

STAT 513 hw 9

1. Suppose a study of the efficacy of a treatment results in the following table of outcomes:

	Successes	Failures	Total
Treatment	20	5	25
Control	10	16	26
Total	30	21	51

(a) State the null and alternate hypotheses which are likely of interest to the researchers.

We wish to test $H_0: p_1 = p_2$ versus $H_1: p_1 \neq p_2$, where p_1 and p_2 are the probabilities of success in the treatment and control groups. We could also formulate the hypotheses as H_0 : *there is no association between the treatment and success* versus H_1 : *there is an association between treatment and success*.

(b) i. Compute the test statistic for the likelihood ratio test of your hypotheses from part (a).
ii. Give the p -value.

The following R code gives the answer:

```
# create matrices of observed and expected counts
O <- matrix(c(20,5,10,16),2,byrow=TRUE)
E <- apply(O,1,sum) %*% t(apply(O,2,sum))/sum(O)

# get LRT test statistic and p-value
G <- 2*sum(O*log(O/E))
1 - pchisq(G,1)
```

The test statistic value is 9.437829 and the p -value is 0.002125549.

(c) i. Compute the test statistic for Pearson's chi-squared test of your hypotheses from part (a).
ii. Give the p -value.

```
# get Pearson's chi-squared test statistic and p-value
Chisq <- sum((O - E)^2/E)
1 - pchisq(Chisq,1)
```

The test statistic value is 9.079121 and the p -value is 0.002585443.

(d) Give the p -value of Fisher's exact test of your hypotheses in part (a).

The following R code gives the answer:

```

# perform Fisher's exact test
Xobs <- O[1,1]
R1 <- O[1,1] + O[1,2]
C1 <- O[1,1] + O[2,1]
N <- sum(O)

obs.hyper.prob <- dhyper(Xobs, n = N - R1, m = R1, k = C1)
all.hyper.probs <- dhyper(max(0,C1-(N-R1)):min(C1,R1),n = N - R1, m = R1, k = C1)
sum(all.hyper.probs[all.hyper.probs<=obs.hyper.prob])

```

The p -value is 0.004171008.

2. The following knee injury data in women collegiate rugby players is taken from [2]. It is of interest to know whether the types of injuries a player experiences are associated with the position (Forward, Back) she plays.

	Meniscal Tear	MCL Tear	ACL Tear	Other
Forward	13	14	7	4
Back	12	9	14	3

- (a) i. Compute the test statistic for the likelihood ratio test of association.
 ii. Give the p -value (make sure you choose the right degrees of freedom!).

The following R code gives the answer:

```

# create matrices of observed and expected counts
O <- matrix(c(20,5,10,16),2,byrow=TRUE)
E <- apply(O,1,sum) %*% t(apply(O,2,sum))/sum(O)

# get LRT test statistic and p-value
G <- 2*sum(O*log(O/E))
1 - pchisq(G,3)

```

The test statistic value is 3.657628 and the p -value is 0.3008863.

- (b) i. Compute the test statistic for Pearson's chi-squared test of association.
 ii. Give the p -value.

```

# get Pearson's chi-squared test statistic and p-value
Chisq <- sum((O - E)^2/E)
1 - pchisq(Chisq,3)

```

The test statistic value is 3.603147 and the p -value is 0.3076286.

3. Consider the following data taken from [1], which result from looking through a microscope at samples of milk film and counting the number of bacterial colonies within the field of vision. A total of 400 observations were gathered and the number of bacterial colonies was recorded for each of them:

# Bacterial Colonies	# Microscopic fields
0	56
1	104
2	80
3	62
4	42
5	27
6	9
7	9
8	5
9	3
10	2
19	1

- (a) Let X_i represent the number of bacterial colonies in microscopic field i , with $i = 1, \dots, 400$. Assume for a moment that the number of bacterial colonies in a microscopic field follows the $\text{Poisson}(\lambda)$ distribution for some value of λ . Compute the maximum likelihood estimator of λ based on the above data.

The maximum likelihood estimator is the sample mean. The following R code computes the mean number of bacterial colonies:

```
numBact <- c(0,1,2,3,4,5,6,7,8,9,10,19)
numMicro <- c(56,104,80,62,42,27,9,9,5,3,2,1)
lambda.hat <- sum(numBact * numMicro)/400
```

The sample mean is 2.44, so the maximum likelihood estimator of λ is $\hat{\lambda} = 2.44$.

- (b) If the bacterial colony counts truly followed a Poisson distribution with mean equal to $\hat{\lambda}$, where $\hat{\lambda}$ is the maximum likelihood estimator of λ computed in part (a), what would be the expected # of microscopic fields corresponding to each number of bacterial counts? That is, in how many microscopic fields out of 400 would we expect to see 0 bacterial colonies, 1 bacterial colony, and so on? Make a table like the table above, but with the numbers in the right-hand column replaced by the expected numbers of microscopic fields. *Hint: the first one is `dpois(0,lambda.hat)*400`.*

We can get the expected numbers of microscopic fields with the R command `round(dpois(numBact,lambda.hat)*400,2)`, which also rounds them to the nearest 1/100.

# Bacterial Colonies	# Microscopic fields
0	34.86
1	85.07
2	103.78
3	84.41
4	51.49
5	25.13
6	10.22
7	3.56
8	1.09
9	0.29
10	0.07
19	0.00

- (c) From looking at these numbers, do you believe the # of bacterial colonies follows a Poisson distribution? Explain your answer.

The expected counts seem quite different from the observed counts, so it seems unlikely that the # of bacterial colonies follows a Poisson distribution.

- (d) Pearson’s chi-squared test is often used to test for what is called the “goodness-of-fit” of a probability distribution to some observed data, as this test statistic provides a useful way to compare observed counts to expected counts. Compute the test statistic of Pearson’s chi-squared test on these data, which is given by

$$\sum_{i=1}^{12} (O_i - E_i)^2 / E_i,$$

where O_1, \dots, O_{12} are observed numbers of microscopic fields and E_1, \dots, E_{12} are the expected numbers of microscopic fields according to the $\text{Poisson}(2.44)$ distribution. *Hint: You get a crazy-huge number.*

The following R code computes the test statistic:

```
O <- numMicro
E <- round(dpois(numBact, lambda.hat)*400, 2)
sum( (O - E)^2/E)
```

We get 152162218, or infinity if we have rounded the expected counts.

- (e) It was naïve of us to compute Pearson’s test statistic in the previous part, because some of the expected counts are very small, almost equal to zero; recall that we require expected counts to be greater than or equal to 5 in order to use Pearson’s chi-squared test. What can we do? Let’s collapse the last few rows of the tables by summing together the rows for which the # of bacterial colonies is greater than or equal to 6, so that we have the following:

# Bacterial Colonies	# Microscopic fields	$\mathbb{E}[\# \text{ Microscopic fields }]$
0	56	34.86
1	104	85.07
2	80	103.78
3	62	84.41
4	42	51.49
5	27	25.13
≥ 6	29	15.23

Recompute the test statistic for Pearson's chi-squared test, this time with

$$\sum_{i=1}^7 (O_i - E_i)^2 / E_i.$$

The following R code computes the test statistic.

```
O.trunc <- c(O[1:6],sum(O[-c(1:6)]))
E.trunc <- c(E[1:6],sum(E[-c(1:6)]))
sum((O.trunc - E.trunc)^2/E.trunc)
```

We get the value 42.75711.

- (f) Under the null hypothesis, the test statistic converges in distribution to a random variable with the χ_6^2 distribution, since we are considering a table with 7 cells in a column, and $7 - 1 = 6$. Use this information to compute a p -value for the test statistic computed in part (e). Use the p -value to make a conclusion about whether or not the # of bacterial colonies follows a Poisson distribution.

We get the p -value $1 - \text{pchisq}(42.75711, 6) = 1.302903 \times 10^{-7}$, so we reject the null hypothesis that the # of bacterial colonies follows a Poisson distribution.

References

- [1] Chester Ittner Bliss and Ronald A Fisher. Fitting the negative binomial distribution to biological data. *Biometrics*, 9(2):176–200, 1953.
- [2] Andrew S. Levy, Merrick J. Wetzler, Marie Lewars, and William Laughlin. Knee injuries in women collegiate rugby players. *The American journal of sports medicine*, 25(3):360–362, 1997.