## THE PROBLEM OF GEOGRAPHIC CONTIGUITY - A MONTE CARLO APPROACH

Stephen F. Gibbens, James A. Tonascia California State Department of Public Health

This paper describes a method for determining whether patterns of counties, census tracts or other types of geographic units depart from a configuration that might have occurred by chance. The method presented is conceptually simple and can be applied with comparative ease through the use of a computer program available through SHARE<sup>1</sup>. This paper describes the problem, the method of solution and gives some useful probability distributions.

Suppose, for example, a rate of some type was calculated for each of the N counties in a particular state and a map was shaded, showing the N/4 counties having the highest rates. This top quartile of counties might distribute themselves in a number of ways - varying from a sin-gle cluster forming a large clump to the opposite situation, where most of the selected counties were isolated, surrounded by counties with low rates. Whatever pattern was formed, most observers would make a judgment on the meaningfulness of the pattern manifested. If the high rate counties tended to be those containing large metropolitan centers, this fact would probably be noted rather quickly. If, however, a judgment was made whether the high rate counties showed a pattern of contiguity - or geographic contagion the problem might become much more difficult. The point at which a pattern departs from simple randomness is obscure indeed, when it must be based upon visual inspection.

A number of criteria could be formulated defining nonrandomness. A tendency for the chosen counties to form a few clusters, a preponderance of chosen counties falling in a single quadrant of the state, a "least squares" solution based on distance from some focal point - all might be developed into probability sampling distributions for departures from randomness. In this paper, two criteria are developed: (1) the largest cluster formed by the selected counties; and (2) the number of pairs of contiguous counties found among those counties selected. These will be referred to as the "clusters" and "pairs" methods.

## Clusters Method

From a shaded map we have a visual picture of geographic contiguity. A nonrandom pattern is characterized by clusterings of shaded counties. However, our data could be more easily quantified if we described the clustering in terms of the number of counties found in the largest cluster formed by the N/4 sample. With this criterion, it is possible to conveniently describe distributions of many samples, and to obtain probabilities of the occurrence of clusters of various sizes. In order to effect this quantification, a listing is required for each state showing each individual county with its contiguous counties.<sup>2</sup>

With the availability of a "county contiguity list" (See Table A), it is possible to select a sample of N/4 counties, from a particular state and note the largest cluster formed by these selected counties.

Briefly described, say we have a state with 100 counties and by use of random numbers select a sample of N/4 counties, notated as A, B, C, . . . Y. We search the county contiguity list for the counties contiguous to A. If counties B and J are contiguous to A, we know we have a cluster of at least three counties. Next, we scan the list of counties contiguous to B and find counties A and F. County A was already counted but F is a new county and so our cluster is known to include at least four counties; A, B, F and J. We proceed similarly to search the lists of counties contiguous to F and J, incrementing the cluster size by one for each sample county found. We continue to county C (the first of the remaining counties sampled that is not already known to be in our first cluster of contiguous counties), repeating the above procedure, and hence on to county Y, the last county included in our sample. When we have finished, we record the number of counties in the largest cluster formed. If we were to repeat this procedure for many samples, accumulating the results, we would have an approximate sampling distribution for the largest cluster size found for that state when N/4 counties are selected by chance. For example, the State of Iowa has 99 counties from which 3,000 independent samples of 25 counties each were selected. In about 5 percent of these 3,000 samples a cluster as large as

<sup>2</sup> It should be noted here that the determination of contiguity or noncontiguity is often quite arbitrary. In this paper, two counties were required to share a common border in order to be considered contiguous. If they were separated by water narrow enough to be bridged, they were considered contiguous. Hence, the Mississippi was not a barrier to contiguity but Lake Michigan was. However, different investigators may want to apply different rules for county connections. For example, when four counties meet at a point, so that their borders form a "cross" (+), a decision must be made as to whether or not the diagonal counties are to be considered contiguous. If "point contiguity" is accepted, i.e., if the diagonal counties are considered contiguous, does this carry over to the situation where the borders describe pairs of obtuse angles, rather than right angles as was the case for the cross? These decisions should be made only after careful consideration of these problems, since alternatives produce significantly different results.

<sup>&</sup>lt;sup>1</sup> The computer program is available through the SHARE users group (identified as GO BC GEOG), with specifications.

10 counties was found; in 1 percent, a cluster as large as 13 counties was found, and in .1 percent, a cluster as large as 16 counties was found.

### Pairs Method<sup>3</sup>

The use of the largest cluster size found is not always an ideal measure of geographic contiguity. If there were several smaller clusters, there might be a departure from randomness even though the largest cluster had only a moderate probability level. Another approach would be to obtain the sampling distribution of the number of pairs of contiguous counties found when selecting many random samples of N/4 counties. A large number of pairs of counties could be found either in the situation where there was a single very large cluster, or, several smaller clusters.

The number of pairs of contiguous counties found in a sample would be ascertained, again by recourse to the county contiguity list, accumulating a tally for every border common to two counties included in the sample. In the above example for the State of Iowa, in 3,000 samples, 19 pairs of counties were found in  $3\frac{1}{2}$  percent of the samples, 21 pairs in 1 percent of the samples, and 24 pairs in .1 percent of the samples.

#### Discussion

The "clusters" method and the "pairs" tend to give similar results when applied to live data. The pairs method is more sensitive to the situation where several clusters of moderate size are formed by the sampled counties, sometimes showing a significant probability level where the clusters method is not significant in terms of the number of contiguous counties in the largest cluster. However, this gain in sensitivity in the pairs method, is in equal measure lost, when contiguity is concentrated in a single cluster. In this situation, the clusters method may show a significant probability level while the pairs method shows a very moderate probability.

Figures 1 and 2 illustrate these different situations. Figure 1 shows the lowest quartile of counties in terms of live birth rate for Iowa in 1962. The largest cluster contains 13 counties. This occurred less than once in every 100 trials when using chance methods. However, among these same counties there were 17 pairs (instances of 2 low quartile counties sharing a common border). This occurred about 12 times in every 100 trials when using chance methods. Figure 2 shows the contrasting situation where the cluster of 8 counties was shown to have occurred over 11 times in every 100 trials whereas the pairs method occurred only twice in every 100 trials.

Since the various states tend to contain different numbers of counties and since the configuration of counties tends to vary as well, clearly,

 $^3$  The authors are indebted to Donald Loveland for suggesting this approach.

each state has its own unique sampling distribution. Even where two states have the same number of counties and the appearance of the counties is similar, their sampling distributions are often distinctly dissimilar when subjected to a Kolmogorov-Smirnov test. The general appearance of the sizes and shapes of the counties in a state gives little indication of the probability pattern that will emerge.

Table 1 shows the Estimated Cumulative Probability Distribution for Cluster Sizes of Contiguous Counties for Selected States based on 3,000 samples drawn from each state. Each Sample included about 25 percent of the counties in the state. Referring again to the State of Iowa, the table shows that in .001 of the samples, the largest cluster formed by the 25 sample counties was 2 counties; that in about half the samples (.492) the largest cluster found was 5 counties, or less. Similarly, in 95.2 percent, the largest cluster formed was 10 counties or less. The largest cluster of contiguous counties found in the 3,000 samples drawn was 16 counties. From a sample of 25 counties in the State of Iowa, it would be rare indeed to find the largest cluster to be as few as 2 counties, or as many as 16 counties.

Since these probability distributions were developed by the Monte Carlo method, they are qualified as estimated distributions. Barring a bias in the pseudo-random number generator used by the computor (BC RNDY), the true probability limits could be estimated by the usual method.

For a 10-county cluster in the State of Iowa which has an estimated cumulative probability of .952, a .05 level confidence interval would be calculated as:

$$952 \pm 1.96 \sqrt{\frac{(.048)(.952)}{3,000}} = .952 \pm .008$$

There appear to be two main problems in the practical application of this methodology. Let us say an investigator is attempting to ascertain whether a particular type of congenital anomaly found in newborn infants shows a nonrandom geographic distribution. Usually the investigator would not know in what percent of the counties the anomaly should show high rates - perhaps in N/4 of the counties or perhaps N/10. This would depend on the distribution of the factors which he hypothesizes gives rise to high incidence of the anomaly. Since each sample size (N/4 or N/10) would have a different sampling distribution, the task of selecting an appropriate sample size becomes critical. The only solution to this would seem to be to rank the counties according to incidence rates and if some of the high incidence counties were tested to show significant departures from a norm - to cut off below these counties and use that number as the sample size. The use of "fortuitous" distributions would not appear to invalidate the ultimate probabilities since these are based on the geographic pattern manifested, rather than on the statistical significance of the rates in the selected counties.

The second problem is more subtle. The probability distribution obtained is valid and useful in an a priori sense. However, usually an investigator has the data on a shaded map and is interested in knowing whether the distribution is nonrandom. When viewed in an a posteriori sense, it may appear invalid to use the tabled probabilities because of gross irregularities in the particular counties included in the shaded portion of the map. For example, if a shaded map of California showed a number of the coastal counties in the high rate quartile, the tabled probabilities would be inappropriate for comparison because the probabilities were obtained from all counties. Since the coastal counties have no contiguous counties on their ocean side, they would have less opportunity to be included in large clusters and no opportunity to contribute pairs for a portion of their borders. Therefore, clusters of counties which included coastal counties would occur less often than suggested by the tabled probabilities. At the other extreme, an unusually large size county would tend to border on more counties than the average size county and would therefore be found in clusters more often than suggested by the tabled probabilities. In atypical situations such as these it might be reasonable to compare the obtained pattern with a large number of samples which included one or more specified counties in each sample selected. The resulting conditional probabilities would reflect this bias introduced by the inclusion of a constant county.

The computer program mentioned in the introduction allows for the inclusion of constant counties for obtaining these biased probability distributions.

Analysis of county contiguity within a state often shows clusters of counties along state boundaries. The question arises as to the results that might have been obtained had the bordering state been included in the analysis. The computer program allows for grouping a number of states into a single analysis and gives probability distributions of clusters and pairs of counties drawn from this larger parameter.

While geographic contiguity is not a common statistical problem, it is frequently encountered in epidemiology and geography. The "shaded map" seems to present a strong stimulus for "closure" or resolution, compelling the observer toward interpretation. In the experience of the authors, several observers may interpret a shaded map as showing various degrees of contagion. The most articulate, or, the senior observer present, wins the argument. The methodologies presented in this paper, when properly applied, provide the investigator with a simple, objective determination of contiguity.

#### REFERENCES

- Geary, R. C. "The Contiguity Ratio and Statistical Mapping", Incorporated Statistician, Volume 5, 1954.
- Ederer, Fred; Myers, Max H. and Mantel, Nathan. A Statistical Problem in Space and Time: Do Leukemia Cases Come in Clusters? 1963 (Unpublished).

# Table A

#### SAMPLE COUNTY CONTIGUITY LISTING

(A portion of the State of Iowa showing counties according to numeric code.)

INDEX COUNTY	COUN	TIES CON	TIGUOUS	TO IN	DEX CO	UNTY
01	15	39	61	88	02	
02	69	15	01	88	87	
03	96	22				
04	93	68	26			
05	83	14	39	15		
06	86	07	10	57	48	
07	38	12	09	10	06	86
08	37	94	40	85	77	25
٠			•	•		•
•	•	•	•	•	•	•
•	•	•	•	•	•	•
99	46	41	40	35	94	





Largest cluster - 13 counties Probability that a cluster as large as 13 counties could have occurred by chance = .008

Number of pairs of counties - 17 Probability that as many as 17 pairs of counties could have occurred by chance =.117

Figure I. A Configuration Favoring The Cluster Method

CRUDE DEATH RATE IN IOWA 1982 (lowest quartile)



Largest cluster - 8 counties Probability that a cluster as large as 8 counties could have occurred by chance =.116

Number of pairs of counties - 20 Probability that as many as 20 pairs of counties could have occurred by chance =.018

Figure 2. A Configuration Favoring The Pairs Method

Table 1
---------

## ESTIMATED CUMULATIVE PROBABILITY DISTRIBUTION FOR CLUSTER SIZES OF CONTIGUOUS COUNTIES SELECTED STATES - 25 PERCENT SAMPLE

(	Based	on	3,00	) tria	is)	
---	-------	----	------	--------	-----	--

NUMBER OF COUNTIES IN LARGEST CLUSTER	ALABAMA <sup>1</sup> (17)	ARKANSAS (19)	CALI- FORNIA (14)	COLO- RADO (16)	10WA (25)	KEN- TUCKY (30)	LOUI- SIANA (16)	MINNE- SOTA (22)	MIS- SOURI (29)	MONTANA (14)	NEBRASKA (23)	PENNL SYLVANIA (17)	SOUTH DAKCTA (17)	TEN- NESSEE (24)	VIR- GINIA (24)	WIS- CONSIN (18)
1 2 3 4 5 6 7 8	.000 .005 .103 .319 .538 .701 .820 .896	.000 .003 .265 .499 .672 .804 .882	.000 .037 .255 .543 .767 .879 .947 .976	.000 .014 .137 .369 .583 .736 .838 .908	.000 .001 .062 .254 .492 .688 .806 .884	.000 .000 .016 .122 .312 .515 .674 .778	.000 .146 .164 .426 .646 .793 .913 .937	.000 .003 .059 .265 .476 .651 .773 .860	.000 .000 .018 .134 .327 .511 .672 .785	.000 .028 .227 .493 .700 .834 .915 .963	.000 .002 .080 .306 .536 .725 .826 .898	.000 .012 .139 .391 .507 .762 .863 .921	.000 .014 .142 .397 .607 .769 .866 .923	.000 .001 .040 .193 .419 .610 .748 .851	.000 .003 .073 .270 .497 .687 .814 .882	.000 .011 .122 .354 .568 .726 .831 .905
9 10	•946 •974	•933 •959	•991 •997	•952 •975	.931 .952	.865 .910	•966 •988	.914 .947	.860 .912	.984 .993	•938 •966	.953 .979	.960 .979	.916 .954	.928 .957	.953 .971
11 12 13 14 15 16 17	•988 •995 •998 •999	•980 •989 •994 •998 •998	.999	•989 •995 •998 •999	.977 .986 .992 .996 .998 .999	.948 .968 .982 .991 .995 .998 .999	•994 •998 •999	.972 .985 .993 .997 .998 .999	.947 .968 .980 .989 .993 .997 .999	•999	•982 •986 •995 •998 •999	.993 .998 .999	.990 .995 .999	973 985 992 995 998 998	.978 .986 .990 .995 .998 .999	.987 .994 .998 .999

1 Number in parentheses is the number of counties used in the 3,000 trials for that particular state and is approximately 25 percent of the total number of counties in the state.

## Table 2

#### ESTIMATED CUMULATIVE PROBABILITY DISTRIBUTION FOR THE NUMBER OF PAIRS OF CONTIGUOUS COUNTIES SELECTED STATES - 25 PERCENT SAMPLE

#### (Based on 3,000 trials)

NUMBER OF PAIRS	ALABAMA 1 (17)	ARKANSAS (19)	CALI- FORNIA (14)	COLORADO (16)	IOWA (25)	KENTUCKY (30)	LOUI- SIANA (16)	MINNE- SOTA (22)	MISSOURI (29)	MONTANA (14)	NEBRASKA (23)	PENN- Sylvania (17)	SOUTH DAKOTA (17)	TENTESSEE (24)	VIPGINIA (24)	WISCONSIN (18)
1 2 3 4 5	.000 .000 .001 .004 .016	.000 .000 .001 .001	.001 .004 .028 .083 .196	.000 .001 .003 .015 .048	.000 .000 .000 .000 .001	.000 .000 .000 .000 .000	.000 .001 .004 .013 .053	.000 .000 .001 .001 .005	.000 .000 .000 .000 .000	.001 .003 .018 .064 .155	.000 .000 .000 .000 .002	.000 .001 .003 .012 .032	.000 .000 .001 .009 .033	.000 .000 .000 .000 .000	.000 .000 .000 .000 .001	.000 .000 .002 .009 .023
6 7 8 9 10 11 12 13 14 14	,052 .122 .234 .368 .523 .664 .779 .868 .930 .967	.016 .050 .109 .214 .344 .495 .637 .746 .835 .901	.353 .523 .685 .815 .894 .948 .976 .991 .996 .998	.104 .216 .361 .517 .657 .778 .863 .926 .962 .982	.004 .012 .031 .074 .144 .238 .354 .493 .632 .752	.000 .000 .002 .004 .012 .031 .060 .111 .183	.137 .260 .415 .566 .721 .822 .896 .941 .967 .983	.011 .031 .064 .139 .231 .338 .478 .597 .715 .805	.000 .001 .005 .009 .024 .052 .096 .161 .256	.298 .470 .629 .758 .860 .925 .964 .982 .993 .997	.010 .027 .069 .147 .245 .359 .501 .631 .756 .851	.080 .170 .310 .464 .615 .745 .842 .912 .952 .977	.089 .191 .320 .471 .625 .754 .852 .923 .958 .978	.001 .004 .012 .031 .073 .137 .227 .345 .467 .592	.003 .007 .017 .045 .094 .164 .272 .404 .521 .641	.061 .128 .245 .371 .528 .674 .8674 .871 .923 .956
16 17 18 19 20 21 22 23 24 25 26 27 28 29 30	.982 .992 .996 .998 .999	•943 •970 -986 •994 •998 •998	.999	•990 •997 •999	.842 .883 .938 .965 .982 .993 .996 .998 .999	.277 .382 .508 .617 .719 .798 .861 .909 .942 .943 .983 .999 .995 .998 .999	.991 .996 .998 .999	.884 .922 .954 .974 .985 .993 .995 .999	.356 .480 .587 .690 .770 .846 .899 .935 .961 .976 .991 .996 .998 .999	<b>.</b> 999	*906 .943 .969 .984 .993 .997 .998 .999	.990 .994 .996 .999	.988 .996 .998 .999	.703 .804 .868 .921 .949 .972 .985 .994 .997 .997	.748 .832 .897 .935 .964 .980 .991 .995 .997 .999	.978 .987 .995 .999

1 Number in parentheses is the number of equaties used in the 3,000 trials for that particular state and is approximately 25 percent of the total number of counties in the state.