Evaluating Spatial Randomness Using the Two-Sample Kolmogorov-Smirnov Test

> Kurt A. Fisher and Russell Fuller Undergraduates, University of Utah

> > April 21, 2013

## 2 Problem: How to detect spatial randomness and non-randomness?

Many disciplines need to determine whether x,y point pairs in a field are randomly or non-randomly distributed in  $\mathbb{R}^2$ .

- Epidemiologists need to know if a new H7N9 outbreak is localized to one neighborhood, or if the flu is spreading randomly across a city.
- Geologists need to know whether mineral samples spread across the surface are randomly distributed or if the samples' distribution are evidence of a sub-surface deposit.
- Criminologists need to know if high crime in a neighborhood reflects a real trend in order to efficiently allocate police resources.
- Biologists need to know if the spatial distribution of pine beetles in a drought-stricken forest is random in order to efficiently allocate biological controls.

3 Problem: Detecting spatial randomness - an abstract example

Are these points randomly or non-randomly distributed?



# 4 Problem: Detecting spatial randomness - a cholera outbreak in London in 1854

Are these cholera cases randomly or non-randomly distributed?



Figure : A variant of the original map drawn by Dr. John Snow, a British physician, showing location of cases of cholera during the London epidemics of 1854. Credit: Wikipedia. Spatial Analysis. http://en.wikipedia.org/wiki/Spatial\_analysis. Accessed Apr. 16, 2013.

# 5 Problem: Detecting spatial randomness - nickel ore deposits in Western, Australia in 2010

Are these nickel deposits randomly or non-randomly distributed within the parent komatiite ore bodies? Should a company open a new mine anywhere in the parent ore body, and should a developer expect to extract a similar grade of nickel ore?



Figure : Fig. 2 in: Mamusel, A. et al. 2010. Spatial statistical analysis of the distribution of komatiite-hosted nickel sulfide deposits in the Kalgoorlie terrane, Western Australia: clustered or not? *Economic Geology*. 105(1):229-242. doi: 10.2113 gsecongeo.105.1.229

6 One solution: Perform the Kolomogrov-Smirnov test on the 2nd nearest neighbor distance between points.

- ► Using the Kolomogrov-Smirnov test, we will demonstrate how to test the distribution of point pairs in ℝ<sup>2</sup> for randomness where there is no underlying regression relationship. See Bain and Engelhardt's Introduction to Probability and Mathematical Statistics at pp. 460-461.
- We will illustrate one limitation of the Kolomogrov-Smirnov test to evaluate spatial randomness: the result of the Kolomogrov-Smirnov test is affected by the shape of the measuring frame.

### 7| Roadmap

- Kurt 20 minutes.
  - The distribution of the kth nearest neighbor (kthnn) distance is a test statistic for spatial randomness. Slides 9 to 16.
  - The distribution of kthnn distances of random points measured within a finite observing frame intrinsically has a non-normal distribution. Since kthnn distance distributions are non-normal, the usual parametric tests cannot be applied. Slides 17 to 19.
  - The Kolomogrov-Smirnov test provides an alternative test by comparing two non-normal distributions. The non-random distribution of a simulated "random" kthnn distances in an observing field is compared to the distribution of kthnn distances in an actual data observing frame. Slides 20 to 24.
  - Case study: Are 890 very small craters in an absolute modeling age (AMA) observing field near lunar crater Hell Q spatially distributed randomly or is the 890 crater sample spatially distributed non-randomly contaminated by secondary impacts? Slides 25 to 38.

### 8 Roadmap - continued

- Russell 10 minutes.
  - What effect does varying the height and width of the observing frame have on the distribution of 2nd nearest neighbor distances?
- Kurt and Russell Optional wrap-up topics 10 minutes.
  - Future work. Slide 39.
  - Other tests for spatial randomness. Slide 40-43.
  - Suggested reading on spatial analysis. Slide 48.
  - R packages for spatial analysis. Slide 49.
  - The inverse problem finding clusters. No slides.
    - The importance of the finding cluster problem in modern e-commerce.
    - Cluster finding to define market segmentation.
    - Cluster finding with kthnn to solve the retail clerk problem.
    - The retail clerk problem automated as "You would like" recommendations on Amazon, Gmail, Yahoo Mail, and other websites.
  - Question and answer period.

## 9 Randomness - In Math 3070, $\Delta \mu$ gave a $\mathbb{R}^1$ test statistic.

In Math 3070, we learned how to detect whether two samples taken from a randomly distributioned population are:

- ► Distinguishable, *i.e.* the distance between the sample means, spatially distributed along a R<sup>1</sup> number line, cannot be attributed to randomness; or,
- Indistinguishable, *i.e.* the distance between the sample means can be attributed to randomness.



Figure : Exercise 8.30 in Using R. Is the difference between the distribution of ages of mothers and fathers in the babies database statistically significant?

10 Randomness - In Math 3070 for two sample tests, the t statistic collapsed the problem of distinguishing between random and non-random differences to  $\mathbb{R}^1$ .

$$t = \frac{\bar{x} - \bar{x} - \Delta \mu}{\sqrt{\frac{s_1^2}{m} + \frac{s_2^2}{n}}}$$
  
with  $df = v$ .

# 11 Randomness - In Math 3080, spatial randomness was associated with a f(x).

In Math 3080 and using regression modeling, we learned how to determine if point pairs distributed in  $\mathbb{R}^2$  could be associated with various non-random patterns of functions, *e.g.* – linear, quadratic, polynomial, exponential, power, or mixed function distributions.



Figure : A quadratic polynomial best-fit model for Math 3080, Final Project No. 4, that associates the distribution of x,y point pairs, total annual cost of college (x) and student debt (y) in  $\mathbb{R}^2$ 

# 12 Spatial non-randomness - In Math 3080, the difference $\hat{\beta}_1 - \beta_1$ provided a $\mathbb{R}^1$ test statistic.

In Math 3080 and using regression modeling, the the difference  $\hat{\beta}_1 - \beta_1$  solved the problem of spatially associating points that may be random in  $\mathbb{R}^2$  with a non-random mathematical function. Testing the difference  $\hat{\beta}_1 - \beta_1$  reduced the problem to  $\mathbb{R}^1$ .



**Figure** : Simulations of  $\hat{\beta}_1$ . Fig. 12.13 from Devore (2012).

13 Spatial non-randomness - In Math 3080 for linear regression, the t statistic collapsed the problem of distinguishing between random and non-random differences in  $\hat{\beta}_1 - \beta_1$  to  $\mathbb{R}^1$ .



## 14 Spatial randomness - But if there is no f(x), what is a viable $\mathbb{R}^1$ test statistic?

If there is no underlying mathematical function to associate with the distribution of the point pairs in  $\mathbb{R}^2$ , then how do you measure and test spatial randomness? The point pairs shown in the figure below have order (by angled geologic strata and by the division between gaseous, solid and liquid matter), they are not spatially distributed randomly, but there is no f(x) to regress on.



Figure : utah.Pictures.com. 2013. Sundial Peak, Big Cottonwood Canyon, Utah. Accessed Apr. 17, 2013. What is a mathematical description of this field's non-randomness?

15 Spatial randomness - If there is no f(x), the kth nearest neighbor distance (kthnnd) provides a  $\mathbb{R}^1$  test statistic. How to find the 2nd nn distance:

- Take any point  $P(x_1, y_2)$ .
- ▶ Find all Euclidean distances between P(x<sub>1</sub>, y<sub>2</sub>) and P(x<sub>i</sub>, y<sub>i</sub>) for i = 1 to n, and store the distances in array D.
- Sort array D.
- ► Choose D<sub>3</sub>. This is the 2nd nearest neighbor P(x<sub>2nn</sub>, y<sub>2nn</sub>) to P(x<sub>1</sub>, y<sub>2</sub>). (D<sub>1</sub> is always zero, the distance between P(x<sub>1</sub>, y<sub>2</sub>) and itself.)



Figure : Finding the Euclidean 2nd nearest neighbor distance.

# 16 Spatial randomness - 2nd nearest neighbor distances provide only a relative indication of spatial distribution.



The 2nnd histogram (b - right) only indicates relative distribution, and then, only if the distribution of 2knnd's in the data sample is compared to a second sample. The 2nnd histogram does not indicate distribution with respect to a normal random distribution.

## 17 | Spatial randomness - The 2nnd distribution is intrinsically non-normal.



(c) Histogram of 2nnds for 347 (d) QQ plot of 2nnds for 347 point points pairs pairs

The 2nnd distribution is intrinsically non-normal, and thus, the usual parametric tests cannot be applied.

18 Spatial randomness - The 2nnd distribution is intrinsically non-normal because of information loss caused by imposing a finite observing frame on an infinite random point field.



Figure : Information loss from a finite measuring frame. Credit: Baddelely (2011) at p. 117.

Although the distribution of 2knn distances of a true infinite spatially random field is normal, imposing an finite observing frame causes information loss of the true 2knn distances for points near the edge of the finite observing frame. 19 An aside - How to represent an infinite random field as two  $\mathbb{R}^1$  cdf vectors that enclose finite observing frames.



Figure : Two finite observing frames within an infinite random field.

Any  $\mathbb{R}^2$  observing frame containing x,y point pairs can be scaled to fit within this real number infinite probability space.

20 How do you test for spatial randomness when the data and reference distributions are intrinsically non-normal? An intrinsically non-normal 2nnd distribution can still be evaluated for randomness by comparing the data field distribution to a non-normal distribution of 2nn distances from simulated random fields.



(a) 2nnds - Simulated (b) Cumulative plot of (c) ECFD plot of suband actual counts subfigure (a) figure (a)

The two-sample Kolmogorov-Smirnov test provides a method for comparing any two non-normal distributions, after conversion of each cumulative distribution to an empirical continuous function distribution (ECFD).

# 21| Evaluating spatial randomness using the two-sample Kolmogorov-Smirnov test.

The two-sample Kolmogorov-Smirnov (KS) test provides a method for determing if any two non-normal distributions are equal.



(d) Two-sample Kolmogorov- (e) Andrej Smirnov (KS) test. The lines are Nikolajewitsch the ECFD of two non-normal distri- Kolmogorov butions. The arrow is the KS test (b.1903-d.1987) statistic D, discussed below.

Credits: Wikipedia. http://en.wikipedia.org/wiki/Kolmogorov-Smirnov\_test, http://en.wikipedia.org/wiki/Andrey\_Kolmogorov. Accessed Apr. 17, 2013.

#### 22| The two-sample Kolmogorov-Smirnov test

The two-sample Kolmogorov-Smirnov (KS) test provides a method for comparing any two non-normal distributions by finding a width  $\pm D$  around  $F(x)_1$  that contains  $F(x)_2$  with probability  $1 - \alpha$ .



Figure : Illustration of two empirical CDFs compared using the Kolmogorov-Smirnov test. After Bain (2009), Figure 13.1. D is the critical characteristic, and D is the maximum y-distance between the two ecdf's.

23 The two-sample Kolmogorov-Smirnov test statistic

 $F_{n_1}(x)$  and  $F_{n_2}(x)$  are both ECFDs.

Source: Zillwinger (2003a). The KS statistic usually is denoted as D, but here KS is used in order to avoid confusion with the use of D for crater diameter, below.

A demonstration computation of the Kolmogorov-Smirnov method with prediction intervals can be found in the *REA Statistics Problem Solver* (1978) at Problem 20-10.

### 24 The two-sample Kolmogorov-Smirnov critical distribution

For n > 40: p-value Critical value of KS  $\alpha = 0.20: KS >= 1.07 \sqrt{\frac{n_1 + n_2}{n_1 n_2}}.$  $\alpha = 0.10: KS >= 1.22 \sqrt{\frac{n_1 + n_2}{n_1 n_2}}.$  $\alpha = 0.05: KS >= 1.36 \sqrt{\frac{n_1 + n_2}{n_1 n_2}}.$  $\alpha = 0.01 : KS >= 1.63 \sqrt{\frac{n_1 + n_2}{n_1 n_2}}.$  $\alpha = 0.001: KS >= 1.96 \sqrt{\frac{n_1 + n_2}{n_1 n_2}}.$ 

Note: Zillwinger (2003b). For critical values for  $n \le 40$  with small sample size corrections, *see id*.

### 25 Case Study: Are 890 very small craters near lunar crater Hell Q spatially distributed randomly?

- H<sub>a</sub>: 890 very small lunar craters in a 16km<sup>2</sup> observing field near lunar crater Hell Q are not spatially distributed randomly.
- ► H<sub>o</sub>: 890 very small lunar craters in a 16km<sup>2</sup> observing field near lunar crater Hell Q are spatially distributed randomly.
- **Test**: Two-sample Kolmogorov-Smirnov to  $\alpha = 0.10$ .
- Purpose: Before an absolute modeling age (AMA) for Hell Q can be estimated, a precondition of modeling is that craters in the observation field must be spatially distributed randomly. If craters in the field are not spatially distributed randomly, the craters in the field are re-examined for secondary impact contamination, the secondary impact craters are excluded, and then the spatial randomness test is repeated. The initial count of 890 craters was preliminarily evaluated for spatial randomness before more labor was expended to measure an additioanl 600-800 craters.

### 26 Case study: background.

In southeast quadrant of large lunar crater Deslandres, the 4 km diameter crater Hell Q is flanked by Cassini's Bright Spot (CSB). In the 1670s, italian astronomer M. Cassini speculated that the bright spot might be an atmospheric cloud moving across the lunar surface.



Figure : Source: Fisher, K. 2012, Nov. 7. LPOD. http:/lpod.wikispaces.com/November+7, +2012 from NASA, JPL, and Ariz. State. Univ. high-resolution LRO imagery. Scale: Hell Q is 4km in dia.

### 27 Case study: Background - the observing frame.

890 small craters were counted in an observing frame north of Hell Q.



Figure : Source: Per above. Scale: Hell Q is 4 km in diameter.

#### 28 Case study: Background - craters within the frame.

Some of the 890 craters counted. Dots are proportional to crater dia.



Figure : Source: Per above. Scale: The observing frame is approx. 2.2 km wide.

### 29 Case study: Background - graphs of craters by diameter bin.



Figure : Source: Per above. Scale: The observing frame is approx. 2.2 km wide and cover 16.2 km<sup>2</sup>.

30 Case study: Counts and 2nn distances for 890 very small craters north of Hell Q.

#### Binned Crater Counts with 2nd Nearest Neighbor Distance Characteristics (N=890)

$D_{low}$	$D_{upp}$	n	$\mu$ k2nnd	$\sigma^2$ k2nnd	$\sigma$ k2nnd
4	6	347	114.8	3879.0	62.3
6	8	283	122.4	5171.6	71.9
8	12	187	163.4	7817.6	88.4
12	20	58	293.1	28720.2	169.45
20	36	14	646.0	110674.0	332.7
36	68	1	na	na	na

Note: All distances are in meters. D = min. Feret's diameter of crater. The bins are open at the top of the interval, that is for bin 4-6, the diameters of craters in the bin are greater than or equal to 4 meters and less than 6 meters.

### 31| Case study: 2nn distances for 890 very small craters north of Hell Q.



Figure : Histograms of distributions of k2nn crater distances in 4 of 5 observed field bins.

32 Case study: Cumulative actual and simulated 2nn distances by binned crater diameter



Figure : Cumulative count distributions of actual and simulated fields for 4 of 5 bins.

### 33 Case study: ECFDs of actual and simulated 2nn distances by binned crater diameters.



Figure : Cumulative Empirical Continuous Function Distributions of actual and simulated fields for 4 of 5 bins.

34 Case study: ECFDs of actual and simulated 2nn distances for the largest 5th bin.



Figure : Cumulative Empirical Continuous Function Distributions of actual and simulated fields.

35 Case study: Results of the two-sample Kolmogorov-Smirnov test.

Kolmogorov-Smirnov two sample test results of 2nn distances binned by crater diameter (N=890) at significance level  $\alpha = 0.10$ 

Bin	$D_{low}$	$D_{upp}$	n	KS stat.	KS p-value	KS c.v.	Random
1	4	6	347	0.814	0.20	0.087	Y
2	6	8	283	0.0812	0.16	0.096	Y
3	8	12	187	0.0599	0.33	0.118	Y
4	12	20	58	0.091	0.003	0.21	N†
5	20	36	14	0.142	< 0.0001	0.314‡	N†
6	36	68	1	na	na	na	

## 36 Case study: Results of the two-sample Kolmogorov-Smirnov test.

Notes to results table:  $\alpha = 0.10$ , two-sided. KS stat. = the Kolmogorov-Smirnov test statistic from comparing the distribution of observed field k2nn distances with a simulated random field. The KS statistic usually is denoted as D, but here KS is used in order to avoid confusion with the use of D for diameter. KS p-value = the p-value of corresponding to the KS statistic for a field comparision. KS c.v = the KS critical value at the  $\alpha = 0.10$ significance level. Random = Conclusion of whether bin-field is spatially random. Y=Yes; N=No.  $\dagger$  = R's ks.test function returned a rejection p-value inconsistent with the critical value in Zillwinger 2003a and 2003b. If the critical values used in Zillwinger 2003a and 2003b were used, these fields would be random. ‡ critical values for  $n \le 40$  per Zillwinger 2003b. All distances are in meters. D = min. Feret's diameter of crater.

37 Case study: Conclusions from the two-sample Kolmogorov-Smirnov test.

- Bins 1-3, dia.s 4-12 m, n=817: Craters are spatially distributed randomly.
- Bin 4, dia.s 12-20 m, n=58: Craters are not spatially distributed randomly, but the non-random distribution may be the result of secondary crater impacts contaminating the field. This diameter bin will be retested after secondary impact craters are identified and removed from this subsample of 58 craters.

38 Case study: Conclusions from the two-sample Kolmogorov-Smirnov test (continued).

Bin 5, dia.s 20-36 m, n=14: Computer tabulated Kolmogorov-Smirnov test results for this small 14 crater field were inconclusive because both R and *Mathematica* returned p-values and Kolmogorov-Smirnov critical values for alpha = 0.10 that were inconsistent with the small sample critical values in Zillwinger (2003b). The computer tabulated results indicate that both R and Mathematica may not incorporate small sample size corrections shown in Zillwinger (2003b). R documentation states that its two sample KS test results may be inaccurate for small sample sizes, but *Mathematica* is silent on that issue. Manual recomputation of the two-sample Kolmogorov-Smirnov test will be undertaken in future work.

### 39 Case study: Future work.

- Plotting techniques in R will be improved to show the two-sample Kolmogorov-Smirnov predictive interval envelope for a given *alpha*. The second ECF distribution should lie within the prediction interval of the first ECF distribution.
- For bin 4, secondary impact craters will be identified and removed from the sample, and then the two-sample Kolmogorov-Smirnov test will be rerun.
- ► For bin 5, manual recomputation of the two-sample Kolmogorov-Smirnov test will be undertaken.
- ► For all bins, each bin will also be tested using the G test implemented in the R-package *spatstat* and or the mean 2nd closest neighbor test. The two-sample Kolmogorov-Smirnov test is strongest for detecting non-random 2nn distances that are large, but the G test is more robust at detecting non-randomness from small clusters of points.

# 40 Other tests - The weakness of the two-sample Kolmogorov-Smirnov test.

- Comprehensive evaluation of the spatial randomness of a field usually requires two test approaches. In the first approach, a test that is sensitive to larger 2nnd distances is used, and in the second approach, a test that is more robust at detecting small clusters of points is applied.
- The two-sample Kolmogorov-Smirnov test is relatively less sensitive to detecting small clusters of points.



#### 41 Other tests - The G test

Another randomness test that is more sensitive to small clusters is the G test (Bivand, Gómez-Rubio, and Pebesma (2008) at 161). The G test is based on the complete spatial randomness theorem, and the G test defines the expected distance between an event relative to an arbitrarily selected event as the G function:

$$G(r) = 1 - \exp(-\lambda \pi r^2)$$

where  $\lambda$  is the intensity of events, that is the mean number of events per unit area, r is the mean distance between all events in a unit area, and the observed G statistic is:

$$\hat{G}(r) = \frac{\#(d_i : d_i \leq r, \forall i)}{n}$$

#### 42 Other tests - The G test

G and G-like tests have a common algorithmic characteristic. The observing field is divided in sub-blocks of random size by seeding the observing frame with random points. A fiducial point is selected within each sub-block, *e.g.* - the geometric center of the sub-block or the actual point closest the geometric center.



Figure : Illustration of dividing an observing frame into random sized sub-areas using random seed points.

The 1st or 2n nearest neighbor distances with respect to the fiducial point are found within each sub-block. The distribution of means or standard deviations of the kthnn distances within each sub-block creates a new distribution. The process is repeated for simulated random fields, and the resulting two distributions are compared.

### 44 End of Part 1 - Kurt

#### Evaluating Spatial Randomness Using the Two-Sample Kolmogorov-Smirnov Test

In Part 2, Russ will talk about how changing the shape of the observing frame affects the result of the two-sample Kolmogorov-Smirnov test.



Figure : Two finite observing frames within an infinite random field.

45 How to obtain the presentation archive copy

#### Evaluating Spatial Randomness Using the Two-Sample Kolmogorov-Smirnov Test

This presentation's archive link (active to May 5, 2013): http://fisherka.csolutionshosting.net/ out/Present.html

#### 46 References

Baddelely, A. 2011. Analyzing spatial point patterns in R. (Workshop Notes). Perth, Australia:CSIRO and University of Western Australia. Retrieved Apr. 6, 2013, from http://www.csiro.au/resources/pf16h.

Bain, L. J., & Engelhardt, M. 2009. *Introduction to probability and mathematical statistics.* Belmont, CA: Brooks/Cole Cengage Learning. At p. 460-461

Bivand, R. S., Gómez-Rubio, V., & Pebesma, E. J. 2008. *Applied spatial data analysis with R.* New York: Springer.

Devore, J. L. 2012. 8th ed. *Probability and statistics for engineers and scientists*. Boston, Mass.:Brooks-Cole.

### 47 References

Research and Education Association. (1978). *The statistics problem solver*. New York: REA. (At Prob. 20-10 demonstrating Kolmogrov-Simirnov test computation).

Zwillinger, D. (ed.) 2003a. (31 Ed.) Table 7.14.9. Critical values, two sample Kolmogorov-Smirnov Test. Section 7.14. In *CRC Standard Mathematical Tables and Formulae*. Boca Raton, La.:Chapman and Hall-CRC.

 2003b. Table 7.14.8. Critical values, Kolmogorov-Smirnov Test.
Section 7.14. In *CRC Standard Mathematical Tables and Formulae*. Boca Raton, La.:Chapman and Hall-CRC.

#### 48 Suggested Reading

Baddelely, A. 2011. Analyzing spatial point patterns in R. (Workshop Notes). Perth, Australia:CSIRO and University of Western Australia. Retrieved Apr. 6, 2013, from http://www.csiro.au/resources/pf16h.

Bivand, R. S., Gómez-Rubio, V., & Pebesma, E. J. 2008. *Applied spatial data analysis with R.* New York: Springer. Available as Marriott Library online holding.

Diggle, P. 2003. *Statistical analysis of spatial point patterns.* London: Arnold. Marriott Library QH323.5 .D6 2003.

Gelfand, A. E. 2010. *Handbook of spatial statistics*. Boca Raton: CRC Press. Available as Marriott Library online holding.

#### 49 R packages for spatial analysis

Baddeley, A. and Turner, R. et al. 2013, Mar. 1. Spatstat, an R Package. Computer software. (Version 1.31-1). Retrieved Apr. 3, 2013 from http://cran.r-project.org/web/packages/ spatstat/index.html. See also http://www.spatstat.org/.

Bivard, R., Lewin-Koh, N., and Pebesma, E. et al. 2013, Feb. 12. Maptools, an R Package. Computer software. (Version 0.8-23). Retrieved Apr. 17, 2013 from http: //cran.r-project.org/web/packages/maptools/index.html (Maptools is needed to manage spatial data used by spatstat.)