Please adhere to the homework rules as given in the Syllabus.

**1.    Hypothesis Testing for the Mean.** Your favorite brand of cookies claims that on average each cookie contains 10 chocolate chips. Being the skeptic that you are, you decide to test this claim.

a) Write out the Null and Alternative Hypotheses for each of the following situations.

    i) You seek to find evidence that the average number of chips per cookie is not 10.
    ii) You seek to find evidence that the average number of chips per cookie is less than 10.
    iii) You seek to find evidence that the average number of chips per cookie is more than 10.

b) Suppose you purchase $n = 9$ cookies and compute $\bar{x} = 9$ and $s = 3$. Calculate the test statistic.

c) Compute the $p$-value for each of the three tests above and make a conclusion. Your conclusion should be in terms of the problem. For each case use $\alpha = 0.05$.

d) Repeat parts $b$ and $c$ for $n = 16$, $n = 25$ and $n = 100$.

**2.** **Data Analysis**. The SENIC dataset contains data from a random sample of 144 hospitals across the US. The dataset can be found at

<div align="center">

`math.unm.edu/~knrumsey/classes/fall16/MiniProjects/senic.csv`

</div>

Extract the column named "MeanLengthOfStay" Test the claim: "the average length of stay at hospitals in America is less than 10 days". Specify which significance level you plan to use **before** looking at the data. Interpret your results.

**3.** **Type I and Type II Errors.** Harry has a single observation $X \sim Exp(1/\mu)$ and he wants to consider the following test. *Note: Remember that in this parameterization $\mu$ is the mean and the CDF is $F(x) = 1 - e^{-x/\mu}$.*

$$H_0 : \mu = 1 \qquad H_a : \mu = 2$$

Harry will use the following "rejection region". $R = \{x | x > 1.5\}$.

a) Find $\alpha$, the probability of Type *I* Error for this test.

b) Find $\beta$, the probability of Type *II* Error for this test. What is the *power* of this test?

c) Harry has decided that he cannot afford to have a Type I error probability which is more than 0.10. Give a new rejection rule which will ensure that $\alpha = 0.10$. What is the power of this new test?

d) Harry is sad, because he has realized that he cannot have high power and low Type I error probability at the same time (for this test at least). He decides that he must collect more data. Now assume that Harry has ten samples $X_1, X_2, \cdots X_{10} \overset{iid}{\sim} Exp(1/\mu)$. Letting $Y = \sum_{i=1}^{10} X_i$, his new rejection rule is: Reject the null in favor of the alternative if $Y > 15.7$. Find the probability of Type I error of this new test. Also find the power of this test. How does it compare to $c$? *Hint: What is the distribution of $Y$? You will need to use an online calculator to help you find the probabilites of error for this problem.*

**4. Challenge Problem.** *Emperical Testing.* Between Janruary 1st, 2016 and April 9th 2018, NBA player James Harden took 1,764 three pointers. Data can be found on the course webpage at

<div align="center">

`math.unm.edu/~knrumsey/classes/fall18/harden.csv`

</div>

The data is recorded sequentially (in order) where a 1 represents a made shot and 0 represents a missed shot. Clearly, we are dealing with Bernoulli trials here. Ideally, we would like to analyze this data by assuming that $Y = \sum_{i=1}^{1764} X_i$ is a Binomial random variable. **However,** the Bernoulli trials have to be independent in order for this result to hold. In this case, it is reasonable to assume that the trials may NOT be indepdent (hot or cold streaks etc.). In general, this is a very difficult thing to test, but we can attempt to do so *empirically*.

**The Hypothesis:** We would like to test,

$$H_0 : X_1, X_2, \cdots X_{1764} \text{ are independent} \qquad H_a : X_1, X_2, \cdots X_{1764} \text{ are not independent}$$

**The Test Statistic:** We need to come up with a statistic that might be able to test this hypothesis in some way. Here is one possible choice.

$$T = \text{the longest streak of made shots}$$

**The Reference Distribution:** Although $T$ can be easily computed from the data, it does no good unless we know the dsitribution of $T$ *given that $H_0$ is true.* This is very hard (and maybe impossible) to find this distribution exactly. But we can get an empirical distribution via simulation to use for reference. *If the Null hypothesis is true*, the $X_1, \cdots X_{91}$ are independent Bernoulli random variables. This data can be simulated in R by using $\mathtt{x} <- \mathtt{rbinom(91, 1,}$ $\mathtt{p)}$ where $p$ should be set to the proportion of Curry's made shots. If you repeat this process a large number of times (say 1000) and calculate $T$ each time, you will have a reference distribution that you can compare the actual observed value of $T$ to.

**Empirical $p$-value** Once you have a reference distribution, you can approximate the $p$-value in the usual way. Think about how you would do this, and ask your instructor for help if necesarry.

**The Problem:** Conduct this Hypothesis test, following the steps above. Provide a histogram showing the reference distribution you simulated, and add a vertical line illustrating where the observed value of the test statistic falls. Give the emprical p-value and interpret your results at the 5% significance level.

**Calculating $T$:** Here is an R function which will calculate the value of $T$ given a binary vector.

```
longest_make_streak <- function(x){
  i <- 0;   m <- 0;
  for(j in 1:length(x)){
    if(x[j] == 1){ i <- i + 1 }
    else{ i <- 0  }
    if(i > m) m <- i
  }
  return(m)
}
```

*Note: You should also look at the longest miss streak which requires a simple change to the function.*