# PGE 338 Data Analytics and Geostatistics

## Lecture 8: Bivariate Distributions

**Lecture outline . . .**

- **Bivariate Statistics**

- **Correlation**

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

Introduction

General Concepts

Univariate

**Bivariate**

Correlation

Regression

Model Checking

Time Series Analysis

Spatial Analysis

Machine Learning
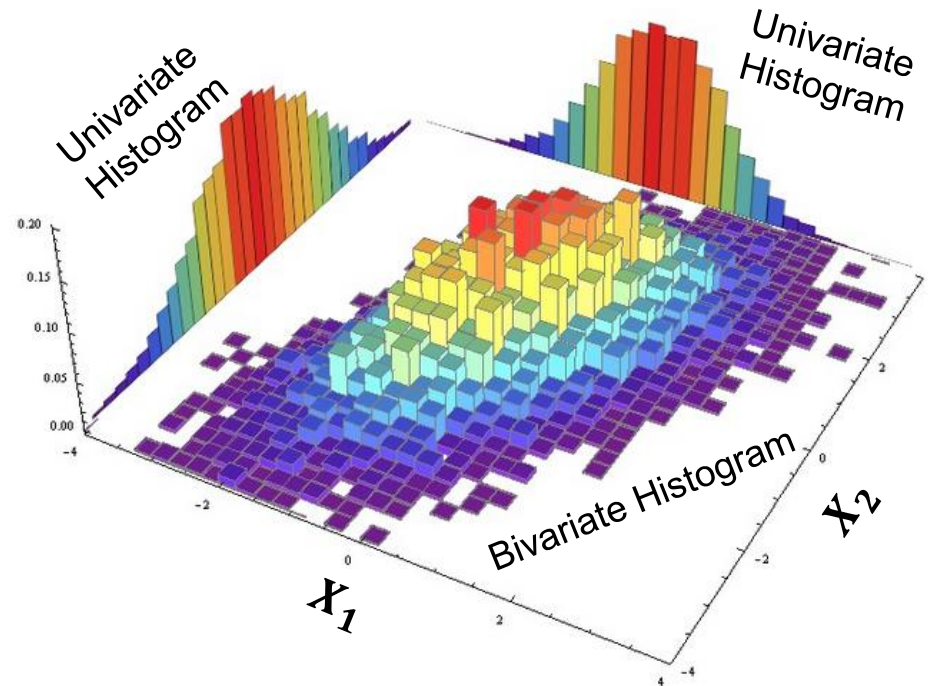
Uncertainty Analysis

# Motivation

We deal with more than one variable or feature?

We must use calculate bivariate statistics and use them in our models.

Note, 'bivariate' is pertaining to 2 variables or features at a time



Univariate and bivariate distributions (modified from ElenaPhys on StackOverflow).

# PGE 338 Data Analytics and Geostatistics
## Lecture 8: Bivariate Distributions

Lecture outline . . .

- **Bivariate Statistics**

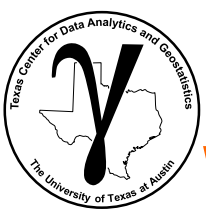| Introduction |
| General Concepts |
| Univariate |
| **Bivariate** |
| Correlation |
| Regression |
| Model Checking |
| Time Series Analysis |
| Spatial Analysis |
| Machine Learning |
| Uncertainty Analysis |

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

# Bivariate Statistics
## What is Bivariate Analysis?


Permeability vs. Porosity
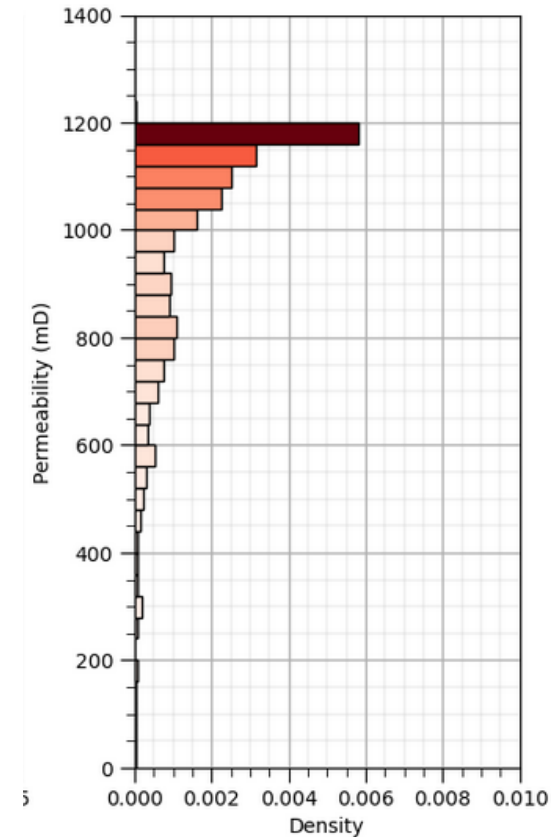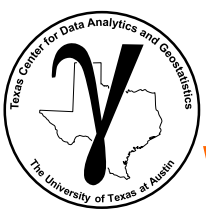
**Bivariate Analysis: Understand and quantify the relationship between two variables**

- Example: Porosity and permeability data

- How can we use this relationship? What would we miss if we only looked at the 2 histograms?

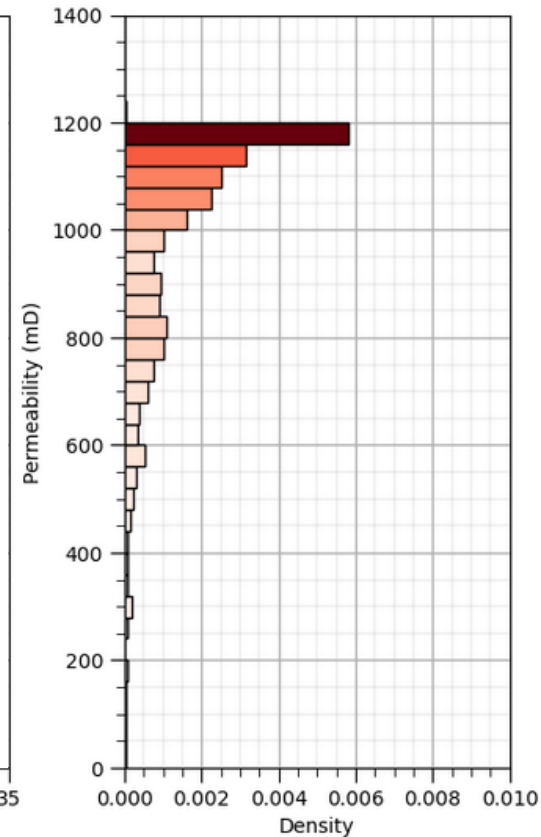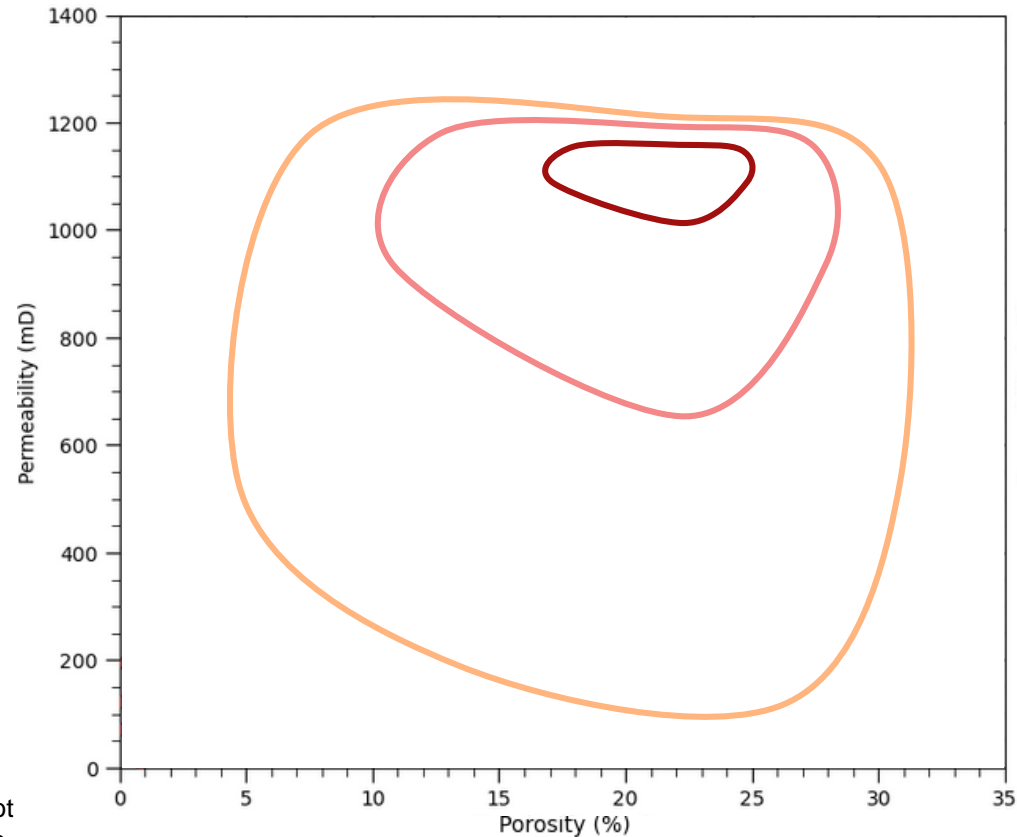- Is there a realizationship between permeability and porosity?

**?**

Porosity and permeability PDFs univariate plot modified from
PythonDataBasics_Bivariate _Visualization.ipynb.

# Bivariate Statistics
## What is Bivariate Analysis?

**Bivariate Analysis: Understand and quantify the relationship between two variables**

- Example: Porosity and permeability data

- How can we use this relationship? What would we miss if we only looked at the 2 histograms?

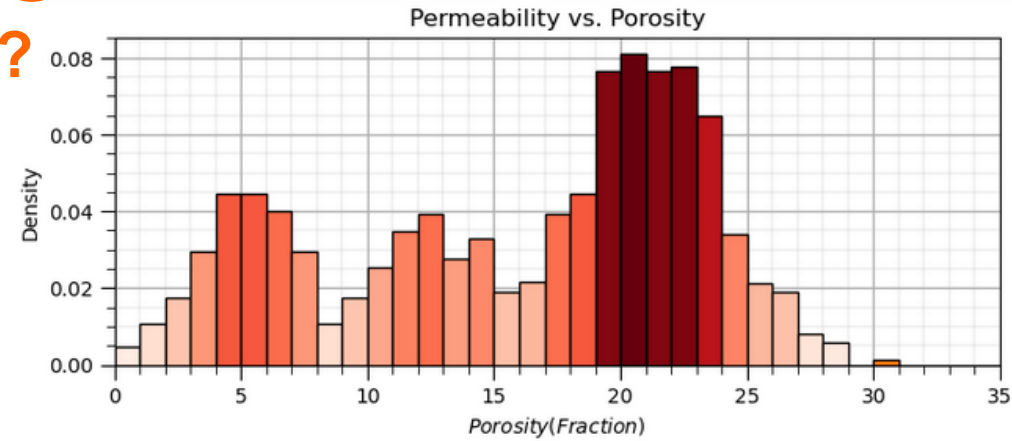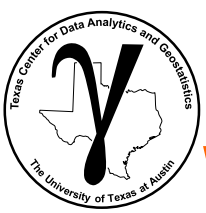- Is there a realizationship between permeability and porosity?
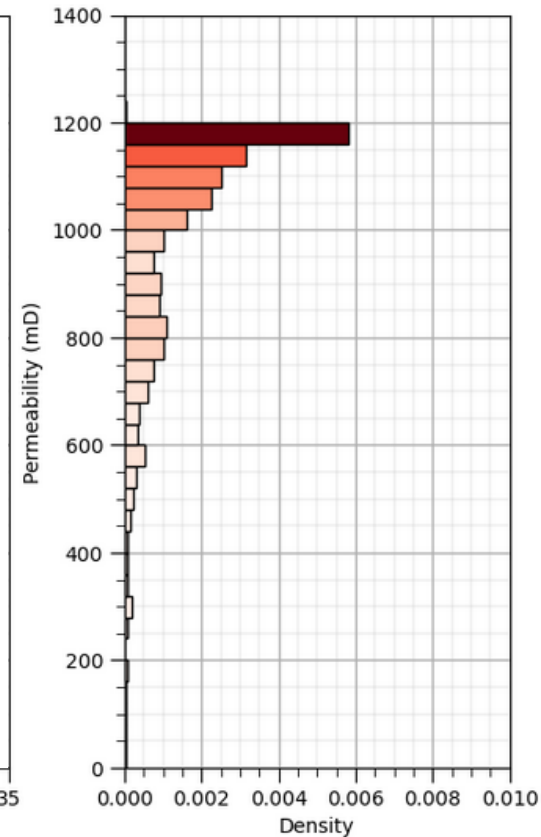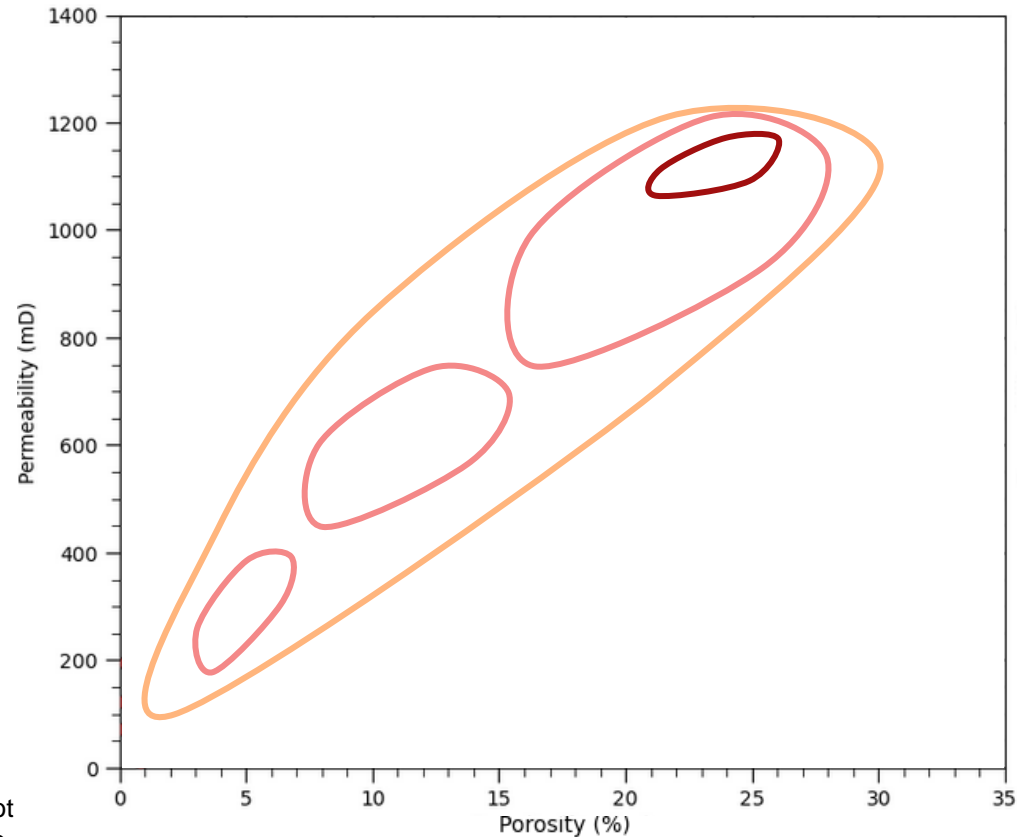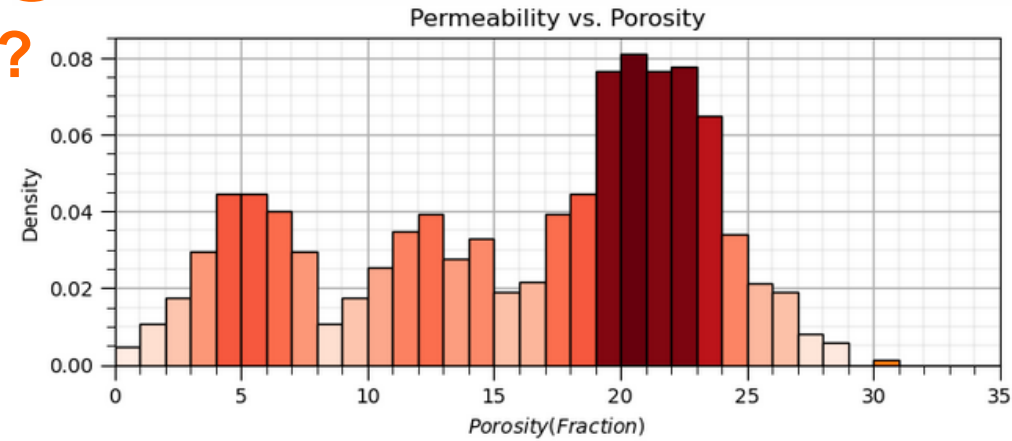
  No Relationship / Independence



Porosity and permeability PDFs univariate and bivariate (schematic) plot modified from PythonDataBasics_Bivariate _Visualization.ipynb.
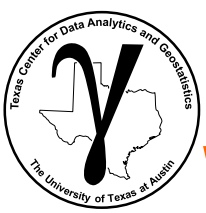
# Bivariate Statistics
## What is Bivariate Analysis?

**Bivariate Analysis: Understand and quantify the relationship between two variables**

- Example: Porosity and permeability data

- How can we use this relationship? What would we miss if we only looked at the 2 histograms?

- Is there a realizationship between permeability and porosity?

Positive Relationship

as Porosity Increases,
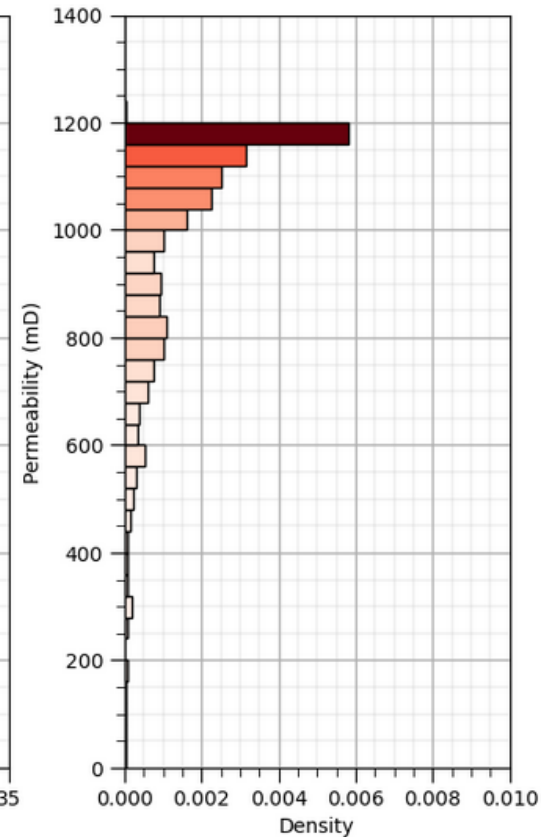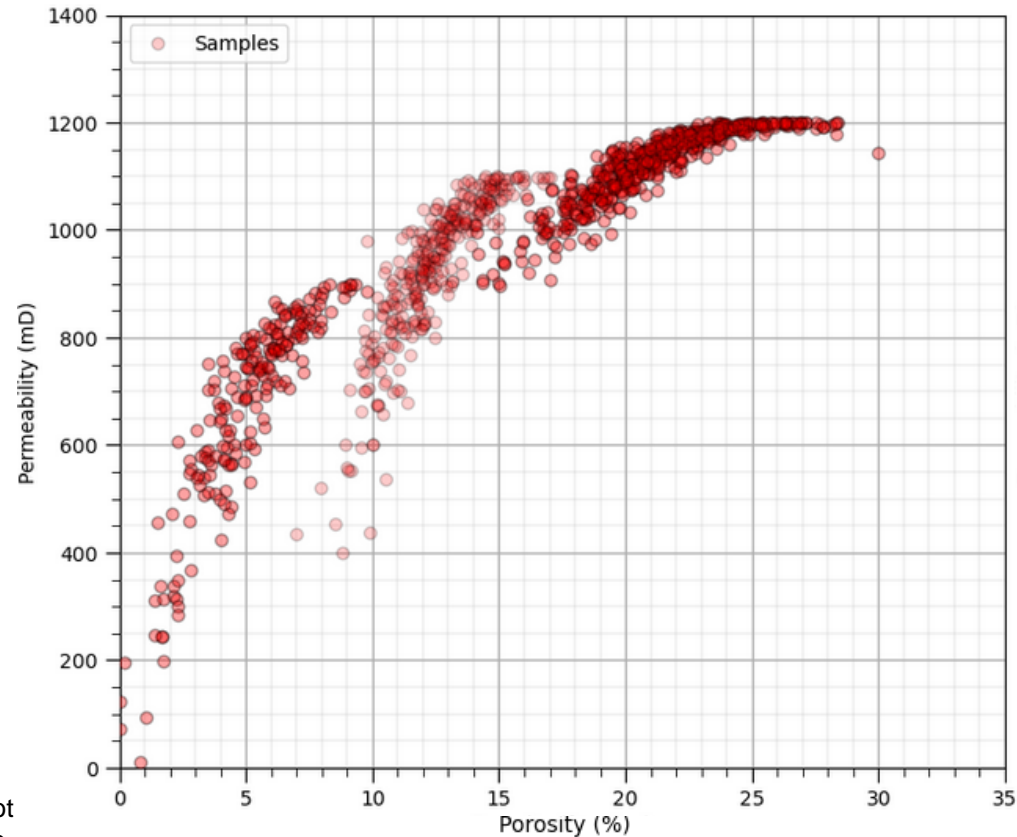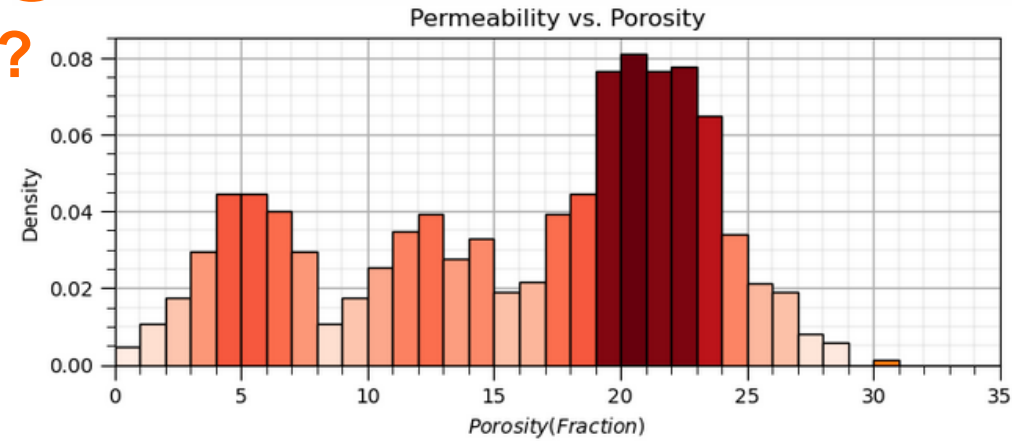
Permeability Increases



Porosity and permeability PDFs univariate and bivariate (schematic) plot modified from PythonDataBasics_Bivariate _Visualization.ipynb.

# Bivariate Statistics
## What is Bivariate Analysis?

**Bivariate Analysis: Understand and quantify the relationship between two variables**

- Example: Porosity and permeability data

- How can we use this relationship? What would we miss if we only looked at the 2 histograms?

- Is there a realizationship between permeability and porosity?

   Now Look at The Data!



Porosity and permeability PDFs univariate and bivariate (schematic) plot modified from PythonDataBasics_Bivariate _Visualization.ipynb.
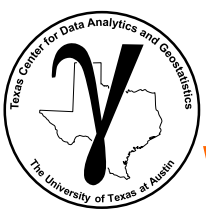
# Bivariate Statistics
## What is Bivariate Analysis?

**Bivariate Analysis: Understand and quantify the relationship between two variables**

- Example: Porosity and permeability data

- How can we use this relationship? What would we miss if we only looked at the 2 histograms?

- Is there a realizationship between permeability and porosity?
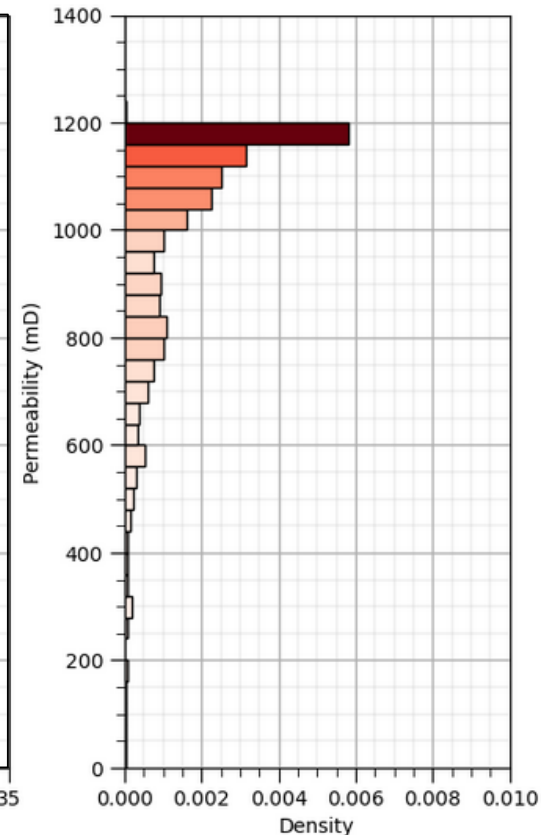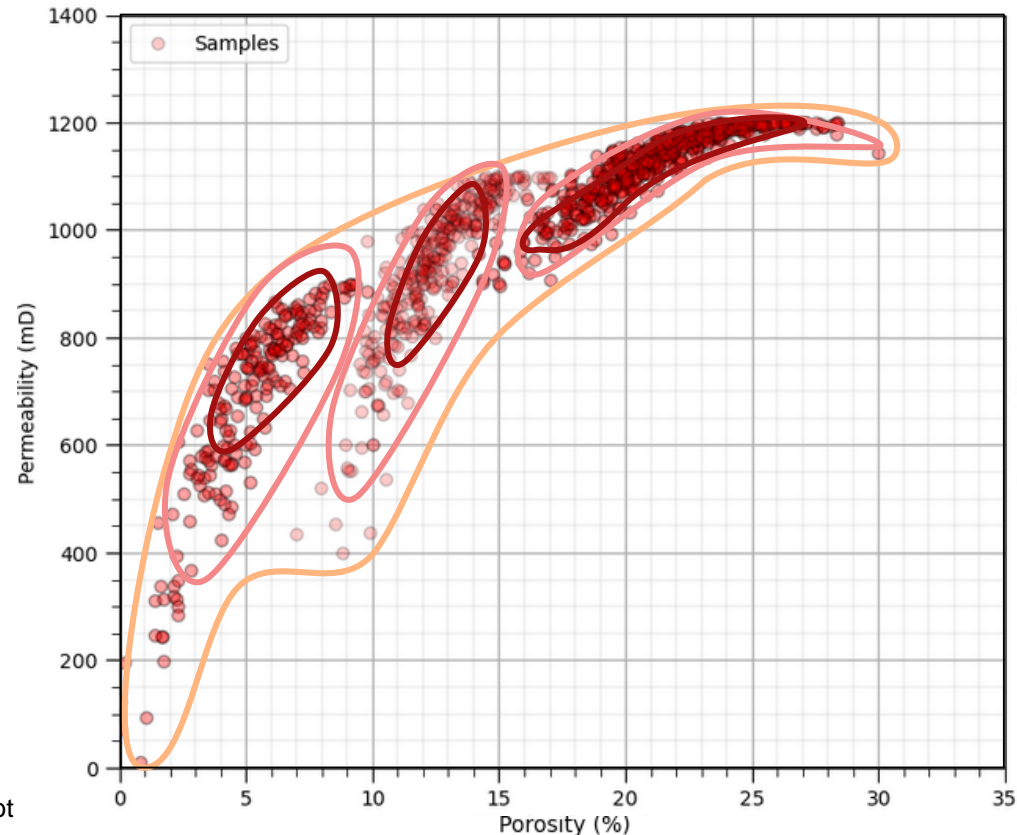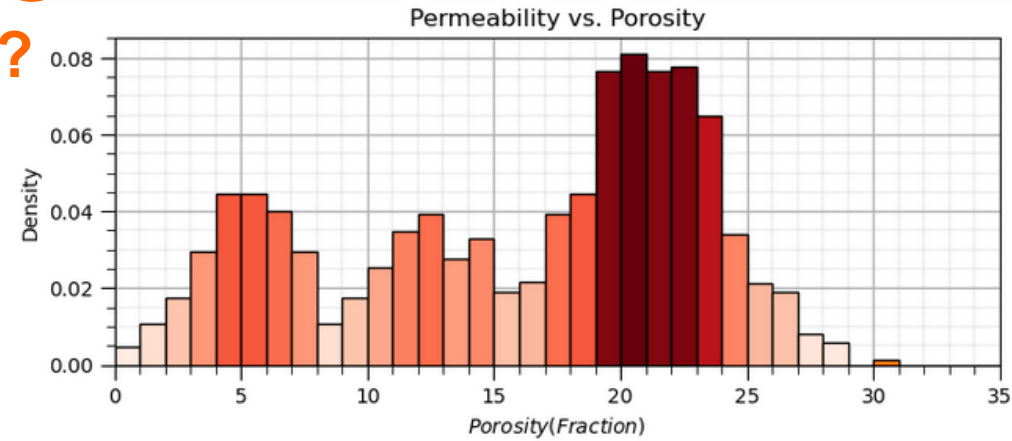
  Now Look at The Data!

  Positive Relationship and…



Porosity and permeability PDFs univariate and bivariate (schematic) plot modified from PythonDataBasics_Bivariate _Visualization.ipynb.
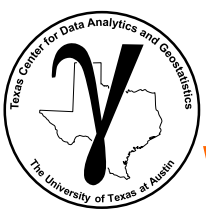
# Bivariate Statistics
## What is Bivariate Analysis?

**Bivariate Analysis: Understand and quantify the relationship between two variables**

- Example: Porosity and permeability data

- How can we use this relationship? What would we miss if we only looked at the 2 histograms?

- Is there a realizationship between permeability and porosity?
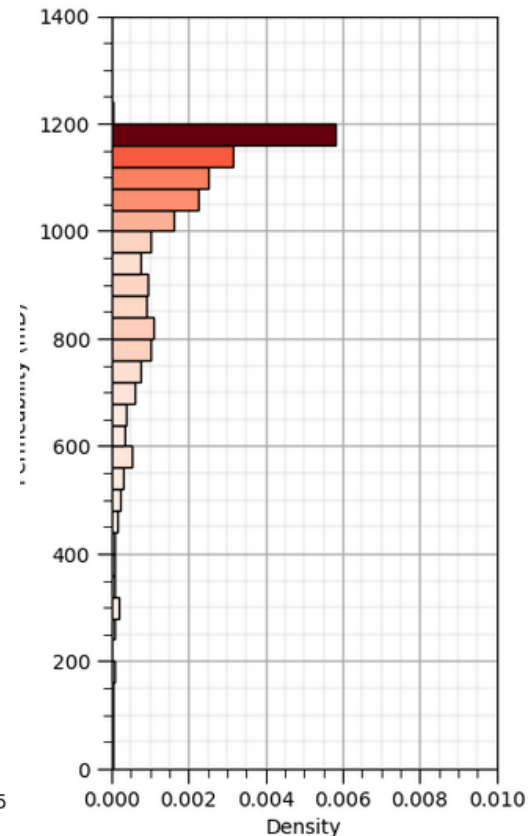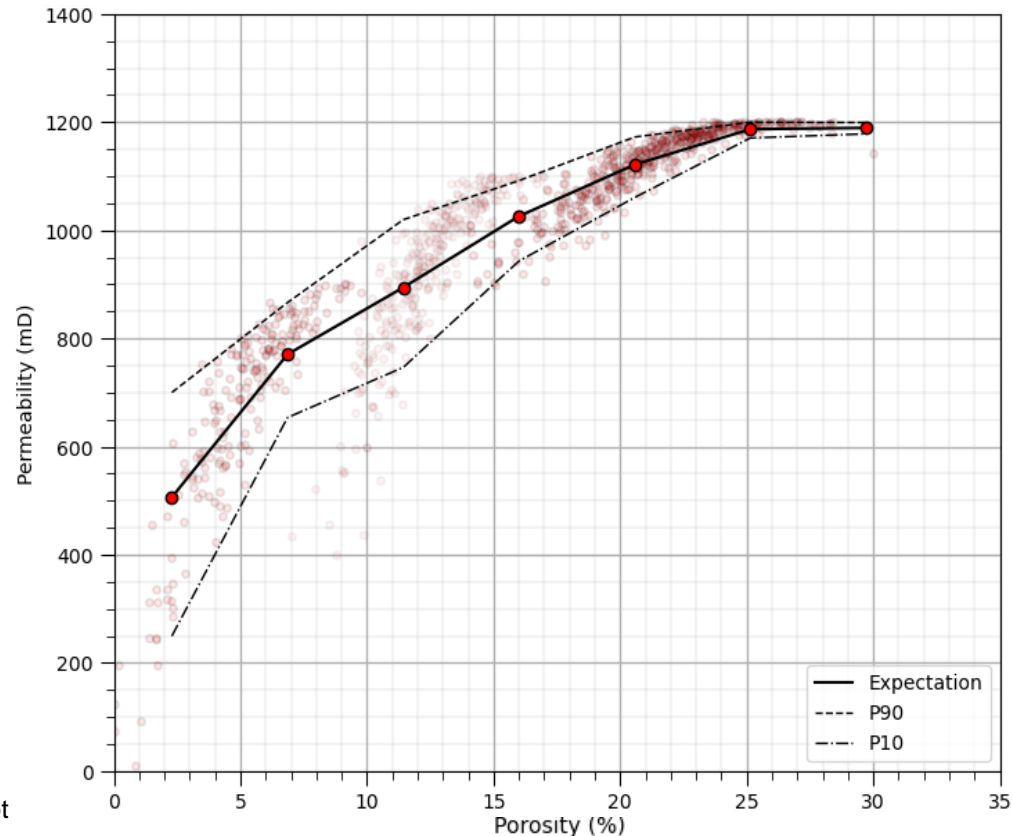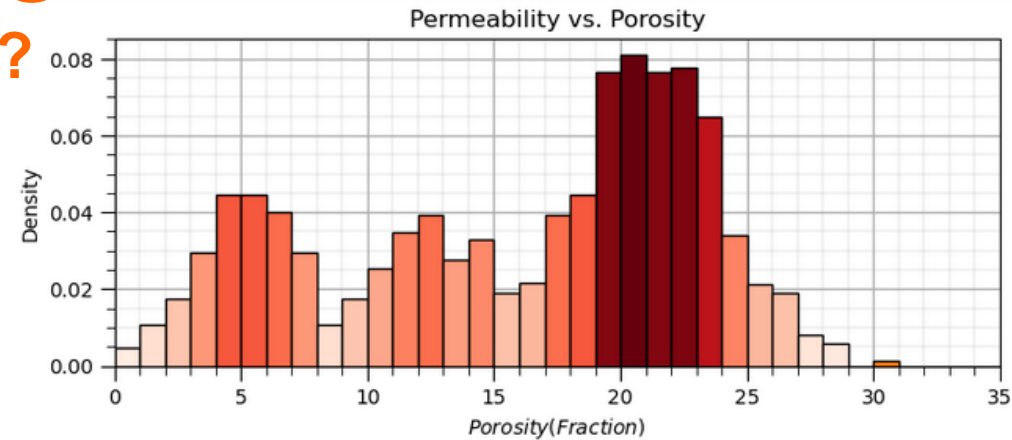
Now Look at The Data!

Positive Relationship and

We Could Model It with Conditional Distributions!



Porosity and permeability PDFs univariate and bivariate (schematic) plot modified from PythonDataBasics_Bivariate _Visualization.ipynb.
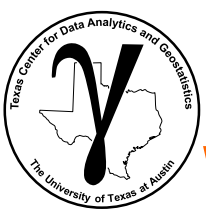
# Bivariate Statistics
## What is Bivariate Analysis?

**Bivariate Analysis: Understand and quantify the relationship between two variables**

- Example: Porosity and permeability data

- How can we use this relationship? What would we miss if we only looked at the 2 histograms?

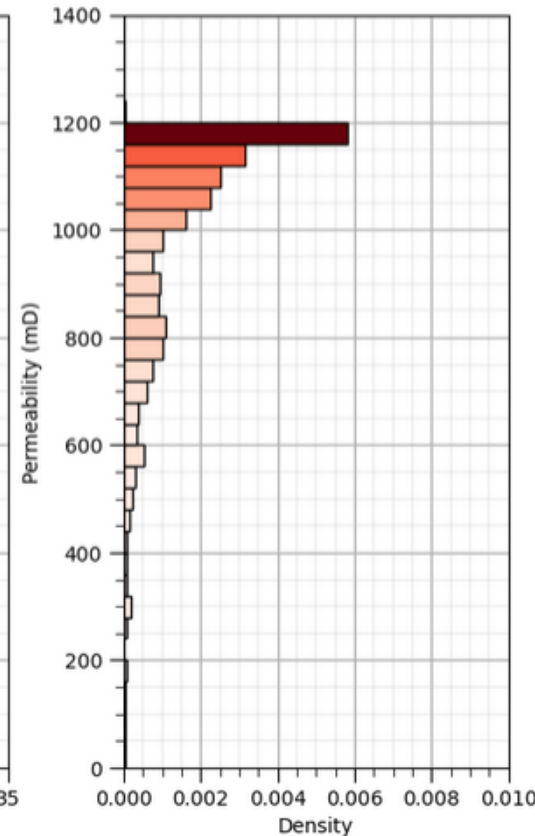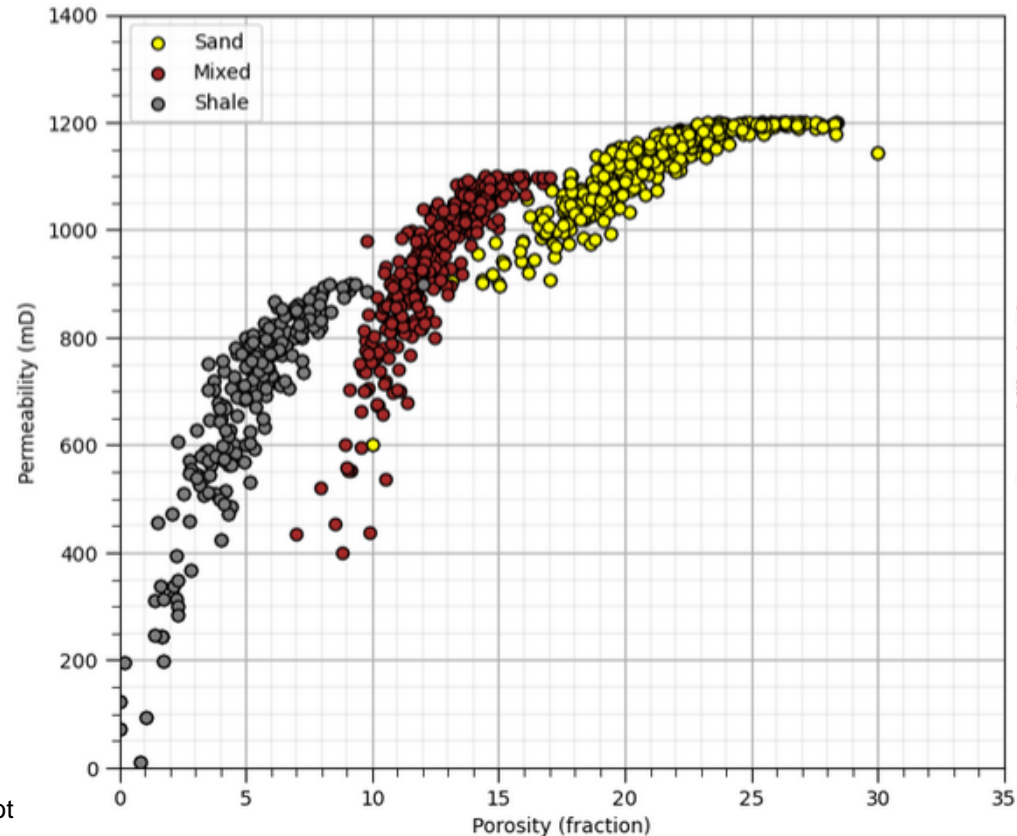- Is there a realizationship between permeability and porosity?
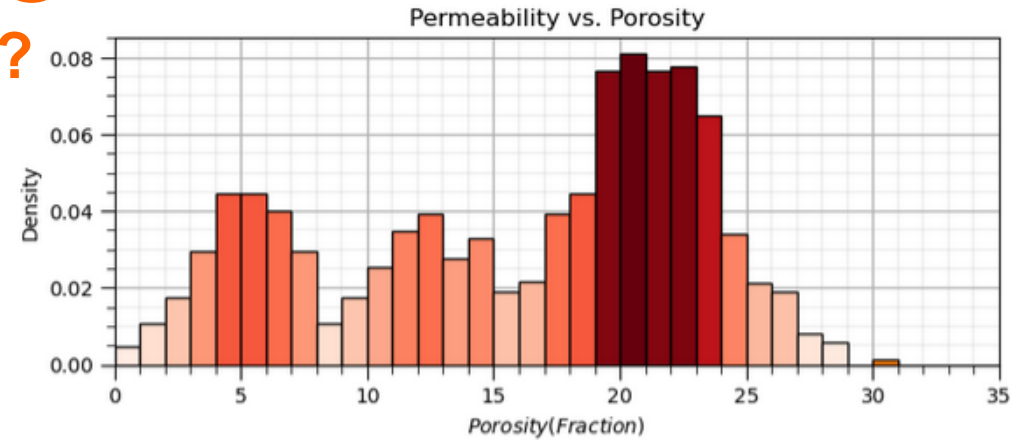
Now Look at The Data!

Positive Relationship and

Multiple Facies / Populations

**We use all of this to build better subsurface models.**



Porosity and permeability PDFs univariate and bivariate (schematic) plot modified from PythonDataBasics_Bivariate _Visualization.ipynb.
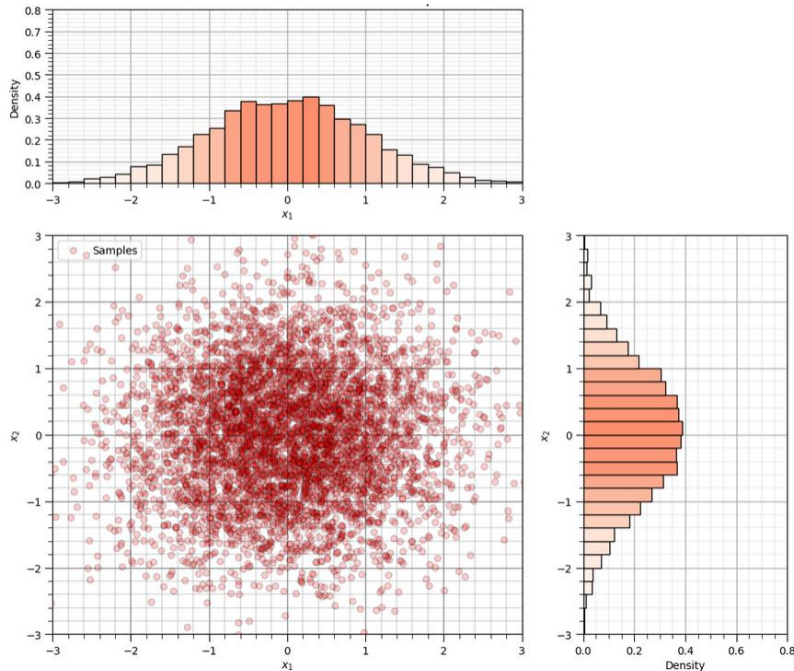
# Bivariate Statistics
## How to Quantify Relationship between two Variables?
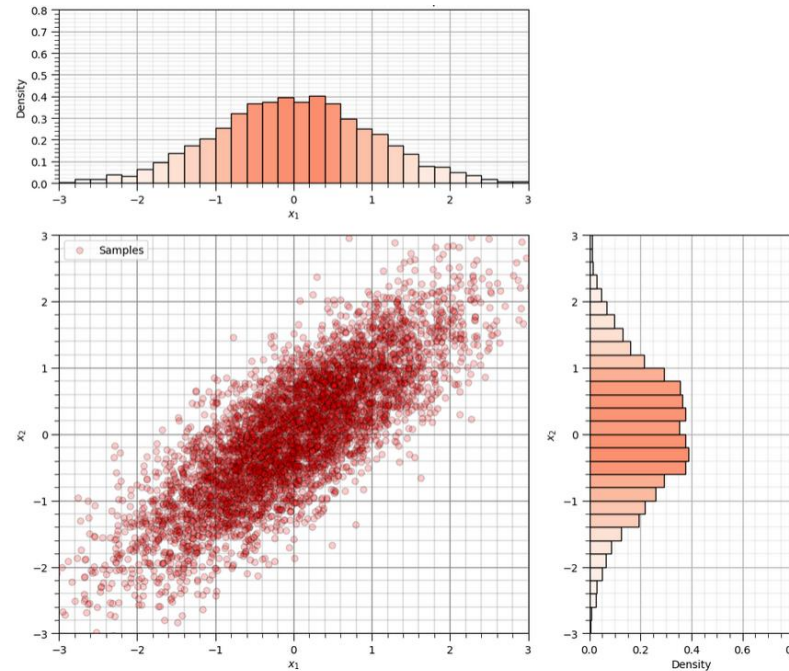
**Bivariate Analysis: Quantifying the strength of the relationship between 2 features**

• Example: porosity and permeability data from a carbonate and a sandstone formation.

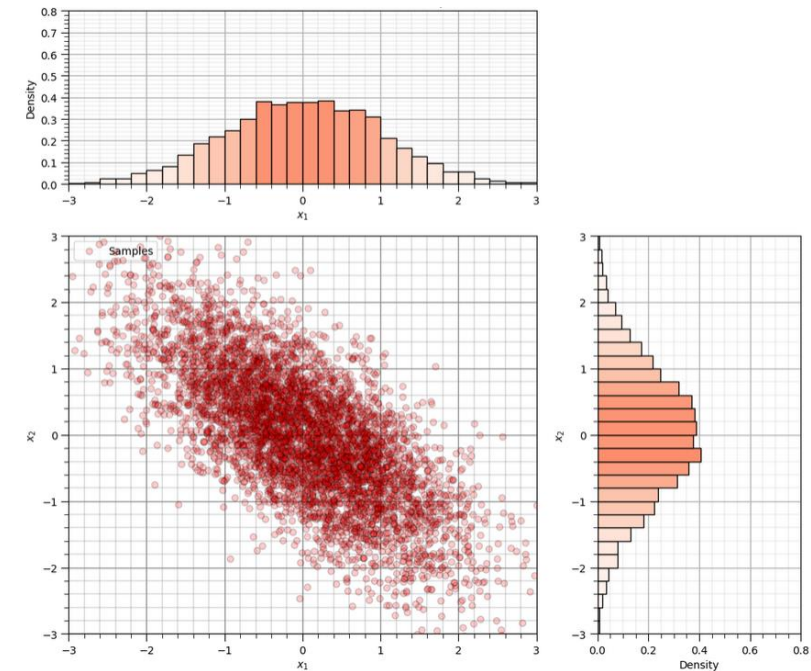• In which case are the features X1 and X2 best correlated? I.e., a stronger relationship?

### Data Set 1        ### Data Set 2        ### Data Set 3



Univariate and bivariate distributions for 3 data sets, generated with Interative_Correlation_Coefficient.ipynb.

# PGE 338 Data Analytics and Geostatistics
## Lecture 8: Bivariate Distributions

Lecture outline . . .

- Correlation

| Introduction |
|---|

| General Concepts |
|---|

| Univariate |
|---|

| **Bivariate** |
|---|

| Correlation |
|---|
| Regression |
| Model Checking |

| Time Series Analysis |
|---|

| Spatial Analysis |
|---|

| Machine Learning |
|---|

| Uncertainty Analysis |
|---|

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

# Bivariate Statistics
## Pearson's Correlation Coefficient

**How Do We Quantify the Relationships Between 2 Features, Bivariate Relationships?**

- We need to go beyond qualitative descriptions, good, bad, strong, weak, none…

# Bivariate Statistics
## Pearson's Correlation Coefficient

**Definition: Pearson's Product-moment Correlation Coefficient ($\rho_{X,Y}$)**

- Provides a measure of the degree of linear relationship.

means of variables
x and y

$$\rho_{X,Y} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{(n-1)\sigma_x\sigma_y}, -1.0 \leq \rho_{xy} \leq 1.0$$

Correlation coefficient of
variables x and y

standard deviation of
variables x and y

number of
data pairs

- Correlation coefficient is a "standardized" covariance. The covariance ($C_{X,Y}$):

Covariance: $\qquad C_{X,Y} = \dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{(n-1)}$ $\qquad$ Correlation Coefficient: $\qquad \rho_{X,Y} = \dfrac{C_{X,Y}}{\sigma_X\sigma_Y}$

Correlation coefficient is covariance standardized by dividing by $\sigma_X\sigma_Y$

**We can see that covariance and variance are related.**

**1. Sample Variance:**

$$s^2 = \frac{1}{n-1} \sum_{i=1}^{N} (x_i - \bar{x})(x_i - \bar{x})$$

A measure of how 1 variable varies with itself.

**2. Sample Covariance:**
- Replace the second term in the square with another feature.

$$C_{xy} = \frac{1}{n-1} \sum_{i=1}^{N} (x_i - \bar{x})(y_i - \bar{y})$$

A measure of how 2 features vary together.

# Bivariate Statistics
## Spearman's Rank Correlation Coefficient

## Definition: Spearman's Rank Correlation Coefficient

- Provides a measure of the degree of monotonic relationship.

means of rank transform of variables x and y

$$\rho_{R_x, R_y} = \frac{\sum_{i=1}^{n}(R_{x_i} - \overline{R_x})(R_{y_i} - \overline{R_y})}{(n-1)\sigma_{R_x}\sigma_{R_y}}, -1.0 \leq \rho_{xy} \leq 1.0$$

Rank correlation coefficient of variables x and y

number of data pairs

standard deviation of Rank transform of variables x and y

- Rank transform, e.g. $R_{x_i}$, sort the data in ascending order and replace the data with the index, $i = 1, \dots, n$.

- Spearman's rank correlation coefficient is more robust in the presence of outliers and some nonlinear features than the Pearson's correlation coefficient

# Bivariate Statistics
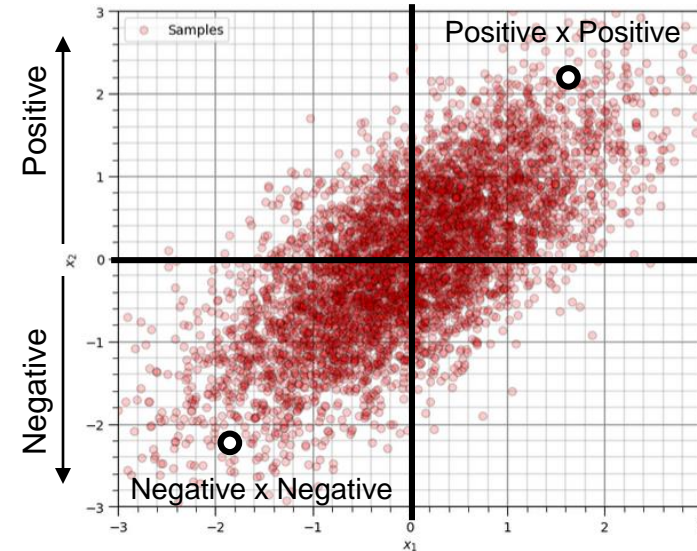## Covariance

**Let's think about covariance.**

- For a thought experiment, consider 2 standard normal variables, N[0,1].

$$C_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{(n-1)}$$
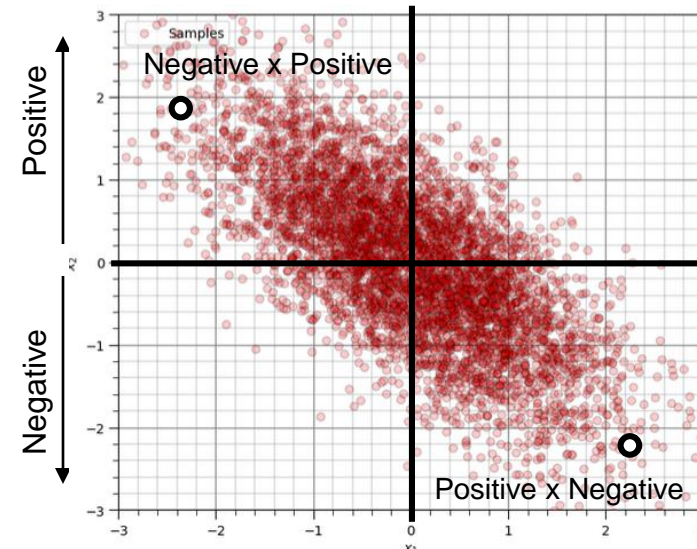
the means, $\bar{x} = \bar{y} = 0$ ∴

$$C_{xy} \sim E\{X\,Y\} \text{ and } \rho_{xy} \sim \frac{E\{X\,Y\}}{\sigma_x \sigma_y}$$

- Positive covariance / correlation if we pair high-high and low-low feature values over the samples.

- Negative covariance / correlation if we pair high-low and low-high feature values over the samples.



$$\rho_{xy} = \frac{C_{xy}}{\sigma_x \sigma_y} > 0.0$$

$$if\ \rho > 0,\ \uparrow C_{xy},\ \sum_{i=1}^{n/2}[\,x^- \times\, y^-] + \sum_{i=n/2}^{n}[\,x^+ \times\, y^+]$$



$$\rho_{xy} = \frac{C_{xy}}{\sigma_x \sigma_y} < 0.0$$

$$if\ \rho < 0,\ \uparrow C_{xy},\ \sum_{i=1}^{n/2}[(\,x^- \times\, y^+] + \sum_{i=n/2}^{n}[(\,x^+ \times\, y^-]$$
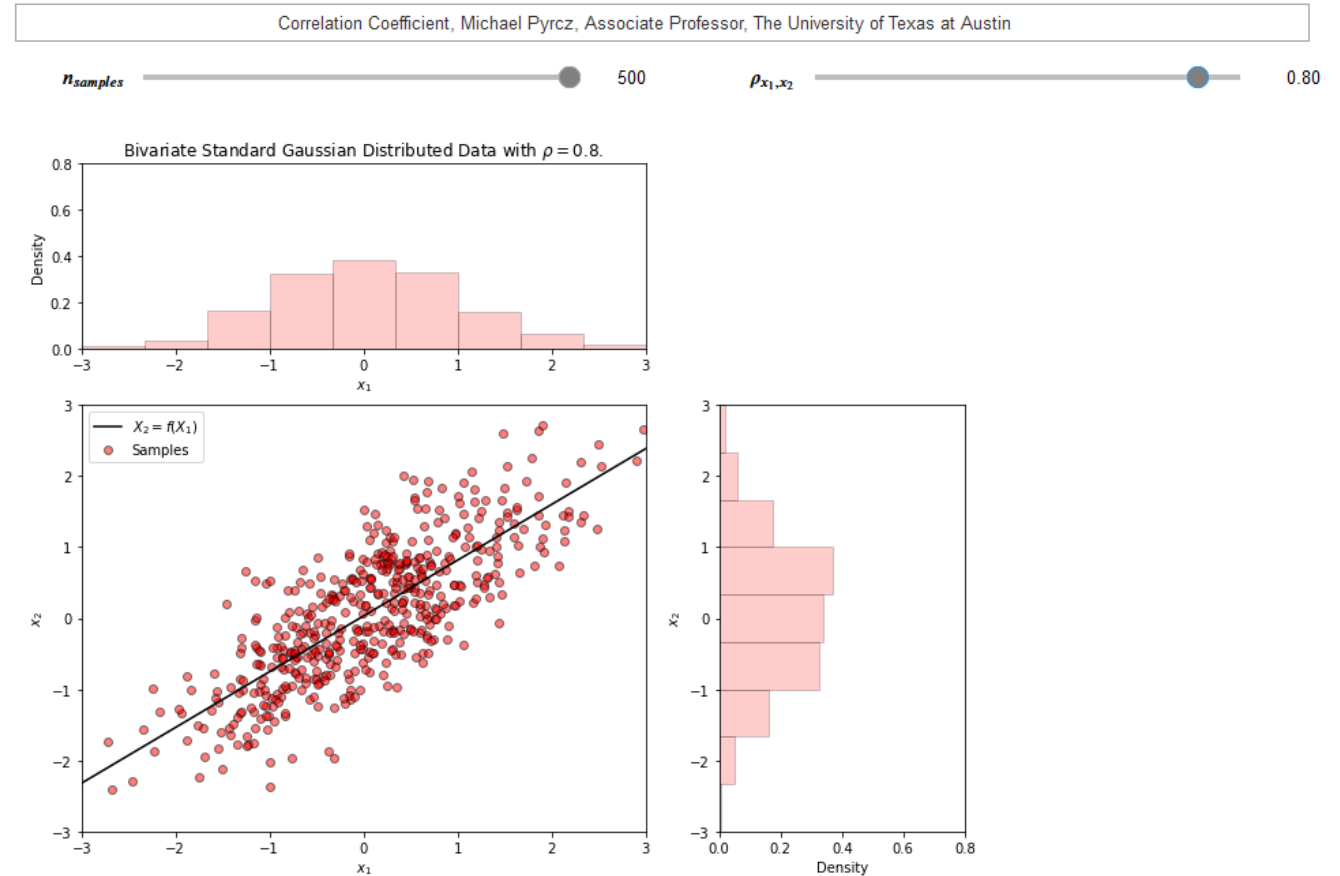
# Interactive Correlation Demonstration

## What does correlation coefficient actually see?

- The correlation coefficient measure linear, homoscedastic relationships between 2 features.

- For the bivariate Gaussian distribution the correlation coefficient completely describes the relationship between features.

$$f_{X,Y}(x,y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho_{x,y}^2}}\exp\left[-\frac{z}{2(1-\rho_{x,y}^2)}\right]$$

where:

$$z = \frac{(x-\mu_x)^2}{\sigma_x^2} - \frac{2\rho_{x,y}(x-\mu_x)(y-\mu_y)}{\sigma_x\sigma_y} + \frac{(y-\mu_y)^2}{\sigma_y^2}$$



Interactive Correlation demonstration, the file is Interactive_Correlation_Coefficient.ipynb.

# Bivariate Statistics
## Correlation and Causation

**Does correlation imply causation?**

- NO!

- We require a "true experiment" where one variable is manipulated, and others are rigorously controlled!

Here's an example to demonstrate a potential issue of assigning causation based on the observation of correlation.



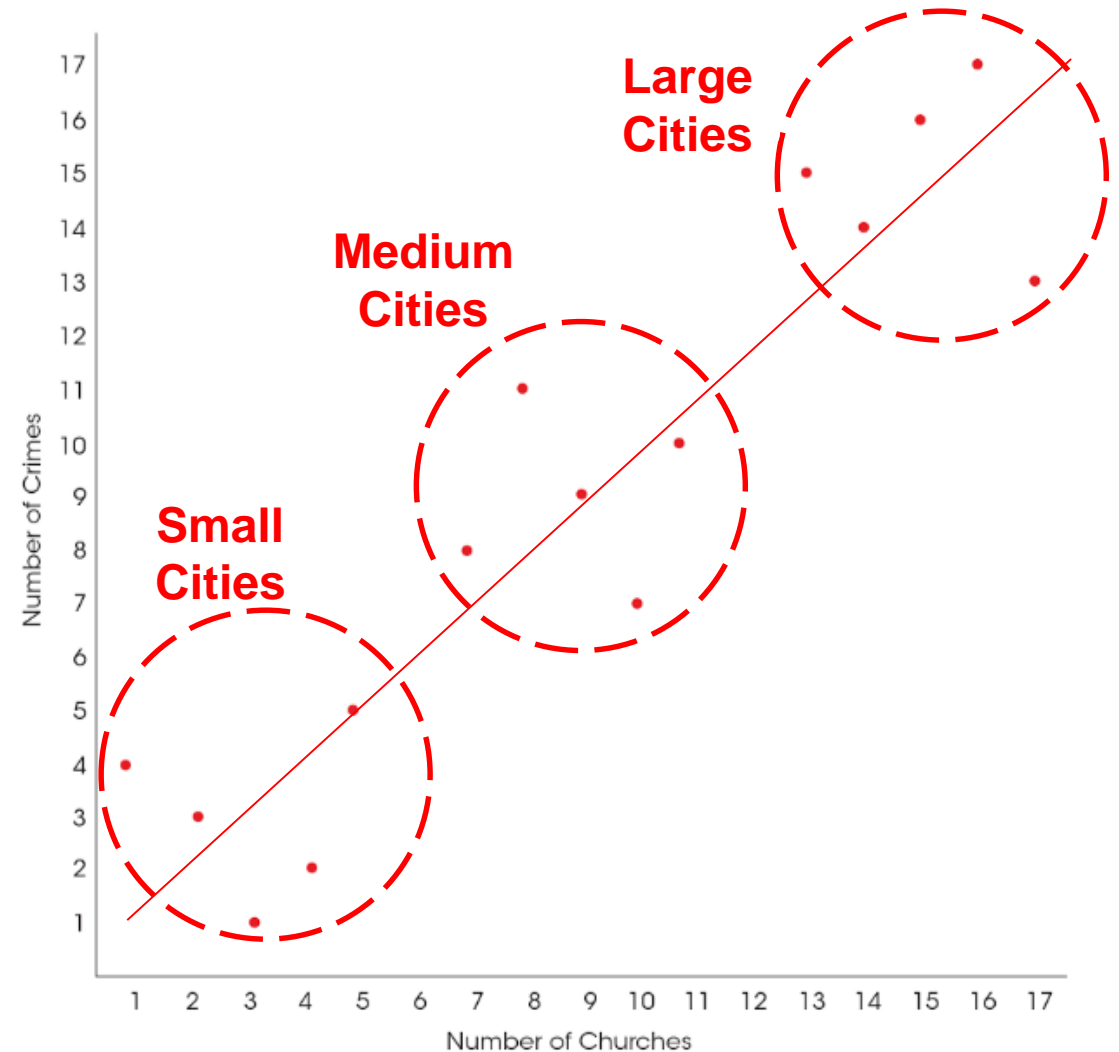Plot of frequency of crimes vs. frequency of churches.

# Bivariate Statistics
## Correlation and Causation

**Does correlation imply causation?**

- NO!

- We require a "true experiment" where one variable is manipulated, and others are rigorously controlled!

Here's an example to demonstrate a potential issue of assigning causation based on the observation of correlation.



Plot of frequency of crimes vs. frequency of churches.

# Bivariate Statistics
## Spurious Correlations

**Can correlation be misleading and misused?**

Spurious Correlation:

A correlation that seems to be causal but is not.

Due to:

- random chance, too few data, data sampling bias
- outliers or artificial truncation of the data
- confounding features

Confounding Features:

feature not part of the study that is related to both the predictor and response feature

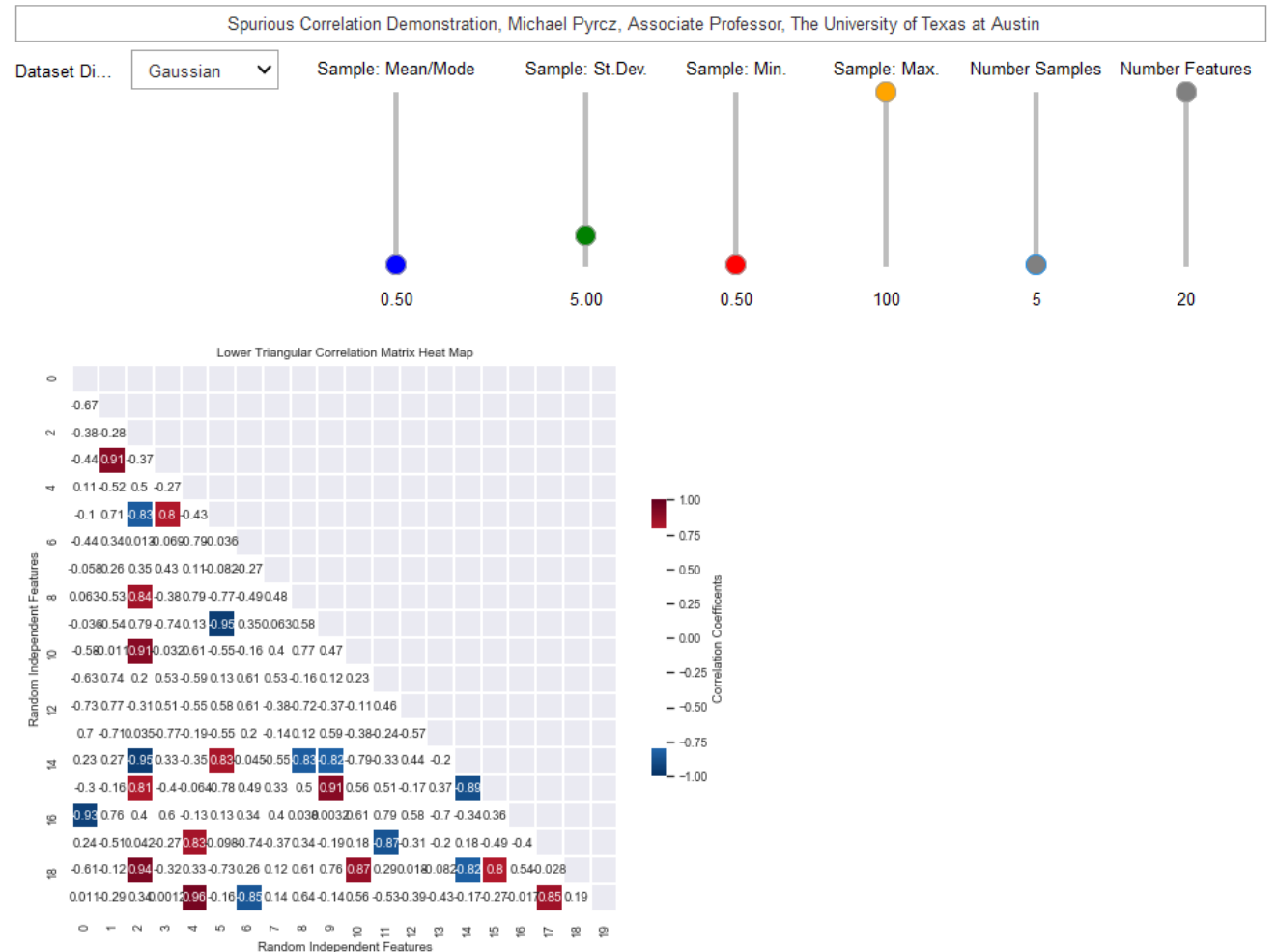- causes spurious correlations



Examples of spurious correlations.

Images from https://www.tylervigen.com/spurious-correlations

# Interactive Spurious Correlations

Let's make uncorrelated random features and check their correlations.

- We can plot all feature pairwise correlations in a matrix color coded by very high or low correlations (>0.8 or <-0.8 respectively)

- What happens has the number of data is small and the number of features is large?



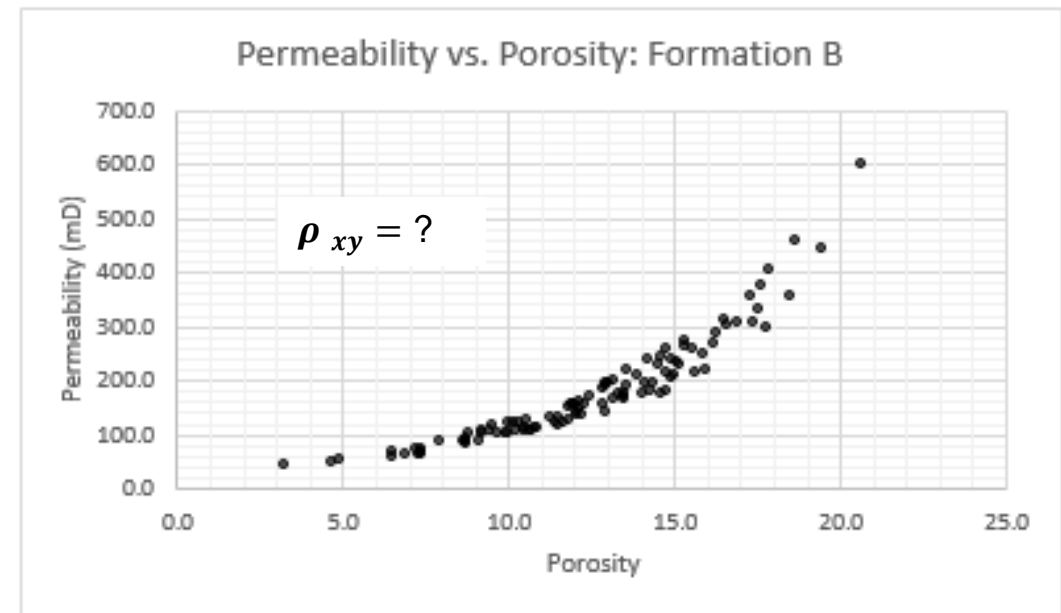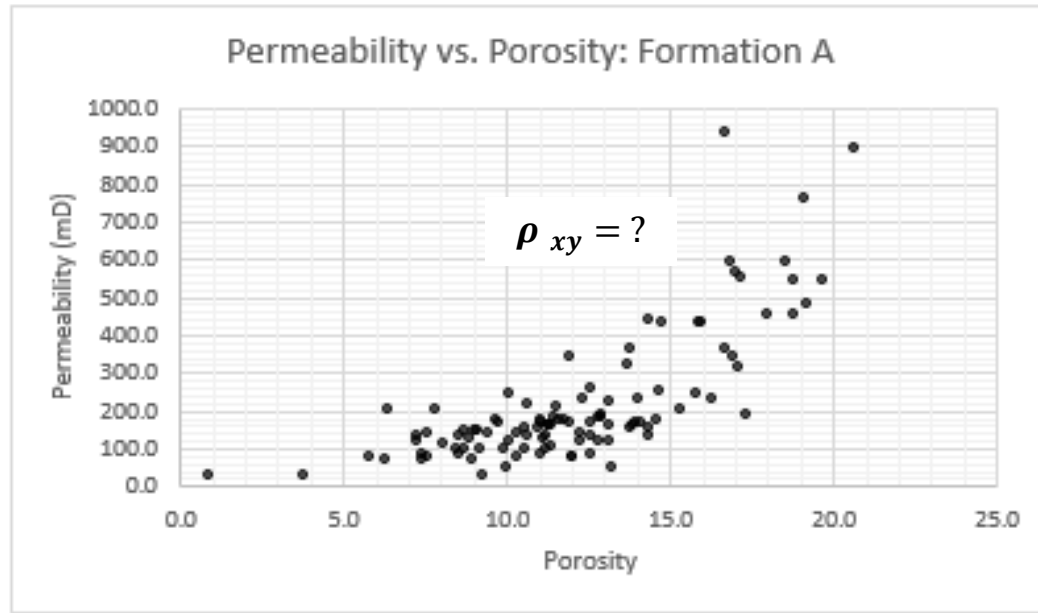Interactive spurious correlations demonstration, the file is Interactive_Spurious_Correlations.ipynb.
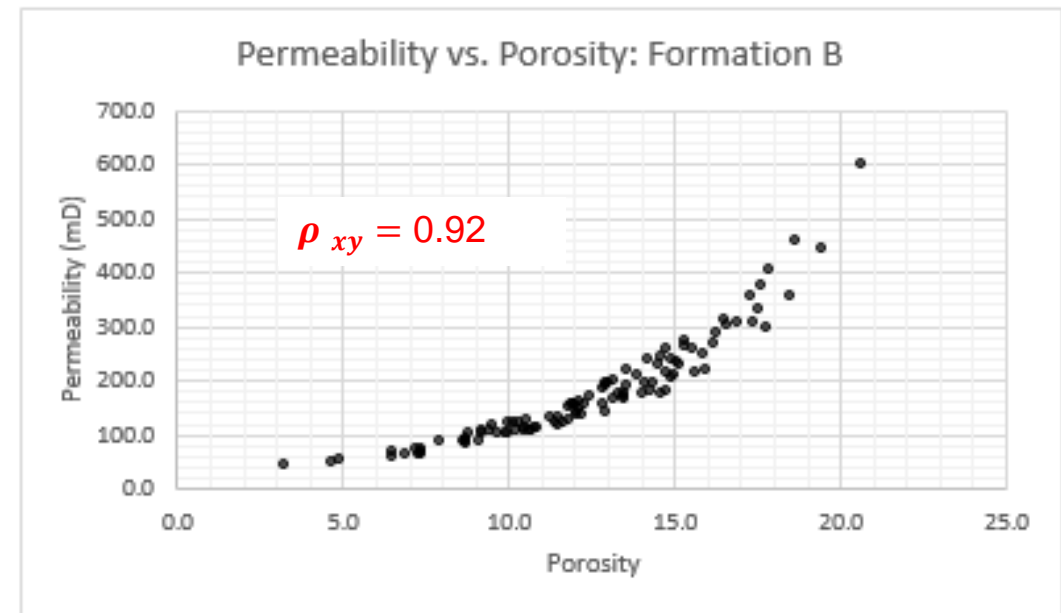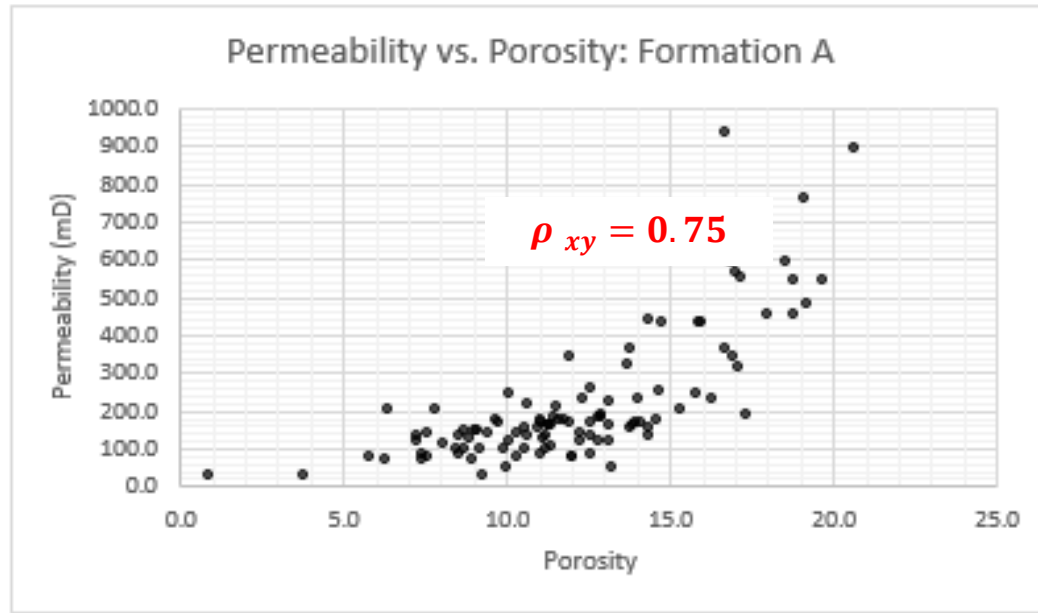
# Bivariate Statistics
## Exercise with Pearson's Correlation Coefficient

Prepare the scatterplot of the core porosity and permeability data for formations A and B provided to you. Estimate correlation coefficient in both cases. How well is porosity and permeability are correlated in these two formations?

**Excel Function Correl(Array1,Array2)**

Can I derive a correlation between porosity and permeability?
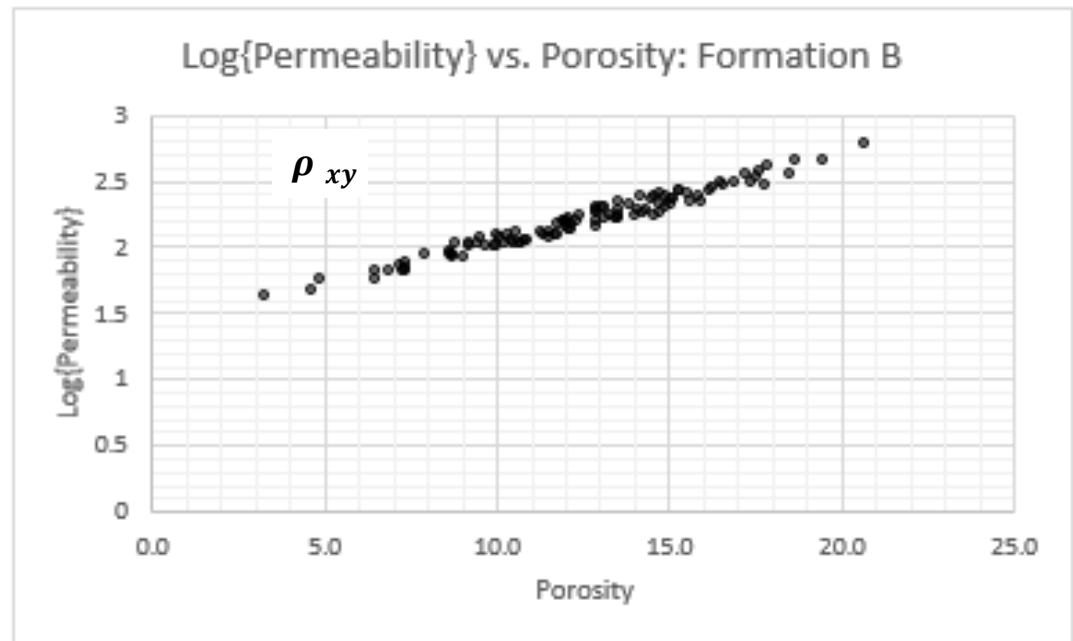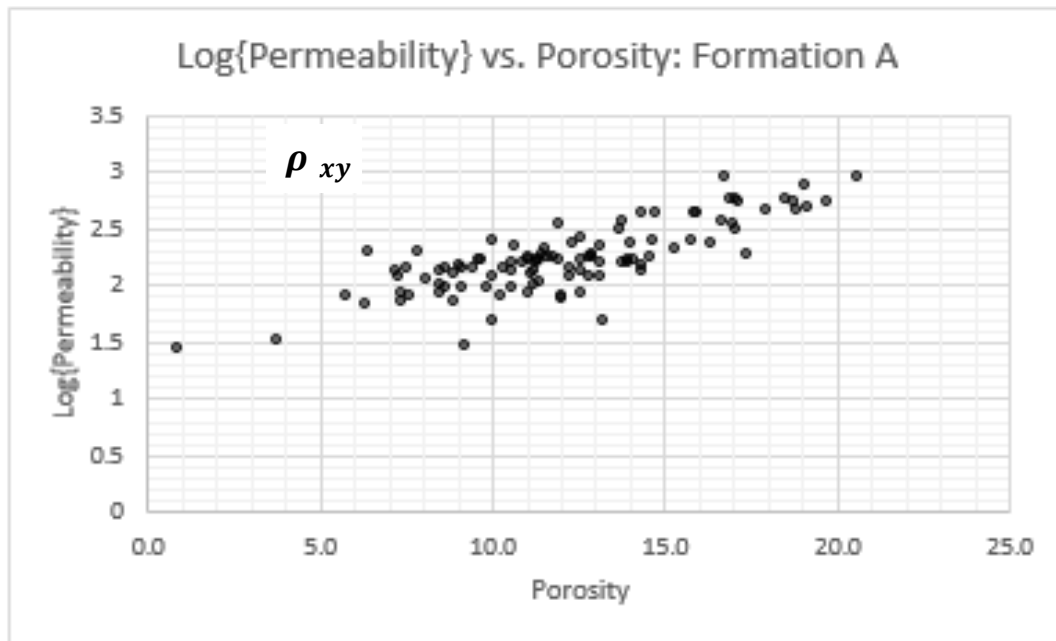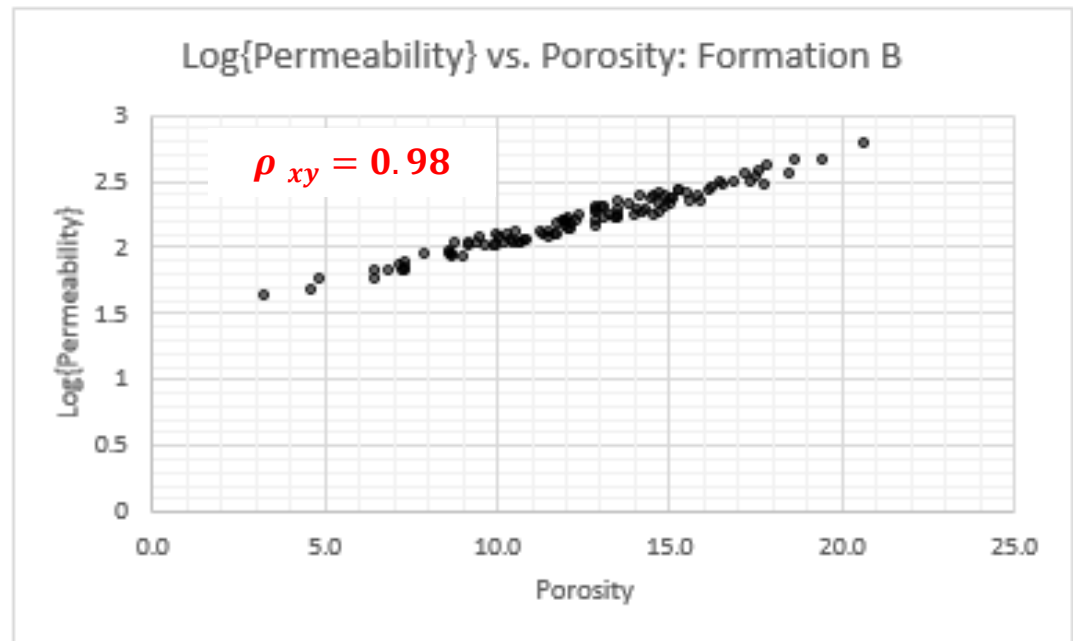How reliable that correlation would be?



Correlation for 2 datasets, file is Univariate_NonLinearPor_Perm.xlsx.

# Bivariate Statistics
## Exercise with Pearson's Correlation Coefficient

Prepare the scatterplot of the core porosity and permeability data for formations A and B provided to you. Estimate correlation coefficient in both cases. How well is porosity and permeability are correlated in these two formations?

**Excel Function Correl(Array1,Array2)**

Can I derive a correlation between porosity and permeability?
How reliable that correlation would be?



Correlation for 2 datasets, file is Univariate_NonLinearPor_Perm.xlsx.

# Bivariate Statistics
## Exercise with Pearson's Correlation Coefficient

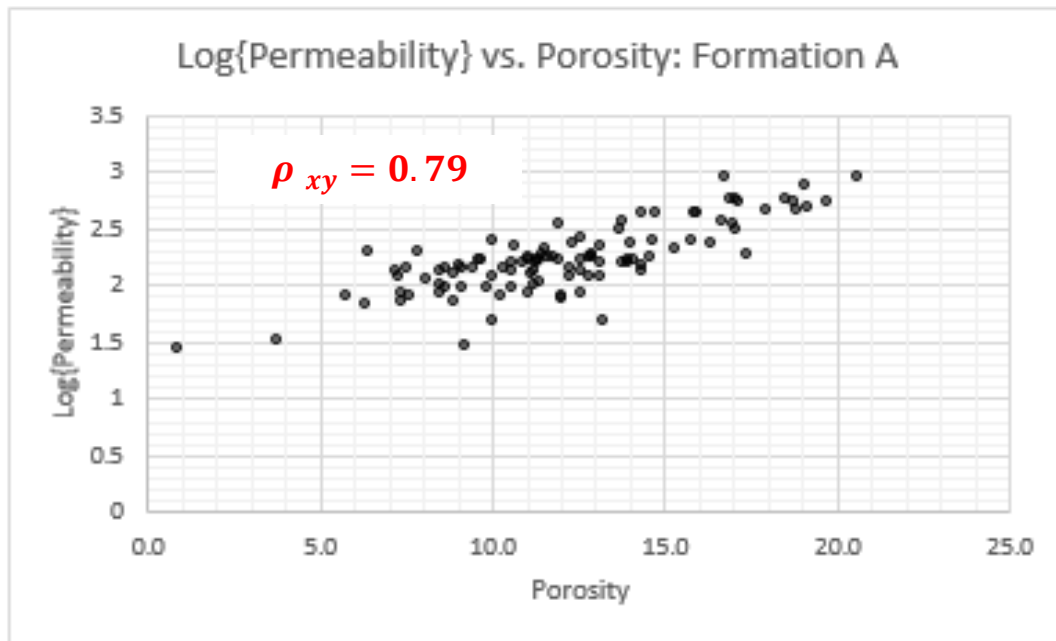Does this nonlinear log transform improve our use / characterization with the correlation coefficient?



Correlation for 2 datasets, file is Univariate_NonLinearPor_Perm.xlsx.
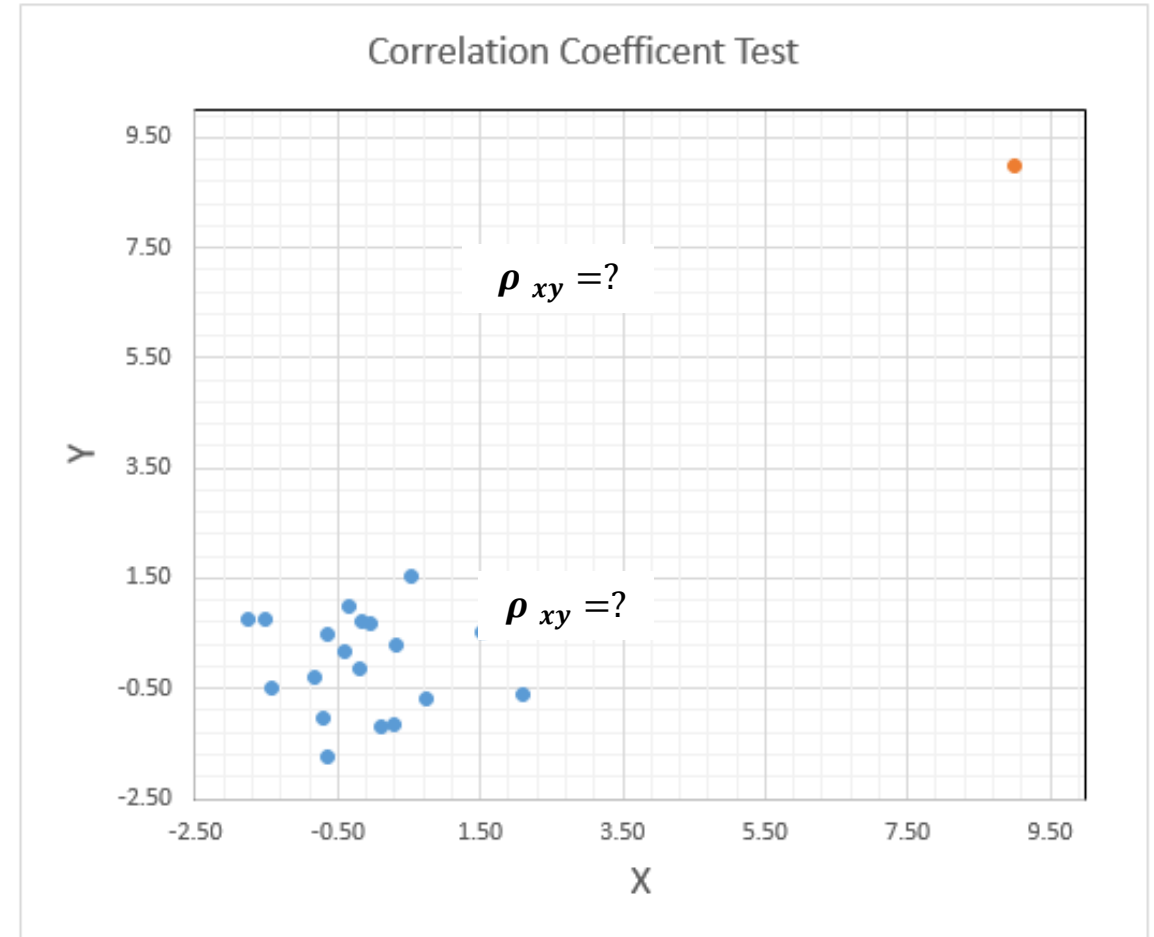
# Bivariate Statistics
## Exercise with Pearson's Correlation Coefficient

Does this nonlinear log transform improve our use / characterization with the correlation coefficient?



Correlation for 2 datasets, file is Univariate_NonLinearPor_Perm.xlsx.
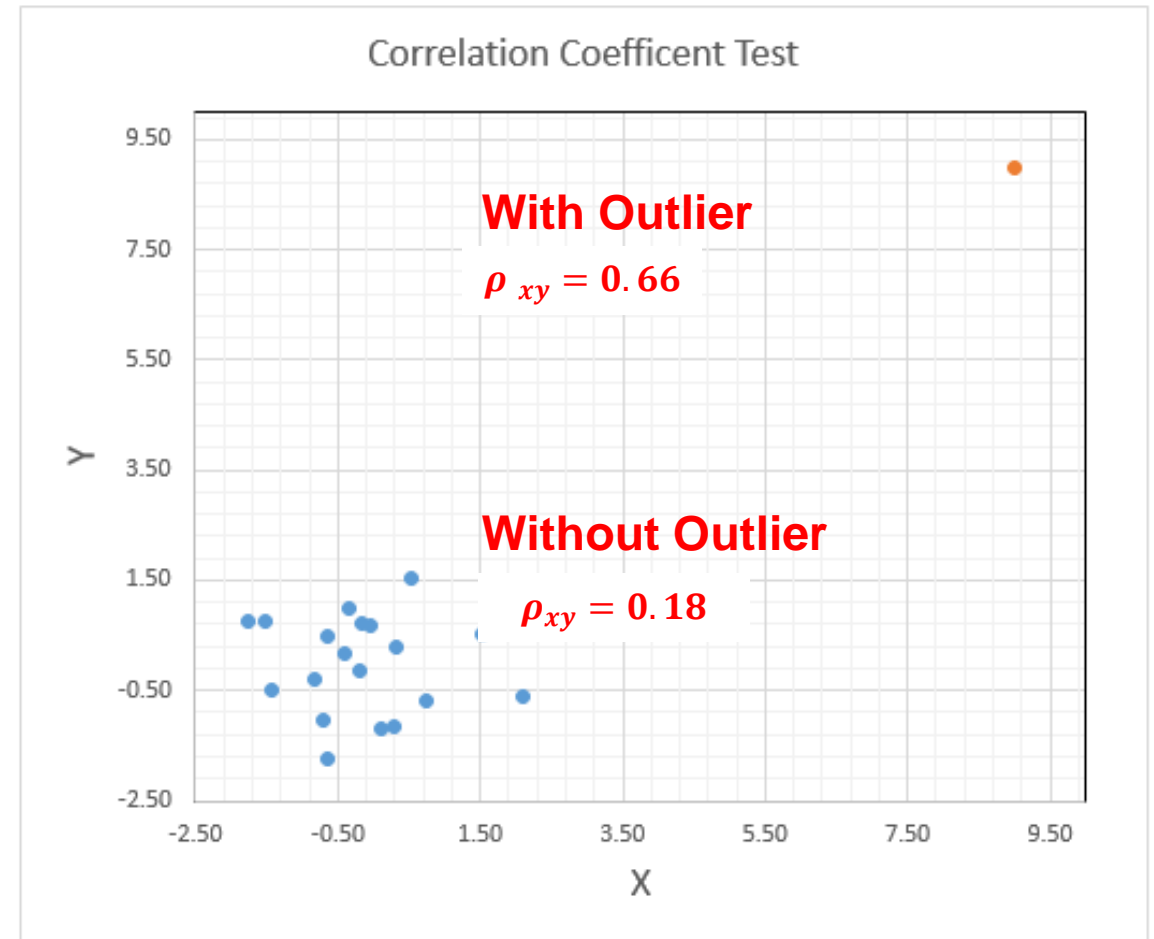
# Bivariate Statistics
## Exercise with Pearson's Correlation Coefficient

Step 1: Generate a random data set of 19 x and y variables and estimate their correlation coefficient (Hint: Rand() in Excel with N[0,1]).

Excel Function NORM.INV(RAND(),0,1)

Step 2: Now add any desired outlier to the data and estimate the correlation coefficient (see example).

How does this outlier affect the correlation coefficient?
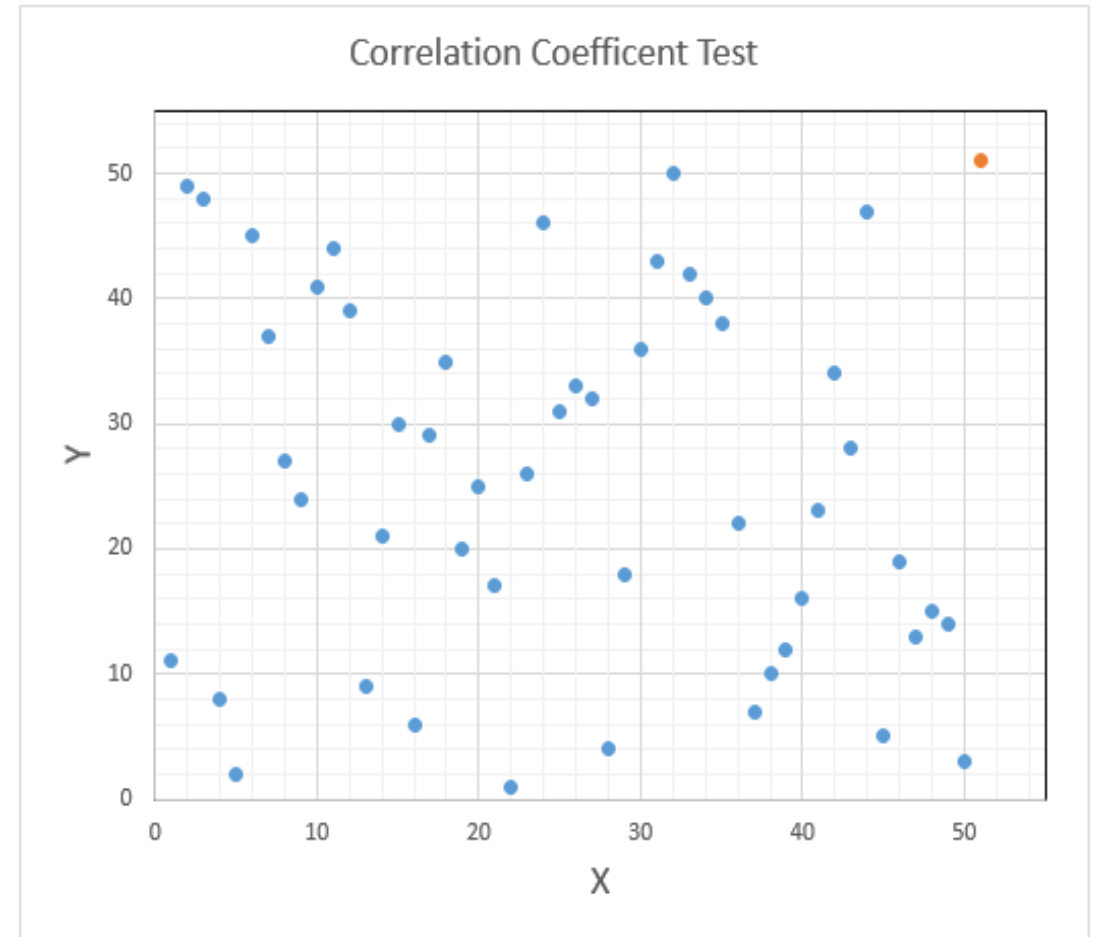


Random data with one bivariate outlier.

# Bivariate Statistics
## Exercise with Pearson's Correlation Coefficient

Step 1: Generate a random data set of 19 x and y variables and estimate their correlation coefficient (Hint: Rand() in Excel with N[0,1]).

Excel Function NORM.INV(RAND(),0,1)

Step 2: Now add any desired outlier to the data and estimate the correlation coefficient (see example).

How does this outlier affect the correlation coefficient?



**With Outlier**

$\rho_{xy} = 0.66$

**Without Outlier**

$\rho_{xy} = 0.18$

Random data with one bivariate outlier.

# Bivariate Statistics
## Exercise with Pearson's Correlation Coefficient

Step 3: Apply the rank transform to the dataset (Hint: 21-Rank.Avg() in Excel).

How does this outlier now affect the correlation coefficient?

This is a more robust form of the correlation coefficient called the rank correlation coefficient.



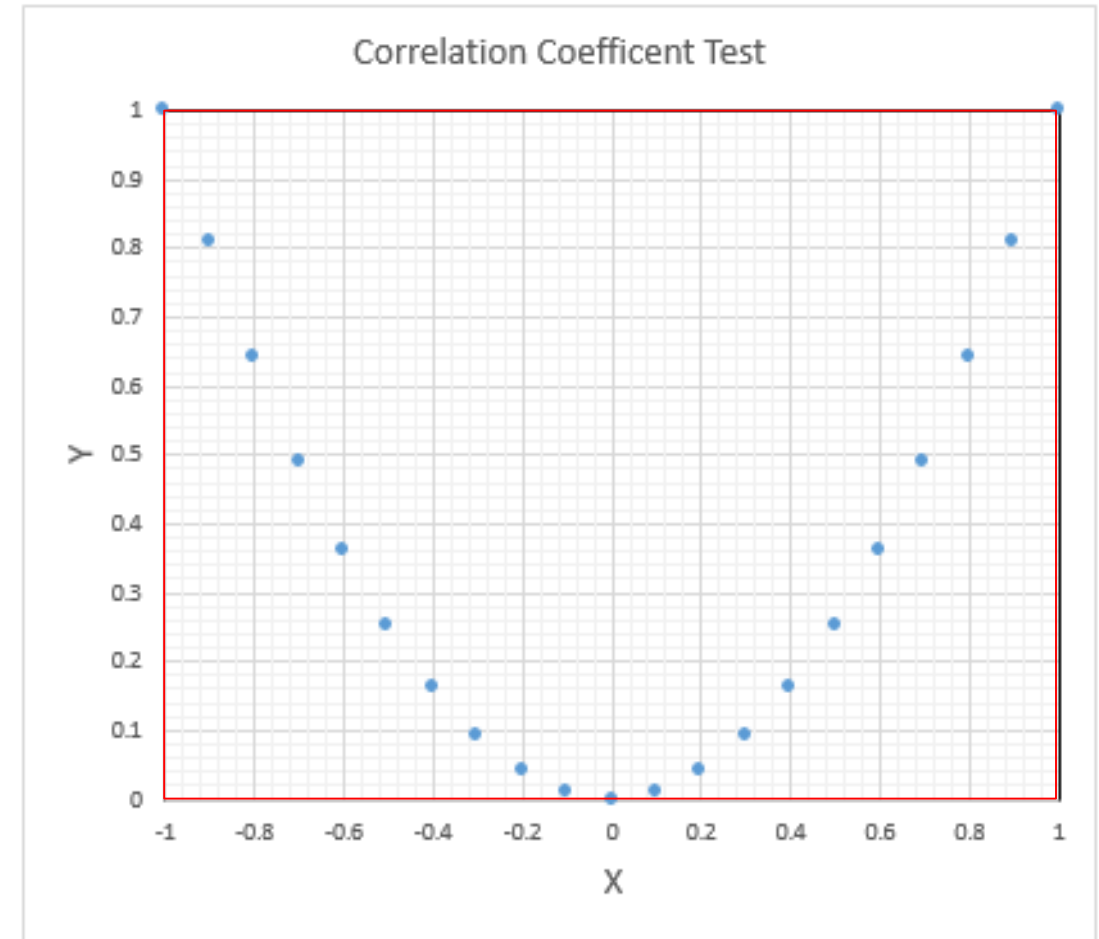Random data rank transformed with one bivariate outlier.

# Bivariate Statistics
## Exercise with Pearson's Correlation Coefficient

Step 3: Apply the rank transform to the dataset (Hint: 21-Rank.Avg() in Excel).

How does this outlier now affect the correlation coefficient?

This is a more robust form of the correlation coefficient called the rank correlation coefficient.



Random data rank transformed with one bivariate outlier.

# Bivariate Statistics
## Measuring Linear Relationships with the Correlation Coefficient

Correlation / Covariance is a measure of linear relationship

- What is the Correlation / Covariance of

$$y = x^2 \text{ over range of } [-1, 1]?$$



Bivariate data that have follow a parabola.

Excel Function Correl(array1,array2)

# Bivariate Statistics
## Measuring Linear Relationships with the Correlation Coefficient

Correlation / Covariance is a measure of linear relationship

- What is the Correlation / Covariance of

$$y = x^2 \text{ over range of } [-1, 1]?$$

- Correlation Coefficient     $\boldsymbol{\rho}_{x,y} = 0.0!$

$$y = x^2 \text{ over range of } [0, 1]?$$

- Correlation Coefficient     $\boldsymbol{\rho}_{x,y} = 0.98!$



Bivariate data that have follow a parabola.

Excel Function Correl(array1,array2)

# Interactive Bivariate Outlier Demonstration

Let's make a random data set, and add and adjust a single outlier

- Watch the Pearson and Spearman correlation coefficients as you increase X and then increase Y for the single outlier.



Interactive Outlier demonstration, the file is Interactive_Correlation_Coefficient_Issues.ipynb.

# Bivariate Statistics
# Mixing Populations

Mixing multiple populations may impact correlation coefficients.



2 populations artificially inflating correlation

2 populations artificially deflating correlation

# Simpson's Paradox

Correlation within multiple groups disappears or reverses when the groups are combined together.



5 groups each with positive correlations.



Combined group with a negative correlation.

Images from Wikipedia gif available at https://commons.wikimedia.org/wiki/File:Simpsons_paradox_-_animation.gif by Pace~svwiki.

# Bivariate Statistics
## Other Issues Correlation Coefficient

For more than two variables make matrix scatterplots:

- By hand in Excel or packages in R and Python.

- Look for linear / nonlinear features

- Look for homoscedasticity (constant conditional variance) and heteroscedasticity (conditional variance changes with value)

- Look for constraints



Matrix scatter plot with a couple interpretations.

# Bivariate Statistics
## Other Issues Correlation Coefficient

Correlation Matrix Plot

- look for high and low correlations

- collinearity, variables with the same correlations with the other variables.



Secondary Variables (26)

Primary Variables (39)

**Correlation matrix plot from McLennon et al., (2006).**

# Bivariate Statistics In Excel

Calculating bivariate statistics in Excel demo

- covariance
- Pearson product-moment correlation coefficient
- Spearman rank correlation coefficient



File is Basic_Statistics_Demo.xlsx.

# Bivariate Statistics In Python

Calculating bivariate statistics in Python demo

- covariance
- Pearson product-moment correlation coefficient
- Spearman rank correlation coefficient

**Data Analytics**

**Basic Bivariate Statistics in Python**

Michael Pyrcz, Associate Professor, University of Texas at Austin

*Twitter* | *GitHub* | *Website* | *GoogleScholar* | *Book* | *YouTube* | *LinkedIn*

**Data Analytics: Basic Bivariate Statistics**

Here's a demonstration of calculation of bivariate statistics in Python. This demonstration is part of the resources that I include for my courses in Spatial / Subsurface Data Analytics and Geostatistics at the Cockrell School of Engineering and Jackson School of Goesciences at the University of Texas at Austin.

We will cover the following statistics:

**Bivariate Statistics**

- Covariances
- Pearson Product Momment Correlation Coefficient
- Spearman Rank Correlation Coefficient

I have a lecture on these bivariate statistics available on YouTube.

File is PythonDataBasics_Bivariate_Statistics.ipynb.

# PGE 338 Data Analytics and Geostatistics

## Lecture 8: Bivariate Distributions

Lecture outline . . .

- **Bivariate Statistics**

- **Correlation**

Introduction

General Concepts

Univariate

Bivariate

Correlation

Regression

Model Checking

Time Series Analysis

Spatial Analysis

Machine Learning

Uncertainty Analysis

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