#### Empirical Research Methods in Information Science

# <u>IS 4800 / CS6350</u>

#### Lecture 14

### Outline

- Reading assessment
- Schedule update
- Finish correlation
- Inferential statistics and power

الاياللارد المساليريان

#### Upcoming

- Assignment I5 due (Friday, EOD)
- Enjoy your spring break!
- Tuesday class review, team up (T1)
- Wednesday class midterm exam
- T1 assignment deadline will be adjusted back a bit ... updates on Piazza



### Comparing r's

 If you want to make statements about how large one correlation is relative to another.

- e.g. one is twice as large as another
- Don't compare r's directly...
- Compare r^2 ("proportionate reduction in error")

## Other measures of association

- Point-biserial [r<sub>pb</sub>]
  - One numeric & one binary (nominal) measure
  - Just dummy code the nominal (0 and 1) and use Pearson correlation.
- Spearman rank correlation coefficient (Spearman's rho) [ρ or r<sub>s</sub>]
  - Two ordinal measures (or for transformed numeric measures if non-linear)
  - Replace each value with its rank order
  - Compute Pearson correlation with ranks
  - Measures degree of monotonicity

### Sample Data for Spearman rho

Monitor Size	21	24	17	19	15
Productivity	3	1	2	5	0



7





## Two meanings of 'correlation'

Correlation <u>statistic</u> vs.

Correlational <u>research model</u>

#### Example: Net Latency & Satisfaction



- NU ITS wants to save money by switching to slower wireless routers, and wants to assess the impact this will have on student satisfaction. You want to know how slow things have to get before students start complaining.
- You have the crew implement a program that randomly chooses a network latency (between 0s and 10s) every time a student logs into NUwave, then adds that latency to every network access from them. After 10 minutes of use a web form pops up asking students to rate their degree of satisfaction with NU ITS (10-item).
- What kind of study is this?
- What statistic would you use to evaluate results?
- Assuming r = -0.8, p=.021, what would you conclude?
- Assuming r = 0.1, p=.342, what would you conclude?



What do you do if your data is clearly not unimodal & symmetric OR there is a clear non-linear relationship?

### Parametric vs. Non-parametric Statistics

- Non-parametric statistics that do not rely on data belonging to any particular distribution
- E.g., Pearson r is a parametric statistic (assumes underlying distributions are normal – can be described using parameters – mean & stddev)
- E.g., Spearman rho is non-parametric



### Power

The "power" of a statistical test is its ability to detect differences in data that are inconsistent with the null hypothesis.

Low South Low South Low Street

p(rejecting H0|H1)

- aka the ability to find a significant result, if your hypotheses are actually true.
- What is it called when this fails (ie, accepting H0 when H1 is true)?
- Why is this a bad situation?

#### Effect size

- والمرجبين الرجسيان وحداثيه الرجد اليجابلي الرجسيان الرجسيان الرجسياني الرجد اليجابلي الرجبين الرجسيان الرجسيان
  - The *amount* of measured difference between study conditions.
  - The greater the effect size, the easier it is to show there is a significant difference in your study (ie, the greater the power).
  - Effect size formula is different for each hypothesis test procedure.
  - Tabulated standard values for "small", "medium", and "large" effect sizes.
  - Only talk about effect size IF significance is established – but then DO present it in your results.





Relationship between alpha, beta, and power. Company and a second for the second for the second s What is the probability of each of these situations occurring? "The Truth" H1 True H1 False

Decide to Reject H( & accept H1

Do not Reject H0 & do not accept H1

0	Correct p = power	Type I err p = $\alpha$
	Type II err $p = \beta$	Correct $p = 1-\alpha$

#### Relationship between power and effect size





#### **Power Analysis**

- Should determine number of subjects you need ahead of time by doing a 'power analysis'
- Standard procedure (part of your study plan):
  - Determine statistic you will use
  - Fix alpha and beta (1-power) (and number of tails if appropriate)
  - Estimate expected effect size from prior studies
  - Then: Determine number of subjects you need
- Note: Power
  - Increases with effect size
  - Increases with sample size
  - Decreases with decreasing alpha

Power analyses are different depending on the statistical test you are using...

First up: t-test for independent means.



 $d = \frac{(\mu_1 - \mu_2)}{(\mu_1 - \mu_2)}$ 

Parameters for population of <u>individuals</u>. (so, use SD-pooled for t-test of indep means)

Cohen: d~0.2 small d~0.5 medium d~0.8 large

#### Power table



TABLE 8–4 Approx for Inde at the .	imate Power fo ependent Mean 05 Significance	r Studies Using the s Testing Hypothe Level	e <i>t</i> Test ses	
humber of Participants	Effect Size			
in Each Group	Small (.20) Medium (.50)		Large (.80)	
One-tailed test				
10	.11	.29	.53	
20	.15	.46	.80	
30	.19	.61	.92	
40	.22	.72	.97	
50	.26	.80	.99	
100	.41	.97	*	
Two-tailed test				
10	.07	.18	.39	
20	.09	.33	.69	
30	.12	.47	.86	
40	.14	.60	.94	
50	.17	.70	.98	
100	.29	.94	*	

#### More Useful and Concise (for practical purposes use a power calculator)

TABLE 8–5Approximate Number of Participants Needed in Each<br/>Group (Assuming Equal Sample Sizes) for 80% Power<br/>for the t Test for Independent Means, Testing<br/>Hypotheses at the .05 Significance Level

	Effect Size		
	Small (.20)	Medium (.50)	Large (.80)
One-tailed	310	50	20
Two-tailed	393	64	26

### Power Analysis Exercise

- Based on related research, we expect that there will be a medium effect size in our study of an LED sign in the Food Court affecting wait times.
- How many subjects do we need for a two-group, two-tailed test at α=0.05, 80% power?

#### More Useful and Concise (for practical purposes use a power calculator)

TABLE 8–5Approximate Number of Participants Needed in Each<br/>Group (Assuming Equal Sample Sizes) for 80% Power<br/>for the t Test for Independent Means, Testing<br/>Hypotheses at the .05 Significance Level

	Effect Size		
	Small (.20)	Medium (.50)	Large (.80)
One-tailed	310	50	20
Two-tailed	393	64	26

#### But, I can't study 786 subjects!

- Increase effect size
  - Increase difference in population means (change manipulation)
  - Decrease population variance (better measures, control more extraneous vars)
  - Redesign study to collect many trials of measures per subject
- Relax criteria for Type I error
  - Increase  $\alpha$  threshold
  - Change from Two-tailed => one-tailed test
  - Decreases credibility of your findings
- Decrease power
  - Decreases likelihood of getting a significant result
- Use a different statistic
  - If possible, maybe consult a statistician
- Practically
  - Usually, redesign experiment so that we have increased effect size or better measures for decreased variance
  - OR, call it a "pilot study"

Interpreting results: Significance & Effect Size

- Significance
  - Just indicates that it is likely there is a nonzero difference between populations.
  - Says nothing about how big the difference is.
- Effect Size
  - Only meaningful if result is significant.
  - Indicates how big the difference is (usually normalized to number of std-deviations)

### Interpreting results: Significance & Effect Size

- Significant & small effect => ?
  - Real difference, but slight.
  - Probably not of practical importance.
- Significant & large effect => ?
  - Real difference, likely meaningful.
- Significant & small sample => ?
  - Significant & possibly important.
- Non-significant & small sample => ?
  - Inconclusive
- Non-significant & large sample => ?
  - Evidence there really is no difference

#### **Group Exercise**

- white the second for the second for
  - Compute effect size for the study
  - Characterize as small/medium/large
  - You are now going to do a follow-up study using similar interventions and measures (ie assume same effect size)
  - Do a power analysis to determine how many subjects you would need for a two-group between-subjects experiment with 80% power, alpha=0.05, two-tailed test.

#### Power & Effect Size for Correlation

- Effect size = |r|
- Power, see table 11-7, pg 465 Aron
  - Usually, given
    - Expected effect size
    - Test criteria
      - Desired significance level (usually 0.05)
      - Desired power (usually 0.8)
      - Directionality of test

			Effect Size	
		Small (r = .10)	Medium $(r \approx .30)$	Large (r = .50)
wo-tailed				
fotal N:	10	.06	.13	.33
	20	.07	.25	.64
•	30	.08	.37	.83
	40	.09	.48	.92
	50		.57	.97
	100	.17	.86	
One-tailed				
Total N:	10	.08	.22	.46
	20	.11	.37	.75
	30	.13	.50	.90
	40	.15	.60	.96
	50	.17	.69	.98
	100	90	92	

and a second second



Approximate number of participants needed for 80% power for a study using the correlation coefficient (r) for testing a hypothesis at the .05 significance level

Effect size			
Small (r=0.1)	Medium (r=0.3)	Large (r=0.5)	
783	85	28	

### Effect Size & Power for X<sup>2</sup> test for independence

- Completely different formulas than for Pearson r or t-test.
- Dependent on df.
- For 2x2, effect size = "phi"


## Effect Size & Power for X<sup>2</sup>



 Table 13-10
 Approximate Total Number of Participants Needed for 80%

 Power for the Chi-Square Test for Independence for Testing

 Hypotheses at the .05 Significance Level

		Effect Size				
Total df		Small	Medium	Large		
+		785	87		26	
2		964	107	· .	.39	
3	· · .	1,090	121		- 44	
4	•	1,194	133		48	

		Effect Size		
Total <i>df</i>	Total N	Small	Medium	Large
1	25	.08	.32	.70
	50	11	.56	.94
	100	.17	.85	*
	200	.29	.99	
2	25	.07	.25	.60
	50	.09	.46	.90
	100	.13	.77	
	200	.23	.97	*
3	25	.07	.21	.54
	50	.08	.40	.86
	100	.12	.71	.99
	200	.19	.96	
4	25	.06	.19	.50
	50	.08	.36	.82
	100	.11	.66	• .99
	200	.17	.94	1997 - <b>1</b> 997 - <b>1</b> 99

.



## Computing effect size

- Some thoughtless authors do not include means & stddevs (per group) in their article...
- R package 'compute.es' contains a variety of methods for computing effect size given other info (e.g., t score, N1, N2)
- Morale: Always include means & stddevs
- Better: Report effect sizes yourself!









#### t-test for independent means

- Tests association between binomial IV and numeric DV.
- Examples:
  - WizziWord vs. Word => wpm
  - Small vs. Large Monitors => wpd
  - Wait time sign vs. none => satisfaction

#### Understanding numeric measures

- Sources of variance
  - IV
  - Other uncontrolled factors ("error variance")

#### Example: Call of Duty vs. Halo





- What variables might affect Satisfaction?
- Typically, one subject's Satisfaction score = TrueSatisfaction + var1 + var2 + var3 + ...
- A sum of random variables.

### Central Limit Theorem

 If (many) independent, random variables with the same distribution are added, the result is approximately a normal curve



Suppose random variable X has distribution

$$X = \begin{cases} 1 & \text{with probability } 1/3, \\ 2 & \text{with probability } 1/3, \\ 3 & \text{with probability } 1/3. \end{cases}$$



0	0	0
1	2	3

From wikipedia

#### Why be normal? Example

Now, consider the distribution of X+X

$$\begin{cases} 1+1 = 2\\ 1+2 = 3\\ 1+3 = 4\\ 2+1 = 3\\ 2+2 = 4\\ 2+3 = 5\\ 3+1 = 4\\ 3+2 = 5\\ 3+3 = 6 \end{cases} = \begin{cases} 2 & \text{with probability } 1/9\\ 3 & \text{with probability } 2/9\\ 4 & \text{with probability } 3/9\\ 5 & \text{with probability } 2/9\\ 6 & \text{with probability } 1/9 \end{cases}$$

## Why be normal? Example

Now, consider the distribution of X+X+X



#### Central Limit Theorem

The sum of independent random variables with the same underlying distribution will be approximately normal



# Estimating probabilities using the normal curve

- p(x>sample>y)
  - Ex: know age of NU students is M=25, SD=3
  - What's the likelihood of a random student being between 25-28?



#### Given what we know so far.. We can already do some hypothesis testing!



#### Example

- You have 100 admins in your company.
- They all use Word.
- You want to consider changing to WizziWord.
- Hypothesize it will increase their net productivity, measured as word per minute typed, averaged over an entire day.

























#### Don't try this at home

- You would never do a study this way.
- Why?
  - Can't control extraneous variables through randomization.
  - Usually don't know population statistics.
  - Can't generalize from an individual.





Single sample test, with sample size>1, and known comparison population...





Comparison Distribution = Distribution of Means

- N = 10
- $\mu = 149, \sigma = 12, \sigma^2 = 144$
- $\mu_{M} = 149$ ,  $\sigma_{M}^{2} = 144/10 = 14.4$
- $Z = (170 149) / \sqrt{14.4} = 5.5$


## Don't try this at home

- You would never do a study this way.
- Why?
  - Can't control extraneous variables through randomization.
  - Usually don't know population statistics.

Can't generalize from an individual.