

A quick derivation of geostatistical Kriging

William S. Harlan

July 2013

INTRODUCTION

A casual encounter with geostatistics can be baffling because of some non-standard terminology. As it turns out, “Kriging,” the core numerical method of geostatistics, can be derived on a napkin, if you are already familiar with some standard least-squares methods.

If you want to appreciate geostatistical applications, then try the popular book “An introduction to Applied Geostatistics,” by Isaaks and Srivastava [3]. If you want a clean mathematical notation, then try “Multivariate Geostatistics: An introduction with Applications,” by Wackernagel [5];

Kriging is just the least-squares solution of a purely under-determined linear inverse problem. Once you see this equivalence, you can see some simple ways to improve distance-weighted interpolation methods.

POSING THE PROBLEM

Geostatistics poses a useful problem that you are unlikely to have encountered elsewhere in least-squares. The solution looks very familiar, but it has unique elements that make it very useful.

Assume that you have a continuous function $v(\mathbf{x})$ of a spatial vector \mathbf{x} . Usually \mathbf{x} has two dimensions, but nothing about the derivation limits us to two dimensions.

The actual function $v(\mathbf{x})$ is unknown and needs to be reconstructed from a collection of samples $\{v_i\}$ at n arbitrarily chosen locations.

$$v_i \equiv v(\mathbf{x}_i) \text{ for } i=1, n. \tag{1}$$

We will interpolate a value v_0 at an unsampled location \mathbf{x}_0 as a linearly

weighted sum of the n known values $\{v_i\}$, at sampled locations $\{\mathbf{x}_i\}$:

$$\begin{aligned} v(\mathbf{x}_0) &= \sum_{i=1}^n w_i \cdot v(\mathbf{x}_i), \\ v_0 &= \sum_{i=1}^n w_i v_i, \text{ or} \\ v_0 &= \mathbf{w}^\top \cdot \mathbf{v}, \\ \text{where } \mathbf{w}^\top &\equiv [w_1, w_2, \dots, w_n], \text{ and } \mathbf{v} \equiv [v_1, v_2, \dots, v_n]^\top. \end{aligned} \quad (2)$$

In the last version we switch to vector notation.

Our problem: how do we optimally estimate weights \mathbf{w} given what we know about the distribution of \mathbf{x} ?

ASSUMPTIONS

We will make a few more assumptions appropriate to a linear least-squares solution.

We will treat all our values $\{v_i\}$ as random variables. Without loss of generality (changing variables if necessary), we will assume these variables have zero expected mean:

$$E(v_i) \equiv 0, \text{ for } i=0, n. \quad (3)$$

Assume a constant stationary variance σ_v^2 for all samples

$$\sigma_v^2 \equiv E(v_i v_i), \text{ for } i=0, n. \quad (4)$$

Notice that we have included the unmeasured value v_0 in our assumptions.

We will find it convenient to construct a covariance \mathbf{C} matrix for our measured values:

$$\begin{aligned} C_{ij} &\equiv E(v_i v_j), \text{ for } i=1, n \text{ and } j=1, n, \text{ or} \\ \mathbf{C} &\equiv E(\mathbf{v}\mathbf{v}^\top). \end{aligned} \quad (5)$$

Note that the covariance is symmetric: $C_{ij} = C_{ji}$.

Define a covariance vector \mathbf{c} for our unknown value relative to the measured values:

$$\begin{aligned} c_i &\equiv E(v_0 v_i), \text{ for } i=1, n, \text{ or} \\ \mathbf{c} &\equiv E(v_0 \mathbf{v}). \end{aligned} \quad (6)$$

Much of geostatistics concentrates on the estimation of these covariances. For now we assume they are known.

Notice that none of our assumptions require that our interpolated value share the same physical units as the values we are interpolating from. Co-kriging differs from Simple Kriging largely by combining values with mixed units. (See Wackernagel for a demonstration of their equivalence.)

LEAST-SQUARES SOLUTION

Next we estimate our weights as a least-squares optimization, and take advantage of our assumptions.

Define an estimation error σ_E^2 as the expected error between the true value and our interpolated value.

$$\sigma_E^2 \equiv E[(v_0 - \sum_{i=1}^n w_i v_i)^2] \quad (7)$$

$$\begin{aligned} &= E(v_0 v_0) - 2 \sum_{i=1}^n E(v_i v_0) w_i + \sum_{i=1}^n \sum_{j=1}^n E(v_i v_j) w_i w_j \\ &= \sigma_v^2 - 2 \sum_{i=1}^n c_i w_i + \sum_{i=1}^n \sum_{j=1}^n C_{ij} w_i w_j. \end{aligned} \quad (8)$$

Naturally we want this error to be small. Since this expression is a quadratic of the weights \mathbf{w} , we can minimize σ_E^2 by differentiating with respect to the weights and setting the value to zero.

$$\frac{\partial \sigma_E^2}{\partial w_i} = -2c_i + 2 \sum_{j=1}^n C_{ij} w_j = 0 \quad (9)$$

$$\Rightarrow \sum_{j=1}^n C_{ij} w_j = c_i \quad (10)$$

$$\Rightarrow \mathbf{C} \cdot \mathbf{w} = \mathbf{c}$$

$$\Rightarrow \mathbf{w} = \mathbf{C}^{-1} \cdot \mathbf{c}. \quad (11)$$

This result is entirely equivalent to Simple Kriging. The optimum weights are intuitive when you see them in this form.

Naive weights \mathbf{w} would simply use the covariance vector \mathbf{c} between values at the known and the unknown locations. Because the known locations are also correlated with each other, we must remove that effect by multiplying by their inverse covariance \mathbf{C}^{-1} .

Once we have our estimated weights \mathbf{w} , we can explicitly quantify the error in our estimate. This will allow us to put Gaussian error bars on each interpolated value. Substitute our solution (10) into our estimation error (8) for

$$\begin{aligned} \sigma_E^2 &= \sigma_v^2 - 2 \sum_{i=1}^n c_i w_i + \sum_{i=1}^n c_i w_i \\ &= \sigma_v^2 - \sum_{i=1}^n c_i w_i \\ &= \sigma_v^2 - \mathbf{c}^\top \cdot \mathbf{w}. \end{aligned} \quad (12)$$

Notice that this expected error is independent of correlations between measured samples. If it were not, then the estimate could be improved.

CHOOSING COVARIANCES

Geostatistics usually concentrates on the estimation of covariances (or variograms, discussed later).

Here is the most geostatistical assumption of all. The covariance between values depends only on the vector distance between their sampled positions:

$$C_{ij} \equiv E[v_i v_j] \equiv E[v(\mathbf{x}_i)v(\mathbf{x}_j)] \equiv C(\mathbf{x}_i - \mathbf{x}_j) \quad (13)$$

for $i=0, n$ and $j=0, n$.

In the absence of other information, it would be appealing to choose a scale-invariant version of a covariance. Scale-invariance allows covariances to be independent of the units chosen for spatial distance. Covariance would be similar at any scale, like a fractal.

Unfortunately, this assumption leads to a degenerate solution. Scale-invariance requires covariances to go to infinity at a distance of zero, which should have a finite variance σ_v^2 . As we approach infinity, the covariance matrix \mathbf{C} becomes diagonal like an identity matrix, and so the inverse covariance has no effect on weights.

We see this degeneracy with a popular interpolation known as Shepard's "Inverse Distance Weighting":

$$v(\mathbf{x}_0) = \sum_{i=1}^n w_i \cdot v(\mathbf{x}_i) \text{ where}$$

$$w_i = \|\mathbf{x}_i - \mathbf{x}_0\|^{-2} / \sum_{j=1}^n \|\mathbf{x}_j - \mathbf{x}_0\|^{-2}. \quad (14)$$

This weight decreases with distance squared, like gravity. The denominator is just a normalization over all samples.

If we are assuming a specific correlation between our measured and interpolated values, then we ought to be able to assume the same correlation between measured values.

Unfortunately, this covariance approaches infinity at a zero distance:

$$\begin{aligned} \mathbf{C}(\mathbf{x}_i - \mathbf{x}_j) &\propto \|\mathbf{x}_i - \mathbf{x}_j\|^{-2} \\ \Rightarrow \sigma_v^2 &\approx \infty \end{aligned} \quad (15)$$

$$\begin{aligned} \Rightarrow \mathbf{C} &\propto \mathbf{I} \\ \Rightarrow \mathbf{w} &\approx \mathbf{c}. \end{aligned} \quad (16)$$

Sure enough, this degenerate solution tells us to ignore covariances between our known values.

As it turns out, the degeneracy is easily remedied if we sacrifice scale invariance below the density of our samples. Choose a small distance Δx below which the correlation does not increase any further.

$$\mathcal{C}(\mathbf{x}_i - \mathbf{x}_j) \propto \min(\Delta x^{-2}, \|\mathbf{x}_i - \mathbf{x}_j\|^{-2}). \quad (17)$$

Covariances at any smaller distance will all equal the marginal σ_v^2 .

Let's see how this assumption works with a simple example.

A SIMPLE EXAMPLE

Let's say we have two measured values to interpolate along a one-dimensional dimension x .

One value v_2 is twice as far as v_1 from the interpolated point v_0 . Specifically $\|x_1 - x_0\| = 1$ and $\|x_2 - x_0\| = 2$.

If we are using Shepard's interpolation, then we already know enough to estimate the weights as $\mathbf{w}^\top = [0.8, 0.2]$.

However, it should matter very much if the point we are interpolating is between the two measured points, or on the same side of both.

For a non-degenerate covariance, we will assume $\Delta x = 0.5$ in our covariance (17).

On opposite sides of the interpolated point, the two measured points are weakly correlated with $\|x_2 - x_1\| = 3$. We get adjusted weights of $\mathbf{w}^\top = [0.82, 0.18]$, which is very close to the Shepard's interpolation weights.

On the same side of the interpolated point, the two measured points are strongly correlated with $\|x_2 - x_1\| = 1$. The value v_2 is actually better correlated with v_1 than it is with the interpolated v_0 . We expect the nearer value v_1 to override and diminish the contribution of v_2 . And sure enough, the adjusted weights are actually $\mathbf{w}^\top = [1, 0]$, which ignores the second point entirely.

VARIOGRAMS

Most geostatistics literature prefers “variograms” to covariances. This is one of the biggest and most unnecessary obstacles to newcomers. The two quantities have a simple equivalence:

$$\gamma_v(\mathbf{x}_i, \mathbf{x}_j) \equiv \frac{1}{2} E\{[v(\mathbf{x}_i) - v(\mathbf{x}_j)]^2\} = \sigma_v^2 - C_{ij}. \quad (18)$$

Variograms have the minor convenience that values no longer need a zero mean.

The geostatistical assumption (13) can be rewritten

$$\gamma_v(\mathbf{x}_i - \mathbf{x}_j) = \sigma_v^2 - C(\mathbf{x}_i - \mathbf{x}_j) \quad (19)$$

Usually variograms are defined in terms of an offset \mathbf{h} :

$$\begin{aligned} \mathbf{h}_{ij} &= \mathbf{x}_i - \mathbf{x}_j; \\ \gamma_v(\mathbf{h}_{ij}) &= \sigma_v^2 - C(\mathbf{h}_{ij}). \end{aligned} \quad (20)$$

Considering only offsets and ignoring sampling, we can write the equivalence this way:

$$\gamma(\mathbf{h}) = C(\mathbf{0}) - C(\mathbf{h}). \quad (21)$$

If we assume values are uncorrelated at infinite distance, we find

$$\begin{aligned} C(\infty) &= 0 \\ \Rightarrow \gamma(\infty) &= C(\mathbf{0}) \\ \Rightarrow C(\mathbf{h}) &= \gamma(\infty) - \gamma(\mathbf{h}). \end{aligned} \quad (22)$$

The most popular of standard analytic variograms is an exponential:

$$\gamma_{\text{exp}}(\mathbf{h}) \propto 1 - \exp(-3\|\mathbf{h}\|/a), \quad \text{or} \quad (23)$$

$$C_{\text{exp}}(\mathbf{h}) \propto \exp(-3\|\mathbf{h}\|/a). \quad (24)$$

The constant a is the distance over which about 95% of the correlation occurs. Unlike our distance weighting, this function is not scale invariant over any range of offsets.

A user must prepare scatterplots and histograms to scan possible values for the constant a , then choose a best fitting value before kriging can begin. For this sort of analysis, variograms have proven more popular than covariances.

KRIGING AS A POSTERIOR ESTIMATE

It is also possible to recognize Kriging as a special case of a least-squares inverse problem, without postulating the specific weighted solution (2). I thank Dave Hale [1] for explaining this equivalence to me. Kriging is a purely under-determined inverse problem, which makes it a bit unusual and less familiar.

We assume that we have a model vector \mathbf{m} and a data vector \mathbf{d} . The data are assumed to be a linear transform (matrix multiplication) \mathbf{F} of the model, plus a vector \mathbf{n} of noise:

$$\mathbf{d} = \mathbf{F}\mathbf{m} + \mathbf{n}. \quad (25)$$

We assume that our model and noise are selected from stationary correlated Gaussian random processes. The expected value of a model given the data (a posterior estimate) should minimize the following objective function:

$$S(\mathbf{m}) = (\mathbf{d} - \mathbf{F}\mathbf{m})^\top \mathbf{C}_n^{-1} (\mathbf{d} - \mathbf{F}\mathbf{m}) + \mathbf{m}^\top \mathbf{C}_m^{-1} \mathbf{m}. \quad (26)$$

The covariance \mathbf{C}_m is the Gaussian covariance of the model \mathbf{m} , and covariance \mathbf{C}_n is the Gaussian covariance of the noise. We also assume that the noise is uncorrelated to the model. The definition (25) determines the data covariance \mathbf{C}_d :

$$\mathbf{C}_d = \mathbf{F}\mathbf{C}_m\mathbf{F}^\top + \mathbf{C}_n. \quad (27)$$

Since this objective function is a quadratic in \mathbf{m} , we can find an optimum $\hat{\mathbf{m}}$ by setting the derivative to zero:

$$\frac{\partial S}{\partial \mathbf{m}^\top}(\hat{\mathbf{m}}) = -\mathbf{F}^\top \mathbf{C}_n^{-1} (\mathbf{d} - \mathbf{F}\hat{\mathbf{m}}) + \mathbf{C}_m^{-1} \hat{\mathbf{m}} = 0; \quad (28)$$

$$\Rightarrow \hat{\mathbf{m}} = (\mathbf{F}^\top \mathbf{C}_n^{-1} \mathbf{F} + \mathbf{C}_m^{-1})^{-1} \mathbf{F}^\top \mathbf{C}_n^{-1} \mathbf{d}; \quad (29)$$

$$\Rightarrow \hat{\mathbf{m}} = \mathbf{C}_m \mathbf{F}^\top (\mathbf{F}\mathbf{C}_m\mathbf{F}^\top + \mathbf{C}_n)^{-1} \mathbf{d} = \mathbf{C}_m \mathbf{F}^\top \mathbf{C}_d^{-1} \mathbf{d}. \quad (30)$$

The first solution (29) is directly solved from the derivative (28). The second solution (30) was shown by Tarantola [4] (Appendix 6.30) to be entirely equivalent. If (29) is equivalent to (30), then

$$(\mathbf{F}^\top \mathbf{C}_n^{-1} \mathbf{F} + \mathbf{C}_m^{-1})^{-1} \mathbf{F}^\top \mathbf{C}_n^{-1} = \mathbf{C}_m \mathbf{F}^\top (\mathbf{F}\mathbf{C}_m\mathbf{F}^\top + \mathbf{C}_n)^{-1}; \quad (31)$$

$$\Rightarrow \mathbf{F}^\top \mathbf{C}_n^{-1} (\mathbf{F}\mathbf{C}_m\mathbf{F}^\top + \mathbf{C}_n) = (\mathbf{F}^\top \mathbf{C}_n^{-1} \mathbf{F} + \mathbf{C}_m^{-1}) \mathbf{C}_m \mathbf{F}^\top \quad (32)$$

$$= \mathbf{F}^\top \mathbf{C}_n^{-1} \mathbf{F}\mathbf{C}_m\mathbf{F}^\top + \mathbf{F}^\top. \quad (33)$$

Both sides of (32) are trivially equal to (33).

This alternate solution (30) is less often used because it requires a matrix inversion in the data space rather than in the model space. Most typically, an overdetermined problem has a greater number of data values than model values. In the case of Kriging, however, the situation is reversed.

Others such as Hansen et al [2] have observed that this alternate solution (30) is equivalent to Kriging. We will elaborate below.

In the case of geostatistical Kriging, we assume that our model \mathbf{m} is a large dense grid of regularly distributed values. The forward operator \mathbf{F} merely subsamples that grid for the values \mathbf{d} . The dimensionality of \mathbf{d} is possibly orders of magnitude smaller than that of the model \mathbf{m} . The transform \mathbf{F} is a rectangular matrix that has many more columns than rows. Each row contains mostly zeros and a single column with a value of 1. Most columns do not have a value of 1 on any row.

Applying the transposed operator \mathbf{F}^\top to subsampled data \mathbf{d} merely repopulates the sampled values of \mathbf{m} and leaves all other values set to 0. This is also the right inverse of the original transform, so that $\mathbf{F}\mathbf{F}^\top = \mathbf{I}$.

If we have sampled the model directly, there is no need to assume any noise at all. If we allow $\underline{\mathbf{C}}_n = \epsilon \mathbf{I}; \epsilon \rightarrow 0$, then the usual estimate (29) becomes degenerate, and $\hat{\mathbf{m}} \rightarrow \mathbf{F}^T \mathbf{d}$. Rather than interpolating any unsampled values, this estimate sets them all to zero.

The alternate form (30) lets us set the noise variance directly to 0. We can then find the following covariances of our sampled data $\underline{\mathbf{C}}_d$:

$$\text{Assume } \underline{\mathbf{C}}_n \equiv 0; \quad (34)$$

$$\Rightarrow \underline{\mathbf{C}}_d = \mathbf{F} \underline{\mathbf{C}}_m \mathbf{F}^T \quad (35)$$

and then rewrite our alternate estimate (30) as

$$\hat{\mathbf{m}} = (\underline{\mathbf{C}}_m \mathbf{F}^T) (\mathbf{F} \underline{\mathbf{C}}_m \mathbf{F}^T)^{-1} \mathbf{d} \quad (36)$$

$$= (\mathbf{F} \underline{\mathbf{C}}_m)^T \underline{\mathbf{C}}_d^{-1} \mathbf{d} \quad (37)$$

$$= \mathbf{W}^T \mathbf{d}; \quad (38)$$

$$\text{where } \mathbf{W} \equiv \underline{\mathbf{C}}_d^{-1} (\mathbf{F} \underline{\mathbf{C}}_m) \quad (39)$$

These weights (39) for all samples are entirely equivalent to those for a single sample (11). Each measured sample of \mathbf{d} is weighted by its covariance to an interpolated point, then these weights are adjusted by the inverse covariance of the data samples. We now have the advantage of directly relating data covariances to that of the model. This expression is also much easier to generalize to data that do more than merely subsample the model.

REFERENCES

- [1] Dave Hale. Personal communication, 2013.
- [2] Thomas Mejer Hansen, Andre Journel, Albert Tarantola, and Klaus Mosegaard. Linear inverse gaussian theory and geostatistics. *Geophysics*, pages 101–111, 2006.
- [3] Edward H. Isaaks and Mohan R. Srivastava. *An Introduction to Applied Geostatistics*. Oxford University Press, USA, January 1990.
- [4] Albert Tarantola. *Inverse Problem Theory and Methods for Model Parameter Estimation*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2004.
- [5] H. Wackernagel. *Multivariate geostatistics: an introduction with applications*. Springer-Verlag, 2nd edition, 2003.