# EVALUATING THE ACCURACY OF LINEAR AND GEOSTATISTICAL INTERPOLATION METHODS IN SUBSURFACE MAPPING

Daniel Hunter

The Pennsylvania State University

December 2015

PENNSTATE

1855

# Contents

# Abstract

The methods by which we model the Earth's subsurface will always necessitate some form of interpolation. Further, inaccurate interpolation of subsurface geology can lead to wasted money and resources. This study compares the results of both linear and geostatistical interpolation methods utilizing a large sampling of boreholes drilled for a subsurface rock investigation at our study site in coastal Central America.

One way to determine the accuracy of an interpolated surface is to compare the values from the surface to additional values collected in the field. In this study, we divide a total population of nearly 500 borings into two parts; a random sampling of 75% of the borings are used as an input to each of the interpolated surfaces, and the remaining 25% are used to assess the surface's accuracy. The linear interpolation method takes the larger 75% sampling of points, generates a triangulated irregular network (TIN), and converts the TIN to a raster. The same 75% sampling are also used to develop a surface through kriging interpolation, a geostatistical method. We then compare each interpolated surface to the values from the remaining 25% sample not used to generate the surface.

The accuracy of each surface was determined through the use of a three-dimensional root mean square error (RMSE) method. This workflow was used to create multiple iterations of each surface using a different random sampling for each scenario and allowing summary statistics to be evaluated across the study. Our results concluded that there was not a statistically significant difference in RMSE values when comparing the linear and geostatistical interpolation surfaces. We were able to achieve nearly identical results from both interpolation methods.

# Introduction

The utilization of spatial statistics and modern day computing in subsurface mapping has introduced advancements in the way we analyze, explore and ultimately interpolate a surface. Traditional linear interpolation methods will always have a practical application in subsurface mapping. They are exact interpolators, easily understood and have an application to a wide variety of industries and use cases. In recent decades, however, geostatistical interpolations have found their way into modern geographic information systems (GIS) and statistical software packages. These alternate interpolations methods are not without challenge and their implementation requires a thorough understanding of the spatial

distribution of one's data. Our study aims to show that geostatistical methods are a viable alternative to traditional linear interpolation methods by quantitatively comparing interpolated and actual values of a subsurface geologic layer.

GIS software packages have become so increasingly user friendly over the past decade that the end user no longer requires an in-depth knowledge base to perform many common analyses.  Using kriging as an example, there are numerous types of kriging methods; ordinary, universal and simple just to name a few (O'Sullivan, 2010, p310).   There are dozens of different models for fitting a semi-variogram, and each model has an infinite range of parameters to select from.  Yet launching a typical commercial package for this interpolation type will yield with one click a kriging surface after every input parameter has been populated with a default value.  The truth is, the integration of well-designed GUI's, simplified tool kits and the population of default parameters to make spatial processes execute without failure has grown such that anyone can perform simple to complex spatial operations with ease.  Building from the fundamental concept of interpolation this paper illustrates that geostatistical interpolation is a viable alternative to linear interpolation and when thoughtful consideration is given to input parameter selection they can yield similar results.

## Background

After a review of numerous source documents its apparent that no one interpolation method fits all scenarios nor is one necessarily better than the other.  Rather it is a thorough understanding of the input dataset and the ability to select the most appropriate method for interpolation for a particular case.  Each method has its pros and cons; however selection of appropriate parameters and an understanding of each for justification of its selection is the single most important aspect of spatial interpolation.  (Chang, 2009).

An infinite combination of interpolation parameters and methods exists to deduce a surface from a set of sampled data.  The question of which interpolation method is most appropriate has been debated as long as more than one method has existed.  In this study we seek to compare two separate methods of interpolation and quantify the differences and errors between them.  Linear interpolation through the use of a Triangular Irregular Network (TIN) was compared to the kriging method.  This study compared the two methods to each other for subsurface rock layers found on a small site in Central America.

The first surface method to be explored is linear interpolation from a TIN.  The TIN interpolation utilizes the Delauney triangulation method by drawing straight lines between data points to create a framework of non-overlapping triangles covering the study area.  Typically linear interpolation is then used to generate a surface of a site (Yalmiz, 2007, p1352).  The TIN algorithm is the most simplistic of all surface generation methods.  It is an exact interpolator and is best used when the data are evenly distributed over the project area (Yalmiz, 2007, p1353).  Interpolation techniques such as inverse distance weighting and basic linear interpolation have a place in approximating a subsurface, but when additional criteria can be inferred from a given data set and particularly where the data have high natural variability (Virdee, 2009, p370) the use of more advanced interpolation techniques should be considered.  One of these approaches is the kriging interpolation method.  For our purposes we refer to the kriging algorithm as the more advanced interpolation method because its mathematics are more complex, it allows for the input of a greater number variables and has multiple sub-methods of interpolation that are not found with the linear interpolation method.  The application of kriging for subsurface interpolation and mapping has been used for decades and was first developed in the 1960's within the gold mining industry of South Africa.  Since this time it has further evolved into its own field of geostatistics and incorporates a distance weighting approach to interpolation along with an expert knowledge of the spatial structure and trends of the study site (O'Sullivan, 2010, p294).   Kriging interpolation works in a similar manner to inverse distance weighting algorithms in the sense that it utilizes the surrounding measured values to interpolate the unknown (Yalmiz, 2007, p1349).  Kriging can be used when data are irregularly spaced and can be either a smooth or exact interpolator.  Also, there are two variations of the kriging method, that of Ordinary and Universal kriging (Yalmiz, 2007, p1349). In order to use the kriging interpolation algorithm correctly one must assess the problem in three separate but sequential steps.  Step one is to evaluate the spatial variation in the sample data; step two is to summarize the spatial variation with a mathematical function; and step three is to use this function to determine the interpolation weights (O'Sullivan, 2010, p.294).  With the nature of the spatial variability defined one can then proceed through the interpolation process.

Today kriging is widely found in the geologic, mining and surface mapping fields with a variety of use cases.  Applications have been linked to subsurface geologic interpretation, mineral investigations, (Virdee, 2009) observation well placement and hydraulic conductivity and transmissivity modeling

(Samui, 2011).  Some advantages to the kriging methodology include its flexibility, the fact the weighting is not selected arbitrarily and also the use of the semi-variogram as the model for characterizing spatial variability or auto-correlation (Samui, 2011, p886).  Kriging has been used reliably to determine rock depth in numerous areas of the world and used widely for many other applications.  While kriging clearly has practical applications to a wide variety of fields both in and out of the geologic industry, its use must be met with caution in the selection of variables and parameters when concluding the trends and spatial structures of a particular data sets (Meyer, 2004, p1).

The background information for both of these surface generation techniques is plentiful.  The literature on TIN algorithms and the simplicity of their implementation was directly applicable to this study.  It allowed for a conceptual understanding of the interpolation method and its practical application for my use.  Alternatively, the kriging topic has a tendency to become immediately detail driven and study specific without a good discussion of the fundamentals of the interpolation method. Thus documentation of kriging fundamentals is more likely found in textbooks than scholarly articles.

## Methodology

One of the most fundamental ways to determine the quality of an output surface is to compare the values from the interpolated surface to other values of the surface collected in the field.  Since logistical challenges, finances and countless other variables all prohibit revisiting the field site to recollect data, an alternative approach is to subdivide your original sampling of data at the onset of your modeling.  The goal is to utilize a larger subset to develop a surface, and a smaller subset to compare and validate the modeled surface (Esri Help, 2015(subset tool)).  In an effort create a reproducible and iterative methodology two nearly identical workflows were developed using the Python programming language and the Esri ArcGIS Geostatistcal Analyst Extension.  These iterative methodologies were used to evaluate the accuracy of both the geostatistical and linear interpolation methods.

### Overview

Beginning with the entire sampling of borings a subset was selected to interpolate each surface.  This study has chosen to use 75% of the sample as an input to each of the interpolated surfaces.  The remaining 25% sample were used as a quality control and validation data set once the surfaces were created.  Utilizing the Esri *subset features* tool the data were divided into two different data sets.  To be

more specific, starting with a sampling of 410 borings, 307 borings (75%) were randomly selected to interpolate a surface using both a kriging geostatistical interpolation method and a linear interpolation from a TIN method. Because of the randomness of the subset process described above a different sampling of borings is created each time the point features are generated. This is advantageous for two reasons. First, the generation of numerous surfaces from random samplings of the same inputs are possible. Second, with these multiple surfaces generated from the different input samplings, the quality control process can draw conclusions from a statistically significant pool of results and not just a comparison of two surfaces.

With the subset of borings defined, the training sampling of points was used to generate a TIN and for simplicity of comparing two surfaces no break lines or additional features were used to supplement the TIN's creation. This TIN was then converted to a raster surface through a linear interpolation method. Then, using the same subset of borings, the training sampling of points was used to generate the surface with kriging.

Those methods found above are part of a workflow that was run 999 times. The objective was to develop a statistical sampling of surfaces. If the comparison of one linear interpolation surface to one kriging surface showed similar results, a number of questions could be asked of the method by which the surfaces were created, the input parameters that were chosen and the significance of those results. However, if the kriging surfaces showed better results more often than the linear interpolation and the results of statistical testing between the RMSE values proved significant it would indicate that it is the preferable interpolation method.

Once created, each surface was subjected to identical quality control methods and the results were aggregated for statistical review. Since the training sampling of the borings were used to generate the input surfaces the test sampling which remain can be used to measure the difference between the interpolated surfaces and measured values at the locations of these test samples. These differences can then be compared against not only the actual measured value but also the differences can be compared against each other to assess which interpolation method generated a more accurate surface.

Accordingly, this study included a comparison of the statistics of all the linear and geostatistical results. A table was structured such that it summarized the residuals of each interpolation method against the

actual value along with maximum, minimum, mean and standard deviation of all points from both surfaces when compared to the actual surface values. These data were representative of the entire study from all simulations. They were then reviewed through visual means such as histograms or other plots. In addition to the summary statistics, an RMSE for each interpolation method was developed. After one simulation, calculate the RMSE for each interpolated surface (once for the linear, once for the geostatistical). After doing this for every run in the study, conclude which RMSE value was more frequently lower. The final statistical validation is to test for statistical significance of the results. A combination of F and T tests were used to validate the statistical significance of the results.

## Data Preparation

The data preparation phase of this study ensured quality surfaces. Since some of the data points were nearly 70 years old with questionable collection methods and others were collected more recently precedence was given to the more current boring if adjacent borings were contradictory. The Esri Geostatistical Analyst Wizard was also used to review the spatial distribution of the dataset. It has a variety of means to evaluate and review the distribution of data. These included a review of histograms and trend analysis to look for anomalies in the data sets. After a cursory review of the data, the primary means for review of the data was a visual inspection along with a semi-variogram plot. The plotting of the values along the semi-variogram identified the spatial auto-correlation of each pair of points and was an indicator of those which were potential outliers. Working through the Esri Geostatistical Analyst Wizard, a graphical review of the semi-variogram, covariance and cross validation plots were used to review the data distribution. In addition to the graphical review the tabular summary from the cross validation check proved to be useful in further identifying outliers. These reviews resulted in the removal of 23 points from my dataset decreasing the original dataset from 433 to 410 borings – which were used as inputs into this study.

## Subset

After completing a review of the spatial distribution of the entire dataset the process of subsetting the data was employed. The data was broken into subsets using the Esri Geostatistical Analyst Toolbox; Utilities tool kit and the Subset Features tool. The subset tool does exactly what it implies and divides the original dataset into two parts. The first is used as an input to the surface and the second is used to cross validate the output surface. For our analysis our dataset was subset at 75% and 25% for the input and validation sets.

PENNSTATE

## Linear Interpolation

To construct a linearly interpolated surface the input points were used to create a Triangulated Irregular Network (TIN). The input points were utilized as mass points such that each point and elevation value represented a single node in the TIN Surface. These points were the only input to the TIN surface and no supplemental break lines or points were included. Additionally, the TIN surface was created with Delaunay constrained triangulation. This means that each segment created between nodes is represented as a single edge, not densified and eliminates the creation of many small triangles (Esri Help). After the TIN generation was complete each TIN surface was converted to a raster surface. It is important to note that the conversion from TIN to raster introduced a surface that no longer exactly passed through the elevation values used to create the TIN. The raster surface was created with a cell size of ten meters square and proved to be reasonable for optimizing file sizes, processing speeds and data coverage. The cell size was also used to preserve the average spacing of the borings in an effort not to over or undersample the density of the data when converting to a raster. The TIN was converted to a raster through linear interpolation. This meant that the TIN triangle faces are all viewed as planes and each raster cell is assigned a value by finding the elevation within the plane which intersects the center of each raster cell. (Esri Help, 2015).

## Geostatistical Interpolation

The geostatistical interpolation in this study required the spatial distribution of the dataset to be evaluated prior to interpolation. In parallel with the original review of the data this distribution was evaluated for input to the geostatistical interpolation. A single model was created from a combination of input and assessment points and then fit to the semi-variogram. To construct the geostatistical interpolation surface the input points were convert to a surface utilizing the optimized model described below and an Ordinary Kriging methodology.

## Parameter Selection

The parameter selection for the geostatistical interpolation consisted of those values shown below. Of primary interest are values for the nugget, range and sill along with measurement error as they are major drivers of the kriging interpolation method.

| Method | Kriging |
|---|---|
| Type | Ordinary |
| Output Type | Prediction |
| **Dataset #** | 1 |
| Trend Type | None |
| **Searching Neightborhood** | **Standard** |
| Neighbors to include | 5 |
| Include at least | 2 |
| Sector type | Four and 45 degree |
| Major semiaxis | 1,025.650158 |
| Minor semiaxis | 1,025.650158 |
| Angle | 0 |
| **Variogram** | **Semivariogram** |
| Number of lags | 12 |
| Lag Size | 128.206269 |
| Nugget | 2.857811 |
| Measurement error % | 0 |
| **Model type** | **Exponential** |
| Range | 1,025.650158 |
| Anisotropy | None |
| Partial Sill | 107.285091 |

Table 1: Input modeling parameters used for the geostatistical interpolation

## Evaluation

The assessment points were used to extract values from both the kriging and linearly interpolated surfaces and compare their interpolated values to the elevation values at each boring not used in the interpolation. The accuracy of each interpolated surface was then determined through the use of a three-dimensional root mean square error (RMSE). Each RMSE statistic was written to a summary table that documented 999 iterations of each interpolation method. Finally, a combination of F and T tests were used to conclude which interpolation method more often resulted in a lower RMSE. Utilizing the Microsoft Excel Data Analysis Add-In, the F-Test Two-Sample for Variances analysis tool is used to compare two sample variances, in this case the RMSE statistic from both the linear and geostatistical interpolation methods. The statistic indicates whether or not the two samples come from distributions

with equal variances (Microsoft Excel Help, 2015).  Knowing the variances of our RMSE values are equal we can further evaluate the data with the use of a T-Test.  The Two-Sample t-Test with Equal Variance analysis tool tests whether or not the means that underlie each sample are equal. (Microsoft Excel Help, 2015).  The result of the T-Test is a p value which is used as a measure of statistical significance for concluding the results.  Utilizing the Python code developed for this project and presented in Appendix A, the entire methodology section described above; subset, linear interpolation, geostatistical interpolation and evaluation were each performed 999 times to then be summarized and reviewed.

## Results

These results discuss the study findings in detail and use specific terminology when referencing each datasets.  An understanding of the terminology and definition of each point and surface relationship will aid the reader in their interpretation of the results.  These relationships are illustrated in the following two figures.



Figure 1:  Graphic of one scenario of borings (100%), input borings (75%) and validation borings (25%)

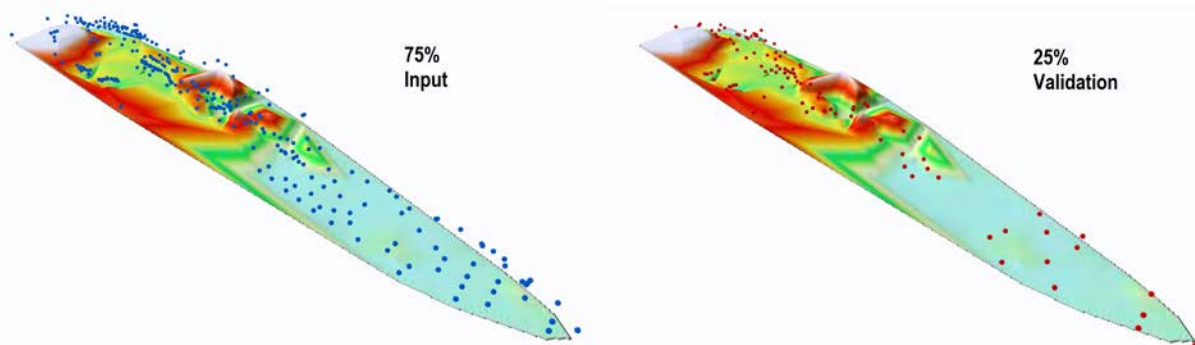

Figure 2:  Graphic of one scenario of input borings (75%) and validation borings (25%) with a linearly interpolated surface.  75% Input on the left was used to create the interpolated surface while the 25% validation on the right is used to compare actual to interpolated values.  The exact same concept was applied for the geostatistical interpolation (not shown).

In an effort to validate our interpolated surfaces we ensured that all surfaces for both the geostatistical and linear interpolation methods conformed to the input boring elevations and concluded that an average RMSE of less than one meter for each surface was achieved. When evaluating against the input borings, the linear interpolation surface averaged a lower RMSE, 0.812 meters, than the geostatistical one, 0.988 meters.

| | Maximum | Minimum | Range | Standard Deviation | Average | Count |
|---|---|---|---|---|---|---|
| **Linear** | 1.693667 | 0.394612 | 1.299055 | 0.166126853 | 0.812565 | 910 |
| **Kriging** | 1.163498 | 0.78811 | 0.375388 | 0.07195292 | 0.988534 | 89 |

Table 2: Summary of RMSE results after comparing interpolated surfaces to the (75%) input points used for creating them.

When testing against the validation points not used for the interpolation our results illustrate that after 999 runs of the model the geostatistical interpolation resulted in a lower root mean standard error (RMSE) more often than the linear interpolation. The geostatistical interpolation showed a lower RMSE 552 times while the linear interpolation resulted in lower RMSE 447 times. We also compiled all the 25% validation points for each of the 999 scenarios and calculated the RMSE using the entire sampling of 96,931 borings. Geostatistical interpolation again returned a lower RMSE, with a value of 3.991 meters versus a linear interpolation RMSE of 4.028 meters.

| | Maximum | Minimum | Range | Standard Deviation | Average | Count |
|---|---|---|---|---|---|---|
| **RMSE Linear** | 6.956974 | 2.48974 | 4.467233 | 0.62566453 | 3.980512 | 447 |
| **RMSE Kriging** | 6.647833 | 2.474804 | 4.173029 | 0.625840782 | 3.94255 | 552 |

Table 3: Summary of RMSE results after comparing interpolated surfaces to the (25%) validation points withheld when creating them.

Figure 3: Histogram plot of RMSE values associated with the geostatistical interpolation



Figure 4: Histogram plot of RMSE values associated with the linear interpolation

In addition to the RMSE evaluation a statistical review of the results was also performed using a combination F and T tests. The F test illustrated equal variance among the RMSE values for each surface type and the T-test evaluated for statistical significant. Among the surfaces being evaluated for RMSE (N=999), there was no statistically significant difference between linear interpolation RMSE (M = 3.9805, SD = 0.62566) and geostatistical interpolation (M = 3.9425, SD = 0.62584), t(1996) = -1.3558 >= .05, CI95

-0.0170, 0.0929.  Therefore, we fail to reject the null hypothesis that there is no difference in RMSE in values between linear and geostatistical interpolation.

| F-Test Two-Sample for Variances | | |
|---|---|---|
| | *RMSE Kriging* | *RMSE Linear* |
| Mean | 3.942550073 | 3.98051151 |
| Variance | 0.391676685 | 0.391456104 |
| Observations | 999 | 999 |
| df | 998 | 998 |
| F | 1.000563489 | |
| P(F<=f) one-tail | 0.496451102 | |
| F Critical one-tail | 1.109804075 | |

Table 4:  Results of F Test

| t-Test: Two-Sample Assuming Equal Variances | | |
|---|---|---|
| | *RMSE Kriging* | *RMSE Linear* |
| Mean | 3.942550073 | 3.98051151 |
| Variance | 0.391676685 | 0.391456104 |
| Observations | 999 | 999 |
| Pooled Variance | 0.391566394 | |
| Hypothesized Mean Difference | 0 | |
| df | 1996 | |
| t Stat | -1.355837612 | |
| P(T<=t) one-tail | 0.087652107 | |
| t Critical one-tail | 1.645617395 | |
| P(T<=t) two-tail | 0.175304214 | |
| t Critical two-tail | 1.961153206 | |

Table 5:  Results of T Test

| Interpolation | n | Mean | SD | t | df | p | 95% Confidence Interval |
|---|---|---|---|---|---|---|---|
| Linear | 999 | 3.9805 | 0.1532 | – | – | – | – |
| Geostatistical | 999 | 3.9426 | 0.1534 | – | – | – | – |
| Total | 1998 | 3.9615 | 0.1533 | -1.3558 | 1996 | 0.1753 | -0.017 - 0.0929 |

Table 6:  Summary results of statistical significance testing

PENNSTATE

# Discussion

There were a number of steps in the methodology in which a decision was made to choose a given parameter for the modeling, some of which may have erred in favor of one method or another. The selection process was not intended to skew results in anyway but rather in many cases it was out of necessity and due to the complexity of the model. The largest concern with the methodology developed was that given the repetitive nature of my method, modeling the semi-variogram repeatedly and automatically defining parameters was not feasible. As such, a single semi-variogram was modeled for the entire dataset and then applied to each geostatistical surface. This possibly gave the geostatistical interpolation an advantage over the linear interpolation since the model to fit the surface was being developed from the entire dataset including those points which are then used for validation.

The selection of model parameters for the geostatistical interpolation could have likewise skewed the results one way or another. The most noticeable parameter to change my results was the measurement error. The measurement error ranges from 0 to 100% and toggling this variable allowed the geostatistical surface to more closely fit the input points (0%) or not as closely (100%). When ensuring the geostatistical surface more closely fit the input points (measurement errors 0%) the RMSE went up for the validation points. This indicates that while geostatistical interpolation can be used to rigidly define the surface it is a better predictor of inexact interpolation or more of a global interpolator. In an effort to model similar scenarios and also approach this from a justification of the surface perspective, I found it more appropriate to ensure the geostatistical surface passed as closely as possible through the points, even at the expense of a larger RMSE.

Yet another unavoidable source of error in the method was the conversion from TIN to raster. By definition the conversion from a TIN to a raster results in a loss of precision. Sampling size and raster cell size can help mitigate this issue but it will never be eliminated. Due to this conversion the RMSE values could have been higher for the linear interpolation method especially if adjacent cells fell on a large transitional area in the TIN.

# Summary

The statistical results showed there was no significant difference between the two interpolation methods.  Accordingly, it cannot be said that one method is better than the other, rather they both have a practical application and the ability to yield highly similar results.

Another variable complicating the outcome is the fact that kriging surfaces are hugely a function of their input parameters which possess a limitless combination of variables.  The selection of kriging parameters is driven by the spatial distribution of the input data and also by the desired outputs of the kriging method.

Furthermore, even though both interpolation methods can be considered exact interpolators an evaluation against the input points illustrates linear interpolation does a better job of conforming to the source data than the geostatistical surface.

An unexpected outcome of this study is documenting the similarities which were created between the two interpolation methods.  Assuming linear interpolation is the simpler method because it is derived from simpler mathematics and minimal input parameters, the fact we matched or exceeded its accuracy with the geostatistical surfaces for nearly every run, is a testament to the use of the geostatistical method and the parameters selected to create it.

# References

Chang, K. (2009). Introduction to Geographic Information Systems. New York: McGraw Hill.
DeMers, M. N. (2003). Fundamentals of Geographic Information Systems (2nd ed.). New York: John Wiley & Sons.

Esri. (2013, August 16). Breaklines in surface modeling. Retrieved December 14, 2013, from http://resources.arcgis.com/en/help/main/10.1/index.html#/Breaklines_in_surface_modeling/00q8000000vq000000/

Esri, (2015). ArcGIS Desktop Help. Retrieved November 2015 from http://desktop.arcgis.com/en/desktop/latest/tools/geostatistical-analyst-toolbox/an-overview-of-the-geostatistical-analyst-toolbox.htm

Kumler, M. P. (1994). An Intensive Comparison of Triangulated Irregular Networks (TINs) and Digital Elevation Models (DEMs). Cartographica , 31 (2), 1-99.

Longley, Paul; Goodchild, Michael; Maguire, David and Rhind, David. (2005) Geographic Information Systems and Science Second Edition. John Wiley and Sons Inc. England.

Microsoft (2015). Microsoft Excel Help. Retrieved November 2015.

Oliver, M. A. (1990). Kriging: A Method of Interpolation for Geographical Information Systems. International Journal of Geographic Information Systems, 4, 313-332.

O'Sullivan, David and Unwin, David (2010). Geographic Information Analysis second edition. John Wiley and Sons Inc. Hoboken, New Jersey.

Samui, Pijush and Thallak G. Sitharam (2011). Application of Geostatistical Models for Estimating Spatial Variability of Rock Depth. Scientific Research, Engineering Volume 3.

Sibson, R. (1981). A Brief Description of Natural Neighbor Interpolation. In V. Barnett, Interpreting Multivariate Data. Chichester: John Wiley.
Tearpock , Daniel J. and Bischke, Richard E. (1991). Applied Subsurface Geologic Mapping. Prentice-Hall Inc. Englewood Cliffs, New Jersey.

Virdee, T.S. and Kottegoda, N. T. (2009). A brief review of kriging and its application to optimal interpolation and observation well selection. Hydrological Sciences Journal. Taylor & Francis, London, England.

Weibel, R., & Heller, M. (1991). Digital Terrain Modeling. In D. Maguire, M. Goodchild, & D. Rhind (Eds.), Geographical Information Systems: Principles and Applications (pp. 269-297). London: Longman.

WordPress (2015). Significance Testing (t-Test). Retrieved September 2015 from https://researchrundowns.wordpress.com/quantitative-methods/significance-testing/

Zoraster, S. (2003). A surface modeling algorithm designed for speed and ease of use withall petroleum industry data. Computers & Geosciences , 29 (9), 1175-1182.

# Appendix 1

## Python Code used for interpolation and comparison of surfaces

```
###################################################################################################################
#   Python code to interpolate and compare a series of surfaces against a subset of points
#   Developed for Penn State Capstone Project GEOG596B
#
#   Author:  Daniel Hunter
#   Contact:  DanielHunterGIS@gmail.com
#
#   Date:  June 2015 – December 2015
###################################################################################################################

# Import arcpy module
import arcpy
import os
import time
arcpy.env.overwriteOutput = "True"

# Check out any necessary licenses
arcpy.CheckOutExtension("GeoStats")
arcpy.CheckOutExtension("3D")
arcpy.CheckOutExtension("spatial")

startTime = time.time()

# Local variables:
srElevPt = r"C:\GIS\PSU\016_GEOG596B\FinalizedResults.gdb\M_SeriesBorings"
inXml = r"C:\GIS\PSU\016_GEOG596B\FinalizedOrdinary.xml"


###################################################################################################################
#   Establish a counter padded with zero's up to 1000 - This will be used for all naming convetions throughout the code...
###################################################################################################################

# A list to store the results for merging
pct25List = []

counter = 1
while counter < 1000:
    if counter < 10:
        counterStr = "000" + str(counter)
        print "Evaluating Scenario: " + counterStr
    elif counter < 100:
        counterStr = "00" + str(counter)
        print "Evaluating Scenario: " + counterStr
    elif counter < 1000:
        counterStr = "0" + str(counter)
        print "Evaluating Scenario: " + counterStr
    elif counter == 1000:
        counterStr = str(counter)
        print "Evaluating Scenario: " + counterStr


###################################################################################################################
#   Subset the points into a 25 and 75% sampling
###################################################################################################################

    # Define 75 and 25 Percent Variables
    sr = arcpy.Describe(srElevPt).spatialReference
    srElevPt75 = arcpy.Describe(srElevPt).path +"\\pct75_pts_" + counterStr
    srElevPt25 = arcpy.Describe(srElevPt).path +"\\pct25_pts_" + counterStr

    # Create a subset of the features
    arcpy.SubsetFeatures_ga(srElevPt, srElevPt75, srElevPt25, "75", "PERCENTAGE_OF_INPUT")
    print "Created subset of Points for scenario " + counterStr
```

```
########################################################################################################################
#  Create the TIN and Linear Interpolation of the Rasters from the 75% sampling
########################################################################################################################

    # Process: Create TIN
    srElevPt75replace = srElevPt75.replace("\\","\\\\")
    srTin75 = str(os.path.dirname(arcpy.Describe(srElevPt).path)) + "\\finalizedTins\\TIN_" + counterStr
    #  This variable contains the field within the points used to create the TIN (could also be shape)
    strTinParameters = srElevPt75replace + " SR_Elevation Mass_Points <None>"
    arcpy.CreateTin_3d(srTin75, sr, strTinParameters, "CONSTRAINED_DELAUNAY")
    print "Created TIN for scenario " + counterStr

    # Process: TIN to Raster
    srTin75_RAS = arcpy.Describe(srElevPt).path +"\\LINT_" + counterStr
    arcpy.TinRaster_3d(srTin75, srTin75_RAS, "FLOAT", "LINEAR", "CELLSIZE 10", "1")
    print "Created the Raster through linear interpolation from the TIN for scenario " + counterStr


########################################################################################################################
#  Create the Raster through a geostatistical interpolation method - Kriging - from the 75% sampling
########################################################################################################################

    # Process: Kriging
    srKrg75_RAS = arcpy.Describe(srElevPt).path +"\\KRIG_" + counterStr
    srKrgVar75_RAS = arcpy.Describe(srElevPt).path +"\\VAR_" + counterStr
    krgLayer = "krgLayer"
    ##  This variable contains the field within the points used to create the TIN (could also be shape)
    strKrgParameter = srElevPt75replace + " X=Shape Y=Shape F1=SR_Elevation"
    arcpy.GACreateGeostatisticalLayer_ga(inXml, strKrgParameter, krgLayer)
    arcpy.GALayerToGrid_ga(krgLayer, srKrg75_RAS, "10", "1", "1")
    ####################################################################

    print "Created the Raster through geostatistical interpolation from the points for scenario " + counterStr


########################################################################################################################
#  Extract raster values from both linear and geostatistical interpolation to the 25% sampling of points
########################################################################################################################

    # Process: Extract Multi Values to Points
    inRasList = [[srTin75_RAS, "LINT"],[srKrg75_RAS,"KRIG"]]
    arcpy.gp.ExtractMultiValuesToPoints_sa(srElevPt25, inRasList, "NONE")
    print "Extracted values from the rasters for the 25% sampling of points for scenario " + counterStr


########################################################################################################################
#  Add and calculate fields for the difference between the actual and the interpolated values along with scenario
########################################################################################################################

    # Process: Add Field
    arcpy.AddField_management(srElevPt25, "dLINT", "DOUBLE", "", "", "", "", "NULLABLE", "NON_REQUIRED", "")
    arcpy.AddField_management(srElevPt25, "dKRIG", "DOUBLE", "", "", "", "", "NULLABLE", "NON_REQUIRED", "")
    arcpy.AddField_management(srElevPt25, "dLINT_sqr", "DOUBLE", "", "", "", "", "NULLABLE", "NON_REQUIRED", "")
    arcpy.AddField_management(srElevPt25, "dKRIG_sqr", "DOUBLE", "", "", "", "", "NULLABLE", "NON_REQUIRED", "")
    arcpy.AddField_management(srElevPt25, "SCENARIO", "TEXT")
    print "Added fields for the 25% sampling of points for scenario " + counterStr

    # Process: Calculate Field
    # RMSE = sqrt(sum((actual-estimated)^2)/count)"
    arcpy.CalculateField_management(srElevPt25, "dLINT", "[SR_Elevation] - [LINT]", "VB", "")
    arcpy.CalculateField_management(srElevPt25, "dKRIG", "[SR_Elevation] - [KRIG]", "VB", "")
    arcpy.CalculateField_management(srElevPt25, "dLINT_sqr", "[dLINT]*[dLINT]")
    arcpy.CalculateField_management(srElevPt25, "dKRIG_sqr", "[dKRIG]*[dKRIG]")
    arcpy.CalculateField_management(srElevPt25, "SCENARIO", "'"+counterStr+"'", "VB", "")
    print "Calculated fields for the 25% sampling of points for scenario " + counterStr

    counter += 1
    currentTime = time.time()

    pct25List.append(srElevPt25)

    print "Elapsed Time = " + str(round((currentTime - startTime),0)) + " seconds"

    print "\n"
```

```
########################################################################################################
#  Merge together all the points to perform a sumary statistic
########################################################################################################

mergedResults = arcpy.Describe(srElevPt).path +"\\mergedResults"
arcpy.Merge_management(pct25List, mergedResults)
print "Successfully merged all features into one feature class: " + mergedResults + "\n"

#  If a null value exists for either the linear or geostatistical interpolation remove the record from the scenaroi
with arcpy.da.UpdateCursor(mergedResults, ["dLint", "dKrig"], ) as cursor:
    for row in cursor:
        if not row[0] or not row[1]:
            cursor.deleteRow()


########################################################################################################
#  Calculate Summary Statistics to conclude results
########################################################################################################

rmseScenario = arcpy.Describe(srElevPt).path +"\\RMSE_Scenario"
rmseCombined = arcpy.Describe(srElevPt).path +"\\RMSE_Combined"
arcpy.Statistics_analysis(mergedResults, rmseScenario,"dLINT_sqr SUM;dKRIG_sqr SUM;dLint MAX;dLint MIN;dKrig MAX;dKrig MIN","SCENARIO")
arcpy.Statistics_analysis(mergedResults, rmseCombined,"dLINT_sqr SUM;dKRIG_sqr SUM;dLint MAX;dLint MIN;dKrig MAX;dKrig MIN", "#")

rmseList = [rmseCombined, rmseScenario]
for rmse in rmseList:
    arcpy.AddField_management(rmse, "dLINT_RMSE", "DOUBLE", "", "", "", "", "NULLABLE", "NON_REQUIRED", "")
    arcpy.AddField_management(rmse, "dKRIG_RMSE", "DOUBLE", "", "", "", "", "NULLABLE", "NON_REQUIRED", "")
    arcpy.AddField_management(rmse, "RMSE", "TEXT")
    arcpy.CalculateField_management(rmse, "dLINT_RMSE", "Sqr([SUM_dLINT_sqr]/[FREQUENCY])", "VB", "")
    arcpy.CalculateField_management(rmse, "dKRIG_RMSE", "Sqr([SUM_dKRIG_sqr]/[FREQUENCY])", "VB", "")
    print "Successfully complete the RMSE calculations for " + rmse

    with arcpy.da.UpdateCursor(rmse, ["dLINT_RMSE", "dKRIG_RMSE", "RMSE"]) as cursor:
        for row in cursor:
            if row[0] > row[1]:
                row[2] = "KRIG"
            elif row[0] < row[1]:
                row[2] = "LINT"
            cursor.updateRow(row)

rmseScenarioResults = arcpy.Describe(srElevPt).path +"\\RMSE_Scenario_Results"
rmseCombinedResults = arcpy.Describe(srElevPt).path +"\\RMSE_Combined_Results"
arcpy.Frequency_analysis(rmseScenario, rmseScenarioResults,"RMSE","#")
arcpy.Frequency_analysis(rmseCombined, rmseCombinedResults,"RMSE","#")

print "Completed the RMSE evaluation"

print "Successfully Completed the code
"
currentTime = time.time()
print "Elapsed Time = " + str(round((currentTime - startTime),0)) + " seconds"
```