

# MATH-493 Applied Biostatistics - Project 3

## Statistical study of Anti-Typhoid Inoculation

Riccardo Cadei

### Introduction

In 1907, Major H.J.M. Buist published an article on BMJ Military Health where he summarized a considerable number of data regarding anti-typhoid inoculation from a census of the army in the British Empire and colonies in those years [1] wondering about the effectiveness of this treatment. His thesis was that the new anti-typhoid vaccine was showing a safe and valuable prophylactic against enteric fever, both in preventing attack and lowering the case mortality. However, the limits of his analysis are, now more than ever, evident: the heterogeneity of the data collected and the poor statistical analysis conducted (mainly comparing percentages without using any statistical result).

Suspiciously, less than 2 years later, G.D. Maynard reconsidered the data collected by Major Buist and tried to answer the same research question using a more sophisticated Discrete Data Analysis. Studying the contingency tables and the coefficients of correlation between inoculations and freedom from attack in several armies dislocated around the British colonies, he found out that the inoculation effectiveness among different location were too uneven to be explained just by randomness effects. So, he focused on locations with a weaker correlation and he hypothesized that there was more than a disease covered by the term *enteric* or *typhoid fever*. In fact, bacteriologically it was known that several varieties of para-typhoid bacilli existed and were the specific cause of fevers clinically indistinguishable from typhoid; and there was no a priori reason for assuming that a vaccine prepared from typhoid bacilli could confer immunity against infection with any of the para-typhoid strains [4]. Taking into account supplementary data approximating the ratio of para- to true typhoid infections, Maynard showed that, in regions where the para-typhoid were more distributed, the correlation between inoculation and freedom from attack correspondingly lowered.

Today, more a 100 years later, I reconsider the same contingency tables studied by G.D. Maynard and I wonder the same research question, trying to answer with modern Discrete Data Analysis tools.

This report is structured as follows: in Section 1 the data are presented together with some graphical representations, in Section 2 different statistical test are used to study independence and homogeneity of the different samples, in Section 3 the results found in previous analysis are summarized and finally in Section 4 a conclusive comparison with the previous works conducted by Major Buist and G.D. Maynard is reported.

# 1 Exploratory Data Analysis

## 1.1 Dataset

The data set that I will consider for the following analysis consists in the same contingency tables analyzed by G.D. Maynard. In particular, the army census around 1906 and 1907 in the following regions are taken into account: Stations Abroad, Indian Stations, 17th Lancers, Coldstream Guards, Transvaal and 7 Large Indian Stations; and for each region the amount of people who have been inoculated, and the amount of people who have been diagnosed attacked by typhoid is considered. All the six 2x2 contingency tables are reported in Table 1.

Location	Inoculation	Typhoid	
		Attacked	Not Attacked
Stations Abroad	Yes	5	786
	No	193	30 757
Indian Stations	Yes	8	2122
	No	770	37 113
17th Lancers	Yes	2	148
	No	58	481
Coldstream Guards	Yes	1	330
	No	13	368
Transvaal Diserticts	Yes	5	219
	No	65	6 690
7 Large Indian Stations	Yes	15	2 192
	No	173	7 940

Table 1: Contingency tables in `typhoidinoc.dat`

Few notes has to be attached with the dataset:

- Stations Abroad and Indian Stations tables are admittedly incomplete and it is not certain that all cases of inoculation were corrected returned as such,
- 17th Lancers and Coldstream Guards tables refer to local outbreaks in two regiments following their removal abroad,
- 7 Large Indian Stations table is stated by Major Buist as the most complete.

In add to this, it is evident that the total number of counts varies a lot among the locations, in fact the 44.23% of total observations comes from the Indian Stations, while only the 0.7% belongs to 17<sup>th</sup> Lancers. This uncontrolled heterogeneity is the main limit to which one must go against when dealing with observational studies and experiment design cannot be formulated *a priori*.

## 1.2 Mosaic Plots

Mosaic plots are one of the best way to quickly visualize the main associations in a contingency table. In Figure 1 I report the mosaic plot of the overall contingency table joining together all the locations considered in this study. Two phenomena are strongly evident:

- The dataset is strongly unbalanced. Only few participants received the inoculation, even less were attacked by Typhoid and just few units both received the inoculation and were attacked.
- The percentage of participants who were attacked by typhoid taking the inoculation is smaller than the percentage of participants who were attacked by typhoid and not inoculated.

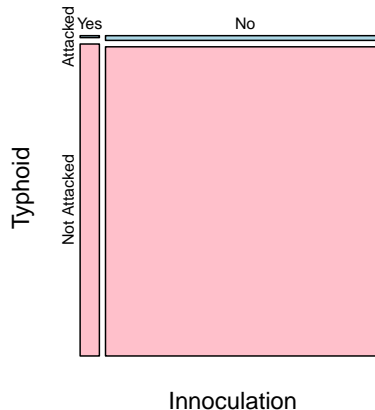


Figure 1: Overall Mosaic Plot

The first observation was already clear from the Table 1, while the second observation is just observed. Stopping here the analysis one could claim that the inoculation reduces the probability to be attacked by typhoid, but as a statistician, I decide to explore more in depth the dataset before to claim such a general statement.

I then report in Figure 2 the Mosaic plot for each location individually (a unique 3-way mosaic plot wouldn't be readable given that the number of total counts varies too much among different locations).

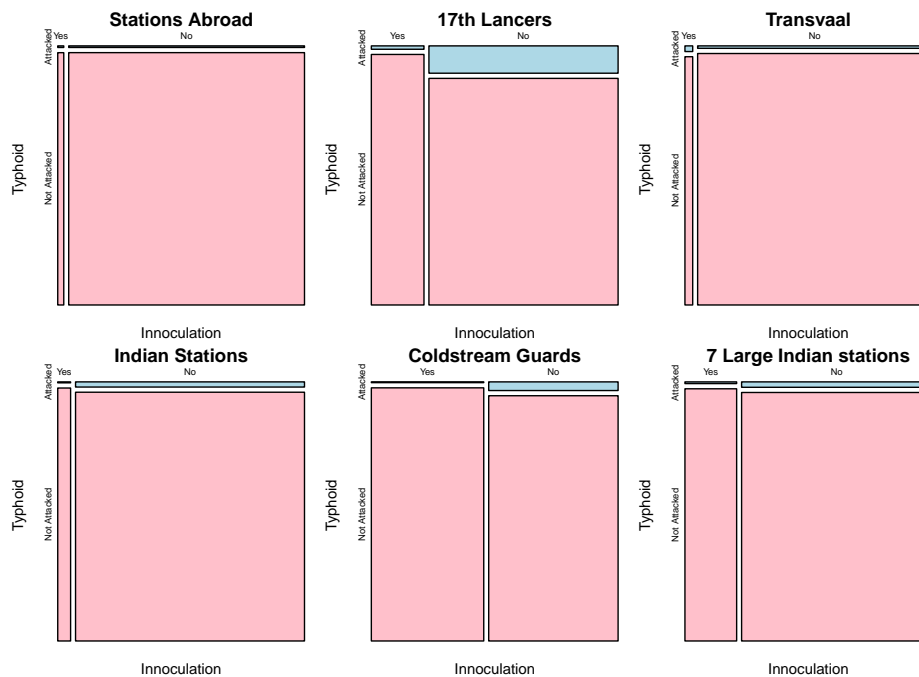


Figure 2: Per Location Mosaic Plots

What is immediately evident from these plots is that the distributions are not all equal among different locations. First of all, the percentages of participants exposed to the inoculation are very different (further limit dealing with an observational study): in Coldstream Guards almost the 50% of the sample were inoculated, in 17<sup>th</sup> Lancers and 7 Large Indian Stations around a fifth, while in Stations Abroad, Indian Stations and Transvaal only a few

percentage points. Then also the associations between inoculation and freedom from attack are different: in Indian Stations, 17<sup>th</sup> Lancers, Coldstream Guards and 7 Large Indian Stations it is positive (the inoculated and attacked are less, in percentage, than the not inoculated and attacked), in Stations Abroad the 2 variables look independent and finally in Tranvaal the association is negative. A third observation is that, even in the same group (treated or no treated) the percentage of attacked vary a lot with the location; in particular the percentage of attacked (among the not inoculated) in 17<sup>th</sup> Lancers is remarkably bigger than in the other stations. These observations are enough to suspect that the global association described by the total mosaic plot is not representing faithfully all the samples considered, and further quantitative analyses are conducted in the following section.

## 2 Discrete Data Analysis

### 2.1 Independence

Two important statistical measurements of the association among 2 nominal values in a 2x2 contingency table are the *Odds Ratio* and the *Relative Risk*. Let  $n_{ij}$  the element in row  $i$  and column  $j$ , then:

$$\text{Odds Ratio} = \frac{n_{11} \cdot n_{22}}{n_{12} \cdot n_{21}} \quad \text{Relative Risk} = \frac{\frac{n_{11}}{n_{11}+n_{12}}}{\frac{n_{21}}{n_{21}+n_{22}}}$$

The first three columns in Table 2 report the *Odds Ratio* and the *Relative Risk* for all the locations considered in this study. As expected from the Mosaic plots, only four locations

Location	<i>Odds Ratio</i>	<i>Relative Risk</i>	Pearson's $\chi^2$ test <i>p</i> -value ( $\chi^2$ )	Fisher's exact test <i>p</i> -value
Stations Abroad	1.01	1.01	0.999 (3.36 · 10 <sup>-25</sup> )	0.8205
Indian Stations	0.18	0.18	1.106 · 10 <sup>-7</sup> (28.17)	1.606 · 10 <sup>-10</sup>
17th Lancers	0.11	0.12	5.436 · 10 <sup>-4</sup> (11.95)	5.659 · 10 <sup>-5</sup>
Coldstream Guards	0.09	0.09	6.718 · 10 <sup>-3</sup> (7.34)	2.253 · 10 <sup>-3</sup>
Transvaal	2.35	2.32	0.1246 (2.35)	0.0734
7 Large Indian stations	0.31	0.32	9.21 · 10 <sup>-6</sup> (19.66)	8.987 · 10 <sup>-7</sup>

Table 2: Statistical measurements (*Odds Ratio*, *Relative Risk*, Pearson  $\chi^2$  test with Yates continuity correction and Fisher exact test) of the association among inoculation and typhoid attack for each location.

have a positive association among inoculation and freedom from attack (*Odds Ratio* < 1), in Stations Abroad seems that the two variable are independent (*Odds Ratio*  $\approx$  1), and in Transvaal the association is opposite (*Odds Ratio* > 1).

An high impact representation of these associations is reported by the FourFold plots in Figure 3 [3]. All the observations about associations deducted by the Mosaic plots are confirmed by the FourFold plots: in Stations Abroad the 95% Confidence Interval rings in adjacent quadrants overlap almost perfectly (independence), in Indian Stations, 17th Lancers, Coldstream Guards and 7 Large Indian stations a strong association among inoculation and freedom from attack is evident, while in Transvaal the trend is almost the opposite but not very conclusive.

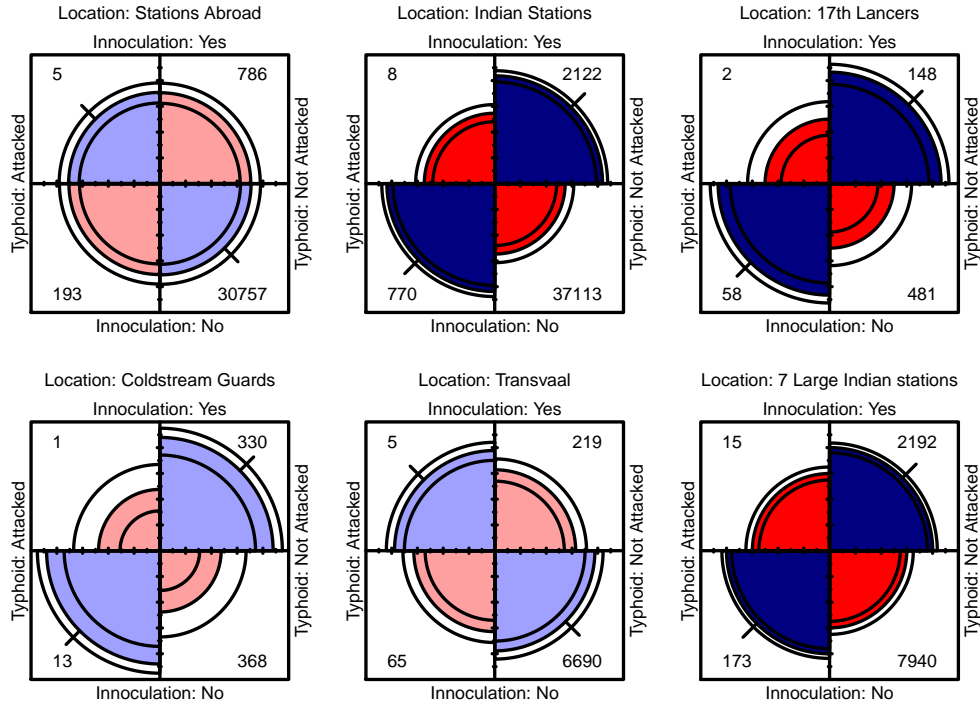


Figure 3: Per Location FourFold Plots

Then, for each location, I verified if the two variables are independent using the Pearson's  $\chi^2$  test. In formula, I tested:

$$\mathcal{H}_0 : \text{Inoculation} \perp\!\!\!\perp \text{Typhoid} \quad \text{vs} \quad \mathcal{H}_1 : \text{Inoculation} \not\perp\!\!\!\perp \text{Typhoid}$$

knowing that under  $\mathcal{H}_0$ , the statistic  $\sum_{i=1}^2 \sum_{j=1}^2 \frac{(n_{ij} - n_{i\bullet}n_{\bullet j}/n_{\bullet\bullet})^2}{n_{i\bullet}n_{\bullet j}/n_{\bullet\bullet}}$  follows a  $\chi_1^2$ . Since in almost all the samples the number of both inoculated and attacked is very small (e.g. 1 in Coldstream Guards), actually I used a slightly modified version of Pearson's  $\chi^2$  test proposed by Yates [5] for contingency tables involving small numbers, considering the corrected statistic  $\sum_{i=1}^2 \sum_{j=1}^2 \frac{(|n_{ij} - n_{i\bullet}n_{\bullet j}/n_{\bullet\bullet}| - 0.5)^2}{n_{i\bullet}n_{\bullet j}/n_{\bullet\bullet}}$ . For completeness I have also considered the Fisher's exact test: an exact test based on hyper-geometric distribution, referring to the same hypotheses of Pearson's test and used when at least a cell in the contingency table is very small. For all the locations the  $p$ -values of both the tests are reported in the last two columns of Figure 2. For the same four locations, both the tests reject the null hypothesis that the treatment and effect are independent ( $p$ -value  $< 0.05$ ), while in Stations Abroad and Transvaal both the test agree that the null hypothesis cannot be discarded ( $p$ -value  $> 0.05$ ).

## 2.2 Homogeneity

The main conclusion from the previous section is that different locations show different associations between the treatment and the effect. Then I want to directly test if there is statistical evidence to claim that, actually, the associations among the 2 variables are equal for all the locations. Using the Woolf's test on the whole dataset on the null hypothesis

that all the *Odds Ratio* are equal I found that, as expected, the null hypothesis of no 3-way association has to be discarded ( $\chi^2 = 29.39$ ,  $p\text{-value} = 1.944 \cdot 10^{-5}$ ). However, what was already suggested both by the Mosaic plots and the Fourfold plots is that maybe the locations could be split in classes where in each class all the data follow the same distribution. So I decide to evaluate also the homogeneity among all the couples of locations to see if there are some clusters. The  $p$ -values of the Woolf's test on all the couples of locations are reported in Table 3.

Indian St.	$2.905 \cdot 10^{-3}$				
17th Lancers	0.0100	0.5497			
Coldstream G.	0.0296	0.4949	0.8331		
Transvaal	0.198	$1.387 \cdot 10^{-5}$	$4.271 \cdot 10^{-4}$	$3.729 \cdot 10^{-3}$	
7 L. Indian St.	0.0266	0.2209	0.1831	0.2273	$2.015 \cdot 10^{-4}$
	St. Abroad	Indian St.	17th Lancers	Coldstream G.	Transvaal

Table 3: Woolf's test per couples of locations (each cell represents the  $p$ -value of the Woolf's test among the row and column locations)

The coloured cells in the table underline all the significant similarity. Two main clusters can be identified:

- **Group A:** Indian Stations, 17th Lancers, Coldstream Guards and 7 L. Indian stations
- **Group B:** Stations Abroad and Transvaal

In fact, all the locations in Group A, as well in Group B, are homogeneous by couples ( $p\text{-value} > 0.05$ ). This result is also confirmed by the Woolf's test on each group (Group A:  $\chi^2 = 3.6473$ ,  $p\text{-value} = 0.3022$ ; Group B:  $\chi^2 = 1.6568$ ,  $p\text{-value} = 0.198$ ). Since the hypothesis of no 3-way association within each group holds, I also propose the Cochran-Mantel-Haenzel's  $\chi^2$  test (CMH's test) within each group to study if the two nominal variable are conditional independent and estimate the common *Odds Ratio*. In formula. I tested:

$$\mathcal{H}_0 : \text{Inoculation} \perp\!\!\!\perp \text{Typhoid} \mid \text{Location} \quad \text{vs} \quad \mathcal{H}_1 : \text{Inoculation} \not\perp\!\!\!\perp \text{Typhoid} \mid \text{Location}$$

knowing that under  $\mathcal{H}_0$ , the CMH's statistic follows a  $\chi^2_{df}$ . For Group A I found strong statistical evidence to claim that the inoculation and typhoid are not conditionally independent ( $df = 3$ ,  $\chi^2 = 67.83$ ,  $p\text{-value} < 2.2 \cdot 10^{-16}$ ), while it is the case for the Group B ( $df = 1$ ,  $\chi^2 = 0.7755$ ,  $p\text{-value} = 0.3785$ ).

### 3 Results

The main result of the whole analysis of `typhoidinoc.dat` dataset is that the six samples considered are different in dimensions and distributions and for this reason, a unique overall association between the inoculation and the freedom of attack is not representative for all of them. Follow that the evaluation of the effectiveness of the treatment cannot be based just on the overall contingency table. However, two main clusters of locations (without 3-way association) can be recognized: Group A (Indian Stations, 17th Lancers, Coldstream Guards and 7 Large Indian stations) and Group B (Stations Abroad and Transvaal). In Figure 4 the

two clusters are represented in different colours, displaying the *Odds Ratio* above discussed with the 95% Confidence Interval, but also the common *Odds Ratio* for each group estimated by CMH's test (represented by a vertical line). Note that the estimated common *Odds Ratio* for the Group B is reported just as a reference, but there is not statistical evidence to say it is different from 1.

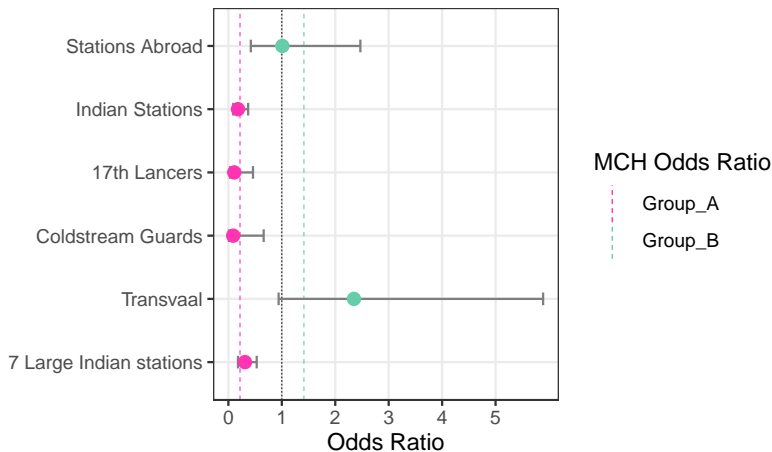


Figure 4: Per Location *Odds Ratio* (95% Confidence Interval)

The divergent nature of these results calls for some explanation to interpret the effectiveness of the treatment. To some extent it may be due to want of homogeneity in the material, owing to different methods of inoculation being employed, and to other causes; but it does not seem probable that the whole effect can be attributed to the results of random sampling.

## 4 Conclusion

In this study I have discussed the effectiveness of anti-typhoid inoculations in 1906-1907 within the army of the British Empire and colonies. Discrete Data Analysis is used to evaluate the independence and associations *within* and *between* the samples considered and two main patterns are discovered. A conclusive explanation of the statistical results is suggested by G.D.Maynard, who first studied this dataset using Discrete Data Analysis (getting results in complete agreement with mine). His hypothesis was that there was more than one disease covered by the term *enteric* or *typhoid fever*. Bacteriologically it was known that several varieties of para-typhoid bacilli exist and they are the specific cause of fevers clinically indistinguishable from typhoid and there is no a priori reason for assuming that a vaccine prepared from typhoid bacilli would confer immunity against infection with any of the para-typhoid strains. It suggests that maybe in the locations where the association between inoculation and freedom from attack is weaker or absent (i.e. Group B) it is due to an high ratio of para- to true typhoid infections. In particular it was shown by Major Statham that in the Trasvaal, actually, the para-typhoid were quite common.

The conclusion is that there is statistical evidence to support the effectiveness of the inoculation against typhoid fever (common *Odds Ratio* of Group A smaller than 1), but it cannot be extended to para-typhoid variants (common *Odds Ratio* of Group B not significantly different from 1).

## References

- [1] HJM Buist. Anti-typhoid inoculation. *BMJ Military Health*, 9(6):613–622, 1907.
- [2] Mr Greenwood and G. Udny Yule. The statistics of anti-typhoid and anti-cholera inoculations, and the interpretation of such statistics in general, 1915.
- [3] Torsten Hothorn and Brian S. Everitt. *A handbook of statistical analyses using R*. CRC press, 2014.
- [4] G.D. Maynard. Statistical study of anti-typhoid inoculation. *Biometrika*, 6(4):366–375, 1909.
- [5] Frank Yates. Contingency tables involving small numbers and the  $\chi^2$  test. *Supplement to the Journal of the Royal Statistical Society*, 1(2):217–235, 1934.