

DAYTUM – SPATIAL DATA ANALYTICS

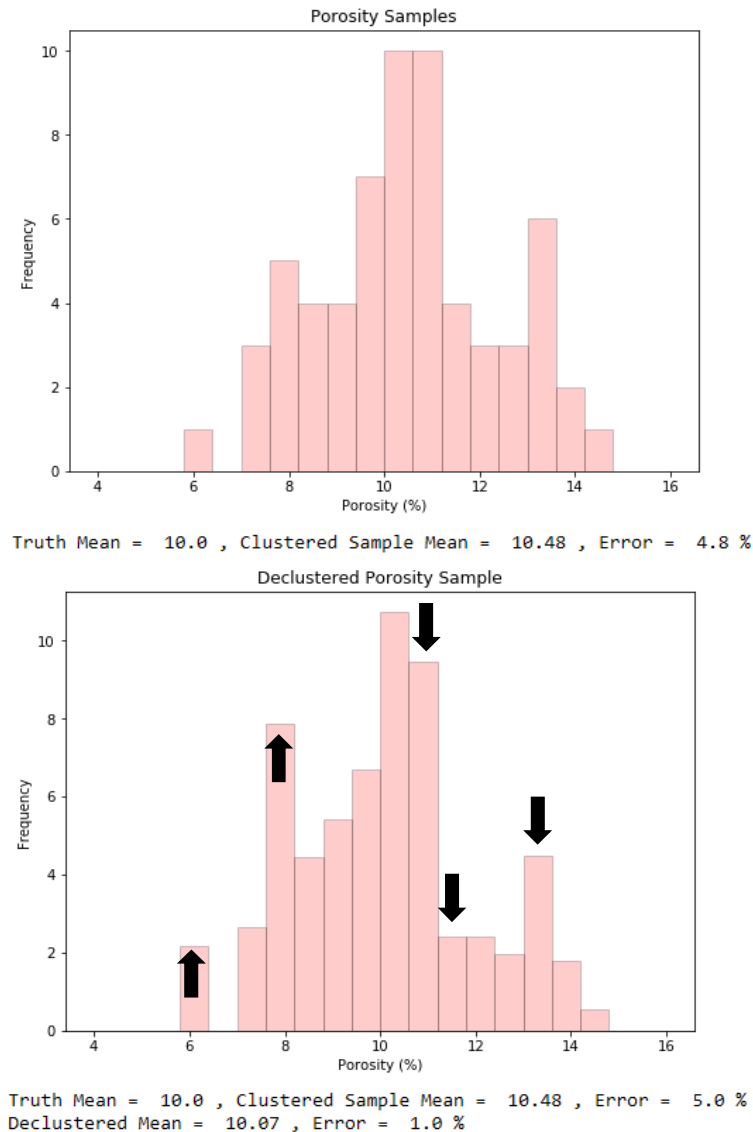
Data Preparation

Lecture outline ...

- ▶ Sampling Limitations
- ▶ Declustering
- ▶ Quantifying Uncertainty

MOTIVATION

- ▶ We work with sparse data
 - The data is sampled in a biased (nonrepresentative) manner
 - Bias in, bias out!
- ▶ There is significant uncertainty in our summary statistics and models
 - We must quantify and account for uncertainty



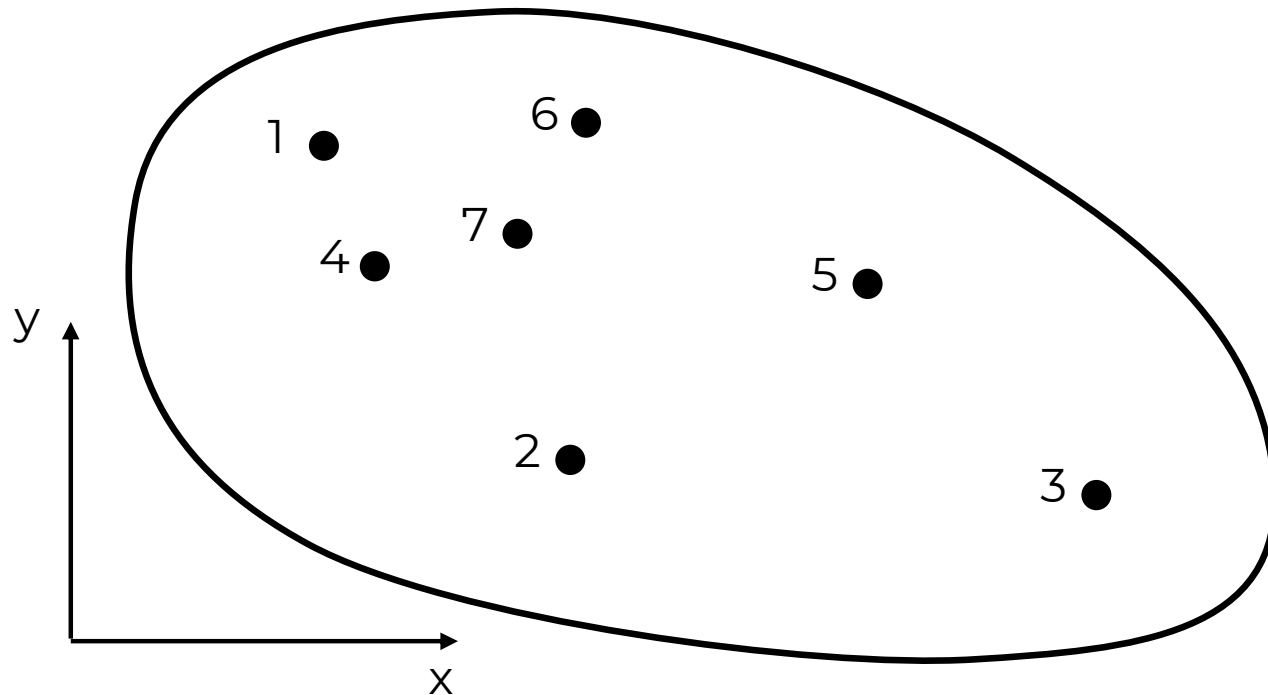
SAMPLING LIMITATIONS

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
 - How far does the reservoir extend? – offset drilling
- ▶ And to maximize NPV directly:
 - Maximize production rates

REPRESENTATIVITY

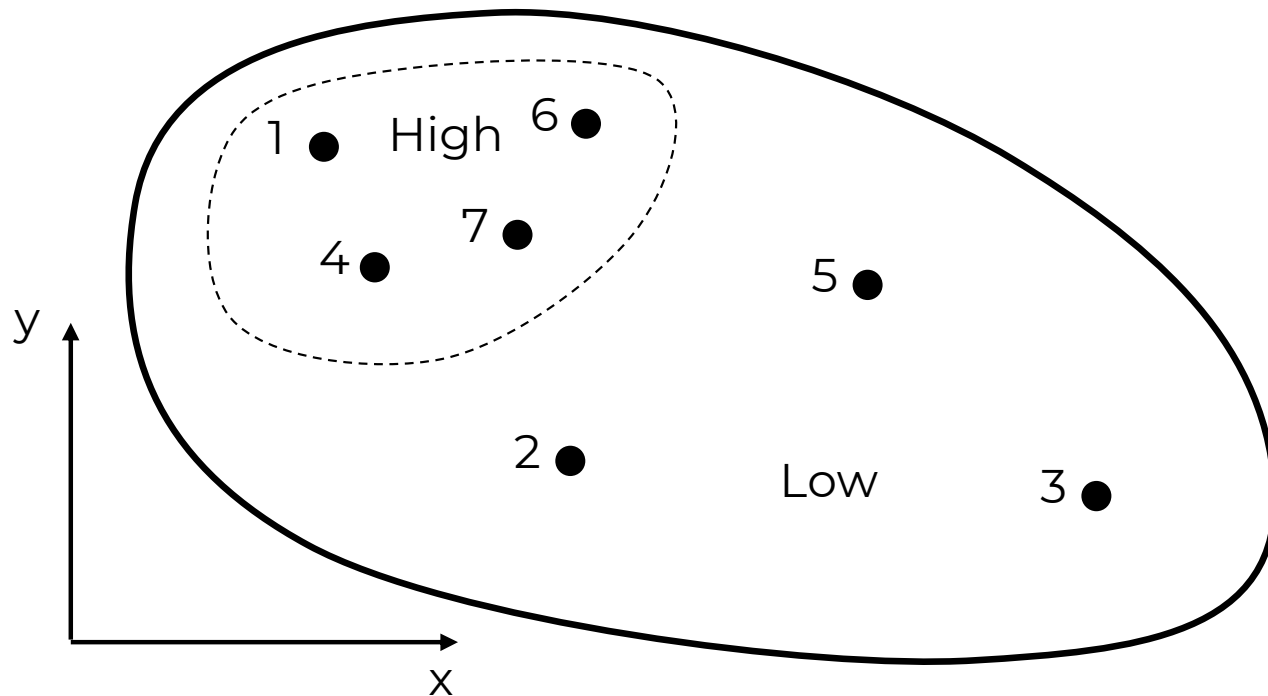
- ▶ The concern is when we attempt to make an estimate:



- ▶ e.g. the average porosity to calculate OIP

REPRESENTATIVITY

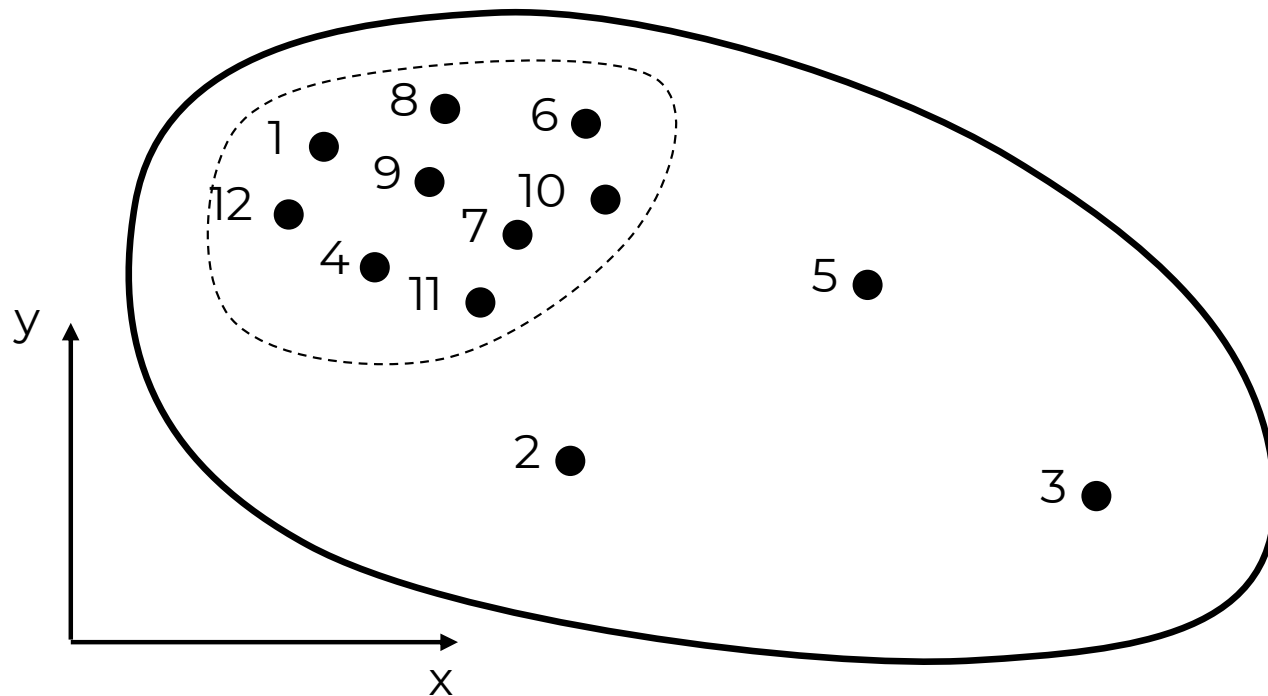
- ▶ The concern is when we attempt to make an estimate:



- ▶ What if we knew from seismic that the reservoir quality is better in the top left area?

REPRESENTATIVITY

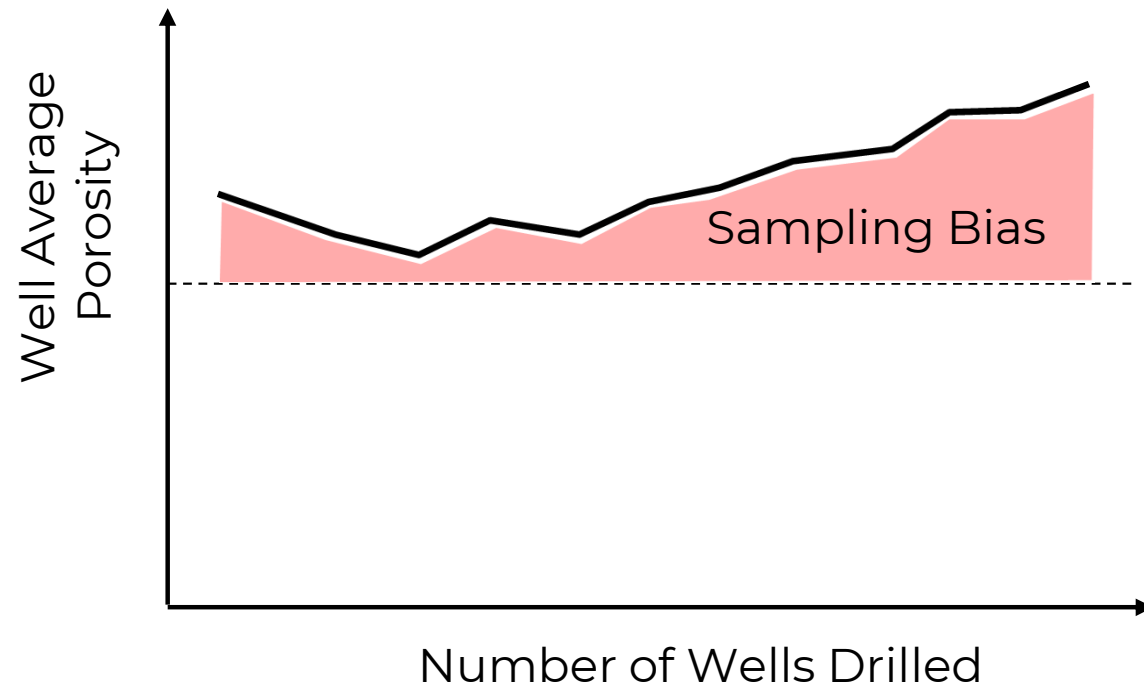
- ▶ The concern is when we attempt to make an estimate:



- ▶ What if we knew from seismic that the reservoir quality is better in the top left area?

REPRESENTATIVITY

- ▶ The concern is when we attempt to make an estimate:



- ▶ We need to mitigate this bias.
- ▶ Drilling representatively? Why do we drill in a biased manner?

(GEO)STATISTICS SAMPLING REPRESENTATIVELY

How Would We Sample for Representativity?

- ▶ **Random Sampling:** when every item in the population has an equal chance of being chosen. Selection of every item is independent of every other selection.
 - Is random sampling sufficient for subsurface? Is it available?
- ▶ **Regular Sampling:** when samples are taken at regular intervals (equally spaced).
 - Less reliable than random sampling.
 - Warning: May resonate with some unsuspected environmental variable.

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.

(GEO)STATISTICS SAMPLING BIAS

Example of Sampling Bias

1. Wells drilled in part of reservoir identified to have the greatest thickness in seismic.
2. Core extracted from the well bore in the location estimated to have the best reservoir.
3. Core plugs extracted from whole cores for porosity / permeability analysis avoiding shales.



Routine core analysis from
https://www.rigzone.com/training/insight.asp?insight_id=325.

(GEO)STATISTICS SAMPLING BIAS

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
- ▶ we can't run permeability evaluation on low permeability rock

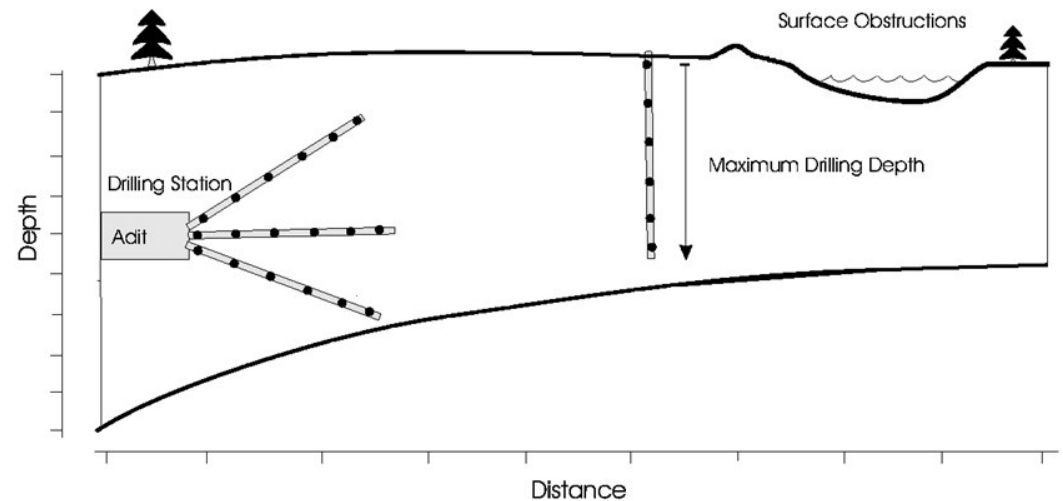


Image from Pyrcz and Deutsch (2003)
<http://gaa.org.au/pdf/DeclusterDebias-CCG.pdf>

SPATIALLY CLUSTERED DATA

- ▶ Data are rarely collected for their statistical representativity:
 - Wells are drilled in areas with the greatest probability of high production
 - Horizontal wells target stratigraphic zones of interest (high pay)
 - Core are taken preferentially from good quality reservoir rock
 - These data collection practices should not be changed:
 - best economics
 - most data in the most important locations

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!

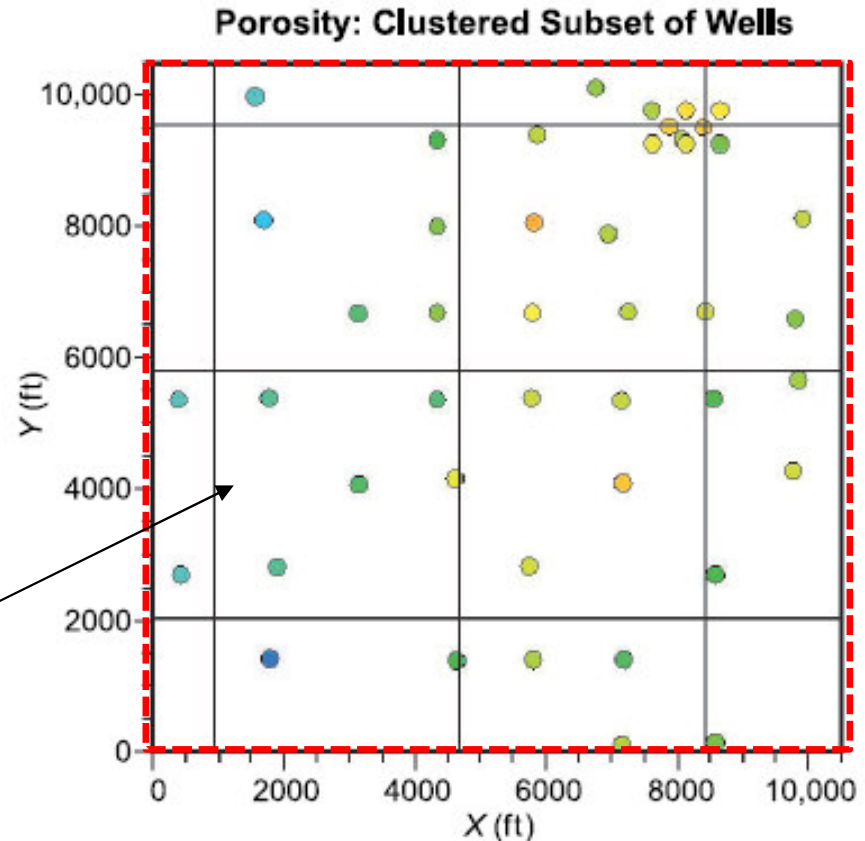
- 1. Mapping:** and summarizing average over map
- 2. Use of Regions:** to pool and use statistical over volumes of high / low reservoir quality (e.g., facies)
- 3. Declustering techniques** assign each datum a weight based on closeness to surrounding data
 - $w_i, i = 1, \dots, n$ (weights are greater than 0 and sum to n)
 - Histogram and cumulative histogram use $w_i, i = 1, \dots, n$ instead of equal weighted, $w_i = 1.0$.
- 4. Debiasing techniques** derive an entirely new distribution based on a secondary data source such as geophysical measurements or expert interpretation

(GEO)STATISTICS GOAL OF SAMPLING AND STATISTICS EXAMPLE

Addressing Bias

- Would it be fair to calculate the average of these wells and to apply that as an average for this area of interest?

What is the average porosity over this reservoir?

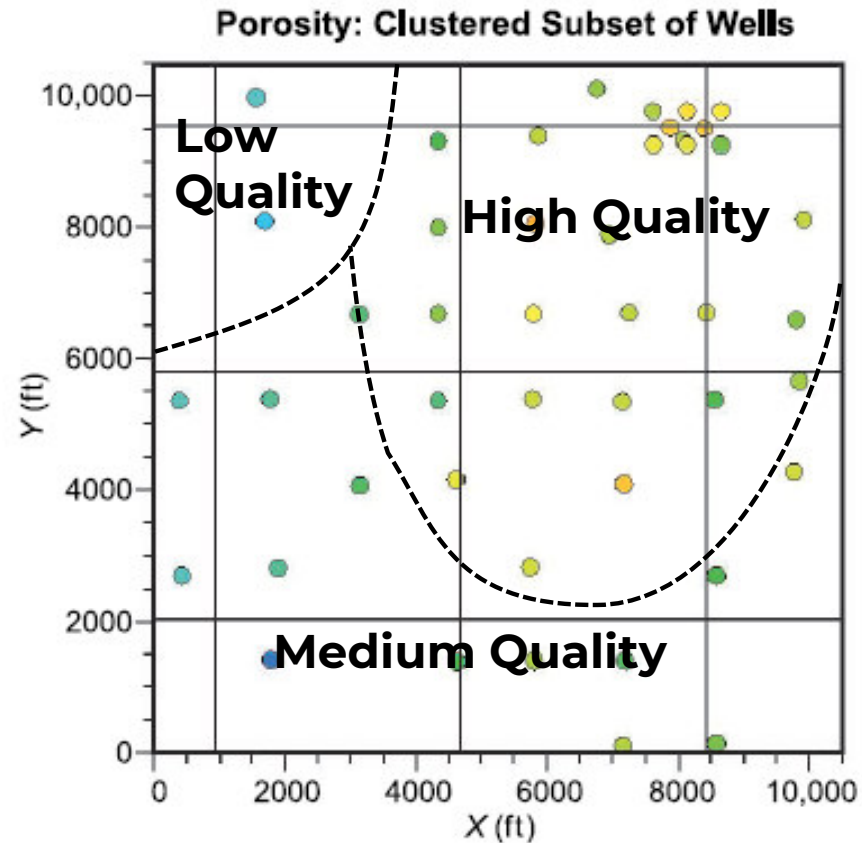


Porosity sample data for an example reservoir (Pyrz and Deutsch, 2014).

(GEO)STATISTICS SAMPLING BIAS

Addressing Bias with Regions

- ▶ Would it be fair to calculate the average of these wells and to apply that as an average for this area of interest?
- ▶ Break model up into subsets.
 - Avoid densely sampled high quality reservoir inflating average over the entire reservoir

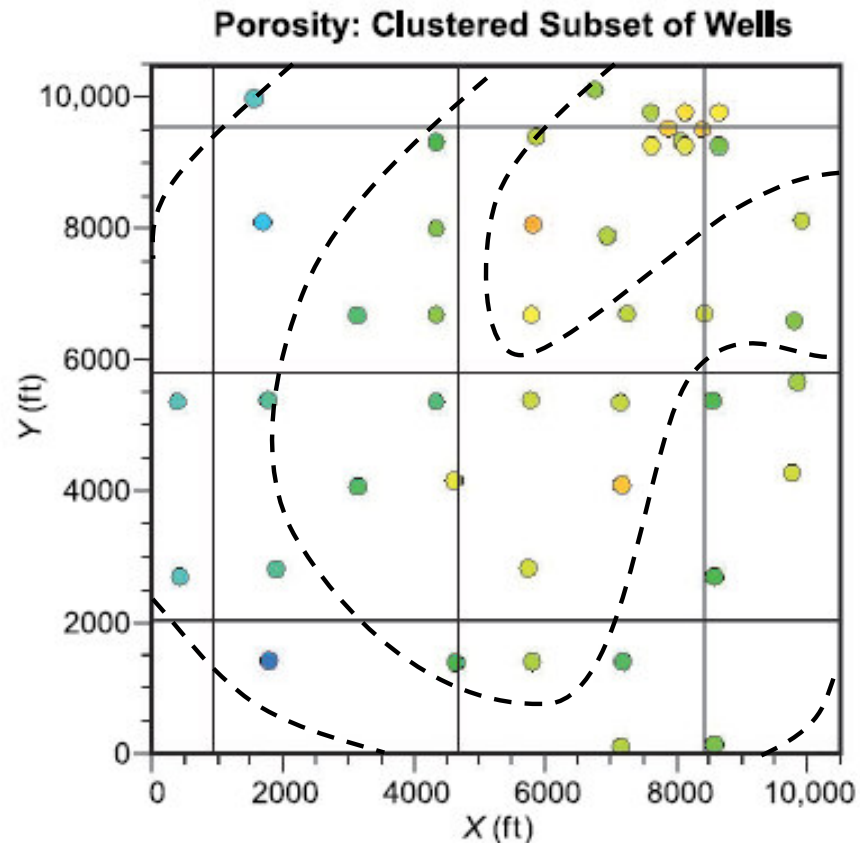


Porosity sample data for an example reservoir (Pyrzcz and Deutsch, 2014).

(GEO)STATISTICS SAMPLING BIAS

Addressing Bias with Geological Mapping

- ▶ Would it be fair to calculate the average of these wells and to apply that as an average for this area of interest?
- ▶ Build a map of the property of interest.
- ▶ Calculate the average of the map
 - Avoid densely sampled high

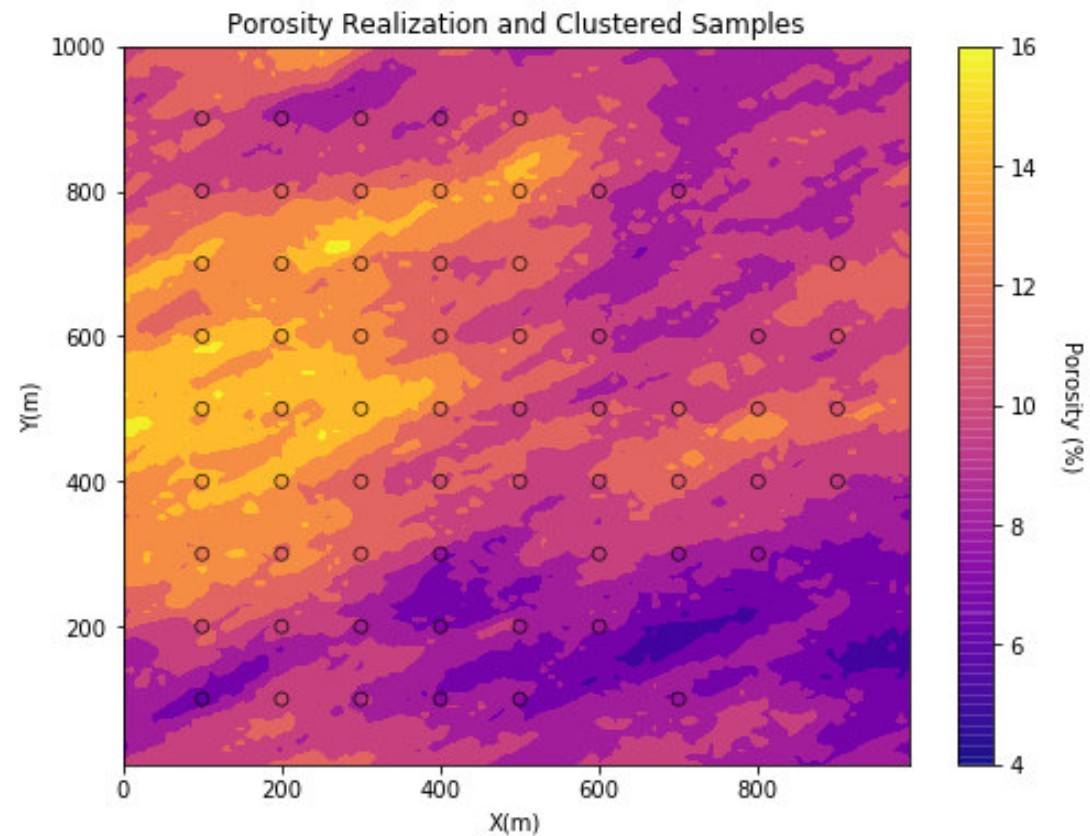


Porosity sample data for an example reservoir (Pyrzcz and Deutsch, 2014).

DECLUSTERING

SPATIALLY CLUSTERED DATA EXAMPLE

What is wrong with this sample set?



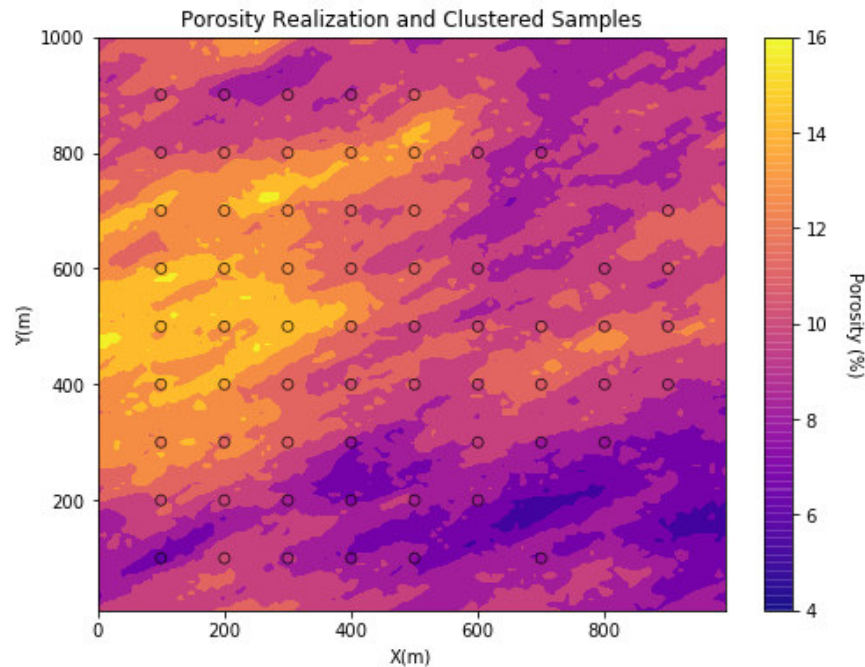
Reservoir population and sample data.

SPATIALLY CLUSTERED DATA

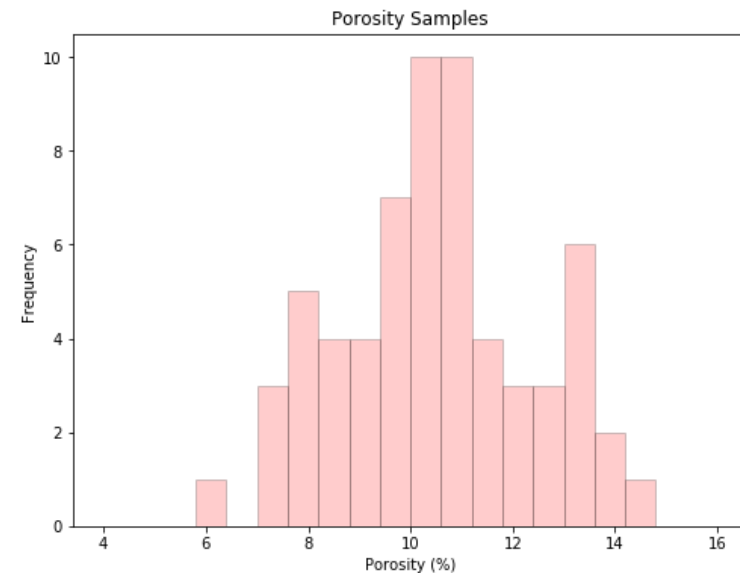
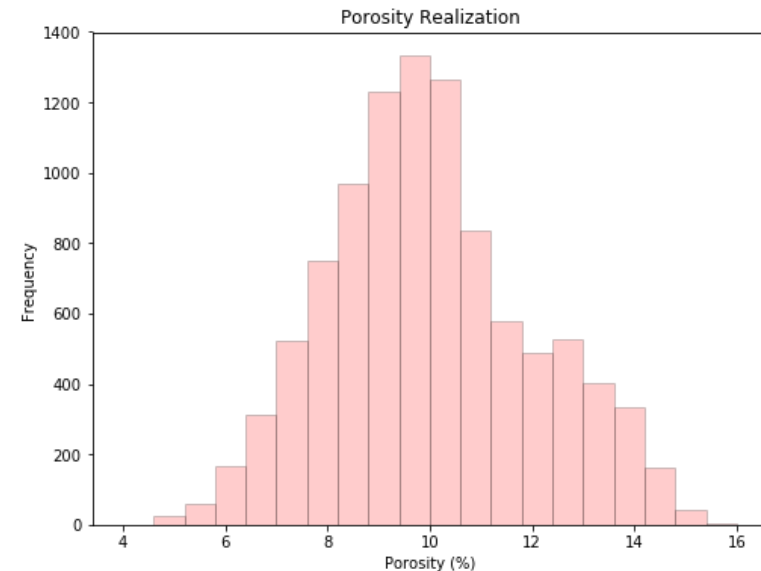
- ▶ There is a need, however, to adjust the histograms and summary statistics to be representative of the entire volume of interest.
- ▶ Declustering techniques assign each datum a weight based on closeness to surrounding data
 - $w_i, i = 1, \dots, n$ (weights are greater than 0 and sum to n)
 - Histogram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

SPATIALLY CLUSTERED DATA

- ▶ Location map of 64 wells with truth model
- ▶ See the error between the samples and the underlying truth model



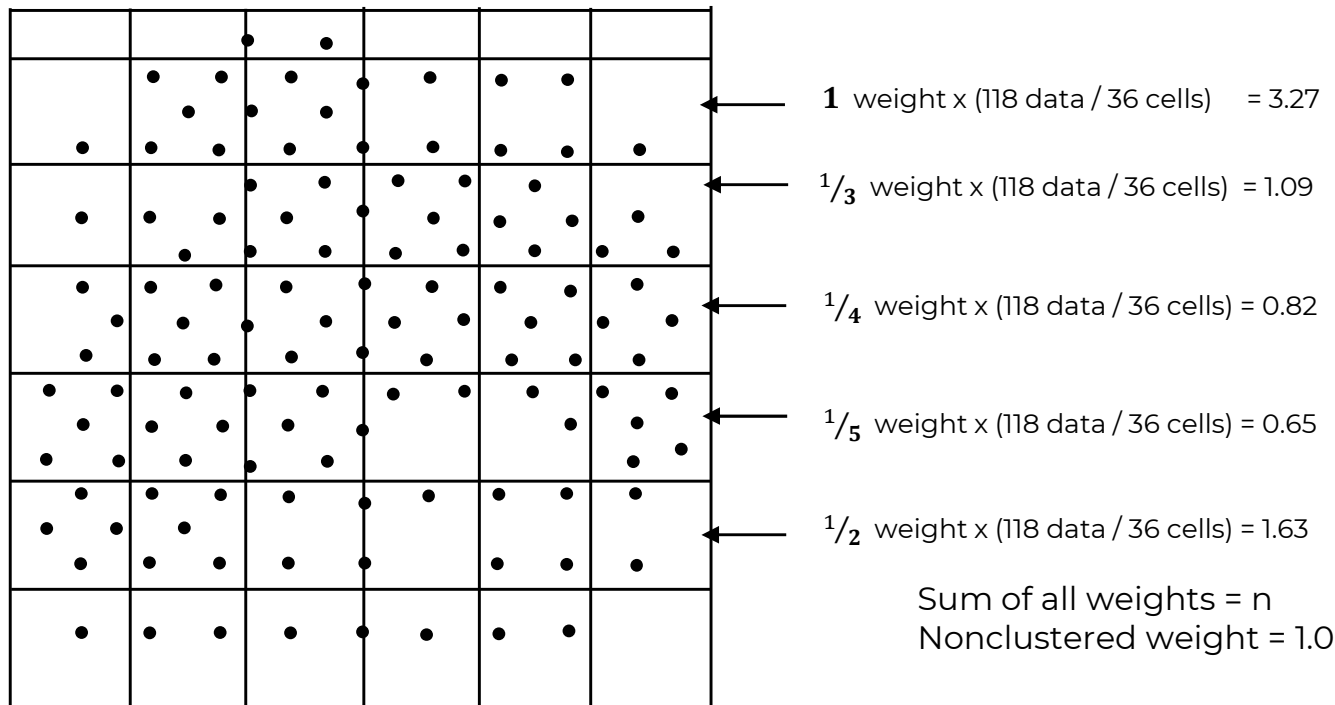
Reservoir population and sample data (above), population distribution (upper right) and sample distribution (lower right).



Truth Mean = 10.0 , Clustered Sample Mean = 10.48 , Error = 4.8 %

CELL DECLUSTERING

- ▶ Cell Declustering, is robust in 3-D and when the limits are poorly defined:
 - divide the volume of interest into a grid of cells $l=1,\dots,L$
 - count the occupied cells L_o and the number in each cell n_l , $l=1,\dots, L_o$
 - weight inversely by number in cell (standardize by L_o)

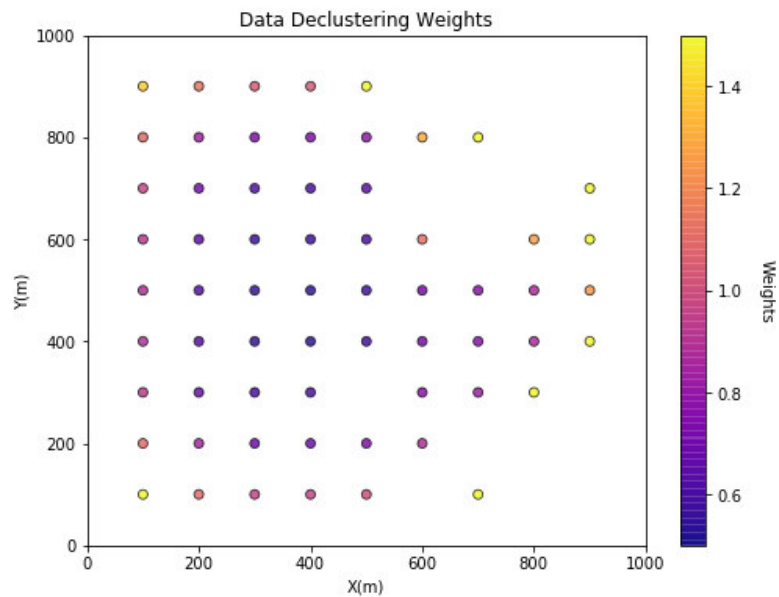
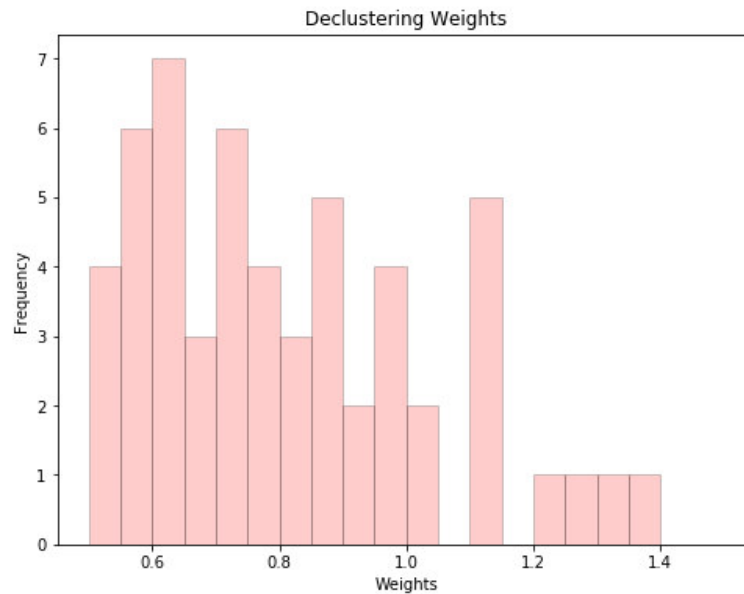


Spatial data and cell-based declustering cells.

- ▶ The issue, of course, is how to choose the cell size...

DECLUSTERING WEIGHTS

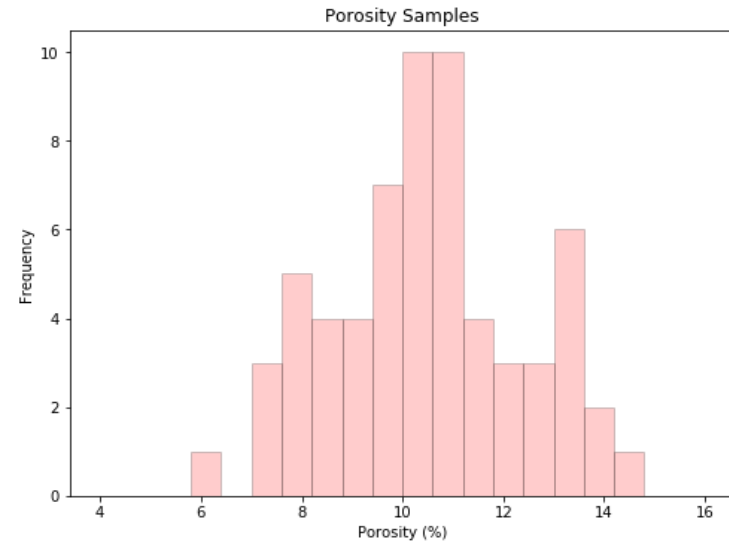
- Declustering weights
 1. 1.0 nominal weight
 2. < 1.0 reduced weight
 3. > 1.0 increased weight



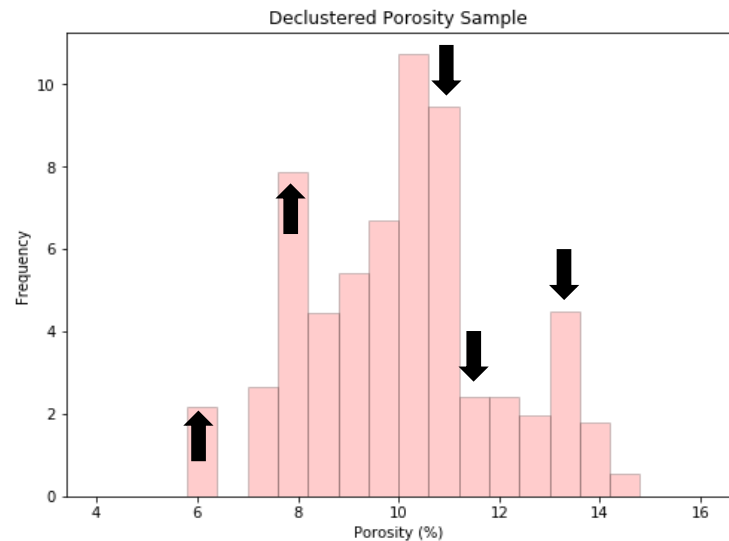
DECLUSTERED DISTRIBUTION

- ▶ Updated distribution with declustering weights
- ▶ Now data file include values and weights based on spatial arrangement.
- ▶ Possible to calculate any weighted statistic.
- ▶ For example, declustered mean:

$$\bar{z} = \frac{\sum_i^n w_i z_i}{\sum_i^n w_i = n}$$



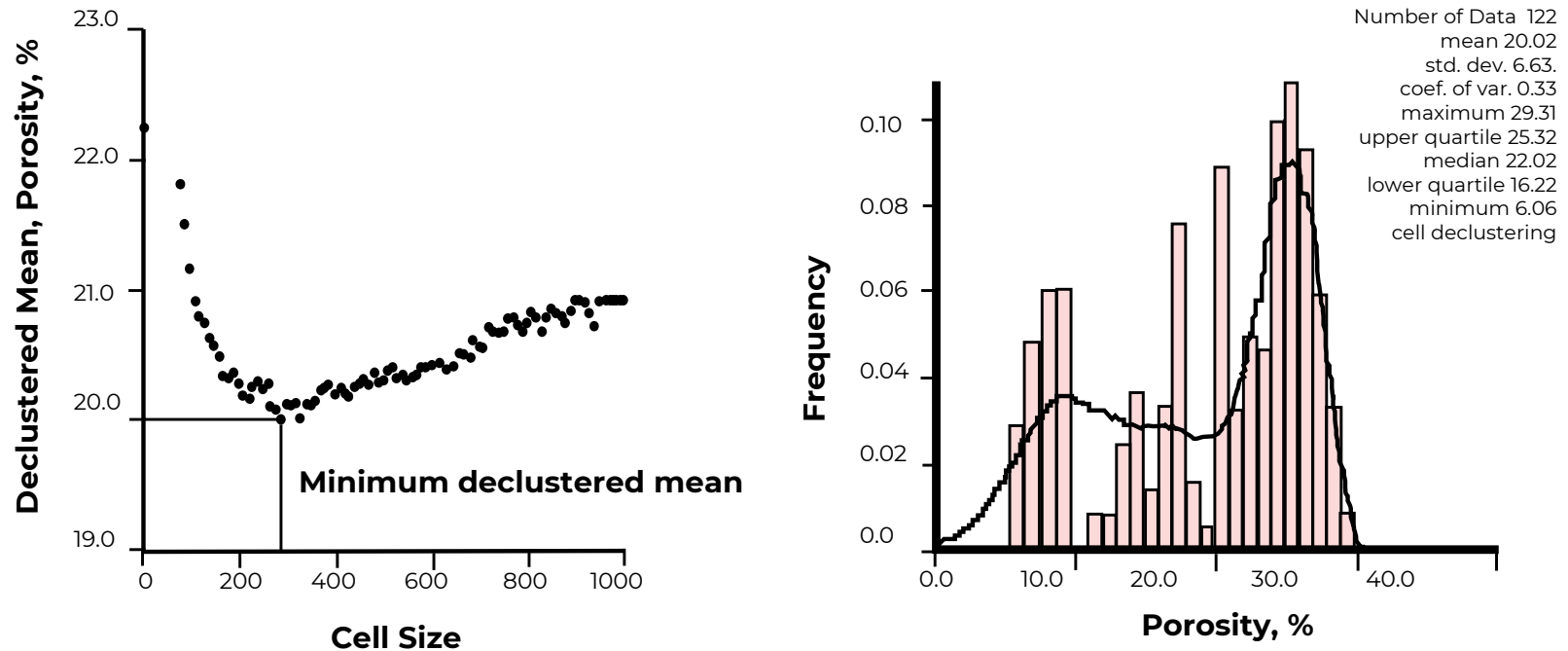
Truth Mean = 10.0 , Clustered Sample Mean = 10.48 , Error = 4.8 %



Truth Mean = 10.0 , Clustered Sample Mean = 10.48 , Error = 5.0 %
Declustered Mean = 10.07 , Error = 1.0 %

THE CELL SIZE

- Plot declustered mean versus the cell size for a range of cell sizes:



- There is no theory that says we are looking for a minimum when the values are clustered in high values or a maximum when clustered in low values – it just seems to make sense
- The result can be very sensitive to large scale trends – it is often better to choose the cell size by visual inspection and some sensitivity studies

PYTHON / GSLIB DECLUSTERING DEMO

Here's a Workflow that Used GeostatsPy

► Things to demonstrate:

1. Load data, visualize
2. Cell-base declustering
3. Visualize the weights
4. Calculated weighted statistics
5. Check multiple cell sizes

Walk through the workflow together.

GeostatsPy: Basic Univariate Statistics and Distribution Representativity for Subsurface Data Analytics in Python

Michael Pyrcz, Associate Professor, University of Texas at Austin

[Twitter](#) | [GitHub](#) | [Website](#) | [GoogleScholar](#) | [Book](#) | [YouTube](#) | [LinkedIn](#)

PGE 383 Exercise: Basic Univariate Summary Statistics and Data Distribution Representativity Plotting in Python with GeostatsPy

Here's a simple workflow with some basic univariate statistics and distribution representativity. This should help you get started data declustering to address spatial sampling bias.

Geostatistical Sampling Representativity

In general, we should assume that all spatial data that we work with is biased.

Source of Spatial Sampling Bias

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

Random Sampling: when every item in the population has a equal chance of being chosen. Selection of every item is independent of every other selection. Is random sampling sufficient for subsurface? Is it available?

- it is not usually available, would not be economic
- data is collected answer questions
 - how large is the reservoir, what is the thickest part of the reservoir
- and wells are located to maximize future production
 - dual purpose appraisal and injection / production wells!

Regular Sampling: when samples are taken at regular intervals (equally spaced).

- less reliable than random sampling.
- Warning: may resonate with some unsuspected environmental variable.

What do we have?

- we usually have biased, opportunity sampling
- we must account for bias (debiasing will be discussed later)

So if we were designing sampling for representativity of the sample set and resulting sample statistics, by theory we have 2 options, random sampling and regular sampling.

Data File is at: <https://git.io/fh0CW> and Jupyter Notebook Workflow is at: <https://git.io/fhgJl>

PROBABILITY AND STATISTICS

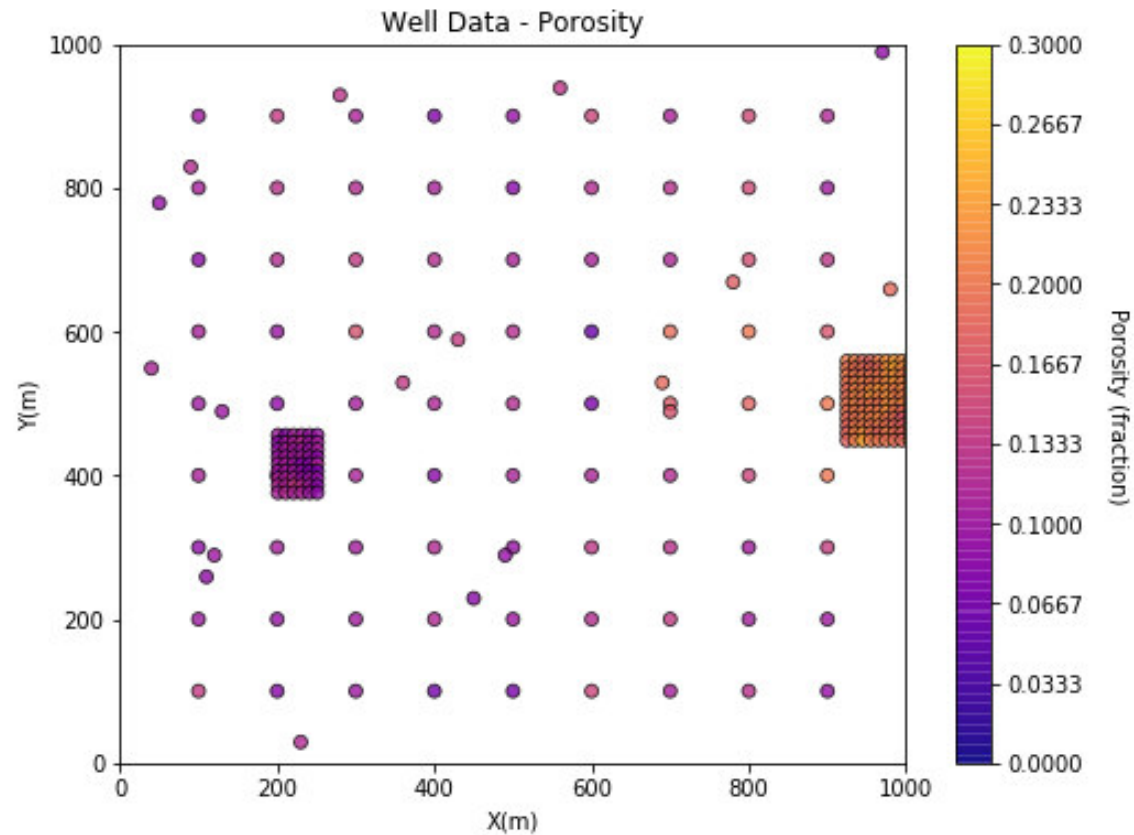
New Tools

| Topic | Application to Subsurface Modeling |
|--------------------------|---|
| Awareness | <p>Every subsurface dataset is sampled to answer questions and add value, not for statistical representativity.</p> <p><i>Assume all data sets are biased, test for bias.</i></p> |
| Cell Declustering | <p>Given the spatial location of the sample data, calculate declustering weights.</p> <p><i>Build representative sample statistics that correct for sampling bias.</i></p> |

QUANTIFYING UNCERTAINTY

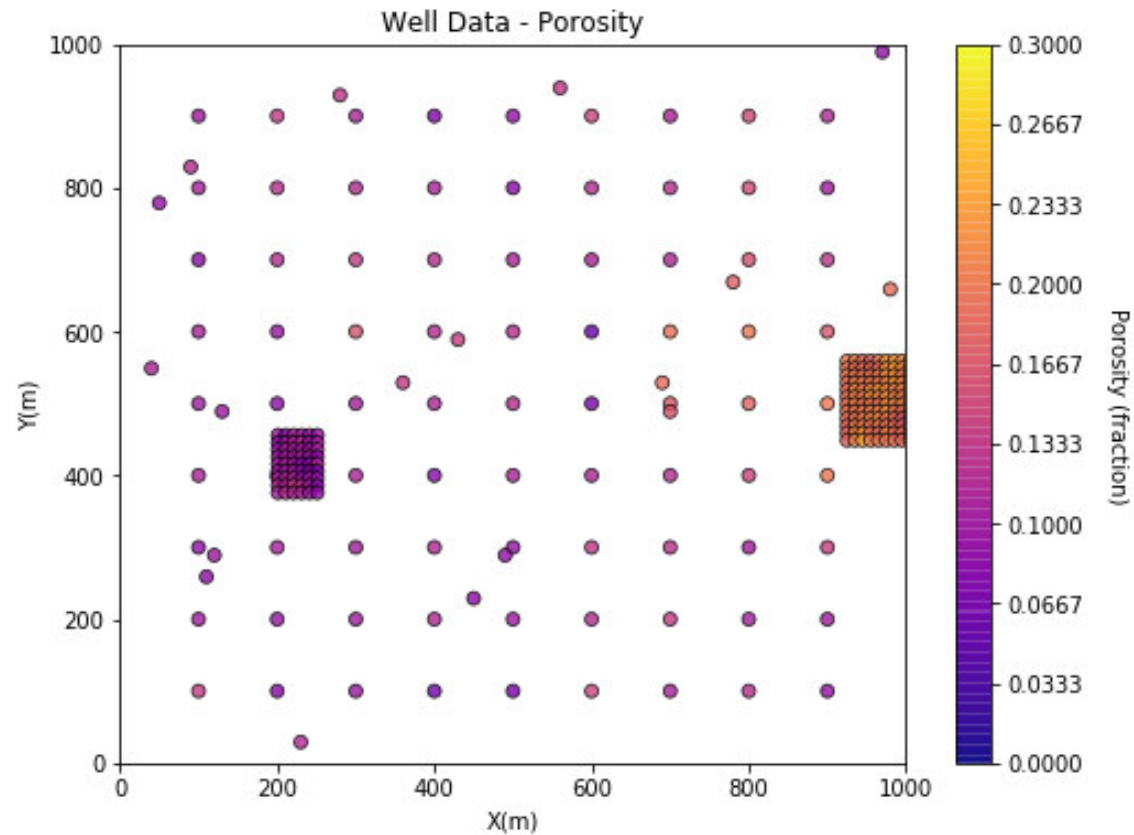
BOOTSTRAP MOTIVATION

- Uncertainty in the Sample Statistics
 - One source of uncertainty is the paucity of data
 - Do these 200 or so wells provide a precise (and accurate estimate) of the mean? standard deviation? skew? P13?



BOOTSTRAP MOTIVATION

- ▶ Would it be Useful to Know the Uncertainty in these Statistics Due to Limited Sampling?
 - What is the impact of uncertainty in the mean porosity e.g. 20%+/-2%?



BOOTSTRAP DEFINITION

▶ Bootstrap

- method to assess the uncertainty in a sample statistic by repeated random sampling with replacement

▶ Assumptions

- sufficient, representative sampling

▶ Limitations

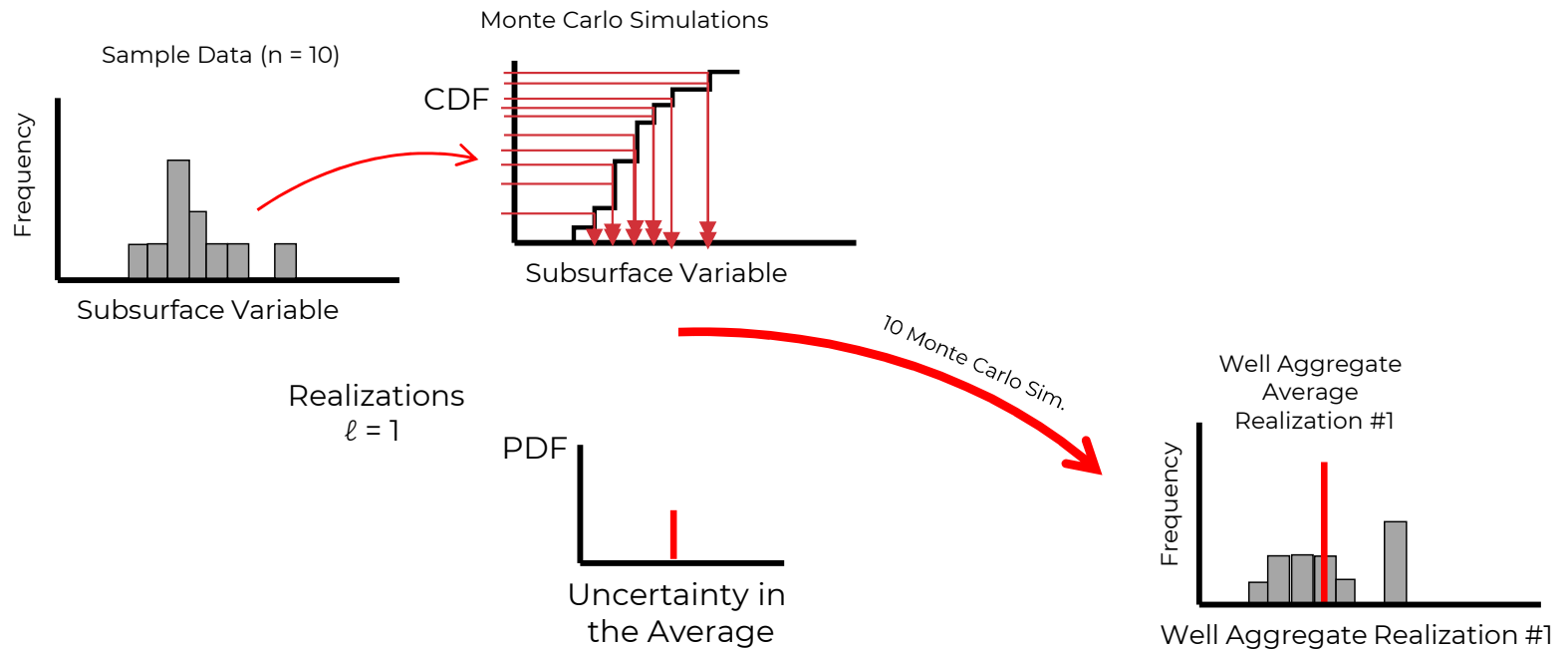
- assumes the samples are representative
- assumes stationarity
- only accounts for uncertainty due to too few samples, e.g. no uncertainty due to changes away from data
- does not account for area of interest
- assumes the samples are independent
- does not account for other local information sources

**No spatial
Context**

UNIVARIATE STATISTICS

Bootstrap

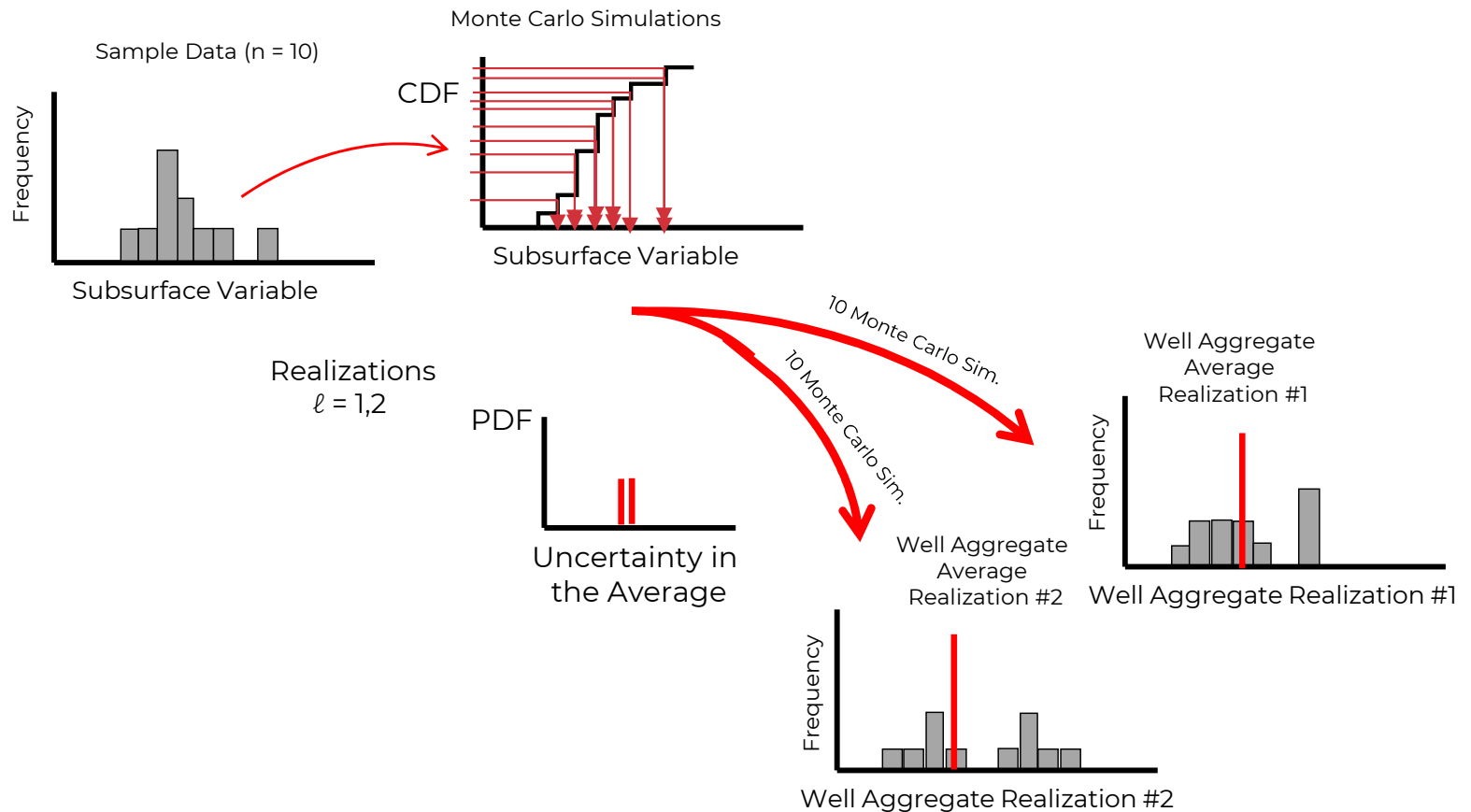
► Bootstrap for Uncertainty in the Mean



UNIVARIATE STATISTICS

Bootstrap

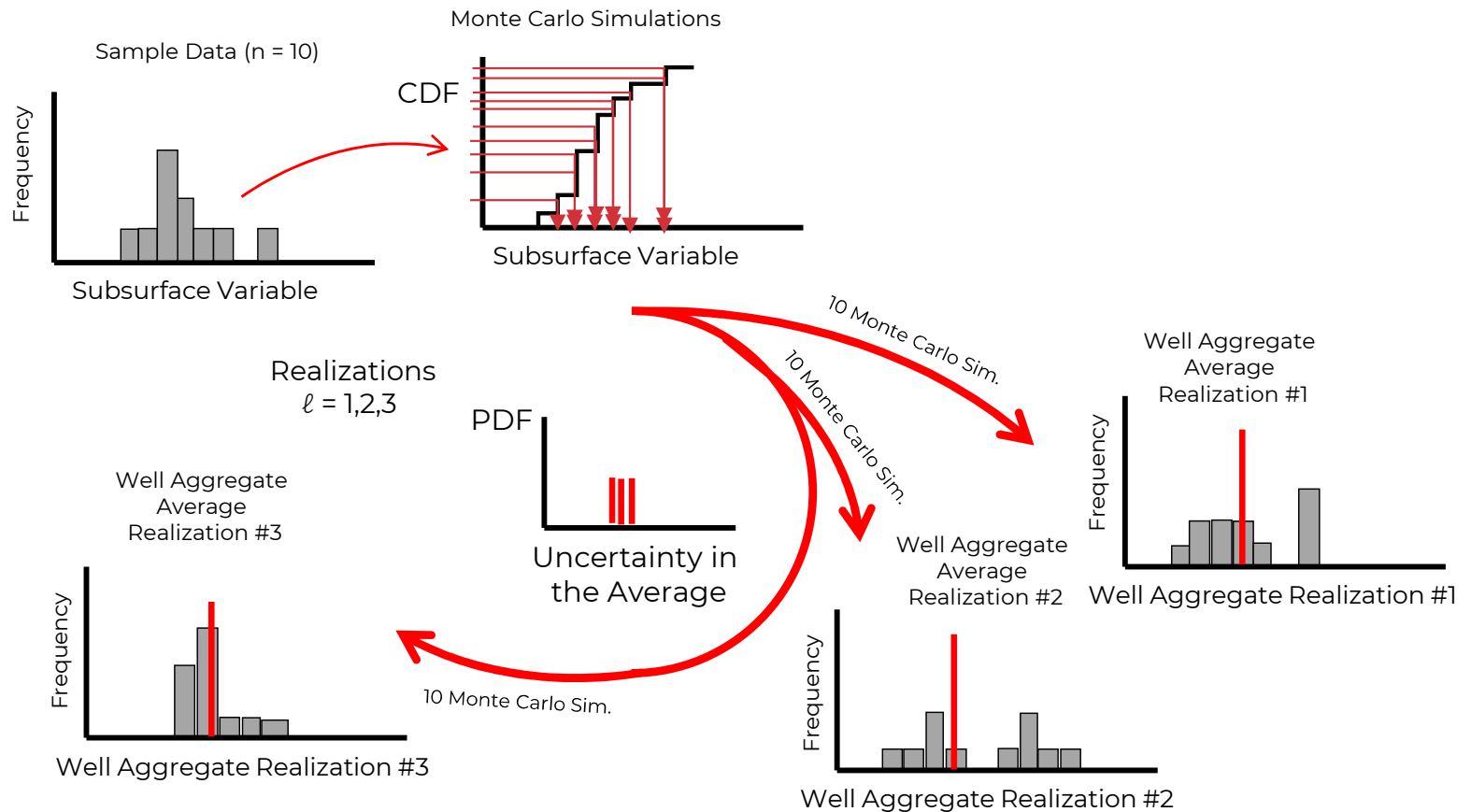
► Bootstrap for Uncertainty in the Mean



UNIVARIATE STATISTICS

Bootstrap

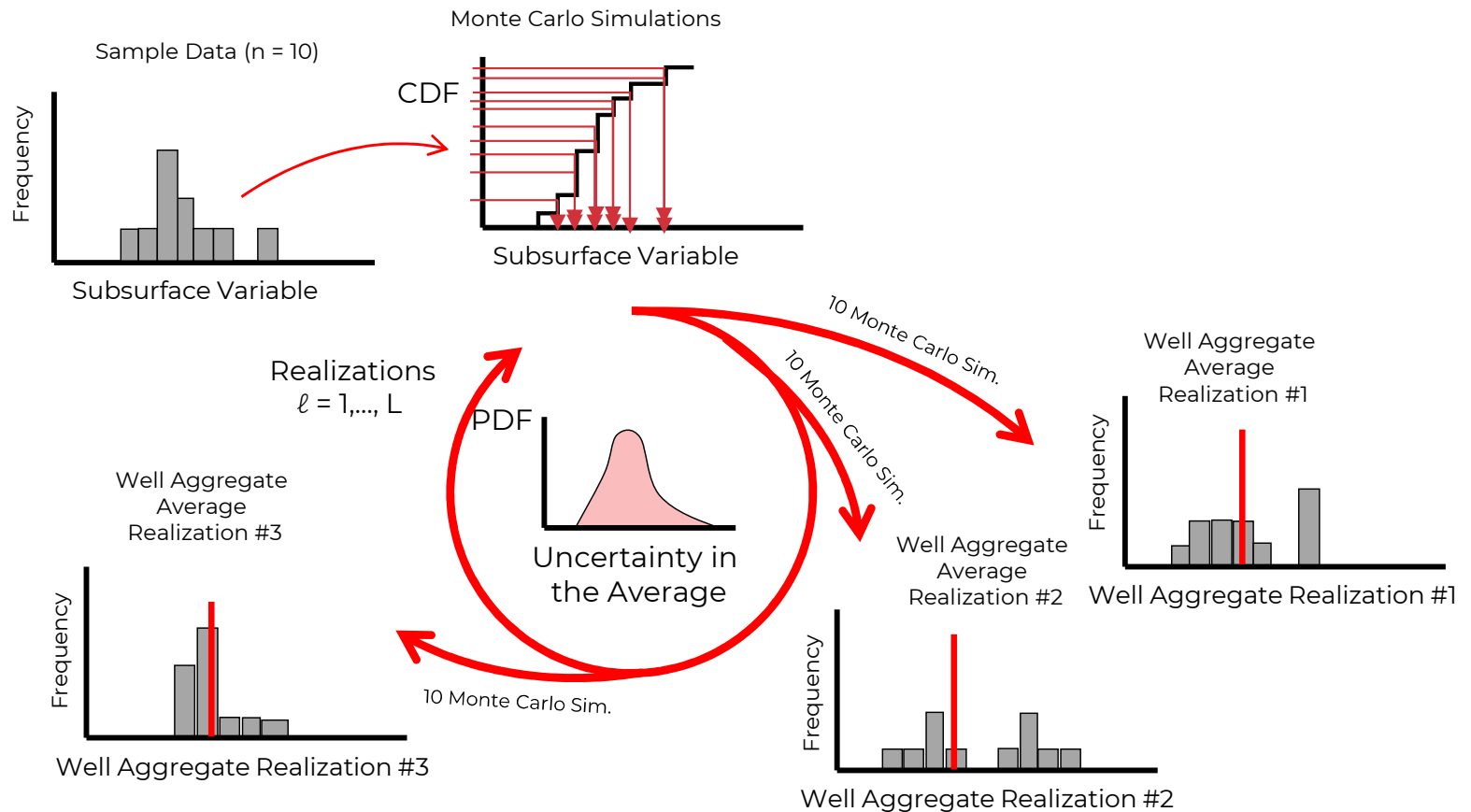
► Bootstrap for Uncertainty in the Mean



UNIVARIATE STATISTICS

Bootstrap

► Bootstrap for Uncertainty in the Mean



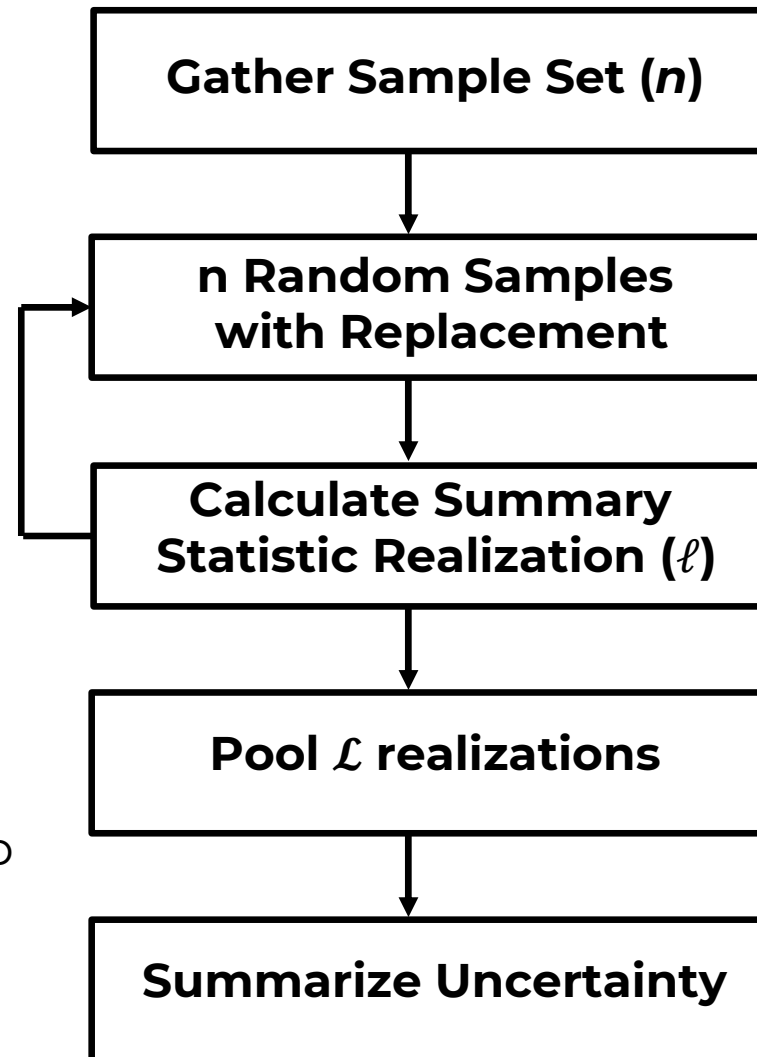
UNIVARIATE STATISTICS

Bootstrap

- ▶ Bootstrap Approach (Efron, 1982)
- ▶ Statistical resampling procedure to calculate uncertainty in a calculated statistic from the data itself
- ▶ For uncertainty in the mean solution is standard error:

$$\sigma_{\bar{x}}^2 = \frac{\sigma_s^2}{n}$$

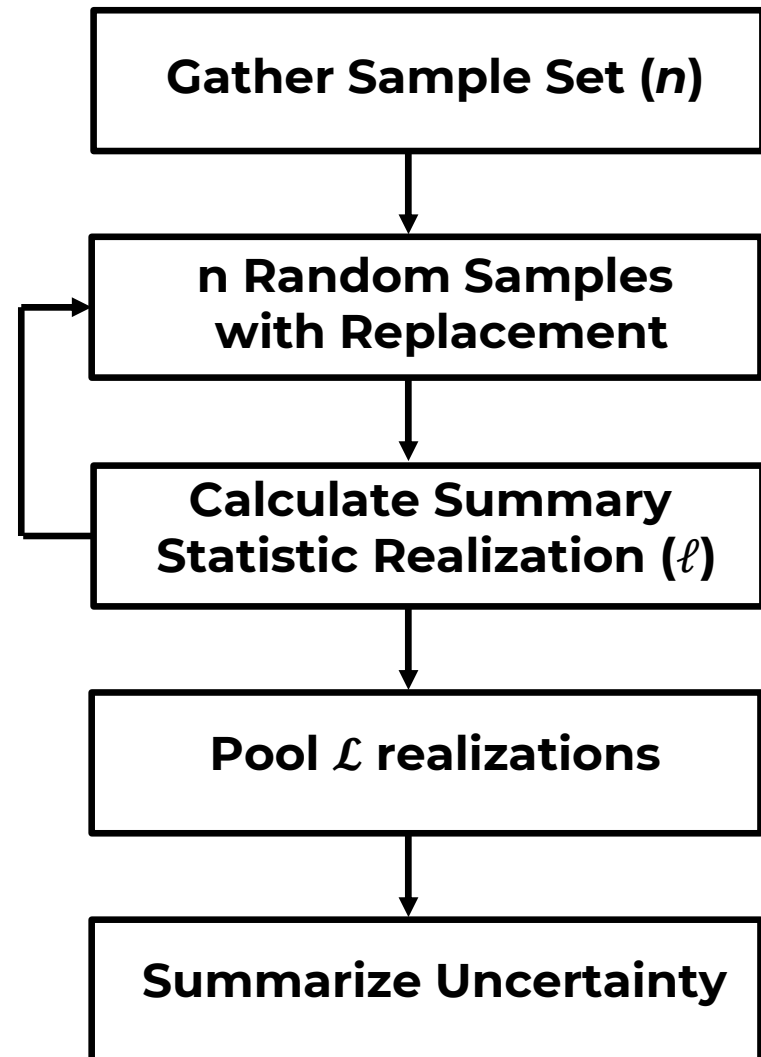
- ▶ Extremely powerful. Could get uncertainty in any statistic! (e.g., P13, skew etc.)
- ▶ Would not be possible without bootstrap
- ▶ Advanced forms account for spatial information and strategy (game theory)



UNIVARIATE STATISTICS

Bootstrap

- ▶ You now know about one of the most powerful tools ever!
- ▶ Caveats:
 - Assumes the sample set is representative
 - Unbiased and covers the full range
 - Assumes all samples are independent if not consider Journel's spatial bootstrap (1993)
- ▶ You can do bootstrap in Excel



UNIVARIATE STATISTICS

Bootstrap Demo

► Things to demonstrate:

1. Load data, visualize
2. Summary statistics
3. Bootstrap Realizations
4. Summarization over Bootstrap Realizations
5. Uncertainty in Average and Standard Deviation

GeostatsPy: Bootstrap for Subsurface Data Analytics in Python

Michael Pyrcz, Associate Professor, University of Texas at Austin

[Twitter](#) | [GitHub](#) | [Website](#) | [Google Scholar](#) | [Book](#) | [YouTube](#) | [LinkedIn](#)

PGE 383 Exercise: Bootstrap for Subsurface Data Analytics in Python

Here's a simple workflow, demonstration of bootstrap for subsurface modeling workflows. This should help you get started with building subsurface models that integrate uncertainty in the sample statistics.

Bootstrap

Uncertainty in the sample statistics

- one source of uncertainty is the paucity of data.
- do 200 or even less wells provide a precise (and accurate estimate) of the mean? standard deviation? skew? P13?

Would it be useful to know the uncertainty in these statistics due to limited sampling?

- what is the impact of uncertainty in the mean porosity e.g. 20%+/-2%?

Bootstrap is a method to assess the uncertainty in a sample statistic by repeated random sampling with replacement.

Assumptions

- sufficient, representative sampling, identical, independent samples

Limitations

1. assumes the samples are representative
2. assumes stationarity
3. only accounts for uncertainty due to too few samples, e.g. no uncertainty due to changes away from data
4. does not account for boundary of area of interest
5. assumes the samples are independent
6. does not account for other local information sources

The Bootstrap Approach (Efron, 1982)

Statistical resampling procedure to calculate uncertainty in a calculated statistic from the data itself.

- Does this work? Prove it to yourself, for uncertainty in the mean solution is standard error:

$$\sigma_{\bar{x}}^2 = \frac{\sigma_x^2}{n}$$

Extremely powerful - could calculate uncertainty in any statistic! e.g. P13, skew etc.

- Would not be possible access general uncertainty in any statistic without bootstrap.
- Advanced forms account for spatial information and sampling strategy (game theory and Journel's spatial bootstrap (1993).

Steps:

1. assemble a sample set, must be representative, reasonable to assume independence between samples
2. optional: build a cumulative distribution function (CDF)

UNIVARIATE STATISTICS

New Tools

| Topic | Application to Subsurface Modeling |
|--|--|
| Awareness of Uncertainty Due to Sparse Sampling | <p>Sample statistics are uncertain due to limited sampling</p> <p><i>Quantify and apply this uncertainty model in subsurface modeling workflows.</i></p> |
| Bootstrap | <p>Resampling with replacement to calculate realizations of statistics</p> <p><i>While aware of the limitations, use them method to calculate uncertainty in e.g., mean porosity and carry through workflow as scenarios</i></p> |

DAYTUM – SPATIAL DATA ANALYTICS

Data Preparation

Lecture Recap

- ▶ Sampling Limitations
- ▶ Declustering
- ▶ Quantifying Uncertainty