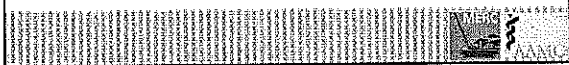


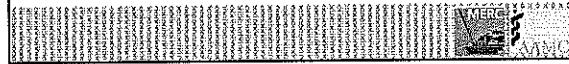
And of course.....

You should consult with a statistician
or research design specialist
while you are *conceptualizing*
the study!



Consulting with a statistician

About what??????



What do you need to ask about?

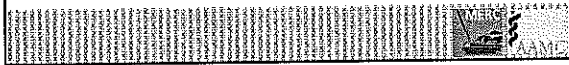
Addressing statistical issues

Comparing sample to population



The Statistical Consultation

Addressing Statistical Issues



The Statistical Consultation

Addressing Statistical Issues

- Statistical significance
- Type I error
- Type II error
- Power
- Effect size
- Appropriate analyses



Statistical Significance

- What level of uncertainty are you willing to tolerate in making statistical decisions?
- What degree of risk are you willing to take of rejecting a null hypothesis that is true?
Called "level of significance" or α



Type I Error

Making an incorrect decision = a *type I error*
(rejecting a null hypothesis that is true)

Set value of α

- 0.05 typical
- 0.01 more conservative
- 0.10 more liberal



Type II Error

- Failing to reject a null hypothesis that is false
- Protecting your study against a Type II error
 - Increase the sample size
 - Look at a *meaningful* effect
 - Increase study's alpha




Power

- Power: the probability of identifying a true hypothesis as true (given the structure of your study)
- Power = (1- Type II error)
- Design studies to have a power of 0.8



Maximizing Power


- Design a meaningful intervention
- Increase sample size
- Limit number of variables in the study



Statistical Significance is Not Everything


Large samples
Very small changes or correlations can be statistically significant but essentially without meaning.

Small samples
Important (meaningful) changes or correlations may not be statistically significant.




Statistical Significance

- How important is statistical significance?
- How do you decide if a finding is significant?
- How do you take chance into account in statistical decision making?
- Can something be statistically significant but not meaningful?
- Can statistics solve all problems?




Effect Size Primer
Effect Size (ES)

- A family of indices that measure the treatment effect's magnitude
- Independent of sample size




Effect Size Primer
Two Simple Measures of Effect Size

- The standardized difference between two means
Cohen's $d = [(M1 - M2) / \text{pooled std dev}]$
- The correlation coefficient
- Rough interpretation guidelines
 - 0.2 - small
 - 0.5 - medium
 - 0.8 - large



The Statistical Consultation

Comparing Sample to Population



Who is in your sample?

- Is the sample representative of the population?
 - Single sample test of proportions
 - Percent female representative
 - Percent African American
 - Single sample test of means
 - Sample MCATs equal to population
 - Sample age equal to population
- Expect results to be statistically non-significant



...and who is not?

- Follow the non-respondents
- Compare respondents to non-respondents
 - Same age (t-test)
 - Same genders (chi-square)
 - Same geographic region (chi-square)



And a Reminder!

When you consult with a statistician while you are *conceptualizing* the study, talk about appropriate analyses for your study data.



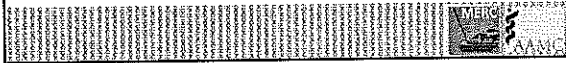
Study 1: Teaching Method Comparison

After midterm, students randomly assigned to

- Classroom lecture plus slides
- Team learning
- CAI

Measures

- Midterm exam
- Final exam



Studies 2 & 3: What Kind of Analysis?

Data available: MCAT total scores, gender, ethnicity, college GPA, % grades for M1 and M2 years, USMLE Step 1 scores

- Study 2: What predicts USMLE score?
- Study 3: What distinguishes students who drop out of medical school from those who do not?



Collect Data



Collect Data Paper Instruments

Advantages

- Good control on distribution and return of instruments
- Responders can indicate where they are unclear
- Anonymity
- Can do double data entry for data verification

Disadvantages

- Entering data takes time and chance for error
- Forms need to be hand reviewed prior to entry
- Comments take time to enter
- Entry expensive
- Anonymity

Web-Based Forms

Advantages

- No data entry
- Easy distribution
- Easy completion

Disadvantages

- Physical appearance has limited appeal
- Requires internet access
- Data enterer cannot be verified
- May be hard to correct entry errors
- Still takes file preparation for analysis

Scannable Forms

Advantages

- Efficient data entry
- Forces forms to be clean and clear
- Can be produced easily
- Software can be relatively inexpensive

Disadvantages

- Forms must be reviewed prior to scanning
- Digital scan requires error checking
- Need to review for scanning problems
- Software can be expensive

Direct Entry

Direct data entry into a computer by someone other than the respondent

- Phone survey
- Consumer survey
- Electronic health record

Set up Data Files

Variable Names

Use short, meaningful name

Add a variable label with more information

Include

- An ID variable (even for anonymous data)
- Grouping variable (if applicable)

Data Naming Example – Study 4

You ask residents in a 3-year training program to

- 1) rate their understanding of 7 categories of complementary or alternative medicine (CAM)
- 2) indicate their gender and age
- 3) write how they describe their ethnicity

- How many variables are there?
- What are their names?

Variable Formats

Variable formats

- Numeric
- String

From previous example, which could be

- numeric?
- string?

Variable Values

- Indicates code for possible responses for a variable:

0 = no

1 = yes

- Missing data codes

How would you code the previous example?

Develop a Codebook

- Variable names, labels and values are listed in a codebook
- Develop codebook before data collection
 - Can be done directly on data collection form
- Codebooks can be generated by computer programs

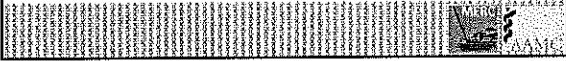
Exercise - Practice Codebook Development

- Using the CAM survey, develop a codebook for the data
- Transfer onto codebook template for your records
 - Variable name
 - Type – string or numeric
 - Label
 - Values – possible values and their meaning

Codebook Template

Variable Name	Type	Label	Values

Enter Data




Software Options

Data entry

- Excel
- Access
- Directly into a statistical package


Data analysis

- Major, common packages: SPSS, SAS
- Less well known: STATA, R



Data Entry

- One line per subject
- Build in data range restrictions
- Double data entry



Merging Data

You can bring together data from two files.

- One variable must have the same name and format in both files.
- Files must be sorted in order by that variable.



Check and Clean Data



Cleaning the Data

- Compute frequencies for each variable
 - Check maximum and minimum values to make sure none are out of range
- Clean/fix missing values
- Compute descriptive statistics for each variable
 - Make sure means and standard deviations make sense



Understanding the Data

For continuously scaled variables

- Construct histograms
- Check skewness



Understanding the Data

For bivariate relationships

- Construct scatterplots
- Examine for bivariate outliers
- Examine for linearity/nonlinearity



Practice Looking at Data

- Review descriptive statistics handout
- Review scatterplots



Example – Study 5

This study had, among others, two variables whose relationship was of interest:

- Scale score composed of combined survey items labeled "understanding others' value"
- A post-intervention score of knowledge labeled "post score"

What to do

- Enter data
- Run frequencies
- Look at distributions
- Start to look at inter-relationships of variables



Statistics from Study 5

Descriptive Statistics

	N		Minimum		Maximum		Mean		Std.		Skewness	
	Statistic	Statistic	Statistic	Statistic	Statistic	Statistic	Statistic	Statistic	Statistic	Statistic	Std. Error	Statistic
Understanding others' value	442	12	72	61.63	8.802	-.497	.116					
Post score	37	7	36	16.19	4.932	1.723	.388					
Valid N (listwise)	36											

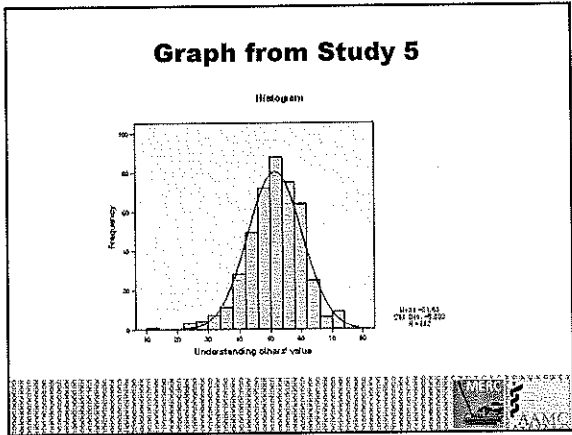


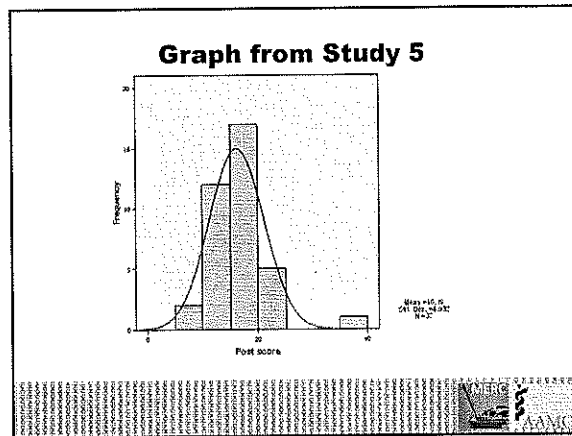
Statistics from Study 5

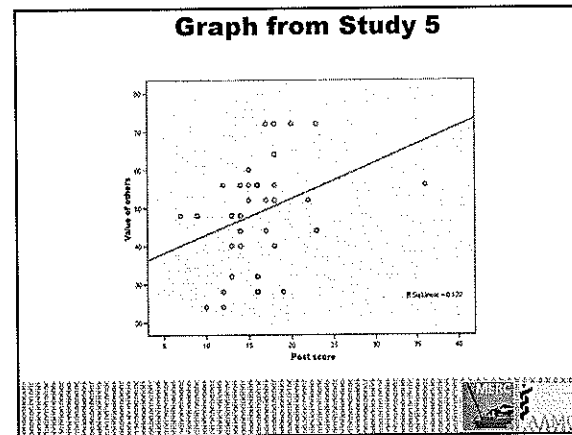
Post score

Valid	Frequency	Percent	Valid Percent	Cumulative Percent
7	1	.2	2.7	2.7
9	1	.2	2.7	5.4
10	1	.2	2.7	8.1
12	3	.7	8.1	16.2
13	4	.9	10.8	27.0
14	4	.9	10.8	37.8
16	4	.9	10.8	48.6
16	4	.9	10.8	59.4
17	3	.7	8.1	67.5
19	6	1.1	16.2	83.7
19	1	.2	2.7	86.4
20	1	.2	2.7	89.1
21	1	.2	2.7	91.8
22	1	.2	2.7	94.5
23	2	.4	5.4	99.9
30	1	.2	2.7	100.0
Total	37	8.2	100.0	
Missing System	412	81.8		
Total	449	100.0		









Now it is your turn

Example 6 - Study set-up

This study was a randomized trial of Oral Case Presentation feedback. About 160 students were randomized to 2 groups – usual feedback and structured intentional feedback.

1. Medical students randomized to the treatment (about 20 a block) were given OCP booklets and asked to gather 1 feedback form a week
2. Data for each student were linked to other data, (e.g., ratings of clinical competence, shelf exam scores)

DO STUDENTS WHO GET STRUCTURED FEEDBACK DO BETTER ON END-OF-CLERKSHIP PRESENTATION?

Your job

1. Make a code book
2. Comb through the data

What to do

- Enter data
- Run frequencies
- Look at distributions
- Link files
- Start to look at inter-relationships of variables

What values do we see?

OCP_M				
	Frequency	Percent	Valid Percent	Cumulative Percent
Valid	1	1.4	1.4	1.4
2	34	4.6	4.6	5.9
3	172	23.2	23.2	29.1
4	294	40.1	40.1	69.2
5	204	27.8	27.8	97.0
6	1	0.1	0.1	97.1
7	134	18.3	18.3	100.0
Total	600	100.0	100.0	
Missing				
System				
Total	600	100.0		



What do we observe?

Time to Observe				
	Frequency	Percent	Valid Percent	Cumulative Percent
Valid	1	0.1	0.1	0.1
2	2	0.3	0.3	0.4
3	4	0.6	0.6	1.0
4	45	7.3	7.3	8.3
5	1	0.2	0.2	8.5
6	2	0.3	0.3	8.8
7	223	37.2	37.2	46.0
8	2	0.3	0.3	46.3
Total	268	100.0	100.0	
Missing				
System				
Total	268	100.0		

CHECK				
	Frequency	Percent	Valid Percent	Cumulative Percent
Valid	563	100.0	100.0	100.0



Move on to distributions...

	N	Descriptive Statistics					
		Range	Minimum	Maximum	Mean	Std. Deviation	Skewness
OCP_M	600	6.00	1.00	5.00	3.7745	1.3668	-.421
OCP_Healthcare	745	5.00	4.00	5.00	2.7345	1.0254	-.340
OCP_SocFam Hx	800	5.00	4.00	5.00	3.7822	1.2671	-.421
OCP_PE	445	4.00	3.00	4.00	2.9776	1.4373	-.713
OCP_Student	427	5.00	4.00	5.00	2.6373	1.0227	-.341
OCP_Accessibility	562	5.00	4.00	5.00	2.6187	1.0316	-.351
OCP_Pwr	340	5.00	4.00	5.00	2.5055	1.0561	-.381
OCP_Organization	362	5.00	4.00	5.00	2.4354	1.0911	-.384
OCP_Stat Space Alloc	381	15.00	3.00	20.00	7.2802	1.9738	2.073
Valid N (listwise)	349						



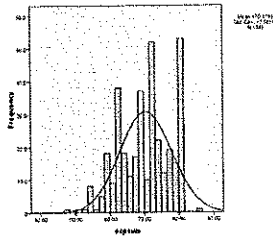
Start to put all of the items together..

Yield	N	%
Yield	319	51.4
Failure	319	24.6
Total	638	100.0

Component A/B's	N of Items
501	5



And then we look at distributions..



And then we aggregate to other scores, etc.

	N	Minimum	Maximum	Mean	Std. Deviation
Yield	319	4.87	8.65	6.9077	.45216
Failure	319	3.35	7.00	6.2545	.33142
Yield+Yield	638	4.18	7.00	6.3121	.36531
Yield+Yield	638	3.00	7.00	6.4214	.26915
Yield	319	87.00	90.00	78.1875	0.84693
Yield	319	8.00	85.25	16.6027	7.19337
Yield	75	5.00	8.00	7.0433	1.63291
Yield	75	4.25	9.00	7.2613	1.08230
Yield	75	6.50	9.00	7.7029	.48042
Yield	75	8.50	9.00	7.7187	.48131
Yield	75	8.00	9.00	7.2121	.66303
Yield	75	8.00	9.00	7.5545	.55415
Yield	75	8.53	9.00	7.4481	.55411
Yield	75	8.00	8.83	7.5192	.63533
Yield	75	8.00	8.83	7.6178	.61101
Yield	75	8.00	8.75	7.9105	.55125
Yield	75	7.00	9.00	7.5251	.49142
Yield+Yield	70				



We've reviewed how to...

- Prepare for your statistical consultation
 - Compare sample to population
 - Address statistical issues
- Collect data
- Set up data files
- Enter data
- Check and clean data



Returning to the Consultant

- What understanding do the researcher and consultant have about their roles with analysis, writing about results, and authorship?
- Recommended (at the very least): a visit to make sure analyses and interpretations are correct!



Questions?



MERC Evaluation Link

Please go to the link below and complete the evaluation:

<http://goo.gl/8pwI6B>