

## **Homework 5, DATA 556: Due Tuesday, 10/31/2018**

**Alexander Van Roijen**

November 6, 2018

Please complete the following:

1. Problem 1 Let  $X$  and  $Y$  be i.i.d.  $\text{Geom}(p)$ , and  $N = X + Y$ .

(a) Find the joint PMF of  $X$ ,  $Y$  and  $N$ .

$$P(X = x, Y = y, N = n) = P(X = x, Y = y, X + Y = n) \quad (1)$$

we note that  $N$  acts as a qualifying condition. We get (2)

$$P(X = x, Y = y, N = n) = \begin{cases} P(X = x, Y = y) = p * (1 - p)^x * p(1 - p)^y & x + y = n \\ 0 & x + y \neq n \end{cases} \quad (3)$$

(b) Find the joint PMF of  $X$  and  $N$ .

$$P(X = x, N = n) = P(X = x, X + Y = n) = P(X = x, Y = n - x) \quad (4)$$

$$= \begin{cases} P(X = x) * P(Y = n - x) = p(1 - p)^x p(1 - p)^{n-x} = p^2(1 - p)^n & x \leq n \\ 0 & x > n \end{cases} \quad (5)$$

(c) Find the conditional PMF of  $X$  given  $N = n$ .

$$P(X = x | N = n) = \frac{P(X = x, N = n)}{P(N = n)} \quad (6)$$

$$P(N = n) = \text{Negative binomial with 2 successes} = \binom{n+2-1}{2-1} p^2(1-p)^n \quad (7)$$

$$P(X = x | N = n) = \begin{cases} \frac{p^2(1-p)^n}{\binom{n+1}{1} p^2(1-p)^n} = \frac{1}{n+1} & n \geq x \\ 0 & n < x \end{cases} \quad (8)$$

2. Let  $X$ ,  $Y$  and  $Z$  be random variables such that  $X \sim N(0, 1)$  and conditional on  $X = x$ ,  $Y$  and  $Z$  are i.i.d.  $N(x, 1)$ .

(a) Find the joint PDF of X, Y and Z.

$$f(x, y, z) = f(y, z|x) * f(x) \text{ and because of independence of Y and Z on x we get} \quad (9)$$

$$f(x, y, z) = f(y|x) * f(z|x) * f(x) \text{ now we shift } Y|x \text{ and } Z|x \text{ to a } N(0,1) \quad (10)$$

$$\text{Note further that shift wont alter the value under the curve, just where the density takes place} \quad (11)$$

$$f(x, y, z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-x)^2}{2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{(z-x)^2}{2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad (12)$$

$$\text{does this need to be simplified?} \quad (13)$$

(b) Find the joint PDF of Y and Z.

$$f(y, z) = \int_{-\infty}^{\infty} f(x, y, z) dx \quad (14)$$

3. Let X and Y be continuous random variables with joint CDF  $F(x, y)$ . Show that the probability that  $(X, Y)$  falls in the rectangle  $(a1, a2) (b1, b2)$  is

$$F(a2, b2) - F(a1, b2) + F(a1, b1) - F(a2, b1). \quad (15)$$

$$f(x, y) = \frac{d}{dx} \frac{d}{dy} F(x, y) \text{ want} \quad (16)$$

$$\int_{a1}^{a2} \int_{b1}^{b2} \frac{d}{dx} \frac{d}{dy} F(x, y) dy dx \quad (17)$$

$$= \int_{a1}^{a2} \frac{d}{dx} (F(x, b2) - F(x, b1)) dx = \int_{a1}^{a2} \frac{d}{dx} F(x, b2) dx - \int_{a1}^{a2} \frac{d}{dx} F(x, b1) dx \quad (18)$$

$$= F(a2, b2) - F(a1, b2) - F(a2, b1) + F(a1, b1) \square \quad (19)$$

4. Let X and Y have joint PDF

$$f(x, y) = x + y \text{ for } 0 < x < 1 \& 0 < y < 1 \quad (20)$$

(a) Check if this is a valid pdf

NTS

$$f(x, y) > 0 \forall x, y \in X, Y \quad (21)$$

$$\text{This is trivial as } 0 < x < 1 \& 0 < y < 1 \Rightarrow 0 < x + y < 2 \quad (22)$$

$$\Rightarrow f(x, y) = x + y > 0 \forall x, y \in X, Y \quad (23)$$

$$\text{And NTS } \int_0^1 \int_0^1 f(x, y) dy dx = 1 \int_0^1 \int_0^1 f(x, y) dy dx = \int_0^1 \int_0^1 (x + y) dy dx \quad (24)$$

$$= \int_0^1 (x + \frac{1}{2}) dx = \int_0^1 (x + \frac{1}{2}) dx = \frac{1}{2} + \frac{1}{2} = 1 \square \quad (25)$$

(b) Find the marginal PDFs of X and Y.

$$f(x) = \int_0^1 f(x, y) dy = \int_0^1 (x + y) dy = yx + \frac{1}{2}y^2 \Big|_0^1 = x + \frac{1}{2} \quad (26)$$

$$f(y) = \int_0^1 f(x, y) dx = \int_0^1 (x + y) dx = yx + \frac{1}{2}x^2 \Big|_0^1 = y + \frac{1}{2} \quad (27)$$

(c) Are X and Y independent?

$$\text{No, as we would need } f(x, y) = f(x) * f(y) \text{ and} \quad (28)$$

$$f(x) * f(y) = (y + \frac{1}{2}) * (x + \frac{1}{2}) \neq x + y \quad (29)$$

(d) Find the conditional PDF of Y given X = x.

$$f(y|X = x) = \frac{f(x, y)}{f(x)} = \frac{x + y}{x + \frac{1}{2}} \quad (30)$$

5. Let X and Y have joint PDF

$$f(x, y) = cxy \text{ for } 0 < x < y < 1 \quad (31)$$

(a) find c to make this a valid pdf

$$f(x, y) > 0 \forall x, y \in X, Y \quad (32)$$

$$\text{This is trivial as } 0 < x < y < 1 \Rightarrow x > 0 \& y > 0 \quad (33)$$

$$\Rightarrow f(x, y) = cxy > 0 \forall x, y \in X, Y \text{ if } c > 0 \quad (34)$$

$$\text{And NTS } \int_0^1 \int_0^y f(x, y) dx dy = 1 \int_0^1 \int_0^y f(x, y) dx dy = \int_0^1 \int_0^y (cxy) dy dx \quad (35)$$

$$= \int_0^1 (\frac{cx^2y}{2} \Big|_0^y) dy = \int_0^1 (\frac{cy^3}{2}) dy = \frac{cy^4}{8} \Big|_0^1 = 1 \quad (36)$$

$$\Rightarrow c = 8 \quad (37)$$

(b) Find the marginal PDFs of X and Y

$$f(x) = \int_x^1 f(x, y) dy = \int_x^1 (8xy) dy = 4xy^2 \Big|_x^1 = 4x - 4x^3 \quad (38)$$

$$f(y) = \int_0^y f(x, y) dy = \int_0^y (8xy) dy = 4x^2 y \Big|_0^y = 4y^3 \quad (39)$$

(c) are X and Y independent?

No, as once again we have (40)

$$f(x) * f(y) = ((4x - 4x^3) * 4y^3) \neq 8xy = f(x, y) \quad (41)$$

(d) Find the conditional PDF of Y given X = x.

$$f(y|x) = \frac{f(x, y)}{f(x)} = \frac{8xy}{4x - 4x^3} = \frac{2y}{1 - x^2} \quad (42)$$

6. Let X and Y be i.i.d. Unif(0, 1).

(a) Use simulations in R (the statistical programming language) to numerically estimate the covariance of X + Y and X - Y.

```
> prob6a = function(n)
+ {
+   x = runif(n,0,1)
+   y = runif(n,0,1)
+   xmy = x-y
+   xpy = x+y
+   return(cov(xmy,xpy))
+ }
> #problem 6a
> print(prob6a(100000))
[1] 0.0005779827
```

(b) Compute the covariance of X + Y and X - Y

$$\text{Cov}[X + Y, X - Y] = E[(X + Y)(X - Y)] - E[X + Y]E[X - Y] = \quad (43)$$

$$E[X^2] - E[Y^2] - (E[X] + E[Y]) * (E[X] - E[Y]) \text{ by independence} \quad (44)$$

$$= \text{Var}(X) + E[X]^2 - \text{Var}(Y) - E[Y]^2 - E[X]^2 + E[Y]^2 = 0 \quad (45)$$

the above is true as X and Y have the same distribution and thus cancel (46)

(c) Are  $X + Y$  and  $X - Y$  independent

We know the following (47)

$$M_{X+Y}(t) = M_X(t)M_Y(t) \text{ \& } M_{X-Y}(t) = M_X(t)M_{-Y}(t) \text{ Due to independence between X and Y} \quad (48)$$

$$\text{Now we want to show: } M_{X+Y, X-Y}(s, t) = M_{X+Y}(s)M_{X-Y}(t) \quad (49)$$

$$M_{X+Y}(s)M_{X-Y}(t) = M_X(t)M_Y(t)M_X(t)M_{-Y}(t) = M_X(t)M_Y(t)M_X(t)M_Y(-t) \quad (50)$$

$$\text{the above is true as you can assign the negative sign} \quad (51)$$

$$\text{from the definition of the MGF to the 't' or R.V.} \quad (52)$$

$$M_{X+Y, X-Y}(s, t) = E[e^{s(X+Y)+t(X-Y)}] = E[e^{X(s+t)+Y(s-t)}] = M_X(s+t)M_Y(s-t) \text{ by indep} \quad (53)$$

$$= M_X(s)M_X(t)M_Y(s)M_Y(-t) \quad (54)$$

$$\text{Thus, since the MGFs are the same, we have sufficiently shown independence between} \quad (55)$$

$$X + Y \text{ and } X - Y \quad (56)$$

7. Let  $X$ ,  $Y$  and  $Z$  be i.i.d.  $N(0, 1)$ . Find the joint MGF of  $(X + 2Y, 3X + 4Z, 5Y + 6Z)$ .

First Notice that we can rewrite this as a Multivariate normal distribution

$$t_1(X + 2Y) + t_2(3X + 4Z) + t_3(5Y + 6Z) = X(t_1 + 3t_2) + Y(2t_1 + 5t_3) + Z(4t_2 + 6t_3) \quad (57)$$

$$\text{We can use some properties to then get the following} \quad (58)$$

$$MGF(X + 2Y, 3X + 4Z, 5Y + 6Z) = E[e^{t_1(X+2Y)+t_2(3X+4Z)+t_3(5Y+6Z)}] \quad (59)$$

$$= e^{t_1 E(X+2Y)+t_2 E(3X+4Z)+t_3 E(5Y+6Z)+\frac{1}{2} \text{Var}(t_1(X+2Y)+t_2(3X+4Z)+t_3(5Y+6Z))} \quad (60)$$

$$\text{by the definition of joint MGF on the multivariate normal} \quad (61)$$

$$= e^{\frac{1}{2} \text{Var}(t_1(X+2Y)+t_2(3X+4Z)+t_3(5Y+6Z))} \text{ since we have } N(0,1) \text{ for } X,Y,Z \quad (62)$$

$$= e^{\frac{1}{2} \text{Var}(X(t_1+3t_2)+Y(2t_1+5t_3)+Z(4t_2+6t_3))} = e^{\frac{1}{2}(t_1+3t_2)^2 \text{Var}(X)+(2t_1+5t_3)^2 \text{Var}(Y)+(4t_2+6t_3)^2 \text{Var}(Z)} \quad (63)$$

$$= e^{\frac{1}{2}(t_1+3t_2)^2+(2t_1+5t_3)^2+(4t_2+6t_3)^2} \quad (64)$$

8. we have the table below

```
> y = matrix(c(0.018, 0.035 , 0.031 ,0.008 , 0.018 ,
+              0.002,0.112 ,0.064,0.032,0.069,
+              0.001,0.066 , 0.094 ,0.032, 0.084,
+              0.001 , 0.018, 0.019, 0.010, 0.051,
+              0.001, 0.029 , 0.032,0.043,0.130), nrow=5 , byrow=TRUE)
> colnames ( y ) = c ( 'farm' , "operatives",'craftsen','sales','professional' )
> rownames ( y ) = colnames ( y )
>
> sum( y )
[1] 1
> print(y)
farm operatives craftsmen sales professional
farm          0.018      0.035    0.031 0.008      0.018
operatives    0.002      0.112    0.064 0.032      0.069
craftsen      0.001      0.066    0.094 0.032      0.084
sales         0.001      0.018    0.019 0.010      0.051
professional 0.001      0.029    0.032 0.043      0.130
```

(a) the marginal probability distribution of a father's occupation

```
sum(y[1,])
#[1] 0.11
sum(y[2,])
#[1] 0.279
sum(y[3,])
#[1] 0.277
sum(y[4,])
#[1] 0.099
sum(y[5,])
#[1] 0.235
```

- (b) the marginal probability distribution of a son's occupation

```
sum(y[,1])
#[1] 0.023
sum(y[,2])
#[1] 0.26
sum(y[,3])
#[1] 0.24
sum(y[,4])
#[1] 0.125
sum(y[,5])
#[1] 0.352
```

- (c) the conditional distribution of a son's occupation, given that the father is a farmer

```
> result = y[1,]/sum(y[1,])
> print(result)
farm      operatives      craftsmen      sales professional
0.16363636  0.31818182  0.28181818  0.07272727  0.16363636
```

- (d) the conditional distribution of a father's occupation, given that the son is a farmer

```
> resultd = y[,1]/sum(y[,1])
> print(resultd)
farm      operatives      craftsmen      sales professional
0.78260870  0.08695652  0.04347826  0.04347826  0.04347826
```

9. You will analyze data from a study of the effects of aspirin on myocardial infarction.

- (a) Calculate the row and columns totals of this table. What is the grand total?

```
#placebo = 28+656=684
#aspirin = 18+658=676
#yes = 28+18 = 46
#no = 656+658 = 1314
placebo = 684
```



```

aspirin = 676
yes = 46
no = 1314
#grandtotal = 1360
grandtotal = yes+no

```

- (b) Calculate the expected cell values under the hypothesis of interaction of Aspirin Use and Myocardial Infarction.

```

> table9 = matrix(c(28,656,18,658), nrow=2 , byrow=TRUE)
> colnames(table9) = c('yes','no')
> rownames(table9) = c('Placebo','aspirin')
> table9
yes  no
Placebo  28 656
aspirin  18 658
> #under interaction we have
> dcell11= table9[1,1]/grandtotal
> dcell12= table9[1,2]/grandtotal
> dcell21= table9[2,1]/grandtotal
> dcell22= table9[2,2]/grandtotal
> dptable = matrix(c(dcell11,dcell12,dcell21,dcell22),nrow=2,byrow=TRUE)
> dpttable = dptable*grandtotal
> colnames(dpttable) = c('yes','no')
> rownames(dpttable) = c('Placebo','aspirin')
> print(dpttable)
yes  no
Placebo  28 656
aspirin  18 658

```

- (c) Calculate the expected cell values under the hypothesis of independence of Aspirin Use and Myocardial Infarction.

```

> #under no interaction, we have a simple set of multiplications
> p1a=sum(table9[1,])/grandtotal
> p2a=sum(table9[2,])/grandtotal
> pa1=sum(table9[,1])/grandtotal
> pa2=sum(table9[,2])/grandtotal
> icell11=p1a*pa1
> icell12=p1a*pa2
> icell21=p2a*pa1
> icell22=p2a*pa2
> # since each cell is a bernoulli RV, our expected value is simply the probability
> iptable = matrix(c(icell11,icell12,icell21,icell22),nrow=2,byrow=TRUE)
> ipttable = iptable*grandtotal
> colnames(ipttable) = c('yes','no')
> rownames(ipttable) = c('Placebo','aspirin')
> print(ipttable)
yes      no
Placebo 23.13529 660.8647
aspirin 22.86471 653.1353

```

- (d) Perform an asymptotic test of independence vs. interaction of Aspirin Use and Myocardial Infarction based on Pearson's chi-square statistic

```

> chiqTest = function(truMatrix,expectedMatrix)
+ {
+   sum =0
+   rcounter = 1
+   ccounter = 1
+   sum = 0
+   while(rcounter <= 2)
+   {
+     ccounter = 1
+     while(ccounter<=2)

```

```

+   {
+       sum = sum + (((truMatrix[rcounter,ccounter]-expectedMatrix[rcounter,ccounter]
+       ccounter= ccounter + 1
+   }
+   rcounter = rcounter +1
+ }
+ return(sum)
+ }

> test1 = (chiqTest(table9,ipptable)) #ipptable is interaction table
> print(test1)
[1] 2.129972

```

- (e) Perform an asymptotic test of independence vs. interaction of Aspirin Use and Myocardial Infarction based on the likelihood ratio statistic  $G^2$

```

> gsquareTest = function(truMatrix,expectedMatrix)
+ {
+   sum =0
+   rcounter = 1
+   ccounter = 1
+   sum = 0
+   while(rcounter <= 2)
+   {
+       ccounter = 1
+       while(ccounter<=2)
+       {
+           sum = sum + (truMatrix[rcounter,ccounter]*log(truMatrix[rcounter,ccounter]/e
+           ccounter= ccounter + 1
+       }
+       rcounter = rcounter +1
+   }
+   return(2*sum)
+ }

```

```

+ }
> test2 = (gsquareTest(table9,ippttable)) #ipptable is interaction table
> print(test2)
[1] 2.147353

```

- (f) Draw conclusions related to the effect of aspirin on the occurrence of myocardial infarction. Summarize your findings in a concise statement

Now, in this scenario where we believe there to be interaction, we have 1 degree of freedom as  $(rows - 1)(cols - 1) = (2 - 1) * (2 - 1) = 1$  Thus we get the following

```

> print(1-pchisq(test1,1))
[1] 0.1444434
> print(1-pchisq(test2,1))
[1] 0.1428159

```

The similarity of the two tests are encouraging and the values they display indicate that we can not reject the null hypothesis assuming an alpha level of 0.05. More plainly, there is not enough evidence to indicate that there is not an independence for myocardial infarctions between aspirin usage and placebos with high confidence.

Have a nice day!