

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_network <- read.table("data/ckm_network.dat")
a <- which(is.na(ckm_nodes$adoption_date))
ckm_network <- ckm_network[-a,-a]
ckm_nodes <- ckm_nodes[-a,]
```

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.

```
doctor = rep(1:125, each = 17)
month = rep(1:17, times = 125)
that_month = rep(ckm_nodes$adoption_date, each = 17)
before_month = (month > that_month)
that_month = (month == that_month)
num_before = function(doc, mon){
```

```

    return(sum(ckm_nodes$adoption_date[which(ckm_network[,doc] == 1)] < mon))
  }
  num_ealier = function(doc, mon){
    return(sum(ckm_nodes$adoption_date[which(ckm_network[,doc] == 1)] <= mon))
  }
  connects_before = mapply(num_before, doctor, month)
  connects_ealier = mapply(num_ealier, doctor, month)
  newdata = data.frame(doctor = doctor, month = month, that_month = that_month, before_mo

```

六列分别是医师人数，月份数，医师是否在该月开始开药，医师是否在该月开始开药，医师在该月开始开药的联系人数量以及人数。在该月及之前开始开处方的医师联系人中的百分比。2125行是 125×17，这是17个月中125位医生的信息。

3. Let

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

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.

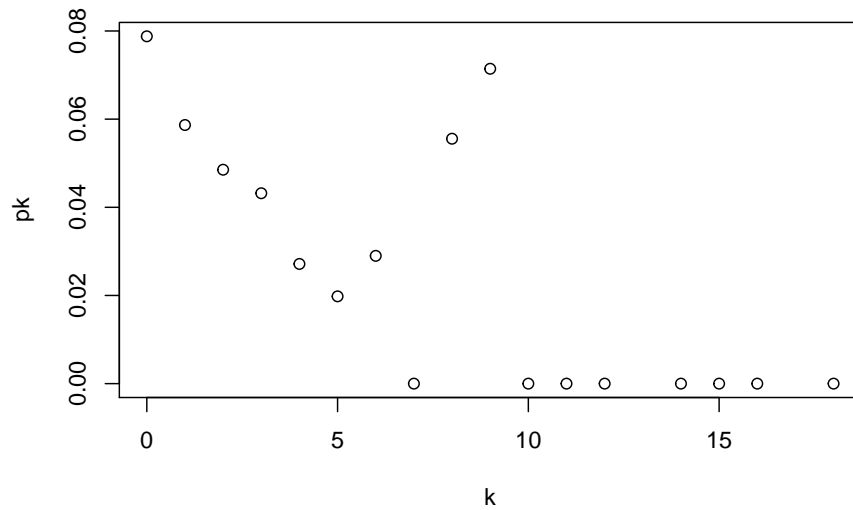
因为数据整理之后，每位医生最多有20位联系人。所以k有21可能(0:20)

b. Create a vector of estimated p_k probabilities, using the data frame from (2). Plot the probabilities against the number of prior-adopter contacts k .

```

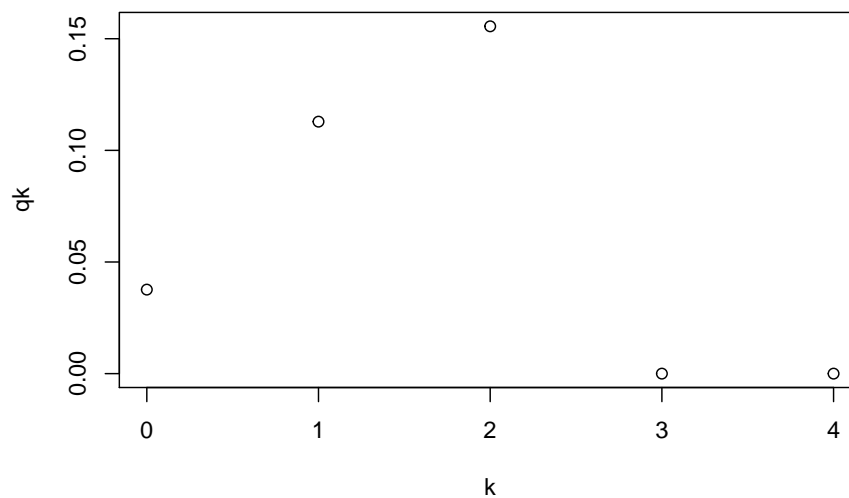
pk = rep(NA, 19)
for (i in 0:18) {
  pk[i+1] = sum(newdata$that_month[which(newdata$connects_before == i)]) / length(which(newdata$that_month == i))
}
plot(0:18, pk, xlab = "k")

```



b. Create a vector of estimated q_k probabilities, using the data frame from (2). Plot the probabilities against the number of prior-or-contemporary-adoptee contacts k .

```
n = max(newdata$connects_ealier-newdata$connects_before)
qk = rep(NA, n+1)
for (i in 0:n) {
  qk[i+1] = sum(newdata$that_month[which(newdata$connects_ealier-newdata$connects_before
}]
plot(0:n, qk, xlab = "k")
```



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.

```
k = c(0:18)
lm.fit = lm(pk~k)
summary(lm.fit)
```

```
##
## Call:
## lm(formula = pk ~ k)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
##	-0.030334	-0.014584	-0.002344	0.005534	0.048694

```
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.0569324  0.0090507   6.290 1.45e-05 ***
## k           -0.0037997  0.0009184  -4.137 0.000877 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.02015 on 15 degrees of freedom
## (2 observations deleted due to missingness)
## Multiple R-squared:  0.533, Adjusted R-squared:  0.5018
## F-statistic: 17.12 on 1 and 15 DF, p-value: 0.0008773
```

估计值: ($a = 0.0569324, b = -0.0037997$)。因为 $\text{pr}(>|t|)$ 小于0.05,所以拒绝 H_0 假设,所以有效。 b. Suppose $p_k = e^{a+bk}/(1 + e^{a+bk})$. Explain, in words, what this model would imply about the impact of adding one more adoptee friend on a given doctor's probability of adoption. (You can suppose that $b > 0$, if that makes it easier.) Estimate the model by least squares, using the values you constructed in (3b).

```
p1 = na.omit(pk)
k1 = which(is.na(pk) == FALSE) - 1
p2 = p1[p1 != 0]
k2 = k1[p1 != 0]
y = log(p2 / (1-p2))
lm.fit1 = lm(y~k2)
summary(lm.fit1)
```

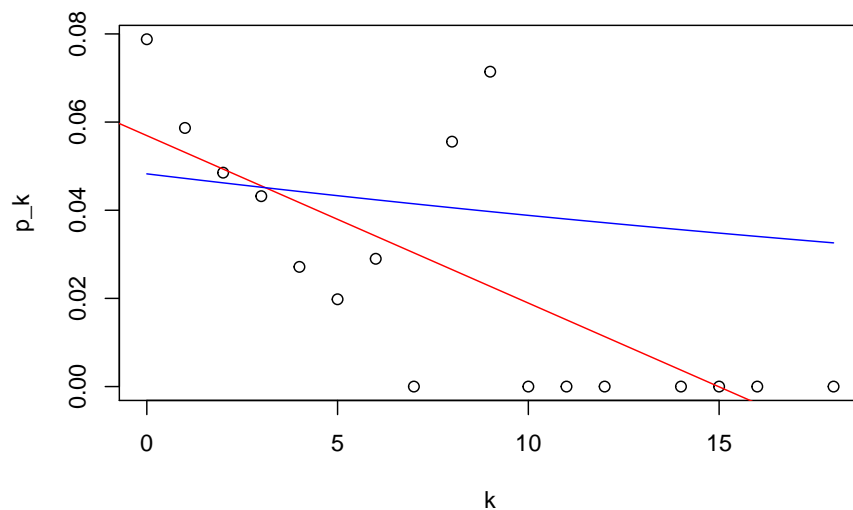
```
##
## Call:
## lm(formula = y ~ k2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.80666 -0.39353  0.05123  0.33021  0.62118
```

```
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2.98180    0.30592  -9.747 2.53e-05 ***
## k2           -0.02270    0.05974  -0.380  0.715
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5193 on 7 degrees of freedom
## Multiple R-squared:  0.02022,    Adjusted R-squared:  -0.1198
## F-statistic: 0.1444 on 1 and 7 DF,  p-value: 0.7152
```

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

```
plot(0:18, pk, xlab = "k", ylab = "p_k", xlim = c(0,18))
abline(lm.fit, col = "red")
```

```
curve(exp(-2.98180 - 0.02270*x) / (1 + exp(-2.98180 - 0.02270*x)), from = 0, to = 18, a
```



我认为(4b)的模型要比(4a)要好，因为非线性比较符合实际而且反应的信息比较多。

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?