stuff you now know ;

Dealing with ambiguity
Random variable
histograms
Comparative boxplots
quantiles
distributions
probability (e.g. from Poisson)
sample mean and variance
distr mean and variance
qqplots
scatterplots
scatterplots correlation
scatterplots correlation regression (multiple, polynomial,)
scatterplots correlation regression (multiple, polynomial,) ANOVA (R^2, s_e ~ RMSE)
scatterplots correlation regression (multiple, polynomial,) ANOVA (R^2, s_e ~ RMSE) overfitting, collinearity, interaction
 scatterplots correlation regression (multiple, polynomial,) ANOVA (R^2, s_e ~ RMSE) overfitting, collinearity, interaction
scatterplots correlation regression (multiple, polynomial,) ANOVA (R^2, s_e ~ RMSE) overfitting, collinearity, interaction sampling distribution
scatterplots correlation regression (multiple, polynomial,) ANOVA (R^2, s_e ~ RMSE) overfitting, collinearity, interaction sampling distribution 1-sample Confidence Interval for
scatterplots correlation regression (multiple, polynomial,) ANOVA (R^2, s_e ~ RMSE) overfitting, collinearity, interaction sampling distribution 1-sample Confidence Interval for 2-sample CI for
scatterplots correlation regression (multiple, polynomial,) ANOVA (R^2, s_e ~ RMSE) overfitting, collinearity, interaction sampling distribution 1-sample Confidence Interval for 2-sample CI for t-distribution

In the recording of today's lecture, I'll first go over the examples in the last lecture. So, please study Them carefully.

Lecture 19 (Ch. 9-end)  
Heve are The CI's we have derived: If 
$$r_{n} = \operatorname{kenown}(or \ln r_{0}^{n} n)$$
  
there are The CI's we have derived: If  $r_{n} = \operatorname{kenown}(or \ln r_{0}^{n} n)$   
 $1 - \operatorname{semple}(C.I. for  $\mathcal{P}_{X}: \quad x \pm z^{*} \ \overline{or} \ h(\overline{n} \ or \ \overline{x} \pm t^{*} \frac{S_{2}}{M}, df = n-1$   
 $1 - \operatorname{semple}(C.I. for  $\mathcal{P}_{X}: \quad p \pm z^{*} \sqrt{\frac{p(1-p)}{n}} \ No \ t^{*} \ |FYI: usetstrap)$   
 $2 - \operatorname{semple}(C.I. for  $\mathcal{P}_{X} - \mathcal{P}_{Z}: \quad \overline{x}, -\overline{x}_{2} \pm z^{*} \sqrt{\frac{p(1-p)}{n}} \ No \ t^{*} \ |friddenterrow \ h(-h)| + \frac{h(1-p)}{n} \ No \ t^{*} \ df = h(h) + \frac{h(1-p)}{n} \ No \ t^{*} \ df = h(h) + \frac{h(1-p)}{n} \ No \ t^{*} \ df = h(h) + \frac{h(1-p)}{n} \ No \ t^{*} \ df = h(h) + \frac{h(1-p)}{n} \ No \ t^{*} \ df = h(h) + \frac{h(1-p)}{n} \ No \ t^{*} \ df = h(h) + \frac{h(1-p)}{n} \ No \ t^{*} \ df = h(h) + \frac{h(1-p)}{n} \ No \ t^{*} \ df = h(h) + \frac{h(1-p)}{n} \ No \ t^{*} \ df = h(h) + \frac{h(1-p)}{n} \ No \ t^{*} \ df = h(h) \ df =$$$$ 

Not-independent Samples

We required The 2 samples (in a 2-sample problem) to be indep. It happened when we wrote  $V[\bar{x}_1 - \bar{x}_2] = V[\bar{x}_1] + V[\bar{x}_2] + 0 = \sigma_1^2 / n_1 + \sigma_2 / n_2$ But There exist problems where The 2 samples are not in dependent. E-2.1 : Does a certain pill increase IQ. > Take 100 peogle, measure Their IQ, w/o pill Z versus -> Take 100 people, measure Their IQ, both before and after pill. E-g.2: Is there a diff in The mean speed of App X and App Y? -> Time App X and App Y on each of 100 computers. -> Such data are called "paired". IQ After You can usually sec /test This by looking at: How do we build a C.I. for M.-M2 from paired data? Make a new column' IQ betore IQ after & C.I. for M. -M2 for paired data:  $\chi_1 \qquad \chi_2 \qquad d = \chi_1 - \chi_2$  $\int dt t = \frac{1}{2} \frac{1}{\sqrt{2}}$ Benefit: paired CIs ave often narrower (more precise). The Math is trivial. "Addressing" pairedness is Complex!

Example Here is how we did The fish example in last lecture: Concentration of zinc in  $\frac{2}{x}$  types of fish. Type I 56 9.15 1.27 3.08 Type II 61 1.71 Are The true/psp. means different? M, = pop. mean zinc in Type I [ Important to define M, Mr. M, = II (The pop. parameters) clearly. 41 95% C.I. for  $M_1 - M_2 = (9.15 - 3.08) \pm 1.96 \sqrt{\frac{(1.27)^2}{56} + \frac{(1.71)^2}{61}}$ 6.07 ± 0.54 = [5.53, 6.61] There are many ways of collecting paired data, e.g. 1) Take one Type I and one Type II, from n lakes. " " " " " " " " with The same weight. 3) In every example, a Type I fish is somehow matched/paired Benefit: In every example, The CI for  $\mu_1 - \mu_2$  will be narrower (ie. more precise) (Once again: If you are analyzing existing data, The squestion Should be "Are the data paired?" I Ffyou are collecting the data yourself do your best to assure factors you "control", The narrower your CI will be.

Jerry Wei (our TA), has found that a negative correlation between the 2 samples can lead to a wider CI (larger p-value, in ch 8) for the paired test! The width of a CI is related to something called power. A wider CI has a lower power - and that's not a good thing.

library(MASS) set.seed(123) s = matrix(c(1, -0.8, -0.8, 1), 2) df = mvrnorm(n=100, mu=c(0, 0.3), Sigma=s, empirical=F) t.test(df[,1], df[, 2], alternative="less", paired=FALSE) t.test(df[,1], df[, 2], alternative="less", paired=TRUE)

https://stats.stackexchange.com/questions/38102/paired-versus-unpaired-t-test The wiki page on "paired difference test" compares the power of the unpaired and paired tests, and indeed shows that a positive correlation is required for higher power. Ie, there \*is\* such a thing as a "bad pairing" which can lead to lower

(FYI) It may be tempting to Think That a paired design does not exist for groportions. But it does. 2 genders (boy/girl), 2 computertypes (Mac(Deld) Here is an example data from an unpaired design: Boy Girl prop. of Give with Mac avop. of Boys E Mac Mac Dell with Mac Mac P2 Dell Dell PI / ! / C.I. for  $m_1 - m_2$ :  $(p_1 - p_2) \pm 2^* / \frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}$ Here is an example data from a paired design: Husband wife > prog. of wives with Mac 200 p. of Husbands ~ Mac Mac with Mac 4 Pi Dell li Pz Mac Dall Dell ( . . However, There is no simple formula for The C.I. of 71-72 when data are paired. (After all, we cannot make a new column of differences, because we cannot subtract "Mac" from "Mac", ic. The data are categorical.

hw lect 19 1 USe t- dicty
Consider the following data on x1 and x2 which was collected in a paired design:
$x_1 = c(-0.27, -0.14, 1.61, 0.09, 0.00, 2.07, 0.56, -1.67, -0.51, -0.54)$
 x2 = c(-0.32, 0.20, 1.93, 0.54, 0.75, 1.77, 0.84, -0.29, -0.33, 0.17)
 a) Compute a 2-sided, 95% CI for the difference between the two true means. You may use R to do simple
their scatterplot:
plot(x1,x2) # I see a linear association
 b) Provide at least one interpretation of the observed CI, AND state the conclusion in English, i.e., the
"Corollary."
Compute an appropriate 95% 2-sided CI.
 $y_1 = c(-0.27, -0.14, 1.61, 0.09, 0.00, 2.07, 0.56, -1.67, -0.51, -0.54)$
 $y_2 = c(0.20, 0.54, -0.33, 1.93, -0.32, 1.77, 0.75, 0.17, -0.29, 0.84)$
d) Provide one interpretation of the observed CI, AND state the conclusion in English, i.e., the "corollary." . e) Which one is narrower?