

Homework 4: Diffusion of Tetracycline

We continue examining the diffusion of tetracycline among doctors in Illinois in the early 1950s, building on our work in lab 6. You will need the data sets `ckm_nodes.csv` and `ckm_network.dat` from the labs.

1. Clean the data to eliminate doctors for whom we have no adoption-date information, as in the labs. Only use this cleaned data in the rest of the assignment.

```
ckm_nodes<-read_csv("data/ckm_nodes.csv")
ckm_nodes$rownum<-c(1:246)
ckm_network<-read.table("data/ckm_network.dat ")
ckm_nodes<-filter(ckm_nodes,!is.na(adoption_date))
ckm_network<-ckm_network[ckm_nodes$rownum,ckm_nodes$rownum]
```

2. Create a new data frame which records, for every doctor, for every month, whether that doctor began prescribing tetracycline that month, whether they had adopted tetracycline before that month, the number of their contacts who began prescribing strictly *before* that month, and the number of their contacts who began prescribing in that month or earlier. Explain why the dataframe should have 6 columns, and 2125 rows. Try not to use any loops.

```
c1<-c()
c2<-c()
c3<-c()
c4<-c()
c5<-c()
c6<-c()
for ( i in c(1:125)) {
  for (j in c(1:17)) {
    c1<-c(c1,ckm_nodes$rownum[i])
    c2<-c(c2,j)
    c3<-c(c3,ifelse(ckm_nodes$adoption_date[i]==j,1,0))
    c4<-c(c4,ifelse(ckm_nodes$adoption_date[i]<j,1,0))
    c5<-c(c5,sum(ifelse(ckm_nodes$adoption_date[ckm_network[i]==1]<j,1,0)))
    c6<-c(c6,sum(ifelse(ckm_nodes$adoption_date[ckm_network[i]==1]<=j,1,0)))
  }
}
newckm <- data.frame(Doc=c1,Mon=c2,began=c3,adopted=c4,ConN1=c5,ConN2=c6)
newckm$began <- factor(c("No", "Yes")[newckm$began+1])
newckm$adopted <- factor(c("No", "Yes")[newckm$adopted+1])
# fail to come up with codes without loops.

# Explain: Because ckm_nodes contains 125 doctors and 17 months ,then 125*17=2125 rows.
# And the tasks require 6 cols.
dim(newckm)
```

```
## [1] 2125    6
```

3. Let

$p_k = \Pr(\text{A doctor starts prescribing tetracycline this month} \mid \text{Number of doctor's contacts prescribing before this month} = k)$

and

$q_k = \Pr(\text{A doctor starts prescribing tetracycline this month} \mid \text{Number of doctor's contacts prescribing this month} = k)$

We suppose that p_k and q_k are the same for all months.

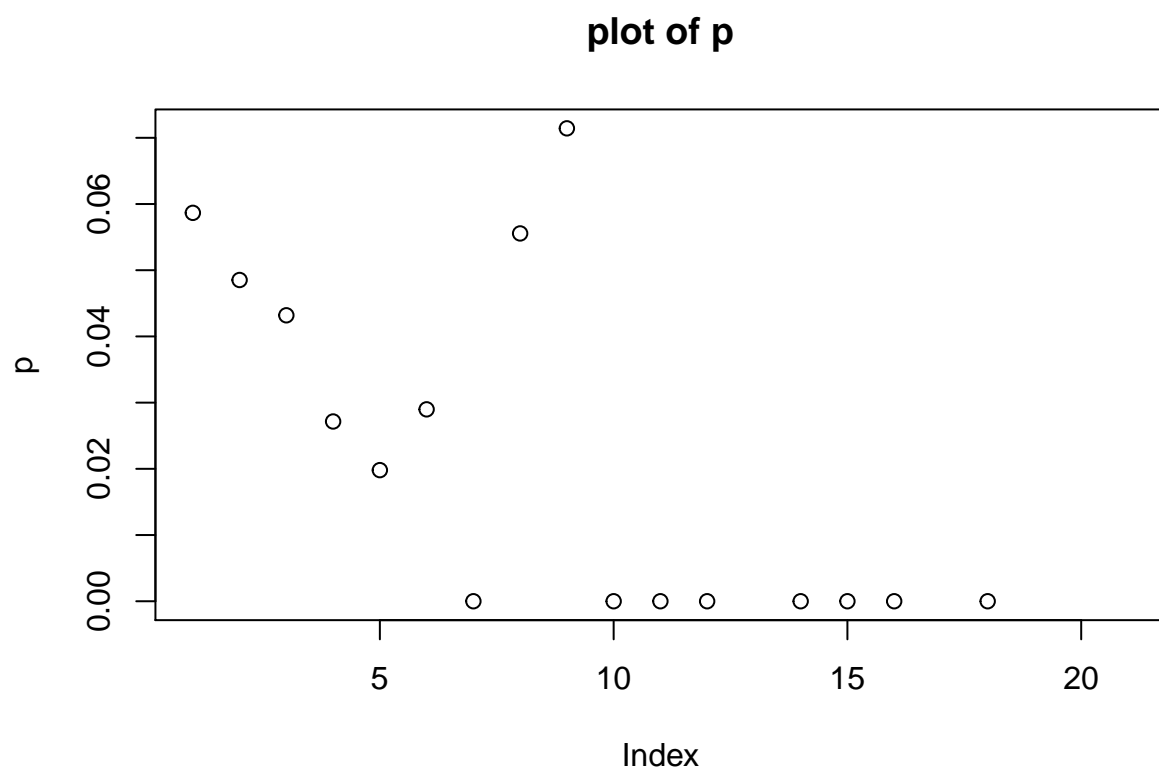
- a. Explain why there should be no more than 21 values of k for which we can estimate p_k and q_k directly from the data.

```
# Explain: Because in our data no doctor contacts more than 21 men totally ,i.e. k is no  
# more than 21 when estimating.  
max(newckm$ConN2)
```

```
## [1] 18
```

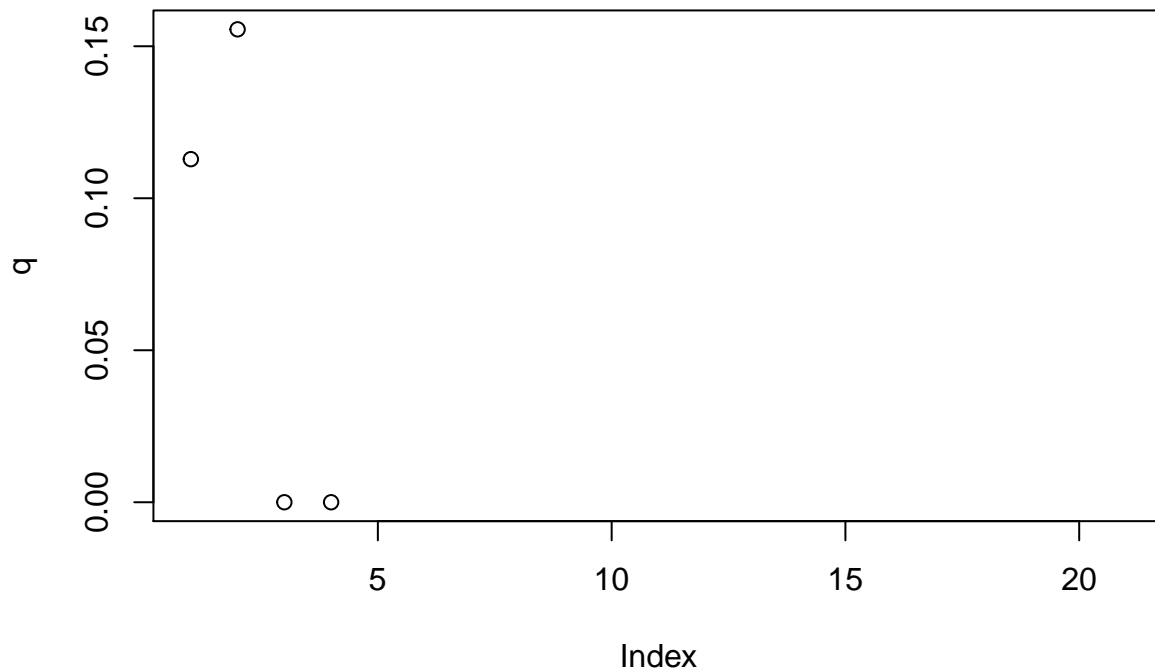
- b. Create a vector of estimated p_k probabilities, using the data frame from (2). Plot the probabilities.
- c. Create a vector of estimated q_k probabilities, using the data frame from (2). Plot the probabilities.

```
p<-c()  
q<-c()  
for (k in c(1:21)) {  
  p<-c(p,nrow(filter(newckm,ConN1==k&began=="Yes"))/nrow(filter(newckm,ConN1==k)))  
  q<-c(q,nrow(filter(newckm,ConN2-ConN1==k&began=="Yes"))/nrow(filter(newckm,ConN2-ConN1==k)))  
}  
plot(p)  
title("plot of p")
```



```
plot(q)
title("plot of q")
```

plot of q



4. Because it only conditions on information from the previous month, p_k is a little easier to interpret than q_k . It is the probability per month that a doctor adopts tetracycline, if they have exactly k contacts who had already adopted tetracycline.
 - a. Suppose $p_k = a + bk$. This would mean that each friend who adopts the new drug increases the probability of adoption by an equal amount. Estimate this model by least squares, using the values you constructed in (3b). Report the parameter estimates.

```
p<-na.omit(p)
ka<-c(1:12,14,15,16,18)
pka<-data.frame(pa=p,k=ka)
modelPa<-coefficients(lm(pa~k,pka))
names(modelPa)<-c("a","b")
modelPa
```

```
##          a          b
## 0.051410875 -0.003328183
```

- b. Suppose $p_k = e^{a+bk} / (1 + e^{a+bk})$. Explain, in words, what this model would imply about the impact of

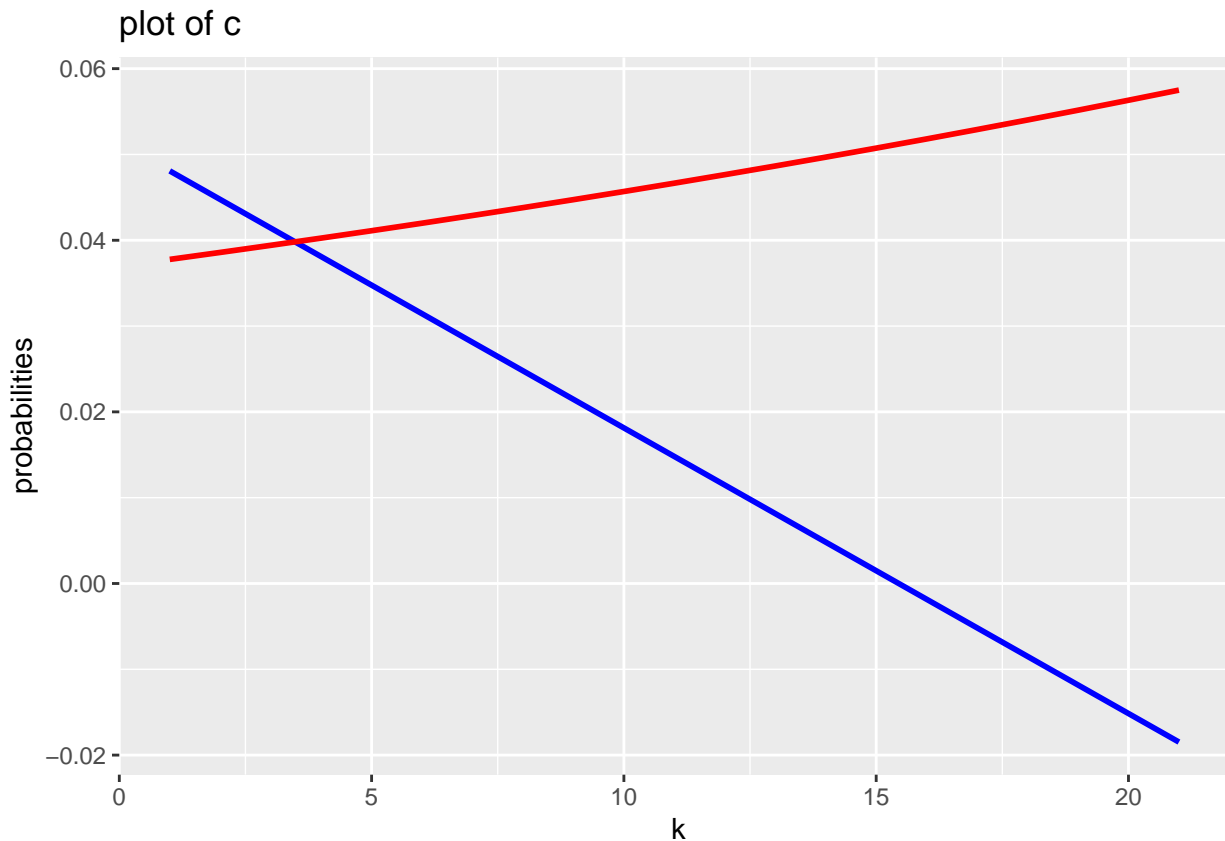
```
#Explain: This would mean that each friend who adopts the new drug increases the
# probability of adoption less and less.
pb<-(log(p/(1-p),base = exp(1)))[(log(p/(1-p),base = exp(1)))!=0]
kb<-ka[(log(p/(1-p),base = exp(1)))!=0]
pkb<-data.frame(pb=pb,kb=kb)
```

```
pkb<-filter(pkb,pb!=--Inf)
modelPb<-coefficients(lm(pb~kb,pkb))
names(modelPb)<-c("a","b")
modelPb
```

```
##          a          b
## -3.2597142  0.0220459
```

c. Plot the values from (3b) along with the estimated curves from (4a) and (4b). (You should have one

```
# curve of 4a is blue ,curve of 4b is red.
c<-c(1:21)
proplot=data.frame(p1=modelPa["a"]+modelPa["b"]*c,p2=1/(1+exp(-modelPb["a"]-modelPb["b"]*c)),c=c)
proplot %>% ggplot() +
  labs(x = " k", y = "probabilities",title = "plot of c")+
  geom_line(aes(x = c, y = p1), col = 'blue', size = 1.0)+
  geom_line(aes(x = c, y = p2), col = 'red', size = 1.0)
```



```
# curve of 4a is blue ,curve of 4b is red ,and I prefer the 4b red one because it"s
# positive correlation.
```

For quibblers, pedants, and idle hands itching for work to do: The p_k values from problem 3 aren't all equally precise, because they come from different numbers of observations. Also, if each doctor with k adoptee contacts is independently deciding whether or not to adopt with probability p_k , then the variance

in the number of adoptees will depend on p_k . Say that the actual proportion who decide to adopt is \hat{p}_k . A little probability (exercise!) shows that in this situation, $\mathbb{E}[\hat{p}_k] = p_k$, but that $\text{Var}[\hat{p}_k] = p_k(1 - p_k)/n_k$, where n_k is the number of doctors in that situation. (We estimate probabilities more precisely when they're really extreme [close to 0 or 1], and/or we have lots of observations.) We can estimate that variance as $\hat{V}_k = \hat{p}_k(1 - \hat{p}_k)/n_k$. Find the \hat{V}_k , and then re-do the estimation in (4a) and (4b) where the squared error for p_k is divided by \hat{V}_k . How much do the parameter estimates change? How much do the plotted curves in (4c) change?

It can be done with flexible use of results in homework 3 and homework 4.