# Experimental design, basic statistics, and sample size determination

Karl W Broman

Department of Biostatistics
Johns Hopkins Bloomberg School of Public Health

http://www.biostat.jhsph.edu/~kbroman

1

# Experimental design

2

# Basic principles

1. Formulate question/goal in advance
2. Comparison/control
3. Replication
4. Randomization
5. Stratification (aka blocking)
6. Factorial experiments

ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

# Example

Question:     Does salted drinking water affect blood
                    pressure (BP) in mice?

Experiment:

1. Provide a mouse with water containing 1% NaCl.
2. Wait 14 days.
3. Measure BP.

ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

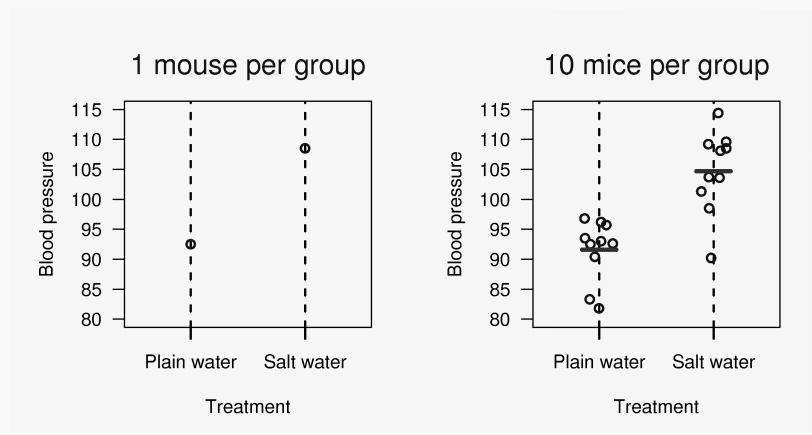# Comparison/control

Good experiments are comparative.

- Compare BP in mice fed salt water to BP in mice fed plain water.
- Compare BP in strain A mice fed salt water to BP in strain B mice fed salt water.

Ideally, the experimental group is compared to concurrent controls (rather than to historical controls).

# Replication

# Why replicate?

- Reduce the effect of uncontrolled variation (i.e., increase precision).

- Quantify uncertainty.

A related point:

An estimate is of no value without some statement of the uncertainty in the estimate.

ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

# Randomization

Experimental subjects ("units") should be assigned to treatment groups at random.

At random does not mean haphazardly.

One needs to explicitly randomize using
- A computer, or
- Coins, dice or cards.

ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

# Why randomize?

- Avoid bias.
  - For example: the first six mice you grab may have intrinsically higher BP.

- Control the role of chance.
  - Randomization allows the later use of probability theory, and so gives a solid foundation for statistical analysis.

# Stratification

- Suppose that some BP measurements will be made in the morning and some in the afternoon.
- If you anticipate a difference between morning and afternoon measurements:
  - Ensure that within each period, there are equal numbers of subjects in each treatment group.
  - Take account of the difference between periods in your analysis.
- This is sometimes called "blocking".

# Example

- 20 male mice and 20 female mice.

- Half to be treated; the other half left untreated.

- Can only work with 4 mice per day.

Question:    How to assign individuals to treatment
            groups and to days?

---

# An extremely bad design

| Week One | | | | | Week Two | | | | |
| M | Tu | W | Th | F | M | Tu | W | Th | F |
| C | C | C | C | C | T | T | T | T | T |
| C | C | C | C | C | T | T | T | T | T |
| C | C | C | C | C | T | T | T | T | T |
| C | C | C | C | C | T | T | T | T | T |

T = treated, C = control, pink = female, blue = male

# Randomized

| Week One | | | | | | Week Two | | | | |
| M | Tu | W | Th | F | | M | Tu | W | Th | F |
|---|----|---|----|---|---|---|----|---|----|---|
| T | T | T | T | T | | C | T | T | C | T |
| C | T | T | T | T | | C | C | C | T | C |
| C | C | C | T | T | | C | C | T | C | C |
| T | C | C | C | C | | C | T | C | T | T |

T = treated, C = control, pink = female, blue = male

13

# A stratified design

| Week One | | | | | | Week Two | | | | |
| M | Tu | W | Th | F | | M | Tu | W | Th | F |
|---|----|---|----|---|---|---|----|---|----|---|
| C | T | T | C | T | | C | C | T | C | T |
| T | T | C | C | C | | T | T | T | C | C |
| C | C | T | T | C | | C | T | C | T | C |
| T | C | C | T | T | | T | C | C | T | T |

T = treated, C = control, pink = female, blue = male

14

# Randomization and stratification

- If you can (and want to), fix a variable.
  - e.g., use only 8 week old male mice from a single strain.
- If you don't fix a variable, stratify it.
  - e.g., use both 8 week and 12 week old male mice, and stratify with respect to age.
- If you can neither fix nor stratify a variable, randomize it.

# Factorial experiments

Suppose we are interested in the effect of both salt water and a high-fat diet on blood pressure.

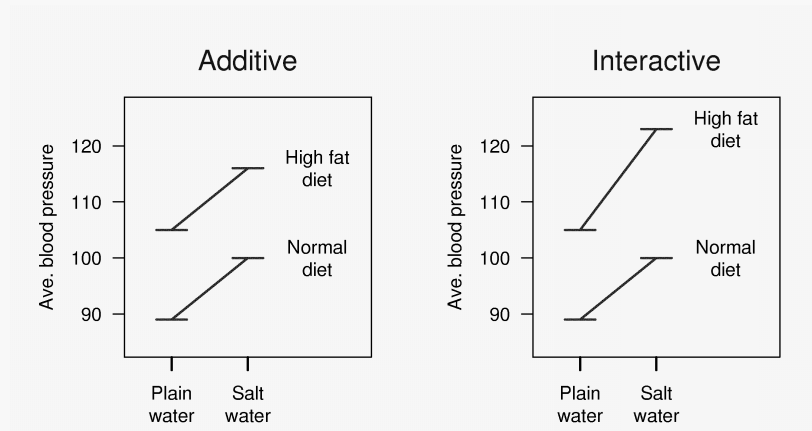Ideally: look at all 4 treatments in one experiment.

|  |  |  |
|---|---|---|
| Plain water | × | Normal diet |
| Salt water | | High-fat diet |

Why?
- We can learn more.
- More efficient than doing all single-factor experiments.

# Interactions

### Additive

Ave. blood pressure

120
110
100
90

High fat diet

Normal diet

Plain water     Salt water

### Interactive

Ave. blood pressure

120
110
100
90

High fat diet

Normal diet

Plain water     Salt water

---

# Other points

- Blinding
  - Measurements made by people can be influenced by unconscious biases.
  - Ideally, dissections and measurements should be made without knowledge of the treatment applied.
- Internal controls
  - It can be useful to use the subjects themselves as their own controls (e.g., consider the response after vs. before treatment).
  - Why?  Increased precision.

# Other points

- Representativeness
  - Are the subjects/tissues you are studying really representative of the population you want to study?
  - Ideally, your study material is a random sample from the population of interest.

# Summary

Characteristics of good experiments:

- Unbiased
  - Randomization
  - Blinding
- High precision
  - Uniform material
  - Replication
  - Stratification
- Simple
  - Protect against mistakes

- Wide range of applicability
  - Deliberate variation
  - Factorial designs
- Able to estimate uncertainty
  - Replication
  - Randomization

# Basic statistics

ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

# What is statistics?

"We may at once admit that any inference from the particular to the general must be attended with some degree of uncertainty, but this is not the same as to admit that such inference cannot be absolutely rigorous, for the nature and degree of the uncertainty may itself be capable of rigorous expression."
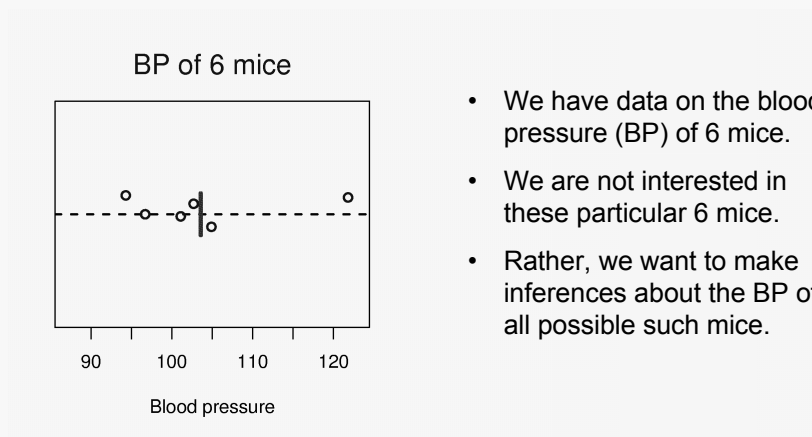
— Sir R. A. Fisher

ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

# What is statistics?

- Data exploration and analysis

- Inductive inference with probability

- Quantification of uncertainty

---

# Example

BP of 6 mice



90    100    110    120

Blood pressure

- We have data on the blood pressure (BP) of 6 mice.

- We are not interested in these particular 6 mice.

- Rather, we want to make inferences about the BP of all possible such mice.

# Sampling

BP of all possible mice

BP of 6 mice

true average

90   100   110   120

Blood pressure

90   100   110   120

Blood pressure

# Several samples

BP of all possible mice

Several possible samples

true average

90   100   110   120

Blood pressure

90   100   110   120

Blood pressure

# Distribution of sample average

BP of all possible mice

average BP of 6 mice



true average

true ave.

90   100   110   120

Blood pressure

90   100   110   120

Blood pressure

27

ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

---

# Confidence intervals

- We observe the BP of 6 mice, with average = 103.6 and standard deviation (SD) = 9.7.

- We assume that BP in the underlying population follows a normal (aka Gaussian) distribution.

- On the basis of these data, we calculate a 95% confidence interval (CI) for the underlying average BP:

$$103.6 \pm 10.2 \quad = \quad (93.4 \text{ to } 113.8)$$
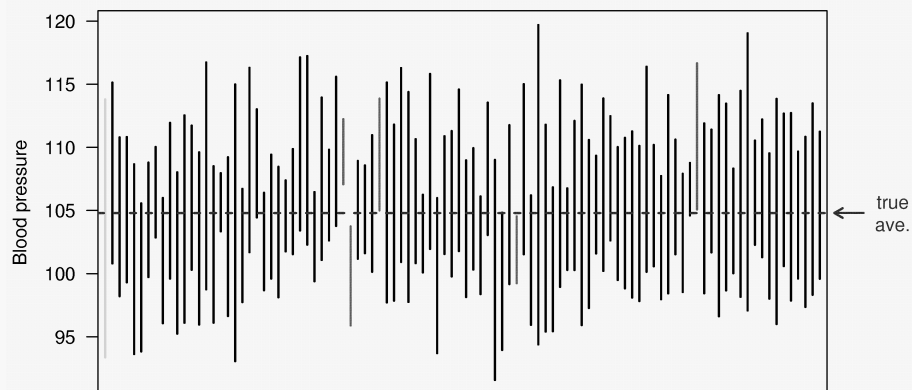
28

ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

# What is a CI?

- The plausible values for the underlying population average BP, given the data on the six mice.

- In advance, there is a 95% chance of obtaining an interval that contains the population average.

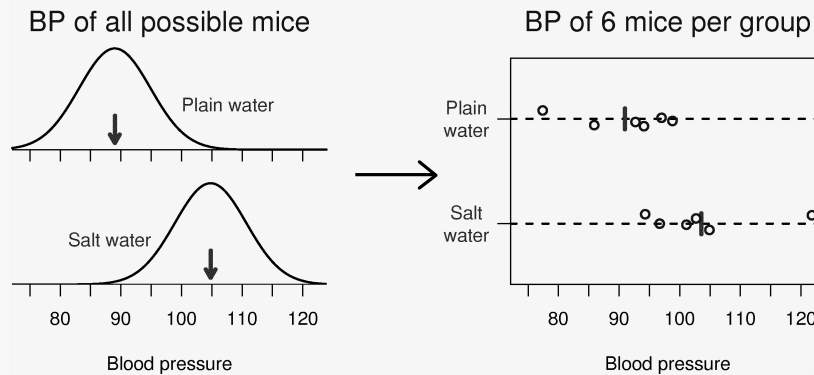29   ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

---

# 100 CIs



30   ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

# CI for difference

BP of all possible mice

BP of 6 mice per group

Plain water

Salt water

Plain water

Salt water

80   90   100   110   120

Blood pressure

80   90   100   110   120

Blood pressure

95% CI for treatment effect = 12.6 ± 11.5

---

# Significance tests

Confidence interval:
    The plausible values for the effect of salt water on BP.

Test of statistical significance:
    Answer the question, "Does salt water have an effect?"

Null hypothesis ($H_0$):        Salt water has no effect on BP.

Alt. hypothesis ($H_a$):        Salt water does have an effect.

# Two possible errors

- Type I error ("false positive")

  Conclude that salt water has an effect on BP when, in fact, it does not have an effect.

- Type II error ("false negative")

  Fail to demonstrate the effect of salt water when salt water really does have an effect on BP.

---

# Type I and II errors

|  | The truth | |
| --- | --- | --- |
| Conclusion | No effect | Has an effect |
| Reject $H_0$ | Type I error | ✔ |
| Fail to reject $H_0$ | ✔ | Type II error |

# Conducting the test

- Calculate a test statistic using the data.
  (For example, we could look at the difference between the average BP in the treated and control groups; let's call this D.)

- If this statistic, D, is large, the treatment appears to have some effect.

- How large is large?
  - We compare the observed statistic to its distribution if the treatment had no effect.
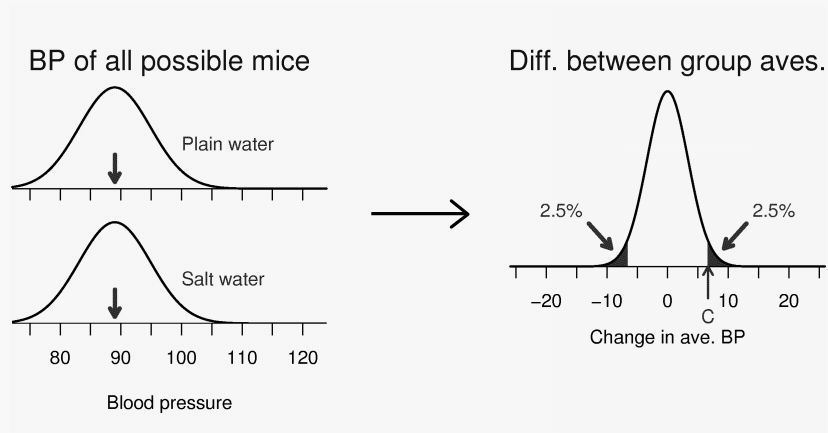
ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

# Significance level

- We seek to control the rate of type I errors.

- Significance level (usually denoted $\alpha$) = chance you reject $H_0$, if $H_0$ is true; usually we take $\alpha = 5\%$.

- We reject $H_0$ when $|D| > C$, for some C.

- C is chosen so that, if $H_0$ is true, the chance that $|D| > C$ is $\alpha$.
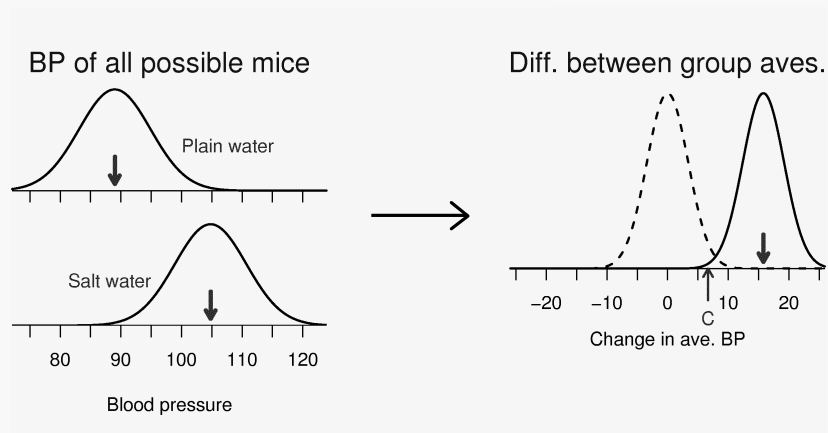
ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

# If salt has no effect

BP of all possible mice

Diff. between group aves.



37

# If salt has an effect

BP of all possible mice

Diff. between group aves.



38

19

# P-values

- A P-value is the probability of obtaining data as extreme as was observed, if the null hypothesis were true (i.e., if the treatment has no effect).

- If your P-value is smaller than your chosen significance level ($\alpha$), you reject the null hypothesis.

- We seek to reject the null hypothesis (we seek to show that there is a treatment effect), and so small P-values are good.
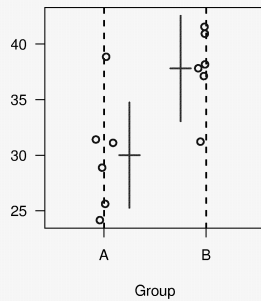
# Summary

- Confidence interval
  - Plausible values for the true population average or treatment effect, given the observed data.

- Test of statistical significance
  - Use the observed data to answer a yes/no question, such as "Does the treatment have an effect?"

- P-value
  - Summarizes the result of the significance test.
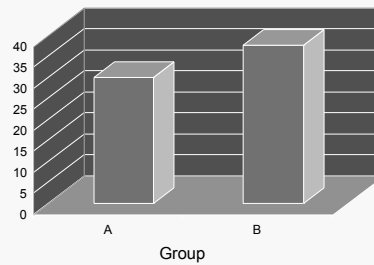  - Small P-value $\rightarrow$ conclude that there is an effect.

  Never cite a P-value without a confidence interval.

# Data presentation

### Good plot



### Bad plot

---

# Data presentation

### Good table

| Treatment | Mean | (SEM) |
|-----------|------|-------|
| A | 11.3 | (0.6) |
| B | 13.5 | (0.8) |
| C | 14.7 | (0.6) |

### Bad table

| Treatment | Mean | (SEM) |
|-----------|------|-------|
| A | 11.2965 | (0.63) |
| B | 13.49 | (0.7913) |
| C | 14.727 | (0.6108) |

# Sample size determination

43

# Fundamental formula

$$n = \frac{\$ \text{ available}}{\$ \text{ per sample}}$$

44

# Listen to the IACUC
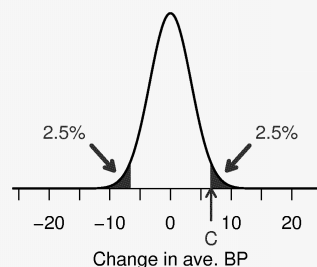
Too few animals    →    a total waste

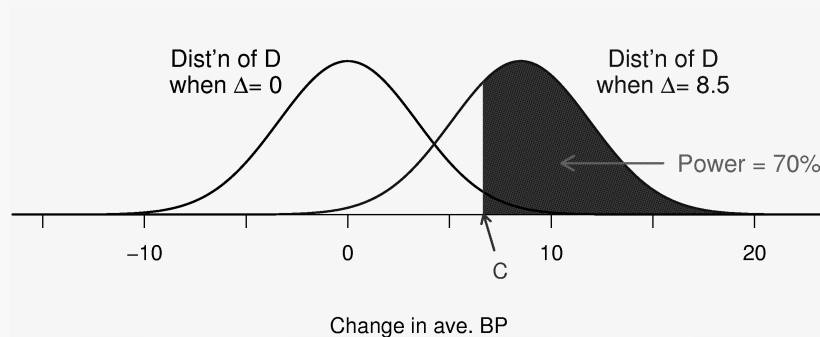Too many animals    →    a partial waste

---

# Significance test

- Compare the BP of 6 mice fed salt water to 6 mice fed plain water.

- $\Delta$ = true difference in average BP (the treatment effect).
- $H_0$: $\Delta = 0$ (i.e., no effect)
- Test statistic, D.
- If $|D| > C$, reject $H_0$.
- C chosen so that the chance you reject $H_0$, if $H_0$ is true, is 5%

Distribution of D
when $\Delta = 0$



2.5%          2.5%

−20   −10   0   C  10   20
Change in ave. BP

# Statistical power

Power = The chance that you reject $H_0$ when $H_0$ is false
(i.e., you [correctly] conclude that there is a treatment
effect when there really is a treatment effect).

Dist'n of D
when $\Delta = 0$

Dist'n of D
when $\Delta = 8.5$

Power = 70%

−10  0  10  20

C

Change in ave. BP

47

---

# Power depends on…

- The structure of the experiment
- The method for analyzing the data
- The size of the true underlying effect
- The variability in the measurements
- The chosen significance level ($\alpha$)
- The sample size

Note: We usually try to determine the sample size to
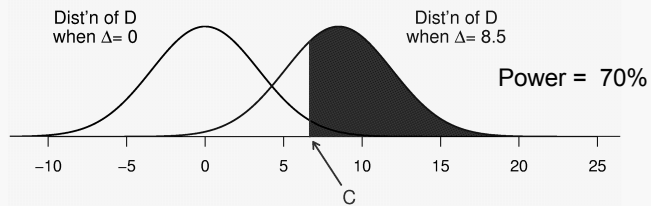give a particular power (often 80%).

48

Johns Hopkins University Center for Alternatives to Animal Testing

**Effect of sample size**

6 per group:

Dist'n of D when Δ= 0    Dist'n of D when Δ= 8.5

Power = 70%

−10   −5   0   5   10   15   20   25

C

12 per group:

Dist'n of D when Δ= 0    Dist'n of D when Δ= 8.5

Power = 94%

−10   −5   0   5   10   15   20   25

C

49    ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH



Johns Hopkins University Center for Alternatives to Animal Testing

**Effect of the effect**

Δ = 8.5:

Dist'n of D when Δ= 0    Dist'n of D when Δ= 8.5

Power = 70%

−10   −5   0   5   10   15   20   25

C

Δ = 12.5:

Dist'n of D when Δ= 0    Dist'n of D when Δ= 12.5

Power = 96%

−10   −5   0   5   10   15   20   25

C

50    ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

# A formula

$$n = \left(\frac{\sigma}{\ }\right)^2 \left[\ \ \phantom{Censored} \ _{1-\beta}\right]^2 \times 2$$

Censored

---

# Various effects

- Desired power ↑ ⇒ sample size ↑

- Stringency of statistical test ↑ ⇒ sample size ↑

- Measurement variability ↓ ⇒ sample size ↓

- Treatment effect ↑ ⇒ sample size ↓

# Determining sample size

The things you need to know:

- Structure of the experiment
- Method for analysis
- Chosen significance level, $\alpha$ (usually 5%)
- Desired power (usually 80%)

- Variability in the measurements
  - if necessary, perform a pilot study, or use data from prior publications

- The smallest meaningful effect

ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

---

# Reducing sample size

- Reduce the number of treatment groups being compared.

- Find a more precise measurement  (e.g., average time to effect rather than proportion sick).

- Decrease the variability in the measurements.
  - Make subjects more homogeneous.
  - Use stratification.
  - Control for other variables (e.g., weight).
  - Average multiple measurements on each subject.

ENHANCING HUMANE SCIENCE IMPROVING ANIMAL RESEARCH

# Final conclusions

- Experiments should be designed.

- Good design and good analysis can lead to reduced sample sizes.

- Consult an expert on both the analysis and the design of your experiment.

---

# Resources

- ML Samuels, JA Witmer (2003) *Statistics for the Life Sciences*, 3rd edition.  Prentice Hall.
  - An excellent introductory text.

- GW Oehlert (2000) *A First Course in Design and Analysis of Experiments*.  WH Freeman & Co.
  - Includes a more advanced treatment of experimental design.

- Course: *Statistics for Laboratory Scientists* (Biostatistics 140.615-616, Johns Hopkins Bloomberg Sch. Pub. Health)
  - Introductory statistics course, intended for experimental scientists.
  - Greatly expands upon the topics presented here.