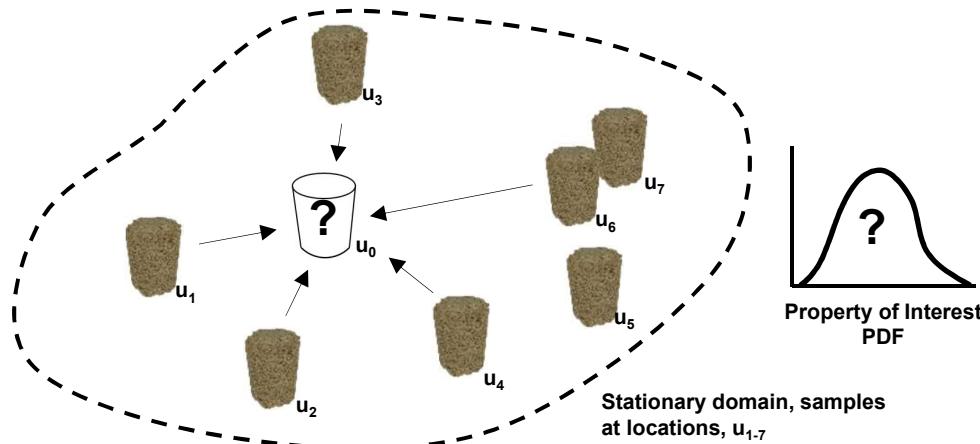


Image from Pyrcz and Deutsch (2014).

**Debias all data sets!**

# Uncertainty

## What is uncertainty?



### Uncertainty is not an intrinsic property of the subsurface.

- At every location ( $u_a$ ) within the volume of interest the true properties could be measured if we had access (facies, porosity etc.).
- **Uncertainty is a function of our ignorance**, our inability to observe and measure the subsurface with the coverage and scale required to support our scientific questions and decision making.

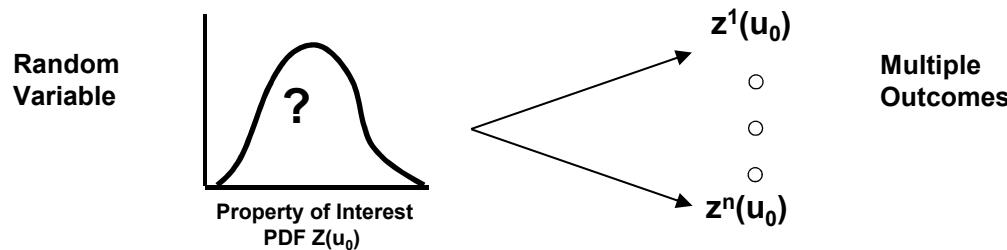
sparsity of sample data + heterogeneity = uncertainty

- If the subsurface was homogeneous, with a few measurements uncertainty would be reduced and estimates resolved to a sufficient degree of exactitude.

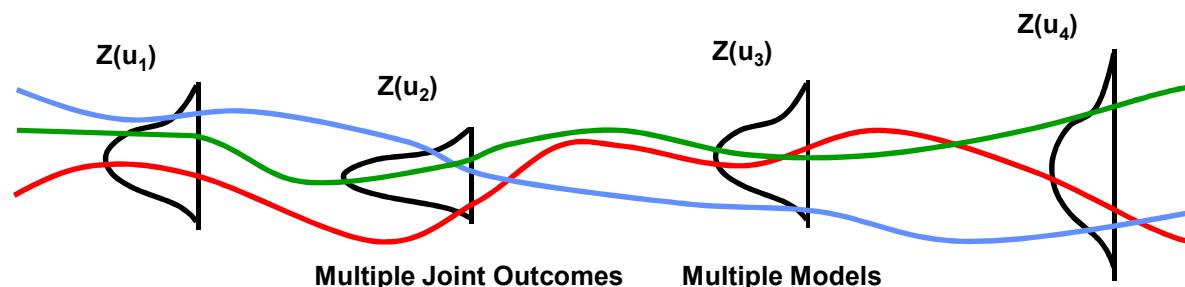
# Uncertainty

## How do we represent uncertainty?

**Random Variables and Functions:** A random variable is a property at a location ( $u_\alpha$ ) that can take on multiple possible outcomes. This is represented by a probability density function (PDF).



If we take a set of random variables at all locations of interest and we impart the correct spatial continuity between them then we have a **random function**. Each outcome from the random function is a potential model of the subsurface.



# Uncertainty

## How do we represent uncertainty?



**Using Multiple Models:** We represent uncertainty with multiple models.

**Scenarios:** when the input decisions and parameters are changed

*Captures interpretation and data uncertainty.*

**Realizations:** when the input decisions and parameters are held constant and only the random number seed is changed

*Captures spatial uncertainty.*

**Working With Multiple Models:** It is generally not appropriate to analyze a single or few scenarios and realizations.

**Use all the models all the time applied to the transfer function**  
(e.g. volumetric calculation, contaminant transport, ore grade scale up, flow simulation etc.).

# Uncertainty

## Comments on uncertainty



**Calculating Uncertainty in a Modeling Parameter:** Use Bayesian methods, spatial bootstrap etc. You must account for the volume of interest, sample data quantity and locations, and spatial continuity.

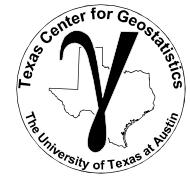
**If You Know It, Put It In.** Use expert geologic knowledge and data to model trends. Any variability captured in a trend model is known and is removed from the unknown, uncertain component of the model. Overfit trend will result in unrealistic certainty.

**Types of Uncertainty:** (1) data measurement, calibration uncertainty, (2) decisions and parameters uncertainty, and (3) spatial uncertainty in estimating away from data. Your job is to hunt for and include all significant sources of uncertainty.

**Be an uncertainty detective! Discover and evaluate various sources.**

# Uncertainty

## Comments on uncertainty



**What about Uncertainty in the Uncertainty?** Don't go there! Use defendable choices in your uncertainty model, be conservative about what you know, document and move on.

**Uncertainty Depends on Scale.** It is much harder to predict a property of tea spoon vs. a house-sized volume at a location ( $u_\alpha$ ) in the subsurface. Ensure that scale and heterogeneity are integrated.

**You Cannot Hide From It.** Ignoring uncertainty assumes certainty and is often a very extreme and dangerous assumption.

**Decision Making with Uncertainty.** Apply all the models to the transfer function to calculate uncertainty in subsurface outcome to support decision making in the presence of uncertainty.

**Ignoring uncertainty is assuming certainty.**

# Facies

## What are the Criteria for Facies?



First some general comments:

1. **Facies / Rock type** is an important decision for subsurface modeling. It should remain a collaborative decision integrating expertise from the project team (Stratigraphers, Reservoir Modelers, Reservoir Engineers, Petrophysicists and Geophysicists).
2. Facies / Rock types **must improve subsurface prediction** away from the data or they do not add value.
3. **Number of facies** is a balancing act between accuracy of geological concepts and statistical inference, and modeling effort
4. Reservoir modeling is **hierarchical**,  
units *contain* depofacies *contain* lithofacies *contain* por/perm
5. **80-90% of flow heterogeneity is captured in the facies models.**

# Facies

## What are the Criteria for Facies?



These are the **criteria for facies** (or any categories in reservoir models).

Criteria	Considerations	Example
<b>Separation of Rock Properties</b>	Facies must divide the properties of interest that impact subsurface environmental and economic performance (e.g. grade, porosity and permeability).	
<b>Identifiable in Data</b>	Facies must be identifiable with the most common data available. e.g. facies identifiable only in cores are not useful if most wells have only logs.	
<b>Map-able Away from Data</b>	Facies must be easier to predict away from data than the rock properties of interest directly, facies improves prediction.	
<b>Sufficient Sampling</b>	There must be enough data to allow for reliable inference of reliable statistics for rock properties for each facies.	

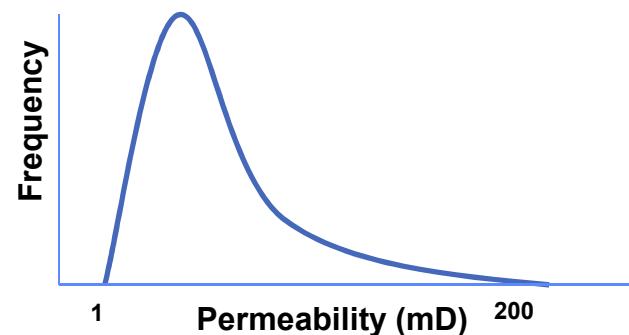
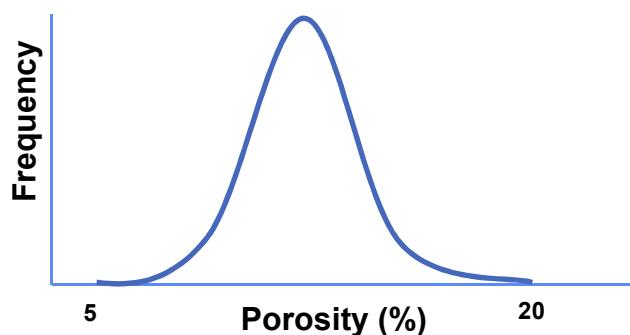
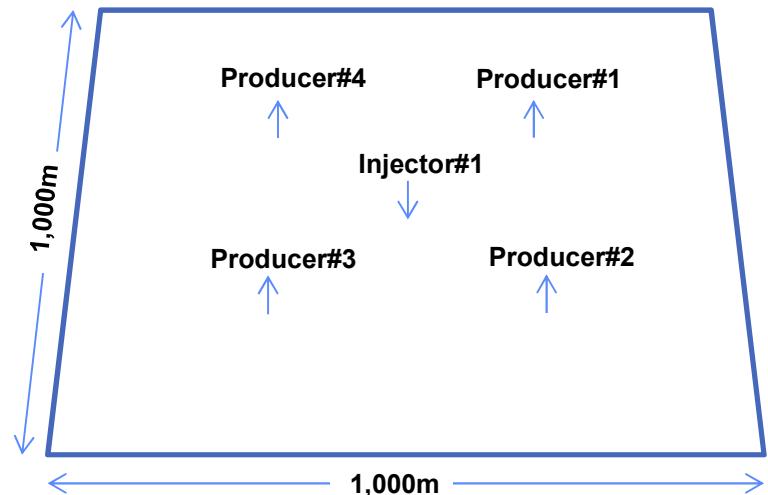
Use these criteria to improve project integration.

# Motivation for Measuring Spatial Continuity



## Simple Example

- Area of interest
- 1 Injector and 4 producers
- Porosity and permeability distributions (held constant for all cases)

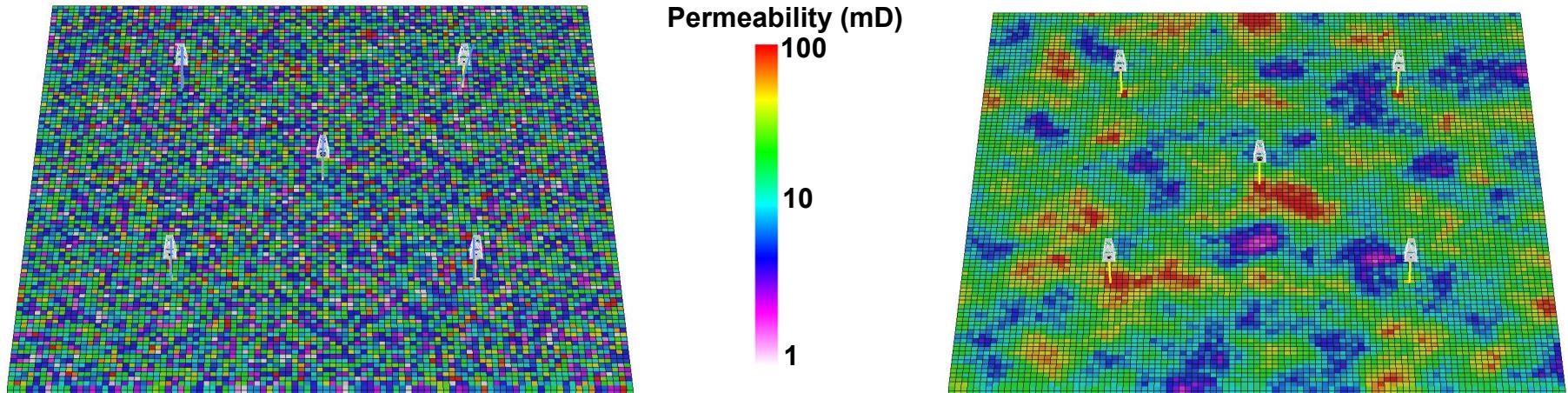


# Motivation for Measuring Spatial Continuity



**Does spatial continuity of reservoir properties matter?**

Consider these models of permeability



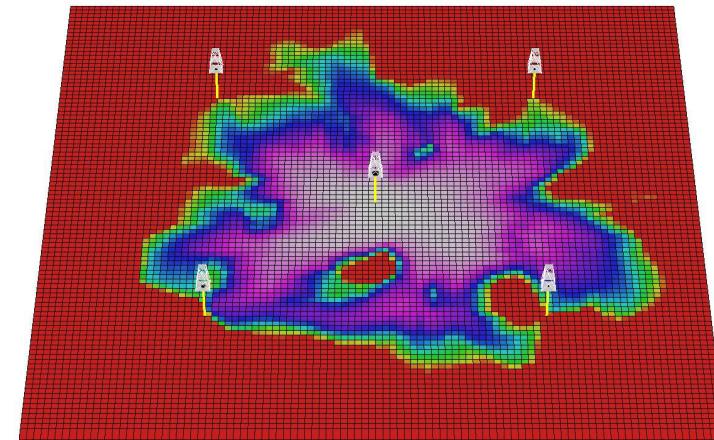
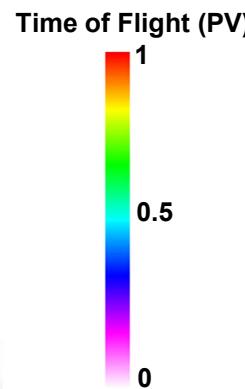
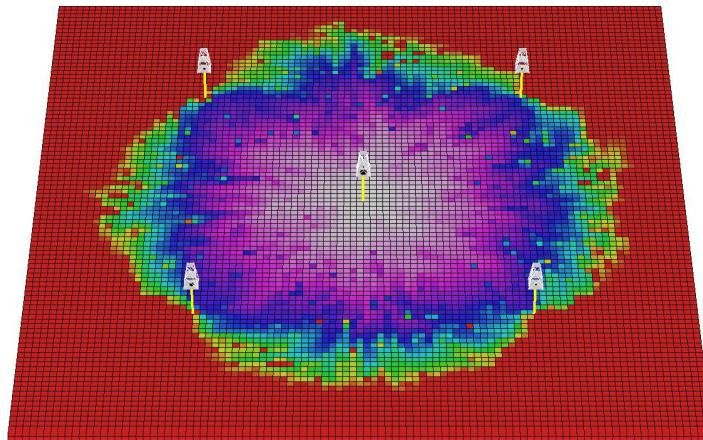
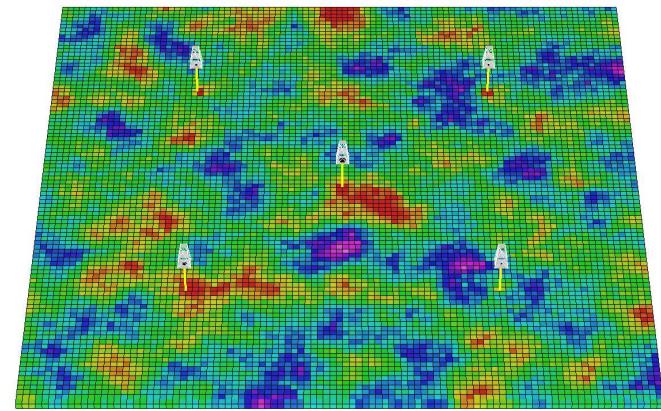
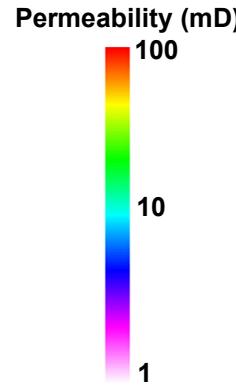
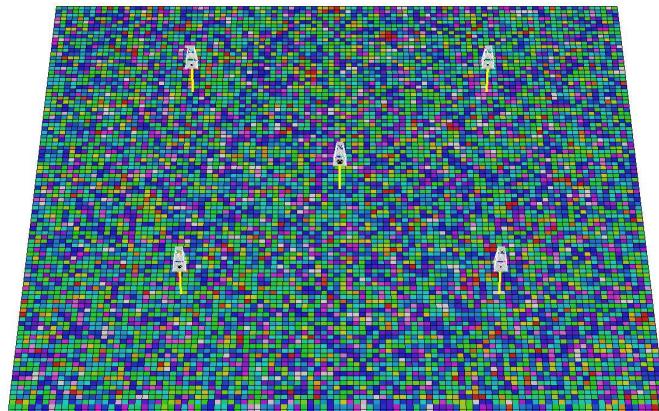
Recall – all models have the same porosity and permeability distributions

- Mean, variance, P10, P90 ...
- Same static oil in place!

# Motivation for Measuring Spatial Continuity



Does spatial continuity of reservoir properties matter?

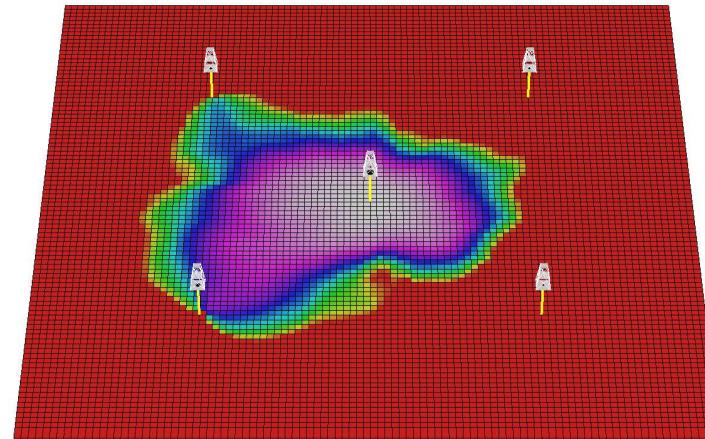
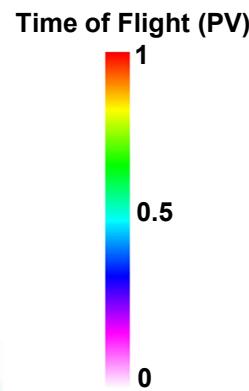
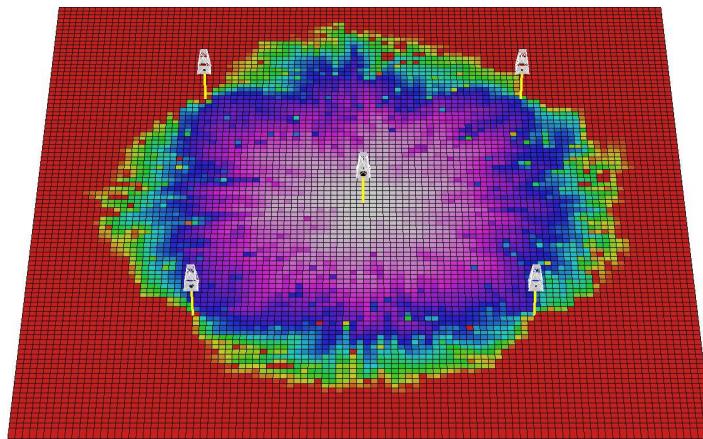
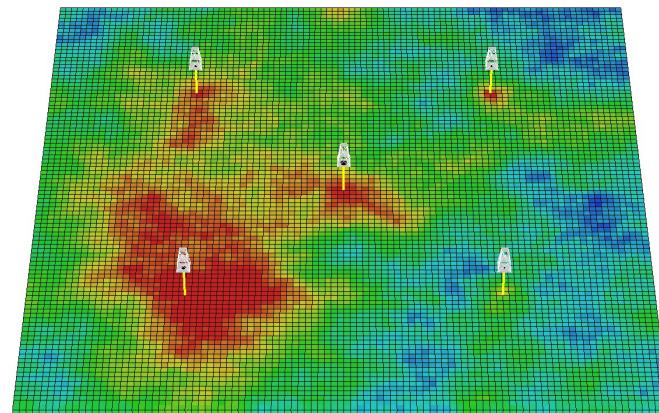
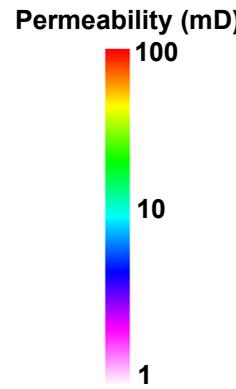
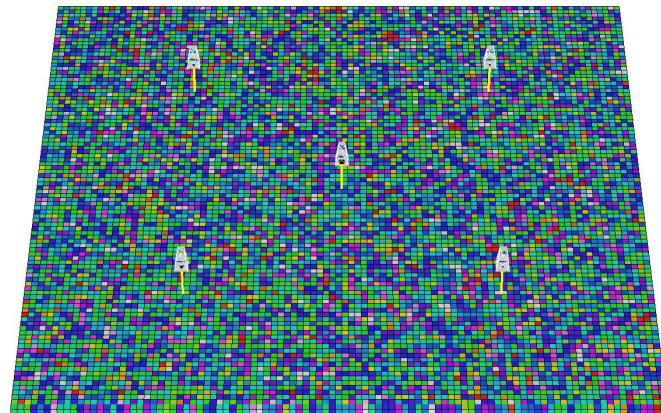


How does heterogeneity impact recovery factor? Well Estimated Ultimate Recovery?

# Motivation for Measuring Spatial Continuity



Does spatial continuity of reservoir properties matter?

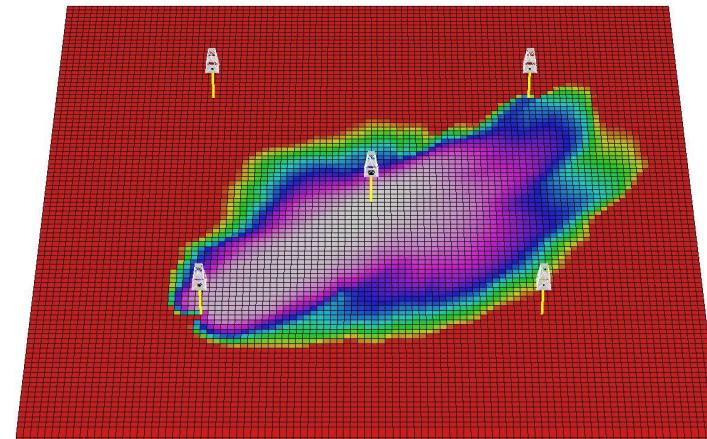
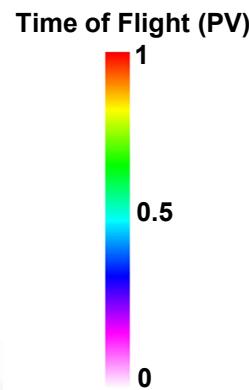
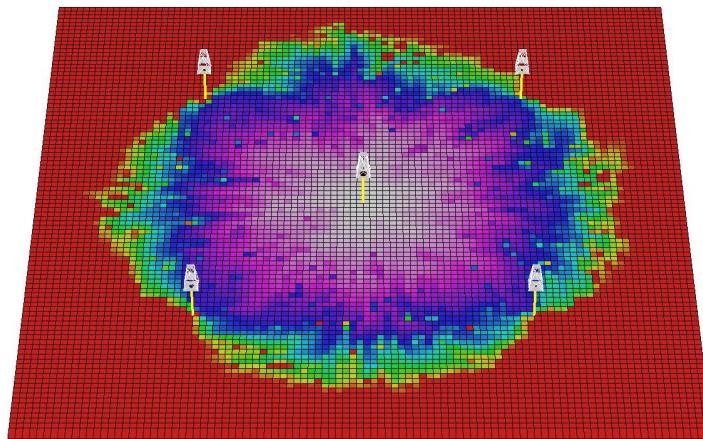
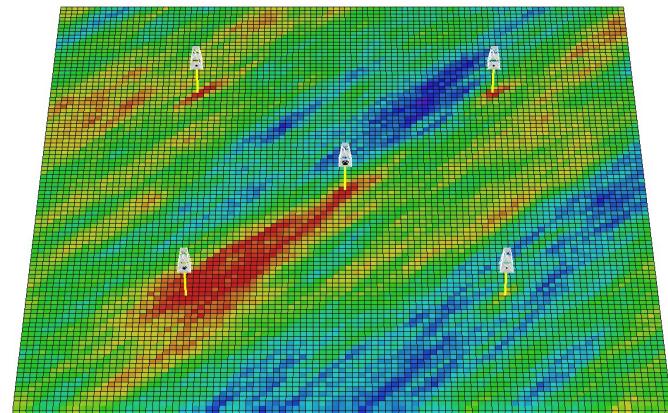
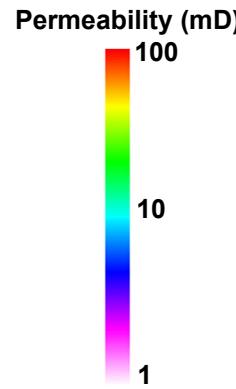
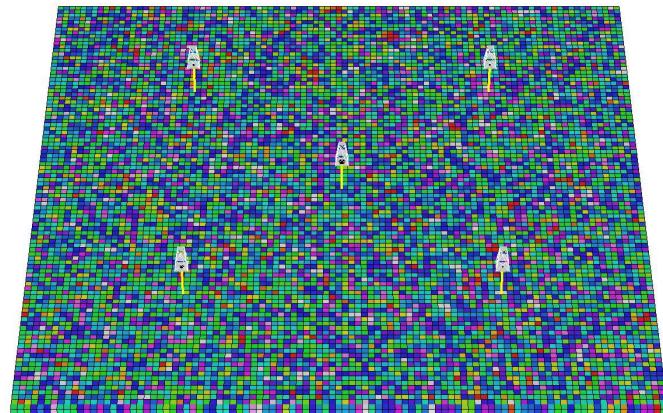


How does heterogeneity impact recovery factor? Well Estimated Ultimate Recovery?

# Motivation for Measuring Spatial Continuity



Does spatial continuity of reservoir properties matter?

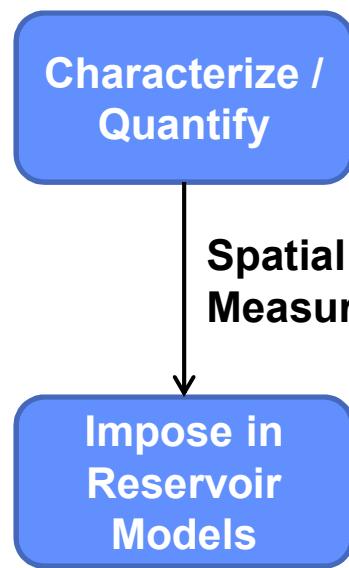


How does heterogeneity impact recovery factor? Well Estimated Ultimate Recovery?

# Motivation for Measuring Spatial Continuity

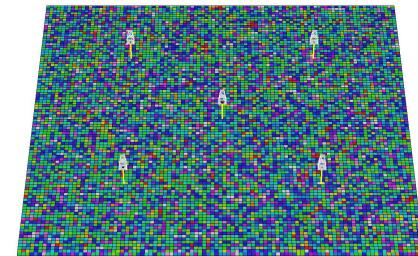


- For the same reservoir property distributions a wide range of spatial continuities are possible.
- Spatial continuity often impacts reservoir forecasts.
- Need to be able to:

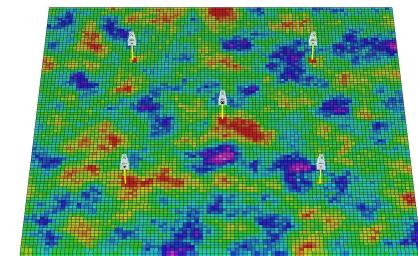


**Spatial Continuity**

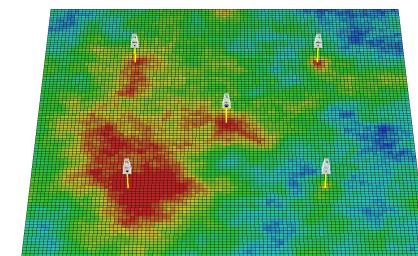
“Very Short”



“Medium”

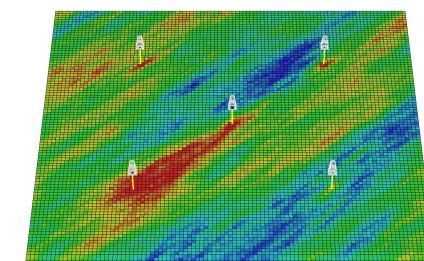


“Long”



“Anisotropic”

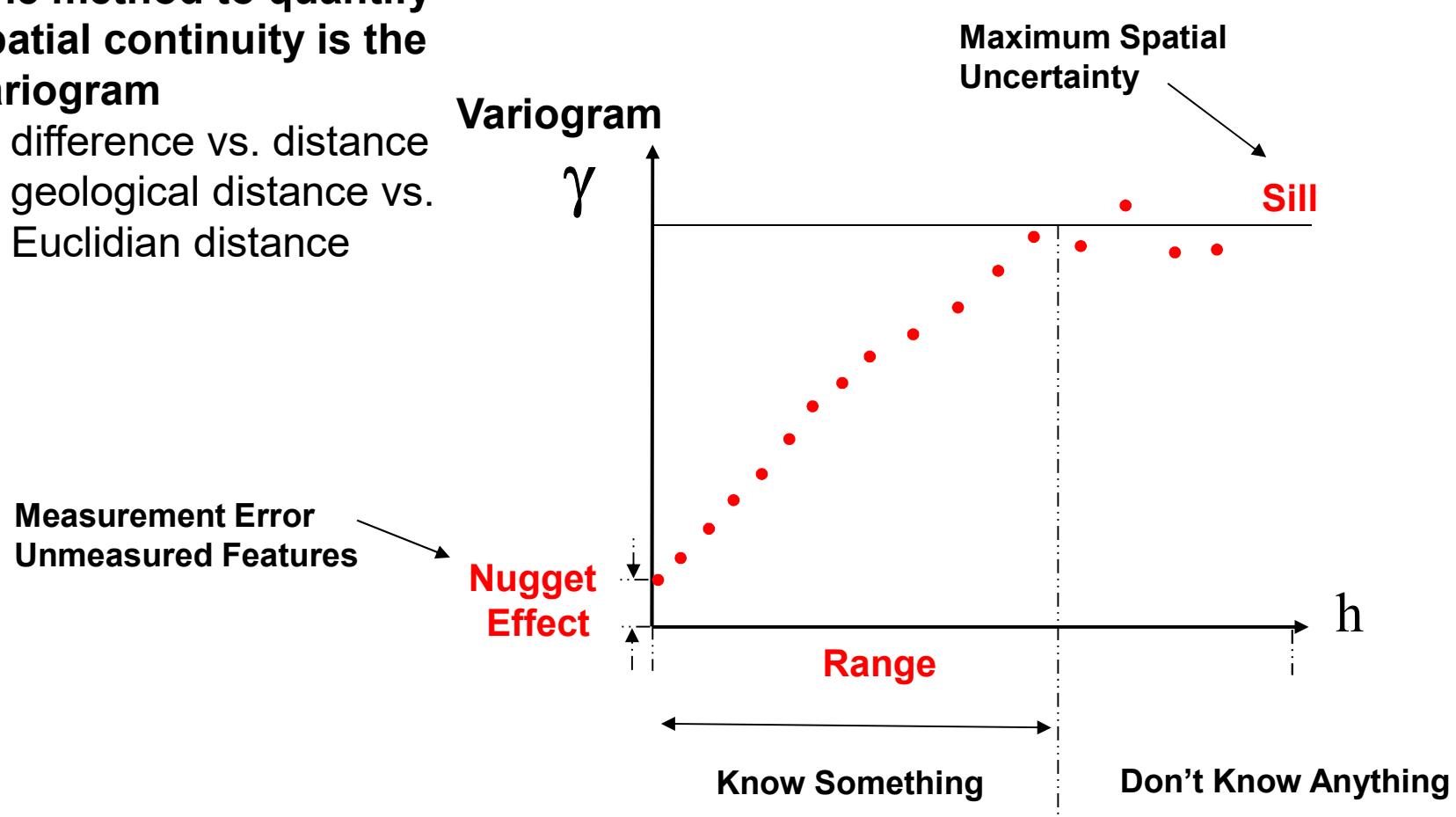
**Continuity matters, but it is complicated.**



# Spatial Continuity

**One method to quantify spatial continuity is the variogram**

- difference vs. distance
- geological distance vs. Euclidian distance



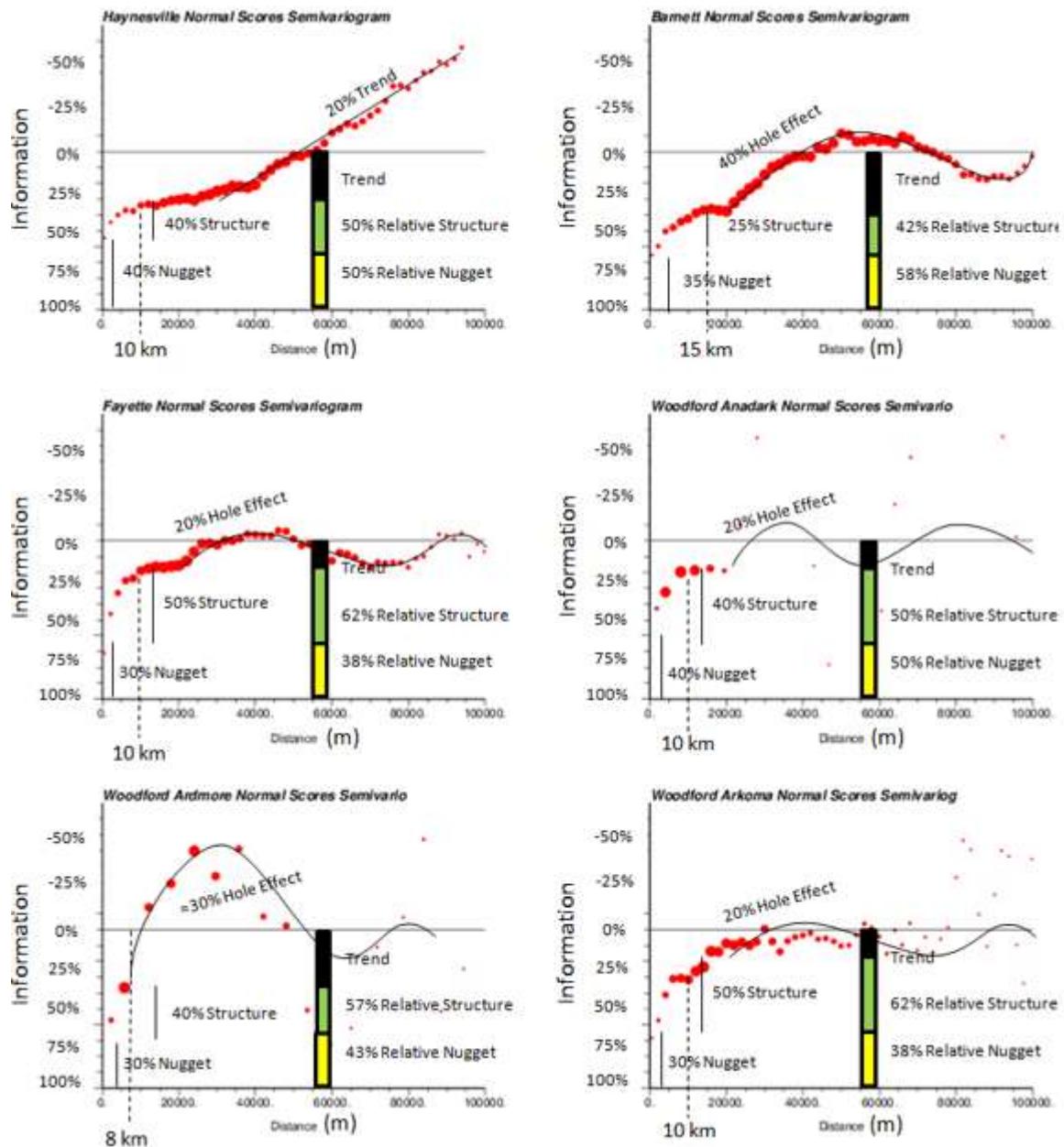
# Spatial Continuity



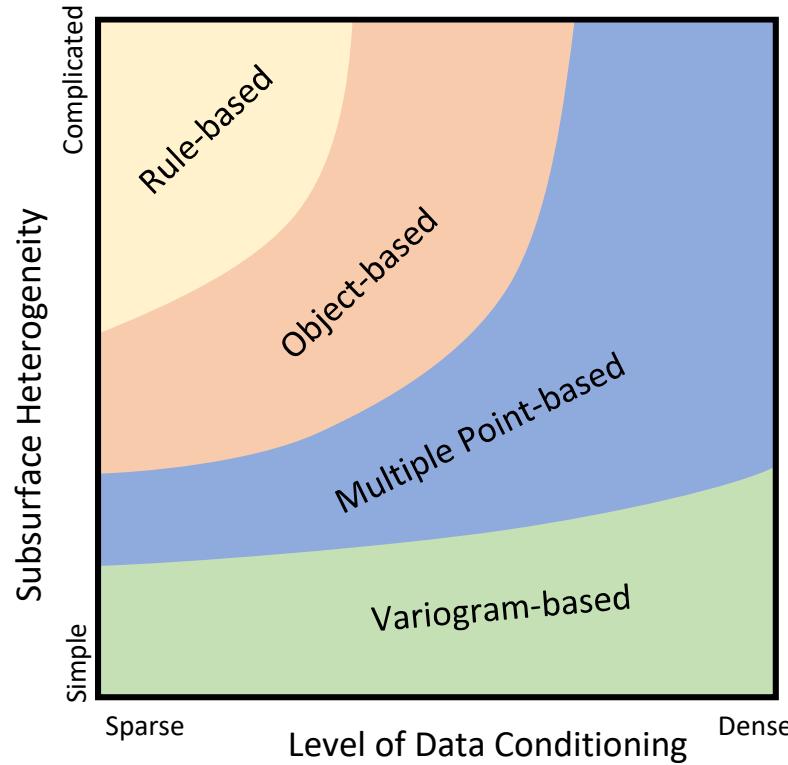
## Example for production from shale assets.

- spatial continuity of unconventional production
- prediction model
- value of well data

Quantification of spatial continuity of shale gas production rates (Pyrcz et al., 2016).



# A Toolbox



Geostatistics includes a tool box of subsurface modeling methods.

- A scheme for selecting between methods, also consider project goals and resources.
- New methods that model the subsurface with data integration, uncertainty, conditional to data. **The tool box is evolving to meet new challenges!**

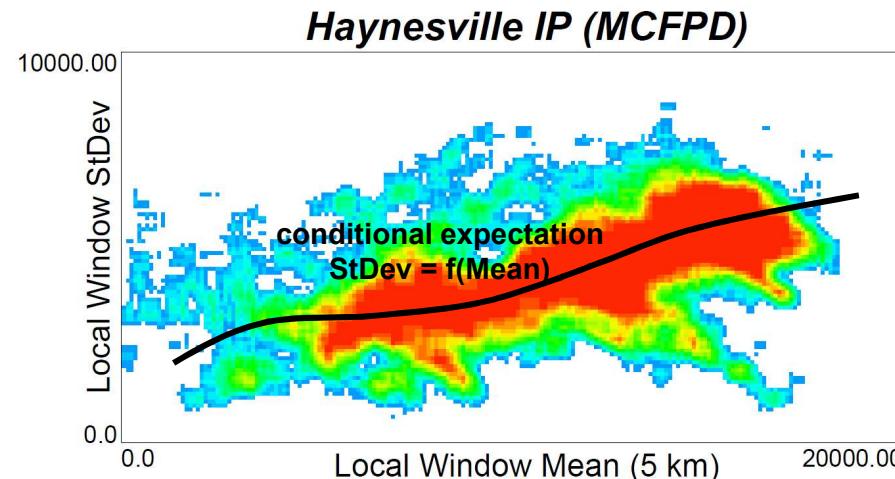
# Motivational (Geo)statistics Examples



## Why Use (Geo)statistics?

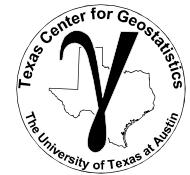
### Quantification / Summarization / Comparison

- abstraction to a small set of parameters allows us to detect features, learn new insights
- robust measure of significance of differences
- seek opportunities for quantification.



Abstraction allows for efficient characterization and leads to insights, (Pyrcz et al., 2016)

# Motivational (Geo)statistics Examples

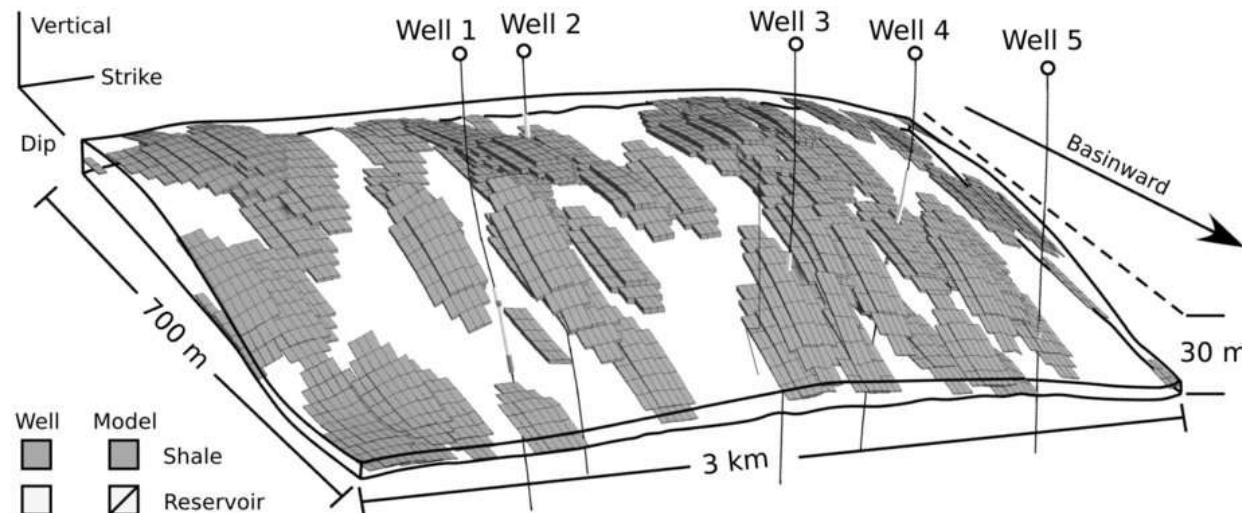


## Why Use (Geo)statistics?

### Model of Uncertainty

sparse sampling + heterogeneity = uncertainty

- if we had enough data and understood the phenomenon perfectly there is no uncertainty, no need for a statistical / stochastic approaches.



Can't know exactly where the shales are from 5 wells and given the shale discontinuity.  
(Pyrcz and Deutsch, 2014)

# Motivational (Geo)statistics Examples



## Why Use (Geo)statistics?

### Too Big / Too Complicated

### Massive Multivariate

- due to the curse of dimensionality we often cannot sample enough to characterize the system
- need to used a statistical multivariate model

Multivariate modeling accounting for complicated relationships (Barnett and Deutsch, 2012)

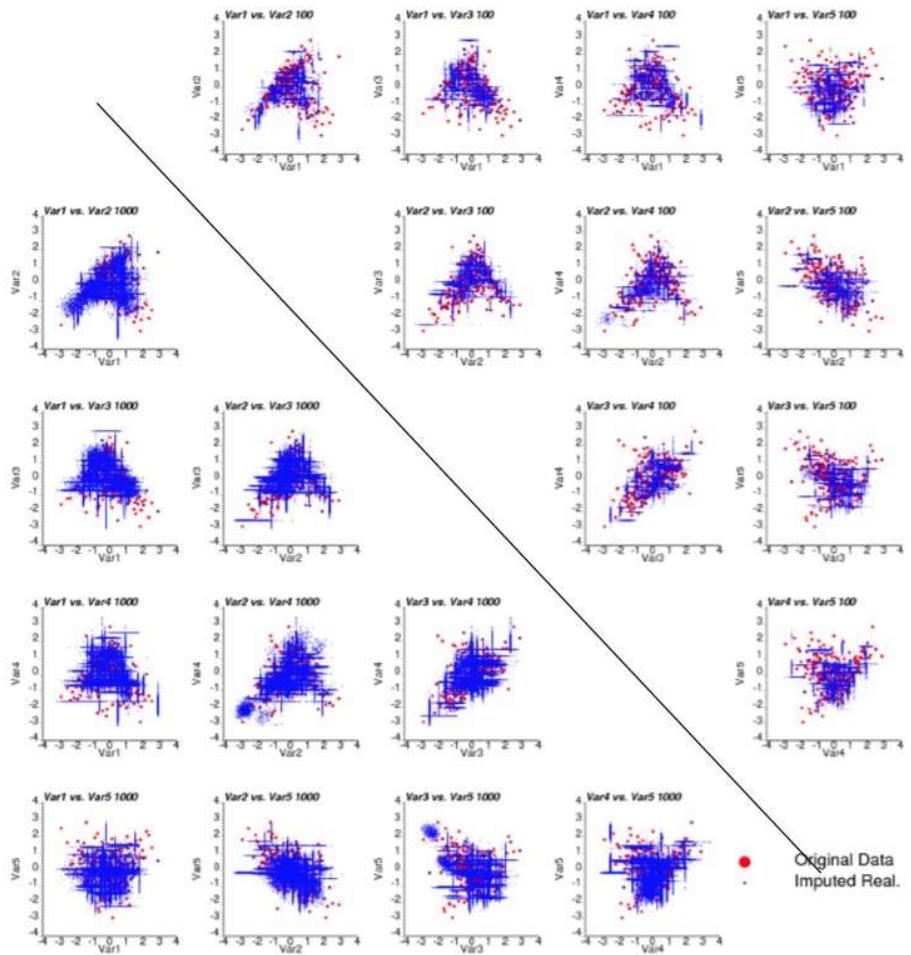
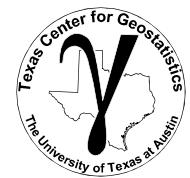


Figure 4: Scatterplots of the sampled observations and imputed observations for 1000 realizations (bottom covariance triangle) and 100 realizations (upper covariance triangle).

# Motivational (Geo)statistics Examples



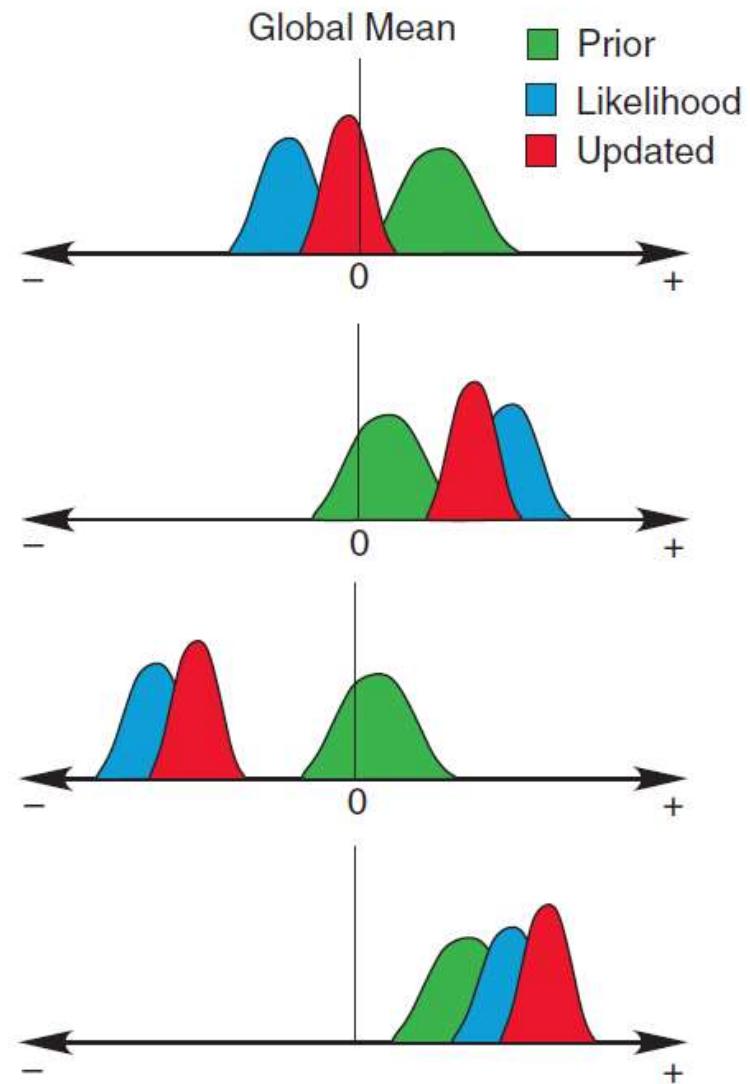
## Why Use (Geo)statistics?

### Combining / Updating with New Information

#### Bayesian Updating

- need statistical models to describe data redundancy

Bayesian updating under the assumption of Gaussianity (Pyrcz and Deutsch, 2012).



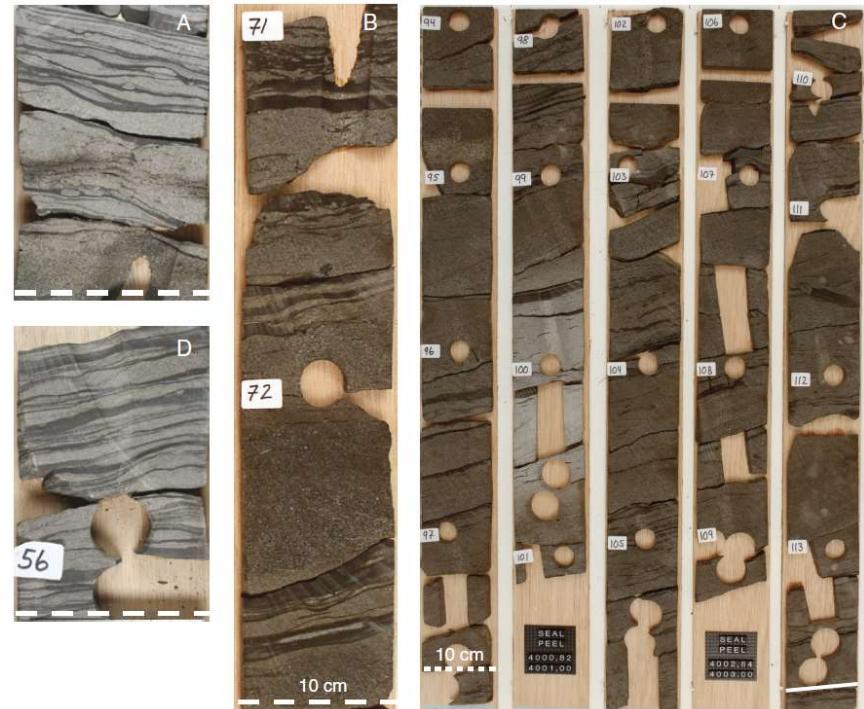
# Motivational (Geo)statistics Examples

Why Use (Geo)statistics?

Accounting for Scale

Pores to Production

- statistical models for change of support size



# Motivational (Geo)statistics Examples



## Why Use (Geo)statistics?

Debias Ourselves

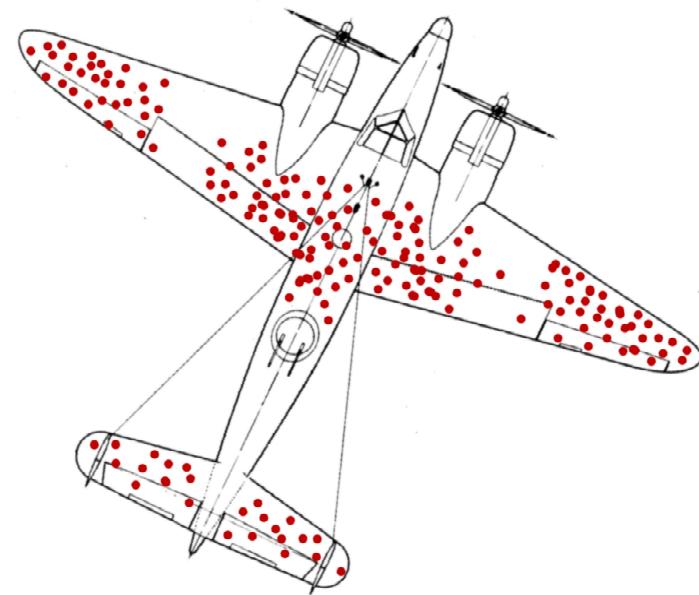
Survivorship Bias

a form of selection bias resulting from selecting samples that “survived” some previous selection process. This often leads to false conclusions.

For example, in WWII the Center for Naval Analyses (@CNA\_org Twitter) compiled a dataset of bomber damage to assess where reinforcement was needed.

Statistician Abraham Wald recognized this was a case of survivorship bias. The planes shot in critical locations did not return to base. Wald suggested reinforcement of locations that were not damaged in planes that safely returned to base!

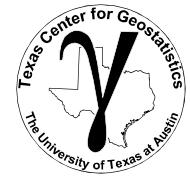
([https://en.wikipedia.org/wiki/Survivorship\\_bias#In\\_the\\_military](https://en.wikipedia.org/wiki/Survivorship_bias#In_the_military))



Hypothetical dataset of aircraft damage for planes that returned to base. Source [https://en.wikipedia.org/wiki/Survivorship\\_bias#/media/File:Survivorship-bias.png](https://en.wikipedia.org/wiki/Survivorship_bias#/media/File:Survivorship-bias.png)

**Is there preselection in our subsurface datasets?**

# Motivational (Geo)statistics Examples

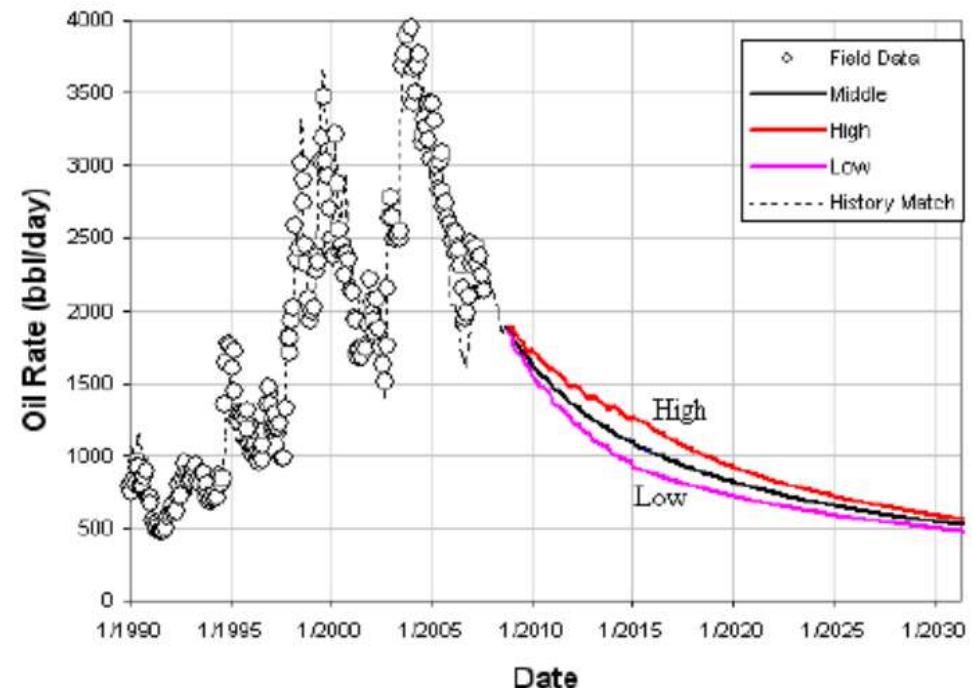


## Why Use (Geo)statistics?

### Forecasting / Decision Making

#### Decision Support

- integrated subsurface models
- to build forecast uncertainty models
- to optimize very expensive project decisions in the presence of subsurface uncertainty



Reservoir forecasting with uncertainty (Yang, 2009).

# Let's Review



**Geostatistics is spatial (big) data analytics**

**We have big data**

**Stationarity is assumed or we are stuck in the well bore**

**Uncertainty is due to our own ignorance**

**Facies must help with spatial prediction**

**We can model spatial continuity**

**Use all of the models**

**Geostatistics helps us integrate, avoid bias, model uncertainty for improved decision making.**



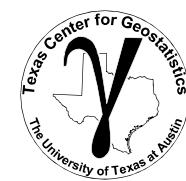
# Work Together?

## Direct Collaboration

- Great opportunity to partner on research for short term and long term deliverables.
- Open to one-on-one partnerships to direct and work with a PhD student, embedded with industry partner's team.

## Professional Development, Teaching, Mentoring

- I have conducted internal training at various companies over the last year on:
  1. Basics of Geostatistics for Geologists, Data Scientists and Engineers
  2. Subsurface Data Analytics and Machine Learning
  3. Advanced Modeling Workflows
  4. Projects Reviews and Feedback
  5. Emerging, Novel Methods
  6. Spatial, Multivariate Modeling
  7. Best Practice
- Lectures, demonstration methods and workflows, hands-on experiential learning



# Work Together?

## Tabular Data with DataFrames in Python for Geoscientists and Geo-engineers

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

Many geoscientists and engineers struggle with getting started with **data analytics**, **machine learning** and **geostatistics** in Python, because they do not know how to handle their data. **DataFrames from the pandas package**, by Wes McKinney and current core team, is designed for high performance, easy to use tabular data structures. Here's a tutorial in Jupyter with Markdown with all the tabular structured data operations needed to build most subsurface modeling workflows. Check it out here: <https://git.io/fNgRW>.

### Jupyter Notebook Tutorial with Markdown

#### Tabular Data Structures / DataFrames in Python for Engineers and Geoscientists

Michael Pyrcz, Associate Professor, University of Texas at Austin

Contacts: [@GeostatsGuy](https://twitter.com/GeostatsGuy) | [www.linkedin.com/in/GeostatsGuy](https://www.linkedin.com/in/GeostatsGuy) | GoogleScholar | Book

This is a tutorial for demonstration of Tabular Data Structures in Python. In Python, the common tool for dealing with Tabular Data Structures is the DataFrame from the pandas Python package.

This tutorial includes the main operations that would commonly be required for Engineers and Geoscientists working with Tabular Data Structures for the purpose:

1. Data Cleaning and Cleaning
2. Data Analysis and Data Analytics
3. Data Analytics / Building Predictive Models with Geostatistics and Machine Learning

#### Table Data Structures

In Python we commonly store our data in three formats, tables and arrays. For sampled data with highly multiple features (e.g.,  $i = 1, \dots, n$ ),  $j = 1, \dots, m$ ) we will work with tables. For exhaustive maps and models usually representing a single feature on a regular grid ( $i = 1, \dots, n$ ,  $j = 1, \dots, m$ ), we will work with arrays.

pandas package provides a convenient DataFrame object for working with data in a table and numpy package provides a convenient ndarray object for working with gridded data. In the following tutorial we will focus on DataFrames although we will utilize ndarray a couple of times. There is another section on Gridded Data Structures that focuses on arrays.

#### Project Goal

Learn the basics for working with Tabular Data Structures in Python.

#### Load the required Libraries

The following code loads the required libraries.

Import as:

Import as np

Import as pd

Set the working directory

If you get a package Import error, you may have to first install some of these packages. This can usually be accomplished by opening a command window on Windows and then typing python -m pip install [package-name]. More assistance is available with the respective package documentation.

Set the working directory

I always like to do this so I don't lose file and to simplify subsequent read and writes (including the full address each time). Also, if you want to save your workspace to place it in the required (see below) file do it in this directory. Then we are done with this tutorial we will write our new dataset back to this directory.

os.chdir(r"\\F:\PESST")

# set the working directory

### Summary Statistics / Data Checking

X	Y	coarseness	id	min	25%	50%	75%	max	
203	293	400000	1113520461	25.000000	1112.000000	2164.000000	3565.000000	3655.000000	
Y		200	1878	150000	4.000000	1081.000000	503.000000	1685.000000	2762.000000
faults									
porosity		203	8	0.43802	0.63200	0.63000	0.63200	0.70100	
permeability		203	28	0.287462	0.447135	0.81680	4.425684	14.50700	46.14000
acoustic_impedance		203	30	1.56220	2.08600	2.47525	2.64665	3.00000	3.80000
penetrator		203	32	14.35000	14.35000	14.35000	14.35000	14.35000	14.35000
poroperm		203	34	33.37000	313.90000	230.00000	19.04669	27.77212	39.78000
permeability0		203	36	25.18000	34.65743	30.00000	47.25848	14.50700	46.14000
permeability100		203	38	1.16000	1.16000	1.16000	1.16000	1.16000	1.16000

### Data Manipulation / Data Formatting

X	Y	factors	porosity	permeability	acoustic impedance	penetrator	coarseness	porosity0	poroperm
0.165	1485.1	0.310	0.170	2.200	0.0001	low	11.64	52.11468	1.16000
0.165	1486.1	0.310	0.275	2.864	0.0001	high	16.64	48.07243	1.16000
0.165	1487.1	0.310	0.375	3.528	0.0001	high	21.64	44.03018	1.16000
0.165	1488.1	0.310	0.475	4.192	0.0001	high	26.64	40.98793	1.16000
0.165	1489.1	0.310	0.575	4.856	0.0001	high	31.64	37.94568	1.16000
0.165	1490.1	0.310	0.675	5.520	0.0001	high	36.64	34.90343	1.16000
0.165	1491.1	0.310	0.775	6.184	0.0001	high	41.64	31.86118	1.16000
0.165	1492.1	0.310	0.875	6.848	0.0001	high	46.64	28.81893	1.16000
0.165	1493.1	0.310	0.975	7.512	0.0001	high	51.64	25.77668	1.16000

X	Y	factors	porosity	permeability	acoustic impedance	penetrator	coarseness	porosity0	poroperm
0.165	1485.1	0.310	0.170	2.200	0.0001	low	11.64	52.11468	1.16000
0.165	1486.1	0.310	0.275	2.864	0.0001	high	16.64	48.07243	1.16000
0.165	1487.1	0.310	0.375	3.528	0.0001	high	21.64	44.03018	1.16000
0.165	1488.1	0.310	0.475	4.192	0.0001	high	26.64	40.98793	1.16000
0.165	1489.1	0.310	0.575	4.856	0.0001	high	31.64	37.94568	1.16000
0.165	1490.1	0.310	0.675	5.520	0.0001	high	36.64	34.90343	1.16000
0.165	1491.1	0.310	0.775	6.184	0.0001	high	41.64	31.86118	1.16000
0.165	1492.1	0.310	0.875	6.848	0.0001	high	46.64	28.81893	1.16000
0.165	1493.1	0.310	0.975	7.512	0.0001	high	51.64	25.77668	1.16000

### Data Extraction

Index	factors	porosity	permeability	acoustic impedance
0	0.1150	4.170	2.000	
1	0.1151	4.170	2.000	
2	0.1152	4.170	2.000	
3	0.1153	4.170	2.000	
4	0.1154	4.170	2.000	

Index	factors	porosity	permeability	acoustic impedance	penetrator	coarseness	porosity0	poroperm
0	0.1150	4.170	2.000		1.16000	1.16000	1.16000	1.16000
1	0.1151	4.170	2.000		1.16000	1.16000	1.16000	1.16000
2	0.1152	4.170	2.000		1.16000	1.16000	1.16000	1.16000
3	0.1153	4.170	2.000		1.16000	1.16000	1.16000	1.16000
4	0.1154	4.170	2.000		1.16000	1.16000	1.16000	1.16000

### Data Transformations

Index	factors	porosity	permeability	acoustic impedance	penetrator	coarseness	porosity0	poroperm
0	0.1150	4.170	2.000		1.16000	1.16000	1.16000	1.16000
1	0.1151	4.170	2.000		1.16000	1.16000	1.16000	1.16000
2	0.1152	4.170	2.000		1.16000	1.16000	1.16000	1.16000
3	0.1153	4.170	2.000		1.16000	1.16000	1.16000	1.16000
4	0.1154	4.170	2.000		1.16000	1.16000	1.16000	1.16000

### Data Visualization



## Gridded Data with ndarrays in Python for Geoscientists and Geo-engineers

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

Many geoscientists and engineers struggle with getting started with **data analytics**, **machine learning** and **geostatistics** in Python, because they do not know how to handle their data. **ndarrays from the NumPy package**, by Jim Hugunin et al., is designed for high performance, easy to use gridded data structures. Here's a tutorial in Jupyter with Markdown with all the gridded data operations needed to build most subsurface modeling workflows.

### Jupyter Notebook Tutorial with Markdown

#### Regular Gridded Data Structures / ndarrays in Python for Engineers and Geoscientists

Michael Pyrcz, Associate Professor, University of Texas at Austin

Contacts: [@GeostatsGuy](https://twitter.com/GeostatsGuy) | [www.linkedin.com/in/GeostatsGuy](https://www.linkedin.com/in/GeostatsGuy) | GoogleScholar | Book

This is a tutorial for demonstration of Regular Gridded Data Structures in Python. In Python, a common tool for dealing with Regular Gridded Data Structures is the NumPy package from the NumPy package (on Kaggle.net).

This tutorial includes the main operations that would commonly be required for Engineers and Geoscientists working with Regular Gridded Data Structures for the purpose:

1. Data Cleaning and Cleaning
2. Data Analysis and Data Analytics
3. Data Modeling

For Data Analytics, Geostatistics, and Machine Learning.

Regular Data Structures

For example, let's say we have sampled data with wells and holes for a geostatistical workflow for subsurface modeling. I think they should be accessible to most geoscientists and engineers. Certainly, there are more advanced, more compact, more efficient methods to accomplish the same tasks. I keep the individual methods that I have found useful for geostatistical workflows for subsurface modeling. These are the typical steps for wells and holes and the grids are the interpreted or simulated models or secondary data structures.

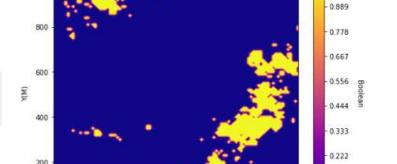
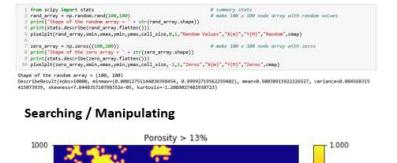
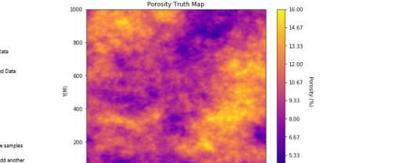
The NumPy package provides a convenient ndarray class for working with regular grids of data with wells and holes and the grids are the interpreted or simulated models or secondary data structures.

The ndarray package is typically simple from wells and holes and the grids are the interpreted or simulated models or secondary data structures.

We can read and write regular grids of data with wells and holes in NumPy.

These are the functions we have found useful:

1. `np.loadtxt` - load CSV/Geo-EAS format regular grid data to or to CSV/HDF5.
2. `np.loadtxt` - write NumPy array to CSV/Geo-EAS regular grid data or 2D.
3. `np.loadtxt` - plot 2D NumPy arrays with same parameters as CSV/EAS project.



Well-documented tutorial workflows for experiential learning.



# Work Together?

## DIRECT: Digital REservoir Characterization Technology Industrial Affiliates Proposal

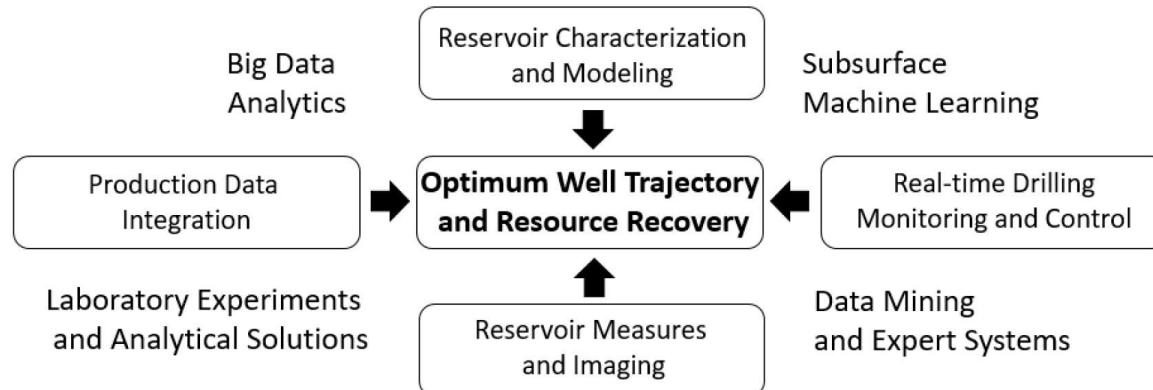
Michael J. Pyrcz<sup>1</sup>, John Foster<sup>1,2</sup>, Carlos Torres-Verdín<sup>1</sup>, and Eric van Oort<sup>1</sup>

1- Hildebrand Department of Petroleum & Geosystems Engineering, the University of Texas at Austin

2 – Institute for Computational Engineering and Science, the University of Texas at Austin

### Opportunity

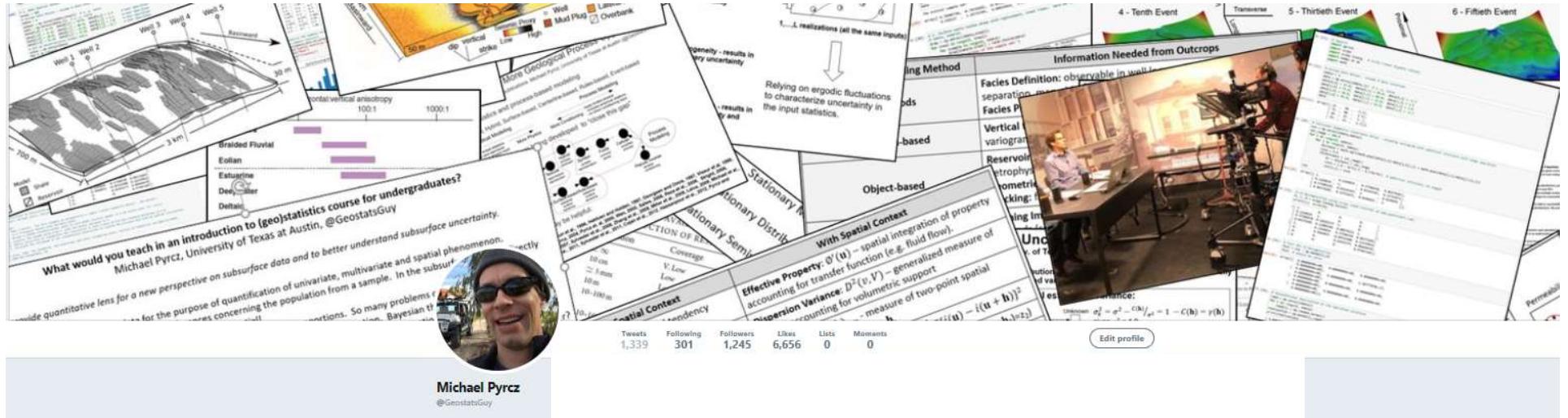
Recent numerical developments and improved computational resources have led to a rapid expansion of big data analytics and machine learning implementations. Oil and gas has a long history with big data from seismic surveys, production monitoring along with various other remote sensing and well-based data, and has developed various physics-based engineering and stochastic statistical workflows. There is an opportunity to combine best-practice and cutting-edge technology in reservoir spatiotemporal characterization and modeling, production data integration, reservoir geophysics and real-time drilling control with big data analytics and machine learning to optimize well trajectory and resource recovery.



*Optimum well trajectory and resource recovery through integration of engineering, data analytics, and machine learning.*

**Kick-off meeting in June, 2019, 3 companies indicated they will join, more interested.**

# Learn More?



For tweets with Subsurface Geostatistical, Data Analytics and Machine Learning resources -

follow @GeostatsGuy



# Learn More?

## GitHub GeostatsGuy

### Demonstration Workflows

#### Excel, R and Python

- Distributions
- Bootstrap
- Cellular Automata
- Hypothesis Testing
- Lorenz Coefficient
- Decision Making
- Bayesian Updating
- Kriging
- Simulation
- Volume-variance



Overview    Repositories 10    Stars 2    Followers 106    Following 8

#### Popular repositories

##### 2DayCourse

2 day short course.

★ 20 ⚡ 2

##### PythonNumericalDemos

A collection of Python demos for geostatistical methods.

● Jupyter Notebook ★ 10 ⚡ 5

##### ExcelNumericalDemos

A set of numerical demonstrations in Excel to assist with teaching / learning concepts in statistics and geostatistics.

★ 9 ⚡ 3

##### 2DayCourse\_Exercises

● Jupyter Notebook ★ 4

##### GeostatsPy

Wrapper / Reimplementation of GSLIB in Python

● Jupyter Notebook ★ 3 ⚡ 1

##### GeostatsLectures

(Geo)statistical course materials released for anyone to use (.pdf format). Enjoy! I'm happy to discuss.

★ 2 ⚡ 2

177 contributions in the last year



# Learn More?



The collage includes:

- A man playing a guitar.
- A document titled "Tabular Data with DataFrames" by Michael Pyrcz.
- A document titled "Analysis in Python for Geoscientists and Geo-engineers" by Michael Pyrcz.
- A heatmap titled "Irregularly Sampled Data Variogram Map".
- A document titled "An explanation of STATIONARITY for geoscientists" by Michael Pyrcz.
- A document titled "Without Replication" with various statistical terms and formulas.
- A video camera filming a group of people in a room.
- A person working at a desk with a laptop.
- A document titled "Geostatistics for Geoscientists" by Michael Pyrcz.
- A document titled "Indicator Semivariogram:  $r_1(\mathbf{h}) = \frac{1}{2} E[(Z(u)-Z(u-\mathbf{h}))^2] - E[Z(u)]^2]$ ".
- A document titled "Multiple Point Statistic:  $P_{ij} = \text{Prob}(A_i|B_{i,j}, C_{i,j})$ ".
- A document titled "Ripley's K Function,  $K(t)$  – multiscale stacking".



GeostatsGuy Lectures

For my lectures check out my YouTube Channel, ‘GeostatsGuy Lectures’.

- ## Example Topics:
- probability theory
  - frequentist vs. Bayesian statistics
  - binomial distribution to model exploration success

**Machine Learning / Statistical Learning**

Ethical Concerns:

- Biased training data
- Rideiro 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/1602.04619.pdf>