### Lecture 1: Fundamentals of Spatial Statistics

#### Description versus Inference

- Description and descriptive statistics
  - Concerned with obtaining summary measures to describe a set of data
- Inference and inferential statistics
  - Concerned with making inferences from samples about populations
  - Concerned with making legitimate inferences about underlying processes from observed patterns

We will be looking at both!



































#### Median and Mean Centers for US Population

#### Median Center:

Intersection of a north/south and an east/west line drawn so half of population lives above and half below the e/w line, and half lives to the left and half to the right of the n/s line

#### Mean Center:

Balancing point of a weightless map, if equal weights placed on it at the residence of every person on census day.











# Lecture 2: Point Pattern

# Point Pattern Analysis

Analysis of spatial properties of the entire body of points rather than the derivation of single summary measures

Two primary approaches:

- Point Density approach using Quadrat Analysis based on observing the frequency distribution or density of points within a set of grid squares.
  - Variance/mean ratio approach
  - Frequency distribution comparison approach
- Point interaction approach using Nearest Neighbor Analysis based on <u>distances</u> of points one from another

Although the above would suggest that the first approach examines *first order* effects and the second approach examines *second order* effects, in practice the two cannot be separated.

### Exhaustive census -used for secondary (e.g census) data



### **Random sampling** --useful in field work



Multiple ways to create quadrats --and results can differ accordingly!





### Frequency counts by Quadrat would be:

	Censu	us Q = 64	Samp	ling Q = 38
Number				
of points				
in				
Quadrat	Count	Proportion	Count	Proportion
0	51	0.797	29	0.763
1	11	0.172	8	0.211
2	2	0.031	1	0.026
3	0	0.000	0	0.000
Q = # of q	uadarts			
P = # of p	oints =	15		

Quadrats don't have to be square --and their size has a big influence

### Quadrat Analysis: Frequency Distribution Comparison

- We can compare <u>observed</u> frequencies in the quadrats (Q= number of quadrats) with <u>expected</u> frequencies that would be generated by
  - a random process (modeled by the Poisson frequency distribution)
  - a clustered process (e.g. one cell with P points, Q-1 cells with 0 points)
  - a uniform process (e.g. each cell has P/Q points)
- The standard Kolmogorov-Smirnov test for comparing two frequency distributions can then be applied – see next slide

A = area of region P = # of points

See Lee and Wong pp. 226-229 for example and further discussion.

# Kolmogorov-Smirnov (K-S) Test

### The test statistic "D" is simply given by:

D = max [Cum Obser. Freq – Cum Expect. Freq]

The largest difference (irrespective of sign) between <u>observed</u> cumulative frequency and expected cumulative frequency

The critical value at the 5% level is given by:

$$D_{(at 5\%)} = \sqrt{\frac{1.36}{Q}}$$
 where Q is the number of quadrats

 Expected frequencies for a random spatial distribution are derived from the Poisson frequency distribution and can be calculated with:

$$p(0) = \mathbf{e}^{-\lambda} = 1 / (2.71828^{P/Q})$$
 and  $p(x) = p(x - 1) * \lambda / x$ 

Where x = number of points in a quadrat and p(x) = the probability of x points

P = total number of points Q = number of quadrats

 $\lambda$  = P/Q (the average number of points per quadrat)

See next slide for worked example for cluster case

Ca	alculatio	on of Poiss	on Freq	uencies for K	olmogoro	v-Smirnov te	st			
Cl	USTER	ED pattern	as used	in lecture						
	А	В	С	D	E	F	G	Н		
		=(	ColA * Co	=Col B / q			!	Col E - Col	G	
N	umber of	Observed			Cumulative	9	Cumulative	Absolute		
P	oints in	Quadrat	Total	Observed	Observed	Poisson	Poisson	Difference		
qu	uadrat	Count	Point	Probability	Probability	Probability	Probability			$\mathbf{D}_{\mathbf{OW}} = 10$
	0	8	0	0.8000	0.8000	0.1353	0.1353	0.6647		KOW IU
	1	0	0	0.0000	0.8000	0.2707	0.4060	0.3940	`	
	2	0	0	0.0000	0.8000	0.2707	0.6767	0.1233		
	3	0	0	0.0000	0.8000	0.1804	0.8571	0.0571		<b>T</b> 1 11
	4	0	0	0.0000	0.8000	0.0902	0.9473	0.1473		The spreadsheet
	5	0	0	0.0000	0.8000	0.0361	0.9834	0.1834		
	6	0	0	0.0000	0.8000	0.0120	0.9955	0.1955		<u>spatstat.xis</u> contains
	7	0	0	0.0000	0.8000	0.0034	0.9989	0.1989		worked avanables for the
	8	0	0	0.0000	0.8000	0.0009	0.9998	0.1998		worked examples for the
	9	0	0	0.0000	0.8000	0.0002	1.0000	0.2000		Uniform/Clustered/
	10	2	20	0.2000	1.0000	0.0000	1.0000	0.0000		Unitorni/ Clusicicu/
										Random data previously
										Random data proviousiy
		e <sup>e</sup> e 0								used, as well as for Lee
	0000	° 8° 0								
										and Wong's data
										$\mathcal{C}$
_		<u> </u>								
Ir	ne Kolmo	ogorov-Smir	nov D tes	st statistic is th	he largest A	Absolute Diffe	rence	0.0047		
_		= largest v	alue in C	olumn h	4.00/			0.6647	0	
CI	ntical Va	lue at 5% fo	or one sa	ample given by	1.36/sqrt((	2) (0.1 + 0.0) (0.1 *		0.4301	Significant	
CI	ntical Va	lue at 5% fo	or two sa	imple given by:	1.36*sqrt((	Q1+Q2)/Q1*	Q2))			
				0	10	(				
nu	imper of	quadrats		Q	10	(sum of colu	IMN B)			
nu	imber of	points		Р	20	(sum of Col	C)			
nu	imper of	points in a	quadrat	X						
			··· (· · )				- >			
pc	usson pr	opapility	p(x) =	p(x-1)"(₽/Q)/x		w iii onward	s)			
if	x=0 ther	p(x) = p(0)	=2.7182	8^P/Q	(Col E, Ro	w 10)				
E	uler's coi	nstant	2.7183							

### Quadrat Analysis: Variance/Mean Ratio (VMR)

 Apply uniform or random grid over area (A) with width of square given by:



- Treat each cell as an observation and count the number of points within it, to create the variable X
- Calculate variance and mean of X, and create the variance to mean ratio: variance / mean
- For an uniform distribution, the variance is zero.
  - □ Therefore, we expect a variance-mean ratio close to 0
- For a random distribution, the variance and mean are the same.
   Therefore, we expect a variance mean ratio around 1
  - □ Therefore, we expect a variance-mean ratio around 1
- For a clustered distribution, the variance is relatively large
  - Therefore, we expect a variance-mean ratio above 1

••	3	1			2	2			0	0	
	5	0			2	2			0	0	
° °	2	1			2	2			10	10	
• •	1	3			2	2			0	0	
• • •	3	1			2	2			0	0	
°°° °		X				X Number				X	
RANDOM		Number of			_	of Points			_	Number of	
• •	Quadrat	Points Per			Quadrat	Per			Quadrat	Points Per	
	#	Quadrat	x^2		#	Quadrat	x^2		#	Quadrat	x^2
· · ·	1	3	9		1	2	4		1	0	0
• •	2	1	1		2	2	4		2	0	0
o 0	3	5	25		3	2	4		3	0	0
0 0	4	0	0		4	2	4		4	0	0
0 0	5	2	4		5	2	4		5	10	100
UNIFORM/	6	1	1		6	2	4		6	10	100
DISPERSED	7	1	1		7	2	4		7	0	0
	8	3	9		8	2	4		8	0	0
	9	3	9		9	2	4		9	0	0
	10	1	1		10	2	4		10	0	0
		20	60			20	40			20	200
	Variance	2.222			Variance	0.000			Variance	17.778	
	Mean	2.000			Mean	2.000			Mean	2.000	
CLUSTERED	Var/Mean	1.111			Var/Mean	0.000			Var/Mean	8.889	
	rand	om			unifo	rm			Clust	ered	
Formulae for	r varianc	e				No	ote:				
$\sum_{i=1}^{n} (\lambda$	$(x_i - \overline{X})^2$	$=\frac{\sum_{i=1}^{n}\lambda_{i}}{\sum_{i=1}^{n}\lambda_{i}}$	$X_{i}^{2}[(\sum X)]$	)²/N	<u></u>	N =	= nun	ıbe	er of <u>(</u>	Quadra	ts = 10
N -	-1		<i>N</i> – 1			Ra	tio =	V	arianc	e/mea	n

## Significance Test for VMR

- A significance test can be conducted based upon the chi-square frequency
- The test statistic is given by: (sum of squared differences)/Mean

$$= \frac{\sum_{i=1}^{n} (X_i - \overline{X})^2}{\overline{X}} = \frac{\sum_{i=1}^{n} X_i^2 [(\sum X)^2 / N]}{\overline{X}}$$

- The test will ascertain if a pattern is significantly more <u>clustered</u> than would be expected by chance (but does <u>not</u> test for a uniformity)
- The values of the test statistics in our cases would be:

$$\frac{\underline{\text{random}}}{2} = 10 \qquad \frac{\underline{\text{uniform}}}{2} = 0 \qquad \frac{\underline{\text{clustered}}}{2} = 80$$

- For degrees of freedom: N 1 = 10 1 = 9, the value of chi-square at the 1% level is 21.666.
- Thus, there is only a 1% chance of obtaining a value of 21.666 or greater if the points had been allocated randomly. Since our test statistic for the clustered pattern is 80, we conclude that there is (considerably) less than a 1% chance that the clustered pattern could have resulted from a random process
- T-test in p.p.234

## Weakness of Quadrat Analysis

- Results may depend on quadrat size and orientation (Modifiable areal unit problem)
  - test different sizes (or orientations) to determine the effects of each test on the results
- Is a measure of dispersion, and not really pattern, because it is based primarily on the density of points, and not their arrangement in relation to one another

For example, quadrat analysis cannot distinguish between these two, obviously different, patterns



Results in a single measure for the entire variations within the region are not recognized (could have clustering locally in some areas, but not overall)

For example, overall pattern here is dispersed, but there are some local clusters



## Nearest-Neighbor Index (NNI)

- uses distances between points as its basis.
- Compares the mean of the distance <u>observed</u> between each point and its nearest neighbor with the <u>expected</u> mean distance that would occur <u>if</u> the distribution were random:

NNI=Observed Aver. Dist / Expected Aver. Dist

For random pattern,NNI = 1For clustered pattern,NNI = 0For dispersed pattern,NNI = 2.149

We can calculate a Z statistic to test if observed pattern is significantly different from random:

Standard Error

if Z is below -1.96 or above +1.96, we are 95% confident that the distribution is <u>not</u> randomly distributed. (If the observed pattern was random, there are less than 5 chances in 100 we would have observed a z value this large.)

(in the example that follows, the fact that the NNI for uniform is 1.96 is coincidence!)



## Nearest Neighbor Formulae

### Index

NN Index: The Nearest Neighbor Index (Uncorrected)

$$NNI = \frac{\overline{d}}{E(\overline{d})}$$

Where:

$$\overline{d} = \frac{\sum_{i=1}^{n} d_i}{n}$$

Exp. Avg.: Expected Average Nearest Neighbor Distance (Uncorrected)

$$E(\overline{d}) = 0.5 \sqrt{\frac{A}{n}}$$

Significance test

**SD:** Standard deviation (Standard error)

$$SD = \sqrt{\left(\frac{1}{4\tan^{-1}1} - \frac{1}{4}\right)\frac{A}{n^2}} = \frac{0.26136}{\sqrt{n^2/A}}$$

Z: Standard z-value (Uncorrected)

$$Z = \frac{\overline{d} - E(\overline{d})}{SD}$$



### RANDOM



#### CLUSTERED



#### UNIFORM

	Nearest			Neares
Point	Neighbor	Distance	Point	Neighbo
1	2	1	1	2
2	3	0.1	2	3
3	2	0.1	3	2
4	5	1	4	5
5	4	1	5	4
6	5	2	6	5
7	6	2.7	7	6
8	10	1	8	9
9	10	1	9	10
10	9	1	10	9
		10.9		
Mean distar	100 1 00		Mean d	listance 0
Aroo of	1.09		Area of	
Aled Ul Dogion	50		Region	ļ
Donaity	0.2		Density	0
Exported	0.2		Expecte	ed
Moon	1 112034		Mean	1.11803
	0.07/026		NNI	0.08944
NNI	0.374320			
Ζ	= -0.15	15	Ζ	= 5.508

	Nearest		
Point	Neighbor	Distance	
1	2	0.1	
2	3	0.1	
3	2	0.1	
4	5	0.1	
5	4	0.1	
6	5	0.1	
7	6	0.1	
8	9	0.1	
9	10	0.1	
10	9	0.1	
		1	
Mean dist	ance 0.1		
rea of			
legion	50		
ensity	0.2		
xpected			
lean	1.118034		
NINI	0.089443		

	Nearest	
Point	Neighbor	Distance
1	3	2.2
2	4	2.2
3	4	2.2
4	5	2.2
5	7	2.2
6	7	2.2
7	8	2.2
8	9	2.2
9	10	2.2
10	9	2.2
		22
Mean dist Area of	ance 2.2	
Region	50	
Density	0.2	
Expected		
Mean	1.118034	
NNI	1.96774	

Ζ = 5.855

Source: Lembro

# Higher-ordered NNI

$$r_{exp} = 0.75\sqrt{A/n}$$

$$SE_{r} = 0.272\sqrt{A/n^{2}}$$

$$r_{exp}(k) = \gamma_{1}(k)\sqrt{A/n} = \frac{k(2k)!}{(2^{k}k!)^{2}}\sqrt{A/n}$$

$$SE_{r}(k) = \gamma_{2}(k)\sqrt{A/n^{2}}$$

$$r_{exp} = 0.5\sqrt{A/n} + (0.0514 + \frac{0.041}{\sqrt{n}})\frac{B}{n}$$

$$SE_{r}^{2} = 0.0683\frac{A}{n^{2}} + 0.037B\sqrt{\frac{A}{n^{5}}}$$
Edge adjustment

## Evaluating the Nearest Neighbor Index

- Advantages
  - NNI takes into account distance
  - No quadrat size problem to be concerned with
- However, NNI not as good as might appear
  - Index <u>highly</u> dependent on the boundary for the area
    - its size and its shape (perimeter)
  - Fundamentally based on <u>only</u> the mean distance
  - Doesn't incorporate local variations (could have clustering locally in some areas, but not overall)
  - Based on point location only and doesn't incorporate magnitude of phenomena at that point
- An "adjustment for edge effects" available but does not solve all the problems
- Some alternatives to the NNI are the G and F functions, based on the <u>entire frequency distribution</u> of nearest neighbor distances, and the K function based on <u>all</u> interpoint distances.
  - See O and U pp. 89-95 for more detail.
  - Note: the G Function and the General/Local G statistic (to be discussed later) are related but not identical to each other





























### Lecture 4: Spatial Autocorrelation (point)

Moran *I* Geary *C* Getis and Ord's *G* 























Name	Eine	Meaunt	NI4	NO	N/2	NI4	NE	NC	N/7	NO	a a
Alahama	1	A	28	13	12	47	CM	IND	N/	INO	Oueens Case
Arizona	4		35	9	40	6	32				Queens cuse
Arkansas		6	22	28	43	47	40	29			a a
California	6	3	4	32	40	4/	40	20			Sparse Configuity
Colorado	8	7	35	1	20	40	31	40	56		<u>spuise</u> contiguity
Connecticut	0	3	44	36	25	40	51	40	50		
Delaware	10	3	24	42	34						Matrix for US
District of Columbia	11	2	51	24							jer eks
Elorida	12	2	13	1							Ctates
Seorgia	13	5	12	45	37	1	47				Sidles
daho	16	6	32	41	56	49	30	53			
Ilinois	17	5	29	21	18	55	19				· Magunt is the
ndiana	18	4	26	21	17	39					
owa	19	6	29	31	17	55	27	46			
Kansas	20	4	40	29	31	8					number of
Kentucky	21	7	47	29	18	39	54	51	17		
Louisiana	22	3	28	48	5						
Maine	23	1	33								neighbors for each
Maryland	24	5	51	10	54	42	11				
Massachusetts	25	5	44	9	36	50	33				
Michigan	26	3	18	39	55						state
Vinnesota	27	4	19	55	46	38					State
Mississippi	28	4	22	5	1	47					$\mathbf{M}$ · $\mathbf{O}$ $\mathbf{M}$ ·
Missouri	29	8	5	40	17	21	47	20	19	31	•Max 1s 8 (Missouri
Montana	30	4	16	56	38	46					
Nebraska	31	6	29	20	8	19	56	46			and Tannasaaa)
Nevada	32	5	6	4	49	16	41				and Tennessee)
New Hampshire	33	3	25	23	50						
New Jersey	34	3	10	36	42						Sum of Moount in
New Mexico	35	5	48	40	8	4	49				•Sum of incount is
New York	36	5	34	9	42	50	25				
North Carolina	37	4	45	13	47	51					218
North Dakota	38	3	46	27	30						210
Ohio	39	5	26	21	54	42	18				
Oklahoma	40	6	5	35	48	29	20	8			•Number of
Oregon	41	4	6	32	16	53					
Pennsylvania	42	6	24	54	10	39	36	34			1 1
Rhode Island	44	2	25	9							common borders
South Carolina	45	2	13	37							common borders
South Dakota	46	6	56	27	19	31	38	30			$(\cdot \cdot \cdot)$
Tennessee	47	8	5	28	1	37	13	51	21	29	(101ns)
Texas	48	4	22	5	35	40					U =
Jtah	49	6	4	8	35	56	32	16			$\Sigma = $
/ermont	50	3	36	25	33						$\rightarrow$ ncount / 2 = 109
Virginia	51	6	47	37	24	54	11	21			
Washington	53	2	41	16							•N1 N2 FIPS codes
West Virginia	54	5	51	21	24	39	42				$-1$ $v_1$ , $1$ $v_2$ FIF S could s
Wisconsin	55	4	26	17	19	27					C · 11
Wyoming	56	6	49	16	31	8	46	30			tor neighbors

#### Weights Based on Distance

- Most common choice is the inverse (reciprocal) of the distance between locations i and j (w<sub>ij</sub> = 1/d<sub>ij</sub>)
  - Linear distance?
  - Distance through a network?
- Other functional forms may be equally valid, such as inverse of squared distance  $(w_{ij} = 1/d_{ij}^2)$ , or negative exponential  $(e^{-d} \text{ or } e^{-d^2})$
- Can use length of shared boundary: w<sub>ii</sub> = length (ij)/length(i)
- Inclusion of distance to all points may make it impossible to solve necessary equations, or may not make theoretical sense (effects may only be 'local')
  - Include distance to only the "nth" nearest neighbors
  - Include distances to locations only within a buffer distance
- For polygons, distances usually measured centroid to centroid, but
  - could be measured from perimeter of one to centroid of other
  - For irregular polygons, could be measured between the two closest boundary points (an adjustment is then necessary for contiguous polygons since distance for these would be zero)































Global indicator or "LISA," statistic.

### Calculating General G

Actual Value for G is given by:

$$G(d) = \frac{\sum \sum w_{ij}(d)x_ix_j}{\sum \sum x_ix_j}$$

Where: *d* is neighborhood distance W<sub>ij</sub> weights matrix has only 1 or 0 1 if j is within *d* distance of i

- 0 if its beyond that distance
- Expected value (if no concentration) for G is given by:

$$E(G) = \frac{W}{n(n-1)}$$
 where  $W = \sum_{i} \sum_{j} w_{ij}(d)$ 

- For the General G, the terms in the numerator (top) are calculated "within a distance bound (d)," and are then expressed relative to totals for the entire region under study.
  - As with all of these measures, if adjacent x terms are both large with the <u>same</u> sign (indicating positive spatial association), the numerator (top) will be large
  - If they are both large with different signs (indicating negative spatial association), the numerator (top) will again be large, but negative













