

Slide 1

Testing Network Hypotheses

Thursday Afternoon

Whether we are using Social Network Analysis as part of a consulting project or in support of academic research, it is important to know if the measures and relationships we see in the data represent a significant phenomenon, or are simply an artifact of the data themselves. But standard methods of calculating significance are inappropriate for network data. In this section we deal with testing hypotheses based on network data.

Objectives:

After completing this module you will be able to:

- Understand and explain why standard statistical methods are inappropriate for network data
- Conceptually explain how the permutation tests for significance work
- Explain the difference between monadic, dyadic, and mixed monadic/dyadic hypotheses
- Explain the concepts of variable and constant homophily
- Run the appropriate hypothesis testing techniques in UCINET

Slide 4

Statistical Issues

- Samples non-random
- Often work with populations
- Observations not independent
- Distributions unknown

Slide 5

Solutions

- Non-independence
 - Model the non-independence explicitly as in HLM
 - Assumes you know all sources of dependence
 - Permutation tests
- Non-random samples/populations
 - Permutation tests

Slide 6

Logic of Permutation Test

- Compute test statistic
 - e.g., correlation or difference in means
 - Correlation between centrality and salary is 0.384 or difference in mean centrality between the boys and the girls is 4.95.
 - Ask what are the chances of getting such a large correlation or such a large difference in means if the variables are actually completely independent?
- Wait! If the variables are independent, why would the correlation or difference in means be anything but zero?
 - Sampling
 - “Combinatorial chance”: if you flip coin 10 times, you expect 5 heads and 5 tails, but what you actually get could be quite different

Slide 7

Logic of Permutation Test

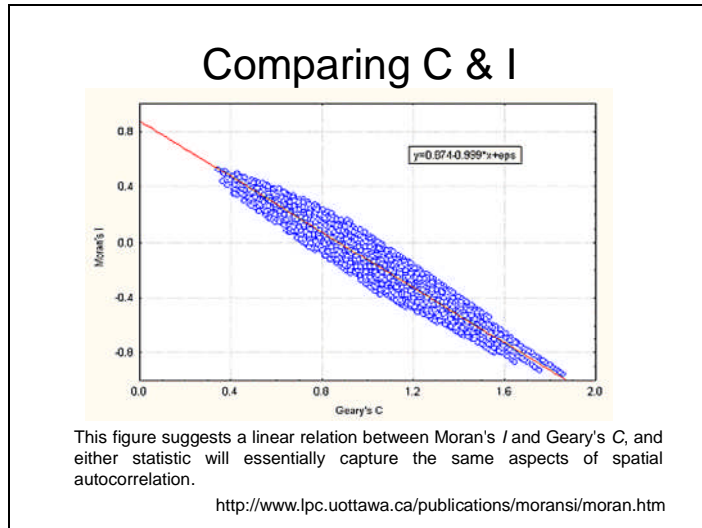
- So to evaluate an observed correlation between two variables of 0.384, we want to
 - correlate thousands of variables similar to the ones we are testing that we know are truly independent of each other, and
 - see how often these independent variables are correlated at a level as large as 0.384
 - The proportion of random correlations as large observed value is the p-value of the test
- How to obtain thousands of independent variables whose values are assigned independently of each other?
 - Fill them with random values
 - But need to match distribution of values
 - Permute values of one with respect to the other

Slide 8

Outline of Permutation Test

- Get observed test statistic
- Construct a distribution of test statistics under null hypothesis
 - Thousands of permutations of actual data
- Count proportion of statistics on permuted data that are as large as the observed
 - This is the p-value of the test

Slide 17



Slide 18

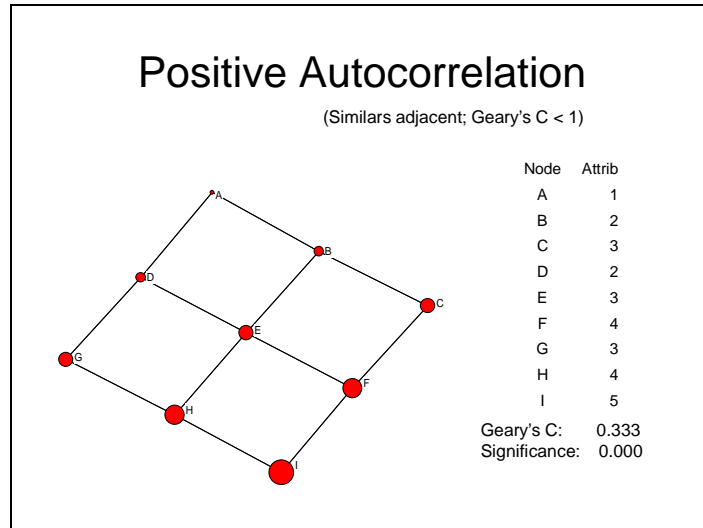
Geary's C

- Let $w_{ij} > 0$ indicate adjacency of nodes i and j , and X_i indicate the score of node i on attribute X (e.g., age)

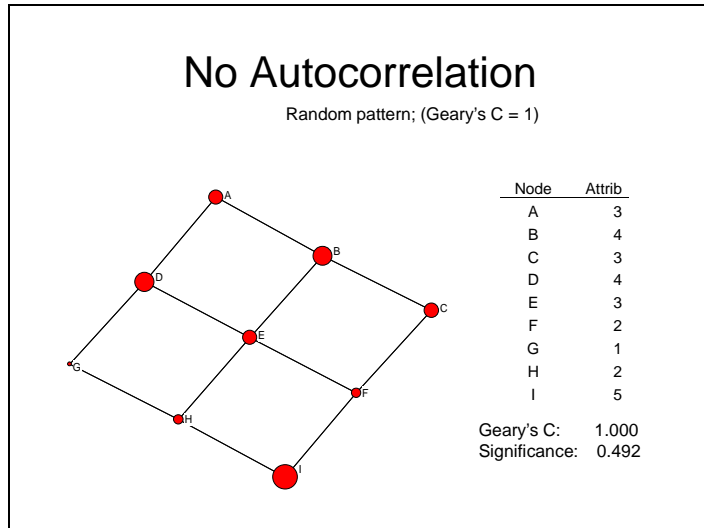
$$C = (n-1) \frac{\sum_i \sum_j w_{ij} (x_i - x_j)^2}{2 \sum_{i,j} w_{ij} \sum_i (x_i - \bar{x})^2}$$

- Range of values: $0 \leq C \leq 2$
 - $C=1$ indicates independence;
 - $C > 1$ indicates negative autocorrelation;
 - $C < 1$ indicates positive autocorrelation (homophily)

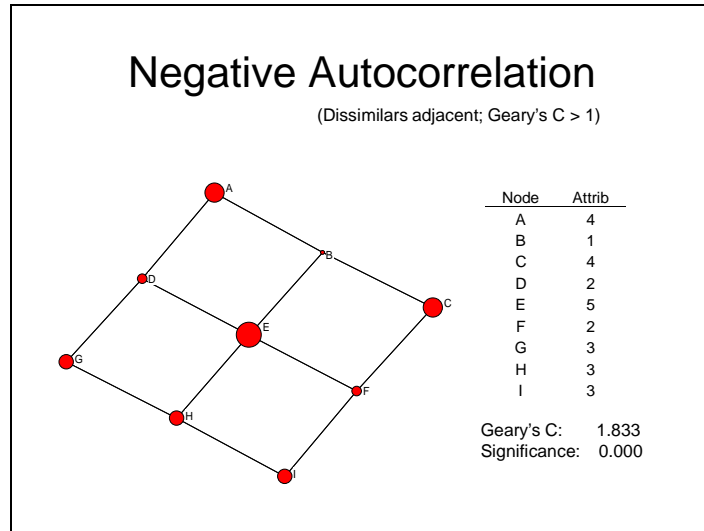
Slide 19



Slide 20



Slide 21



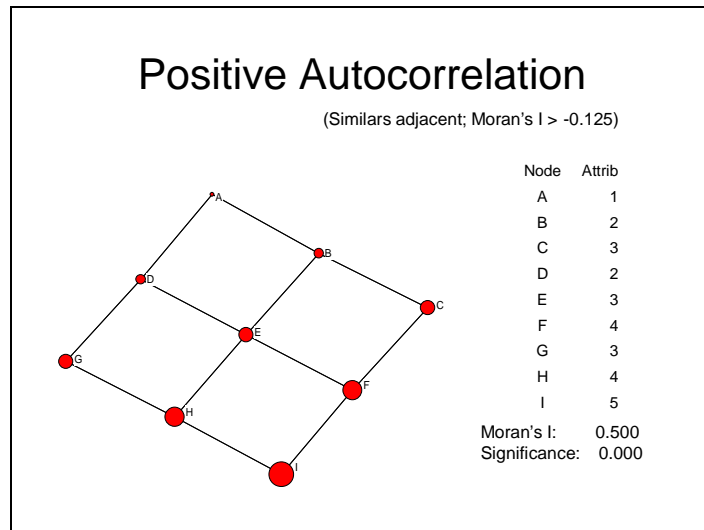
Slide 22

Moran's I

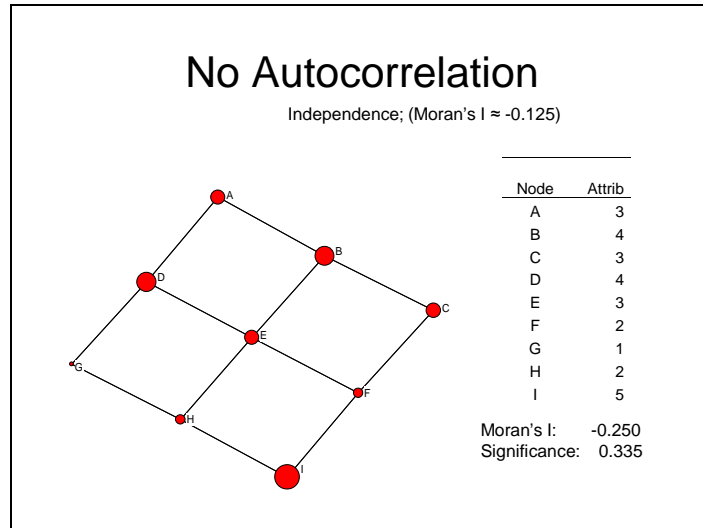
- Ranges between -1 and +1
- Expected value under independence is $-1/(n-1)$
- $I \rightarrow +1$ when positive autocorrelation
- $I \rightarrow -1$ when negative autocorrelation

$$I = n \frac{\sum_{i,j} w_{ij}(x_i - \bar{x})(x_j - \bar{x})}{\sum_{i,j} w_{ij} \sum_i (x_i - \bar{x})^2}$$

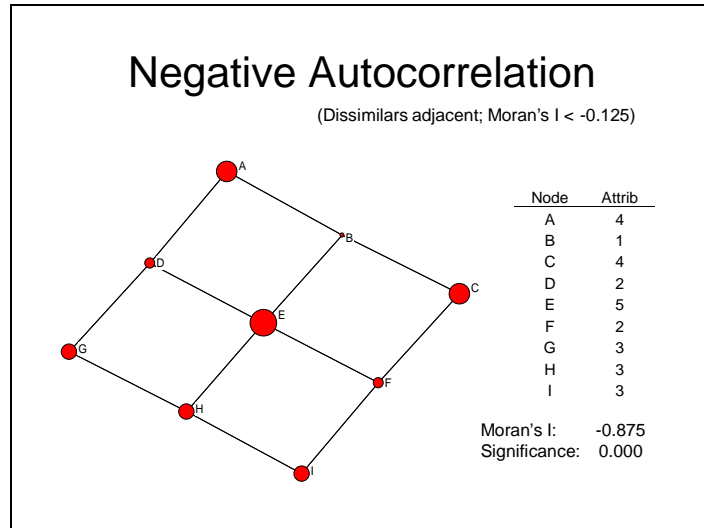
Slide 23



Slide 24



Slide 25



Interpreting Autocorrelation

- With Moran's I
 - A value near +1.0 indicates **clustering** (adjacency tends to accompany similarity along a dimension)
 - A value near -1.0 indicates **dispersion** (adjacency tends to accompany dissimilarity along a dimension)
 - a value near 0 indicates **random** distribution
- For Geary's C
 - just substitute 0, 2, and 1 for 1, -1, and 0 above

Slide 27

Another Approach

- Convert the attribute vector into a matrix
 - Use Data | Attribute to Matrix in UCINET
- QAP this new matrix against the adjacency matrix
 - Significances will be the same because it uses same underlying permutation method
 - Values will follow same pattern (but not same values) as Moran's I

Using QAP for Autocorrelation

Gender		HOL	BRA	CAR	PAM	PAT	JEN	PAU	ANN	MIC	BIL	LEE	DON	JOH	HAR	GER	STE	BER	RUS
HOLLY	1	HOLLY	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
BRAZEY	1	BRAZEY	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
CAROL	1	CAROL	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
PAM	1	PAM	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
PAT	1	PAT	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
JENNIE	1	JENNIE	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
PAULINE	1	PAULINE	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
ANN	1	ANN	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
MICHAEL	1	MICHAEL	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
BILL	2	BILL	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
BILL	2	LEE	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
LEE	2	DON	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
DON	2	JOHN	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
JOHN	2	HARRY	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
HARRY	2	GERY	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
GERY	2	STEVE	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
HARRY	2	BERT	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
GERY	2	RUSS	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
STEVE	2																		
BERT	2																		
RUSS	2																		

This matrix was constructed based on "exact match"
but you can use different transformations

Slide 29

Comparing QAP & Moran's I

Moran's I Output

A value of -0.059 indicates perfect independence. Autocorrelation: 0.667
Significance: 0.001

QAP Output

Independent	Un-stdized Coefficient	Stdized Coefficient	Significance
Intercept	0.056250	0.000000	0.999
CAMPATTR2-MAT	0.251969	0.330131	0.001

Hypothesis Testing Lab

For this lab we will use four datasets:

CAMPNET:

This is a dichotomous adjacency matrix of 18 participants in a qualitative methods class. Ties are directed and represent that the ego indicated that the nominated alter was one of the three people with which s/he spent the most time during the seminar.

ZACKAR & ZACHATTR:

ZACKAR is another stacked dataset, containing a dichotomous adjacency matrix, ZACHE, which represents the simple presence or absence of ties between members of a Karate Club, and ZACHC, which contains valued data counting the number of interactions between actors. ZACHATTR is a rectangular matrix with three columns of attributes for each of the actors from the ZACKAR datasets.

KRACK-HIGH-TEC & HIGH-TEC-ATTRIBUTES:

KRACK-HIGH-TEC is another stacked dataset, containing three dichotomous relations (REPORTS_TO, ADVICE, FRIENDSHIP). HIGH-TEC-ATTRIBUTES contains several attributes about the nodes in KRACK-HIGH-TEC, including Age, Level (CEO, Manager, Staff), Tenure, and Department.

WIRING:

This is a stacked dataset that includes many different files. This is a dichotomous adjacency matrix of 14 employees of the bank wiring room of Western Electric used in the famous Hawthorne Studies. Ties are symmetric and represent participation in games during work breaks. RDGAM records people playing games together, RDCON records conflict between people, RDPOS is positive interactions, RDCON is negative interactions.

2009 LINKS Center Summer SNA Workshop

1) Testing dyadic hypothesis

- a. Run Data | Unpack on ZACKAR (if you have not yet), which will create ZACHE and ZACHC. ZACHE has dichotomous data about the ties and ZACHC has valued data (the strength of ties, a count).
- b. Run Tools | Similarities and use the cross-product measure to compute similarities of ZACHE. (The cross product is a very powerful and common matrix operation that, in this case, will count how many friends each pair of actors have in common.) Call the output FOF (Friends of Friends).
- c. Go to Tools | Testing Hypotheses | Dyadic (QAP) | QAP Correlation and browse to include both ZACHC and FOF to be correlated and click okay. What do the results mean?
- d. Congratulations, you have just statistically demonstrated the first part of Granovetter's famous "strength of weak ties" theory, which states that I have stronger ties (ZACHC) with those people with whom I share more friends in common (FOF).

2) Testing multivariate dyadic hypotheses

- a. Unpack the WIRING dataset if you have not done so yet.
- b. Go to Tools | Testing Hypotheses | Dyadic (QAP) | QAP Regression | Full Partialling. Put RDCON (conflict between members about whether the windows should be open or shut) in as the dependent variable. Put in RDPOS (positive relationships), RDNEG (negative relationships), and RDGAM (playing games together) in as independent variables. Before running it, what do you think would most significantly predict conflict? After running it, are your results what you expected? How would you explain the results?
- c. Record the standardized coefficient and significance for any significant predictor, and run the same procedure two more times (still using the default value of 2000 for the number of permutations) and record the same results. Now, run the same procedure three more times setting the number of random permutations set to 50000. Record the same results. How did

2009 LINKS Center Summer SNA Workshop

the parameter affect the results? Why? When running with 2000 permutations, why did one number change but the other remain constant?

- d. Now run Tools | Testing Hypotheses | Dyadic (QAP) | QAP Regression | Double Dekker (MRQAP) still using the same independent and dependent variables, and setting your permutations to 50000. Compare these results with your previous run. Compare the time required to run it.

3) Testing monadic hypotheses.

- a. You should have already unpacked the KRACK-HIGH-TEC dataset, but if not, do so now. You will get three datasets (REPORTS_TO, ADVICE, FRIENDSHIP). We are going to use the ADVICE dataset. Run Network | Centrality | Degree on this dataset, using the directed version, telling it NOT to treat the data as symmetric, and calling your output ADVISING. Record which column has InDegree centrality. This is a measure of how many people said they sought advice from each person.
- b. Display (D) the HIGH-TEC-ATTRIBUTES dataset to determine which columns the AGE and TENURE attributes are in.
- c. Now, it is common wisdom that people look to the “senior” people for advice, but is unclear in an organizational context whether senior is “older than” or “longer tenured than”. You will test if either of these is supported by the data. Run Tools | Testing Hypotheses | Node-Level | Regression specifying ADVISING for your dependent dataset and the appropriate column, and HIGH-TEC-ATTRIBUTES for your independent dataset and the appropriate columns, and set the number of permutations to 10000. Which meaning of “senior” does the data support?
- d. Why did we use the Regression option of Node-Level instead of T-Test or Anova? When would we use those?

4) Testing Mixed-Dyadic Monadic hypotheses

2009 LINKS Center Summer SNA Workshop

- a. Since it is only fitting that we end where we started, we shall use the campnet data for these final exercise.
 - b. You will run Tools | Testing Hypotheses | Mixed Dyadic/Nodal | Categorical attributes | Anova Density twice. For both, specify CAMPNET as the network matrix, and the gender column of the CAMPATTR2 matrix as the Actor Attribute. For the first run, choose “Constant Homophily” for your model, and for the second, choose “Variable Homophily”. Interpret both sets of results. What do they mean? Is there homophily? Who tends to be more homophilous?
- 5) Using QAP for Mixed Monadic/Dyadic Hypotheses testing.
- a. Using Data | Attribute to matrix, create a matrix of exact matches among the actors in Campnet based on gender.
 - b. View this new matrix (named CAMPATTR2-MAT by default) in Netdraw. What does the diagram show?
 - c. Use Tools | Testing Hypotheses | QAP Regression to regress the Campnet network on this new matrix of gender similarity. What do the results show?
 - d. Do you prefer this approach, ANOVA Density Tables, Moran’s I, or Geary’s C? When might you use each of these separate techniques?