

Instructions: A printed copy of your homework should be handed in at **the start of class** the day it is due. If you have any supplementary electronic files you wish to turn in (e.g. R scripts or wxMaxima files) email them to the instructor prior to class with file name format: `Lastname-hwX.ext`. Each part of each exercise is worth 10 points unless stated otherwise.

Exercise 1: Suppose the results of a Poisson Regression analysis done on transect count data (similar to the one done in class) reveals a best fit model of $y_i \sim \text{Poisson}(\lambda = 0.75 \text{Length}_i)$, where the data do not appear to be overdispersed. Based on the underlying model, what would you estimate to be the average distance between organisms along the transects? Justify your answer.

Exercise 2: See the file [rssmle.R](#) on the course website.

- (a) First, generate a fake data set of 50 points along a parabola (i.e. a quadratic polynomial) with two real roots. Whether or not you add noise is up to you.
- (b) Estimate the parameters r_i in $y = (x - r_1)(x - r_2)$ or $-(x - r_1)(x - r_2)$ based on graphically inspecting the data (e.g. plot it and carefully read off where the roots occur).
- (c) Fit the function $y = a + bx + cx^2$ to your simulated data using two different methods in `optimx()` (e.g., use `method=c("Nelder-Mead", "BFGS")`). Use the quadratic formula to calculate the roots r_i , and compare to part (b).

Exercise 3: Fit the rate parameter of an exponential distribution from a simulated data set (use `set.seed(123); x=rexp(200,rate=pi);`) using an MLE based optimization using `optimize()` and the `dexp()` function in R (i.e., don't use the analytical formula for the MLE; see Exercise 4).

Exercise 4: Next, repeat the above parameter estimation but here find an analytical expression for the MLE \hat{r} by differentiating the negative log likelihood function for the scenario above (i.e., for an exponential distribution with rate r) for a given data set (x_1, \dots, x_N) . Recall

$$f(x) = r \exp(-rx), \quad \text{for } x \geq 0.$$

Exercise 5:

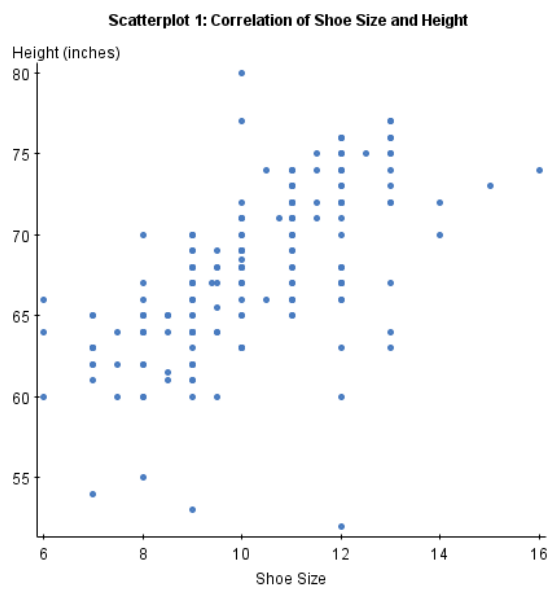
(420 Students Only) Use the `boot` package in R to calculate 95% confidence intervals for the estimate of r used in Exercise 4 above (for example, see <https://www.statmethods.net/advstats/bootstrapping.html>).

(620 Students Only) The **non-parametric bootstrap** is a resampling-based method of constructing a distribution for an estimator that can then be used to create confidence intervals. In short, iteratively sample from your data with replacement a new data set with the same sample size and use it to recalculate your estimate. Repeat this 10,000 or more times and each time store the estimate, then use that set of 10,000+ estimates as a sample from the estimator distribution to construct the confidence interval. Implement this for the estimation procedure in Exercise 3 above, plot a histogram of the bootstrapped estimates, an empirical CDF (using `plot(ecdf(rbootest))` if the vector of bootstrap estimates is stored in `rbootest`), a CDF plot using the estimate \hat{r} (use `curve(pexp(x,r.est),0,2.8,col="red",type="l",lty=2` or overlay it on the empirical CDF by adding the argument `add=TRUE`).

- Does this estimate seemed biased?
- What is the 95% Confidence Interval on r ? (Hint: if the bootstrap estimates are in a vector `rbootest`, the 2.5% and 97.5% percentile of r can be obtained using the R code `quantile(rbootest, c(0.025, 0.975))`).

Exercise 6 (Parametric Bootstrap): Consider this small study of human height versus shoe size at (<https://www.statcrunch.com/5.0/viewreport.php?reportid=35115>), summarized in Figure 1 below, which found the following best fit SLR regression model correlating U.S. shoe size and height in inches: $Height = 50.87 + 1.657 ShoeSize$ with $\hat{\sigma} = 3.78$ which yielded an R^2 value of 0.387. This analysis is repeated in the R script **hw3-shoesize-height.R** using the data in **hw3-sullivan-statistical-survey-data.xlsx**.

- What is the basic definition of R^2 in this context, and how should we interpret R^2 values?
- How expected or unexpected is this level of explanatory power (i.e., the $R^2 = 0.387$ value) under the assumption that this best fit model is true, and given these sample sizes? To answer that question, use the following number of individuals with each shoe size in the data set to simulate synthetic data under their best fit linear model given above. Repeat this data simulation 50,000 times using a for loop (`for(i in 1:50000) { ...}`) and each time store the R^2 value using something like `fit=lm(...); fitsum <- summary(fit); Rsq[i]<-fitsum$adj.r.squared;`. Plot a histogram of these R^2 values, use `abline()` to draw a vertical line indicating the $R^2 = 0.387$, and see Exercise 5 above, and calculate a corresponding p-value.
- Discuss what your results mean in terms of how well this particular data set meets the assumptions of the SLR model.



Simple linear regression results:
 Dependent Variable: Height
 Independent Variable: Foot
 Height = 50.874798 + 1.6565183 Foot
 Sample size: 199
 R (correlation coefficient) = 0.6222
 R-sq = 0.38713068
 Estimate of error standard deviation: 3.7840705

Parameter estimates:

Parameter	Estimate	Std. Err.	Alternative	DF	T-Stat	P-Value
Intercept	50.874798	1.5311509	≠ 0	197	33.22651	<0.0001
Slope	1.6565183	0.14849721	≠ 0	197	11.155215	<0.0001

Analysis of variance table for regression model:

Source	DF	SS	MS	F-stat	P-value
Model	1	1781.8633	1781.8633	124.438835	<0.0001
Error	197	2820.8804	14.31919		
Total	198	4602.7437			

Predicted values:

X value	Pred. Y	s.e.(Pred. y)	95% C.I. for mean	95% P.I. for new
10	67.43998	0.2691875	(66.90912, 67.97084)	(59.958637, 74.921326)

Residuals stored in new column, Residuals.

Figure 1: Source: <https://www.statcrunch.com/5.0/viewreport.php?reportid=35115> based on data from the Sullivan Statistics Survey at <https://www.statcrunch.com/app/index.php?dataid=450205>.