# Solution to Homework 2

***Warning****: This note is only used as a reference solution for the homework, and the solution to each question is not unique. The solution may contain factual and/or typographic errors and comments and criticism are kindly welcomed.*

**Problem 1** *Recall that the Kolomogorov-Smirnov One-Sample Statistic $D_n = \sup\limits_{x}|F_X(x) - \hat{F}_n(x)|$ is distribution free ($\hat{F}_n$ is the empirical CDF). Suppose $n = 100$ and the observed $D_n = 0.04$. Would you reject the null hypothesis $H_0 : X_1, \ldots, X_n \sim F_X$ at level $\alpha = 0.05$. Write a simulation to justify your answer.*

*Solution:*

```r
1  # D: Kolomogorov-Smirnov One-Sample Statistic
2  # n: the number of samples
3  # nrep: the number of samples for simulation
4  kspvalue<-function(D,n,nrep){
5    Dn <- c()
6    for (i in 1:nrep) {
7      x.obs <- rnorm(n,0,1)
8      Fn.hat <- ecdf(x.obs)
9      Dn[i] <- max(abs(Fn.hat(x.obs) - pnorm(x.obs)))
10   }
11   return(sum(Dn > D)/nrep)
12 }
13 # Calculate the p-value of the statistics under the problem setting
```

```
14  D = 0.04; n = 100; nrep = 5000
15  kspvalue(D,n,nrep)
```

The result shows that the $p$-value is $0.983 > 0.05$ and we fail to reject the null hypothesis.    □

**Problem 2**    *Use the following r code to generate 50 random numbers from Cauchy distribution.*

```
1  set.seed (2)
2  d <- rcauchy (50)
```
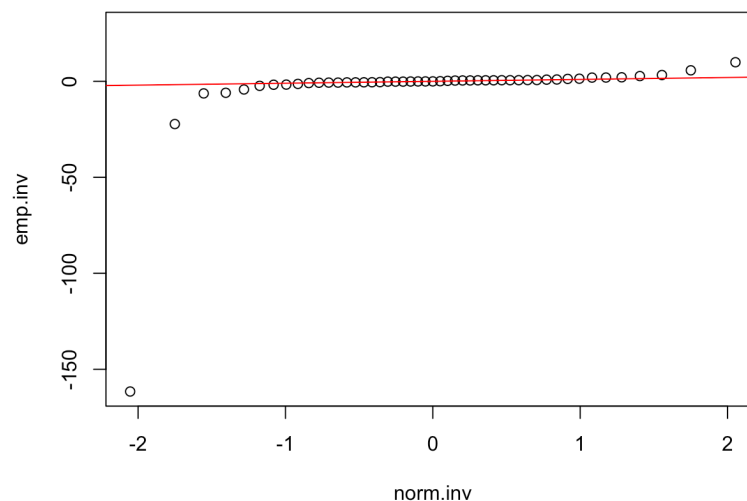
*Use a QQ plot to see how good (or bad) your data fits a normal distribution.*

*Solution:*

```
1  emp.inv <- sort(d)
2  sel <- seq(1, 50, by=1)/50
3  norm.inv <- qnorm(sel)
4  plot(norm.inv, emp.inv)
5  abline(0,1,col="red")
```



where we only compare the data to the standard normal distribution. If you want to show the goodness of fit for the normal family, kindly use 'abline(a=mean(d), b=sd(d), col="red")'.    □

**Problem 3**   *Use a test discussed in class to test whether the following sequence is random at level*
$\alpha = 0.1$,

$$A, B, B, A, A, A, A, A, A, B, A.$$

*Use simulation to justify your answer.*

*Solution: (Solution is not unique and the following code is based on the number of run)*

```r
1  # n: the number of samples
2  # nrep: the number of samples for simulation
3  # sample: the original sample
4  # run: the number of run in the samples
5  randomtest.run<-function(n,nrep,sam,run){
6    runs <- c()
7    for (i in 1:nrep) {
8      count <-1
9      temp <- sample(sam, 11, replace = F)
10     for (j in 1:(n-1)){
11        count <- count + abs(temp[j]-temp[j+1])
12     }
13     runs[i] <- count
14   }
15   return(2*min(sum(runs < run), sum(runs > run))/nrep) # two-sided p-value
16 }
17 # Calculate the p-value of the statistics under the problem setting
18 n=11; nrep = 5000; run=5
19 sam <- c(1, 0, 0, 1, 1, 1, 1, 1, 1, 0, 1)
20 randomtest.run(n,nrep,samp,run)
```

The corresponding $p$-value is $0.489 > 0.1$ and we fail to reject the null.                          □

**Remark**   Please note that randomness implies independency but *does not suggest* sampling with equal probability (i.e. Bernoulli($p$) with $p = 0.5$). Thus, in the simulation of Problem 3, you shall sample locally rather than creat new sequence with equal probability.