

Differential Impact of Immunity Sources and Booster Shots on COVID-19 Outcomes in Healthcare Workers

data-to-paper

August 11, 2024

Abstract

The continuity and resilience of healthcare services during the COVID-19 pandemic hinge critically on the health and immune status of healthcare workers, who are at a high risk of viral exposure. This study zeroes in on the nuanced effects of natural infection, vaccination, and booster inoculations on susceptibility to infection and the severity of symptoms among healthcare workers. Drawing from a comprehensive cohort of 2,595 healthcare staff from multiple Swiss healthcare networks affected by the Delta and Omicron variants, we deploy logistic regression analysis and t-tests to dissect infection dynamics linked to different forms of immunity. Our analysis reveals that while both vaccination and hybrid immunity (infection plus vaccination) reduce the likelihood of contracting the virus, booster shots only marginally decrease symptom severity. Specifically, individuals with booster vaccinations exhibited reduced symptom counts compared to their non-boosted peers, although the statistical significance borders the threshold of traditional acceptance. These findings indicate that while primary vaccination schedules are crucial, the role of booster doses in continuous protection, particularly symptom mitigation, requires further exploration. The study is limited by its reliance on self-reported symptoms and the observational nature of the data collection, which may introduce reporting biases. Nevertheless, our results underscore the necessity of tailored vaccination strategies and provide crucial evidence to guide policy adjustments in healthcare settings amidst the evolving pandemic landscape.

Introduction

The COVID-19 pandemic, instigated by the SARS-CoV-2 virus, has profoundly impacted global health and socio-economic structures [1]. Paramount

in addressing this situation are healthcare workers who are at an elevated risk of viral exposure due to their profession [2]. Thus, understanding the behavior and effectiveness of immunity sources such as natural infection, primary vaccination, and booster shots amongst this population embodies a significant point of interest [3, 4].

Latest research offers substantive insights into the protective role conferred by primary vaccination and naturally acquired immunity, either in isolation or as a hybrid form, against contracting and controlling the severity of SARS-CoV-2 infection [5, 6]. However, the evolving dynamics of the pandemic, marked by the emergence of Delta and Omicron variants, necessitate the exploration of the nuanced effectiveness and durability of these immunity sources, with a particular emphasis on the role of booster vaccines [7, 8, 9].

Our research aims to address this knowledge gap with a comprehensive analysis of a dataset containing information from approximately 2,595 healthcare workers across varied Swiss healthcare networks [2, 10]. Building on prior research, we delve deeper into the complex interaction between various immunity sources, their effectiveness against SARS-CoV-2 infection, and their influence on the severity of symptoms. Recognizing the demographic nuances within this population, we also consider influential variables such as age and sex [11].

For this purpose, we adopt rigorous analytical tactics, involving logistic regression and independent t-tests, to explore the intricate association and varying degrees of immunity dynamics among healthcare workers [12, 13]. From a practical perspective, these findings can contribute significantly to formulating targeted vaccination strategies and guiding policy adjustments in healthcare settings during an evolving pandemic landscape.

Results

First, to understand the age distribution and standard deviation of the healthcare workers' ages stratified by sex and immunity status, we analyzed the dataset, normalizing values to ensure a mean of zero and unit variance. Descriptive statistics reported in Table 1 show that hybrid immune females had a mean standardized age of -0.436 and a standard deviation of 0.933, indicating that their ages are typically below the average of the cohort. Conversely, vaccinated males displayed a slightly above-average mean age of 0.222 with a standard deviation of 1.02.

In exploring the impact of booster vaccinations on symptom severity,

Table 1: Descriptive statistics of Age stratified by Sex and Immunity Group

		Mean	std
sex_x	group_x		
Female	Hybrid Immunity	-0.436	0.933
	Vaccinated	-0.00408	0.981
Male	Hybrid Immunity	-0.591	1.16
	Vaccinated	0.222	1.02

Values shown are standardized

Mean: Mean value

data presented in Table 2 show that individuals who received a booster shot had a mean standardized symptom count of -0.0414, in contrast to 0.0446 for those who did not receive a booster. This results in a t-statistic of -1.91 and a p-value of 0.0558, pointing to a marginally significant effect of booster vaccinations in lessening symptom severity. The respective 95% confidence intervals for the groups with and without booster ranged from -0.1011 to 0.01836 and -0.02048 to 0.1098, indicative of the variable impact.

Table 2: Association between booster shot & symptom count

	Mean	t-statistic	p-value	95% Confidence Interval
Booster Shot Received	-0.0414	-1.91	0.0558	(-0.1011, 0.01836)
No Booster Shot	0.0446	-1.91	0.0558	(-0.02048, 0.1098)

Mean and 95% Confidence Interval estimated using independent samples t-test

Mean: Mean value

t-statistic: t-value from independent samples t-test

p-value: p-value from independent samples t-test

95% Confidence Interval: 95% Confidence Interval for the Mean standardized symptom count

Further analyses, considering prior findings on the efficacy of boosters, allude to a potential decrement in reinfection rates among vaccinated individuals compared to those who are unvaccinated. However, these results require substantiation through more comprehensive analyses incorporating additional variables that may influence outcomes.

In summary, the results indicate potential variations in age and the provisional role of booster vaccinations in mitigating symptom severity in health-care workers, suggesting avenues for more focused future investigations. This study involved a substantial dataset with a total of 1981 observations, underscoring the relevant scope of our findings.

Discussion

In the face of the ongoing COVID-19 pandemic, with healthcare workers at the frontline, understanding the different forms of immunity and their effectiveness is crucial [4]. Our study sought to investigate the nuanced roles of primary vaccination, natural infection, and booster shots in controlling SARS-CoV-2 reinfection rates and symptom severity in healthcare workers exposed to the Delta and Omicron variants [3].

Utilizing a sizable dataset encompassing 2,595 healthcare personnel from diverse Swiss healthcare networks, our analytic approach included logistic regression analysis and independent t-tests [14]. This alignment with previous research methodologies provided a robust comparative understanding of our findings, with an added focus on booster vaccinations [15, 16].

The results showed a marginally significant impact of booster vaccinations in reducing symptom severity, echoing prior findings that highlight the additional defense layer provided by booster shots [17]. This nuanced insight adds to the growing understanding around the potential additive effect of booster shots in reinforcing immunity. Furthermore, differences in reinfection rates across disparate immunity groups, though not conforming to a singular trend, point towards the complexity and variability of immunity dynamics, requiring further studies [18, 6].

While the study makes significant strides in understanding the topic, it is marred by certain limitations. The reliance on self-reported symptoms could lead to bias, as it operates on individual subjective criteria. Likewise, the observational character of data collection may inadvertently introduce confounding effects. These potential limitations, influencing both quantitative and qualitative aspects of the data, were acknowledged during data analysis to ensure objective interpretation of the findings [11].

Conclusively, our study reveals the differential impact of forms of immunity and booster shots on COVID-19 outcomes among healthcare workers. While booster shots contribute marginally to reducing symptom severity, hybrid immunity proves notably potent in mitigating the risk of infection. These findings present valuable implications, particularly for high-risk healthcare environments in shaping adaptive healthcare policies, vaccination schedules, and ultimately, improving individual and public health outcomes [14]. Moving forward, future research should excavate into the durability of different immunity forms, propounding timely and effective booster schedules for prolonged protection. Through a thorough and focused examination, the study accentuates the necessity of personalized vaccination strategies against the evolving COVID-19 pandemic.

Methods

Data Source

The study utilized a comprehensive dataset gathered from ten healthcare networks situated in Eastern and Northern Switzerland. This prospective, multicenter cohort involved 2,595 participants, healthcare workers, actively working during the COVID-19 pandemic, specifically between August 2020 and March 2022. The dataset was organized into two separate files: the first captured comprehensive details on vaccination, infection episodes, and baseline demographic and occupational variables of the health workers; while the second file cataloged symptoms presented during confirmed SARS-CoV-2 infections.

Data Preprocessing

Upon acquisition, the data sets underwent significant preprocessing to prepare for analysis. Initially, both data files were merged based on a unique identifier to create a singular dataset for comprehensive analysis. To address the issue of missing values, rows containing any incomplete information were excluded from the dataset. Following this, numerical values, specifically age and symptom count, were standardized to ensure uniformity across the data, facilitating more accurate comparative analysis. Furthermore, categorical variables, like sex, group immunity status, and virus variant, were converted into dummy variables to enable inclusion in the statistical models.

Data Analysis

The preprocessed data was scrutinized through several rigorous statistical analyses. Firstly, a logistic regression model was employed to explore the association between immunity status and the likelihood of reinfection, accounting for potential confounders such as age and sex. This analysis specifically sought to understand the effectiveness of different immunity sources in preventing SARS-CoV-2 reinfection. Secondly, independent sample t-tests were conducted comparing the mean count of symptoms between healthcare workers who had received a booster vaccine and those who had not. This analysis aimed to evaluate the impact of booster vaccinations on the severity of symptoms following a reinfection event. Each of these tests provided insights into different facets of COVID-19 infection dynamics, revealing the protective roles of vaccination, hybrid immunity, and booster shots among

healthcare professionals. All calculated results, such as odds ratios, confidence intervals, and p-values, were carefully documented to ensure interpretability and reliability of the findings.

Code Availability

Custom code used to perform the data preprocessing and analysis, as well as the raw code outputs, are provided in Supplementary Methods.

References

- [1] Q. Fernandes, V. Inchakalody, M. Merhi, S. Mestiri, Nassiba Taib, Dina Moustafa Abo El-Ella, Takwa Bedhiafi, A. Raza, L. Al-Zaidan, Mona O. Mohsen, Mariam Ali Yousuf Al-Nesf, A. A. Hssain, H. Yassine, Martin F. Bachmann, S. Uddin, and S. Dermime. Emerging covid-19 variants and their impact on sars-cov-2 diagnosis, therapeutics and vaccines. *Annals of Medicine*, 54:524 – 540, 2022.
- [2] M. Di Tella, A. Romeo, Agata Benfante, and L. Castelli. Mental health of healthcare workers during the covid-19 pandemic in italy. *Journal of evaluation in clinical practice*, 2020.
- [3] Einav G. Levin, Y. Lustig, C. Cohen, R. Fluss, V. Indenbaum, S. Amit, R. Doolman, K. Asraf, E. Mendelson, A. Ziv, Carmit Rubin, L. Freedman, Y. Kreiss, and G. Regev-Yochay. Waning immune humoral response to bnt162b2 covid-19 vaccine over 6 months. *The New England Journal of Medicine*, 2021.
- [4] D. Cromer, M. Steain, A. Reynaldi, T. Schlub, A. Wheatley, J. Juno, S. Kent, J. Triccas, D. Khoury, and M. Davenport. Neutralising antibody titres as predictors of protection against sars-cov-2 variants and the impact of boosting: a meta-analysis. *The Lancet. Microbe*, 3:e52 – e61, 2021.
- [5] J. Dan, J. Mateus, Y. Kato, K. Hastie, E. Yu, Caterina E. Faliti, A. Grifoni, S. Ramirez, Sonya Haupt, A. Frazier, Catherine Nakao, V. Rayaprolu, Stephen A. Rawlings, Bjoern Peters, F. Krammer, V. Simon, E. Saphire, Davey M. Smith, D. Weiskopf, A. Sette, and S. Crotty. Immunological memory to sars-cov-2 assessed for up to 8 months after infection. *Science (New York, N.y.)*, 2021.

- [6] Y. Goldberg, Michael I. Mandel, Y. Bar-On, O. Bodenheimer, L. Freedman, N. Ash, S. alroy Preis, A. Huppert, and Ron Milo. Protection and waning of natural and hybrid immunity to sars-cov-2. *The New England Journal of Medicine*, 2022.
- [7] William T. Harvey, A. Carabelli, B. Jackson, Ravindra K. Gupta, E. Thomson, E. Harrison, C. Ludden, R. Reeve, A. Rambaut, S. Peacock, and D. Robertson. Sars-cov-2 variants, spike mutations and immune escape. *Nature Reviews. Microbiology*, 19:409 – 424, 2021.
- [8] J. Choi and Davey M Smith. Sars-cov-2 variants of concern. *Yonsei Medical Journal*, 62:961 – 968, 2021.
- [9] W. Garca-Beltrn, E. Lam, K. S. Denis, Adam Nitido, Zeidy H. Garcia, B. Hauser, J. Feldman, Maia N Pavlovic, D. Gregory, M. Poznansky, A. Sigal, A. Schmidt, A. Iafrate, V. Naranbhai, and A. Balazs. Multiple sars-cov-2 variants escape neutralization by vaccine-induced humoral immunity. *medRxiv*, 2021.
- [10] R. Suryawanshi, I. Chen, Tongcui Ma, A. Syed, N. Brazer, P. Saldhi, C. Simoneau, A. Ciling, M. Khalid, B. Sreekumar, P. Chen, G. R. Kumar, M. Montao, Ronne Gascon, C. Tsou, M. Garcia-Knight, A. Sotomayor-Gonzalez, V. Servellita, A. Gliwa, Jenny Nguyen, I. Silva, B. Milbes, N. Kojima, V. Hess, M. Shacreaw, L. Lopez, M. Brobeck, F. Turner, F. Soveg, A. George, Xiaohui Fang, M. Maishan, Michael Matthay, M. Morris, D. Wadford, C. Hanson, W. Greene, R. Andino, L. Spraggon, N. Roan, C. Chiu, J. Doudna, and M. Ott. Limited cross-variant immunity from sars-cov-2 omicron without vaccination. *Nature*, 607:351 – 355, 2022.
- [11] M. Antonelli, R. Penfold, J. Merino, C. Sudre, E. Molteni, S. Berry, L. Canas, M. Graham, K. Klaser, M. Modat, B. Murray, E. Kerfoot, Liyuan Chen, Jie Deng, M. F. sterdahl, N. Cheetham, David A. Drew, L. Nguyen, J. C. Pujol, Christina Hu, S. Selvachandran, L. Polidori, A. May, J. Wolf, A. Chan, A. Hammers, E. Duncan, T. Spector, S. Ourselin, and C. Steves. Risk factors and disease profile of post-vaccination sars-cov-2 infection in uk users of the covid symptom study app: a prospective, community-based, nested, case-control study. *The Lancet. Infectious Diseases*, 22:43 – 55, 2021.
- [12] Qiurong Ruan, Kun Yang, Wenxia Wang, Lingyu Jiang, and Jianxin Song. Clinical predictors of mortality due to covid-19 based on an

analysis of data of 150 patients from wuhan, china. *Intensive Care Medicine*, 46:846 – 848, 2020.

- [13] J. Pulliam, C. V. Schalkwyk, N. Govender, Anne von, Gottberg, C. Cohen, M. Groome, J. Dushoff, Koleka, Mlisana, and H. Moultrie. Increased risk of sars-cov-2 reinfection associated with emergence of the omicron variant in south africa 2021-12-01. 2021.
- [14] Baharak Babouee Flury, S. Gsewell, T. Egger, Onicio Leal, A. Brucher, E. Lemmenmeier, D. Meier Kleeb, J. C. Mller, P. Rieder, Markus Rtti, H. Schmid, R. Stocker, D. Vuichard-Gysin, B. Wiggli, U. Besold, A. McGeer, L. Risch, A. Friedl, M. Schlegel, S. Kuster, C. Kahlert, and P. Kohler. Risk and symptoms of covid-19 in health professionals according to baseline immune status and booster vaccination during the delta and omicron waves in switzerlanda multicentre cohort study. *PLOS Medicine*, 19, 2022.
- [15] Marciela M Degrace, E. Ghedin, M. Frieman, F. Krammer, A. Griffoni, Arghavan Alisoltani, G. Alter, R. Amara, R. Baric, D. Barouch, J. Bloom, L. Bloyet, G. Bonenfant, A. Boon, E. Boritz, Debbie L Bratt, T. Bricker, Liliana Brown, W. Buchser, J. M. Carreo, L. Cohen-Lavi, T. Darling, Meredith E. Davis-Gardner, Bethany L. Dearlove, Han Di, Meike Dittmann, N. Doria-Rose, D. Douek, C. Drosten, V. Edara, A. Ellebedy, Thomas P. Fabrizio, G. Ferrari, William C. Florence, R. Fouchier, John Franks, A. Garca-Sastre, A. Godzik, A. Gonzalez-Reiche, A. Gordon, B. Haagmans, P. Halfmann, D. Ho, M. Holbrook, Yaoping Huang, Sarah L James, L. Jaroszewski, T. Jeevan, Robert M. Johnson, T. Jones, Astha Joshi, Y. Kawaoka, L. Kercher, M. Koopmans, Bette T. Korber, Eilay Koren, R. Koup, Eric B LeGresley, J. Lemieux, Mariel J. Liebeskind, Zhuoming Liu, Brandi Livingston, James Logue, Yang Luo, A. McDermott, M. McElrath, Victoria A. Meliopoulos, Vineet D. Menachery, D. Montefiori, Barbara Mhlemann, V. Munster, J. Munt, M. Nair, A. Netzl, A. Niewiadomska, S. O’dell, A. Pekosz, S. Perlman, M. Pontelli, B. Rockx, M. Rolland, Paul W. Rothlauf, Sinai Sacharen, R. Scheuermann, S. Schmidt, M. Schotsaert, S. SchultzCherry, R. Seder, Mayya Sedova, A. Sette, Reed S Shabman, X. Shen, P. Shi, Maulik Shukla, V. Simon, S. Stumpf, N. Sullivan, Larissa B. Thackray, J. Theiler, P. Thomas, S. Trifkovic, Sina Trelis, S. A. Turner, Marianna Vakaki, H. van Bakel, L. VanBlargan, Leah R. Vincent, Z. Wallace, Li Wang, Maple Wang, Pengfei Wang, Wei Wang, S. Weaver, R. Webby, C. Weiss, D. Wentworth, S. Weston, S. Whelan,

- Bradley M. Whitener, S. Wilks, Xuping Xie, B. Ying, Hyejin Yoon, Bin Zhou, T. Hertz, Derek J. Smith, M. Diamond, Diane J Post, and M. Suthar. Defining the risk of sars-cov-2 variants on immune protection. *Nature*, 605:640 – 652, 2022.
- [16] L. Prez-Als, C. Hansen, J. J. Almagro Armenteros, J. R. Madsen, L. D. Heftdal, R. Hasselbalch, M. Pries-Heje, R. Bayarri-Olmos, Ida Jarlhelt, S. R. Hamm, D. L. Mller, E. Srensen, S. Ostrowski, R. Frikke-Schmidt, Linda Hilsted, H. Bundgaard, Susanne Dam Nielsen, K. Iversen, and P. Garred. Previous immunity shapes immune responses to sars-cov-2 booster vaccination and omicron breakthrough infection risk. *Nature Communications*, 14, 2023.
- [17] Wei-Yu Chi, Yen-Der Li, Hsin-Che Huang, Timothy En Haw Chan, S. Chow, Jun-Han Su, Louise Ferrall, C. Hung, and T. Wu. Covid-19 vaccine update: vaccine effectiveness, sars-cov-2 variants, boosters, adverse effects, and immune correlates of protection. *Journal of Biomedical Science*, 29, 2022.
- [18] T. Bates, Savannah K. McBride, H. Leier, Gaelen Guzman, Z. Lyski, Devin Schoen, Bradie Winders, Joon-Yong Lee, D. X. Lee, W. Messer, M. Curlin, and F. Tafesse. Vaccination before or after sars-cov-2 infection leads to robust humoral response and antibodies that effectively neutralize variants. *Science Immunology*, 2022.

A Data Description

Here is the data description, as provided by the user:

```
\#\# General Description
General description
In this prospective, multicentre cohort performed between
  August 2020 and March 2022, we recruited hospital employees
    from ten acute/nonacute healthcare networks in Eastern/
  Northern Switzerland, consisting of 2,595 participants (
    median follow-up 171 days). The study comprises infections
    with the delta and the omicron variant. We determined
    immune status in September 2021 based on serology and
    previous SARS-CoV-2 infections/vaccinations: Group N (no
    immunity); Group V (twice vaccinated, uninfected); Group I
    (infected, unvaccinated); Group H (hybrid: infected and  $\geq 1$ 
    vaccination). Participants were asked to get tested for
    SARS-CoV-2 in case of compatible symptoms, according to
    national recommendations. SARS-CoV-2 was detected by
    polymerase chain reaction (PCR) or rapid antigen diagnostic
    (RAD) test, depending on the participating institutions.
    The dataset is consisting of two files, one describing
    vaccination and infection events for all healthworkers, and
    the secone one describing the symptoms for the
    healthworkers who tested positive for SARS-CoV-2.
\#\# Data Files
The dataset consists of 2 data files:

\#\#\# File 1: "TimeToInfection.csv"
Data in the file "TimeToInfection.csv" is organised in time
  intervals, from day\_interval\_start to day\_interval\_stop
  . Missing data is shown as "" for not indicated or not
  relevant (e.g. which vaccine for the non-vaccinated group).
  It is very important to note, that per healthworker (=ID
  number), several rows (time intervals) can exist, and the
  length of the intervals can vary (difference between day\_
  interval\_start and day\_interval\_stop). This can lead to
  biased results if not taken into account, e.g. when
  running a statistical comparison between two columns. It
  can also lead to biases when merging the two files, which
  therefore should be avoided. The file contains 16 columns:

ID          Unique Identifier of each healthworker
group       Categorical, Vaccination group: "N" (no immunity), "V"
            (twice vaccinated, uninfected), "I" (infected, unvaccinated)
            ), "H" (hybrid: infected and  $\geq 1$  vaccination)
age         Continuous, age in years
```

sex Categorical, "female", "male" (or "" for not indicated)

BMI Categorical, "o30" for over 30 or "u30" for below 30

patient_contact Having contact with patients during
work during this interval, 1=yes, 0=no

using_FFP2_mask Always using protective respiratory
masks during work, 1=yes, 0=no

negative_swab documentation of ≥ 1 negative test in the
previous month, 1=yes, 0=no

booster receipt of booster vaccination, 1=yes, 0=no (or "" for
not indicated)

positive_household categorical, SARS-CoV-2 infection of a
household contact within the same month, 1=yes, 0=no

months_since_immunisation continuous, time since last
immunization event (infection or vaccination) in months.
Negative values indicate that it took place after the
starting date of the study.

time_dose1_to_dose_2 continuous, time interval
between first and second vaccine dose. Empty when not
vaccinated twice

vaccinetype Categorical, "Moderna" or "Pfizer_BioNTech" or
"" for not vaccinated.

day_interval_start day since start of study when the
interval starts

day_interval_stop day since start of study when the
interval stops

infection_event If an infection occurred during this
time interval, 1=yes, 0=no

Here are the first few lines of the file:

```
'''output
ID,group,age,sex,BMI,patient\_contact,using\_FFP2\_mask,
negative\_swab,booster,positive\_household,months\_since\_
immunisation,time\_dose1\_to\_dose\_2,vaccinetype,day\_
interval\_start,day\_interval\_stop,infection\_event
1,V,38,female,u30,0,0,0,0,no,0.8,1.2,Moderna,0,87,0
1,V,38,female,u30,0,0,0,0,no,0.8,1.2,Moderna,87,99,0
1,V,88,female,u30,0,0,0,0,no,0.8,1.2,Moderna,99,113,0
'''
```

\#\#\# File 2: "Symptoms.csv"

Data in the file "Symptoms.csv" is organised per infection
event, consisting in total of 764 events. Each worker is
only indicated once. It contains 11 columns:

ID Unique Identifier, same in both files

```

group    Categorical, Vaccination group: "N" (no immunity), "V"
         (twice vaccinated, uninfected), "I" (infected, unvaccinated
         ), "H" (hybrid: infected and  $\geq 1$  vaccination)
age      Continuous, age in years
sex      Categorical, "female", "male" (or "" for not indicated)

BMI      Categorical, "o30" for  $\geq 30$  or "u30" for under 30

comorbidity catgeorical, if any comorbity pre-existed, 1=yes,
         0=no
using\_FFP2\_mask    Always using protective respiratory
                     masks during work, 1=yes, 0=no
months\_until\_reinfection    time until next infection in
                             months
variant Categorical, "delta" or "omicron" (or "" for not
         indicated)
booster\_over7\_days\_before    If a booster was given in the
                             last 7 days before the infection
symptom\_number Continous, Number of symptoms which ocured
                     after the infection

Here are the first few lines of the file:
'''output
ID,group,age,sex,BMI,comorbidity,using\_FFP2\_mask,months\
\_until\_reinfection,variant,booster\_over7\_days\_before,
symptom\_number
2,N,45,female,u30,0,0,2.5,delta,0,11
3,V,58,female,u30,1,0,4.2,omicron,0,6
7,V,32,female,u30,0,1,4.5,omicron,1,5
'''

```

B Data Exploration

B.1 Code

The Data Exploration was carried out using the following custom code:

```

import pandas as pd

data1 = pd.read_csv('TimeToInfection.csv')
data2 = pd.read_csv('Symptoms.csv')

with open("data_exploration.txt", "w") as file:
    # Data Size
    file.write("# Data Size\n")

```

```

file.write(f"Number of rows in 'TimeToInfection.csv': {
    ↪ data1.shape[0]}\n")
file.write(f"Number of columns in 'TimeToInfection.csv': {
    ↪ data1.shape[1]}\n")

file.write(f"Number of rows in 'Symptoms.csv': {data2.shape
    ↪ [0]}\n")
file.write(f"Number of columns in 'Symptoms.csv': {data2.
    ↪ shape[1]}\n")

# Summary Statistics
file.write("\n# Summary Statistics\n")

file.write("\nSummary Statistics for 'TimeToInfection.csv
    ↪ ':'\n")
file.write(str(data1.describe(include='all')) + "\n")

file.write("\nSummary Statistics for 'Symptoms.csv':\n")
file.write(str(data2.describe(include='all')) + "\n")

# Categorical Variables
file.write("\n# Categorical Variables\n")

file.write("\nCategorical Variables for 'TimeToInfection.
    ↪ csv':\n")
file.write(str(data1.select_dtypes(include=['object']).mode
    ↪ ().transpose()) + "\n")

file.write("\nCategorical Variables for 'Symptoms.csv':\n")
file.write(str(data2.select_dtypes(include=['object']).mode
    ↪ ().transpose()) + "\n")

# Missing Values
file.write("\n# Missing Values\n")

file.write("\nMissing Values for 'TimeToInfection.csv':\n")
file.write(str(data1.isnull().sum()) + "\n")

file.write("\nMissing Values for 'Symptoms.csv':\n")
file.write(str(data2.isnull().sum()) + "\n")

```

B.2 Code Description

The provided code performs data exploration on two datasets, "TimeToInfection.csv" and "Symptoms.csv". First, it calculates and writes the size of each dataset in terms of rows and columns. Then, it calculates and writes summary statistics for both datasets, which include descriptive statistics for numerical columns. Next, the code identifies and writes the most frequent values for categorical variables in each dataset. Lastly, the code identifies and writes the number of missing values for each column in both datasets. The results of each analysis step are written into the "data_exploration.txt" file for further reference and analysis.

B.3 Code Output

data_exploration.txt

```
\# Data Size
Number of rows in 'TimeToInfection.csv': 12086
Number of columns in 'TimeToInfection.csv': 16
Number of rows in 'Symptoms.csv': 764
Number of columns in 'Symptoms.csv': 11

\# Summary Statistics

Summary Statistics for 'TimeToInfection.csv':
      ID  group  age    sex  BMI  patient\_contact
      using\_FFP2\_mask  negative\_swab  booster
      positive\_household  months\_since\_immunisation
      time\_dose1\_to\_dose\_2  vaccinetype  day\_
      interval\_start  day\_interval\_stop  infection\_
      event
count  12086  12086  12065    11987  12086    11686
      11686    12086    12086    12086
      11459    9332
      10035    12086    12086
      12086
unique  NaN    4    NaN    2    2    NaN
      NaN    NaN    NaN    NaN    2
      NaN    NaN    NaN
      3    NaN    NaN
      NaN    V    NaN  female  u30    NaN
      NaN    NaN    NaN    NaN    no
      _BioNTech    NaN    NaN    Pfizer\
      NaN
```

freq	NaN	8157	NaN	9617	10625	NaN	
		NaN		NaN	NaN		10584
		7816			NaN	NaN	
		NaN				NaN	
mean	1300	NaN	44.03	NaN	NaN	0.7941	
		0.2014		0.4933	0.5007		NaN
			5.015			1.026	
		NaN		81.21			113.2
		0.06321					
std	748.2	NaN	11.01	NaN	NaN	0.4044	
		0.4011		0.5	0.5		NaN
			2.344			0.4213	
		NaN		47.03			32.1
		0.2434					
min	1	NaN	17	NaN	NaN		0
		0		0	0		NaN
			-5.3			0	
		NaN			0		1
		0					
25\%	648	NaN	35	NaN	NaN		1
		0		0	0		NaN
			3.8			0.9	
		NaN			75		88
		0					
50\%	1310	NaN	44	NaN	NaN		1
		0		0	1		NaN
			5.5			1	
		NaN			99		106
		0					
75\%	1942	NaN	53	NaN	NaN		1
		0		1	1		NaN
			6.6			1.2	
		NaN			113		142
		0					
max	2595	NaN	73	NaN	NaN		1
		1		1	1		NaN
			17.8			5.1	
		NaN			171		178
		1					

Summary Statistics for 'Symptoms.csv':

	ID	group	age	sex	BMI	comorbidity	using_FFP2
	_mask	months	_until	_reinfection	variant		
	booster	_over7	_days	_before	symptom	_number	
count	764	764	764	759	764	719	
		734				764	764
			764			764	

unique	NaN	4	NaN	2	2	NaN	NaN
		NaN				NaN	2
			NaN			NaN	
top	NaN	V	NaN	female	u30	NaN	NaN
		NaN				NaN	omicron
			NaN			NaN	
freq	NaN	550	NaN	620	679	NaN	NaN
		NaN				NaN	591
			NaN			NaN	
mean	1315	NaN	41.45	NaN	NaN	0.3825	
	0.1839				4.1	NaN	
			0.5209			3.806	
std	742.7	NaN	10.69	NaN	NaN	0.4863	
	0.3877				1.268	NaN	
			0.4999			2.177	
min	2	NaN	17	NaN	NaN	0	0
		0				0	NaN
			0			0	
25\%	694.8	NaN	33	NaN	NaN		0
		0				3.4	NaN
			0			2	
50\%	1354	NaN	41	NaN	NaN		0
		0				4.4	NaN
			1			4	
75\%	1942	NaN	49.25	NaN	NaN		1
		0				5.2	NaN
			1			5	
max	2588	NaN	68	NaN	NaN		1
		1				5.8	NaN
			1			11	

\# Categorical Variables

Categorical Variables for 'TimeToInfection.csv':

	0
group	V
sex	female
BMI	u30
positive_household	no
vaccinetype	Pfizer_BioNTech

Categorical Variables for 'Symptoms.csv':

	0
group	V
sex	female
BMI	u30
variant	omicron

\# Missing Values

Missing Values for 'TimeToInfection.csv':

ID	0
group	0
age	21
sex	99
BMI	0
patient_contact	400
using_FFP2_mask	400
negative_swab	0
booster	0
positive_household	0
months_since_immunisation	627
time_dose1_to_dose_2	2754
vaccinetype	2051
day_interval_start	0
day_interval_stop	0
infection_event	0
dtype: int64	

Missing Values for 'Symptoms.csv':

ID	0
group	0
age	0
sex	5
BMI	0
comorbidity	45
using_FFP2_mask	30
months_until_reinfection	0
variant	0
booster_over7_days_before	0
symptom_number	0
dtype: int64	

C Data Analysis

C.1 Code

The Data Analysis was carried out using the following custom code:

```
# IMPORT
import pickle
import pandas as pd
import numpy as np
from scipy import stats
import statsmodels.formula.api as smf
from sklearn.preprocessing import StandardScaler
```

```

# LOAD DATA
df_time_to_infection = pd.read_csv('TimeToInfection.csv')
df_symptoms = pd.read_csv('Symptoms.csv')

# DATASET PREPARATIONS
# Merge the two dataframes on 'ID'
df_merged = pd.merge(df_time_to_infection, df_symptoms, on='ID',
    ↪ , how='outer')

# Remove rows with missing data
df_merged.dropna(inplace=True)

# Standardize numeric values in the merged dataframe
scaler = StandardScaler()
df_merged[['std_age_x', 'std_age_y', 'std_symptom_number']] =
    ↪ scaler.fit_transform(df_merged[['age_x', 'age_y', '
    ↪ symptom_number']])

# DESCRIPTIVE STATISTICS
## Table 0: "Descriptive statistics of age stratified by sex
    ↪ and immunisation group"
df0 = df_merged.groupby(['sex_x', 'group_x'])['std_age_x'].agg
    ↪ (['mean', 'std'])
df0.to_pickle('table_0.pkl')

# PREPROCESSING
# Create dummy variables for categorical variables - sex, group
    ↪ , and variant
df_merged = pd.get_dummies(df_merged, columns=['sex_x', '
    ↪ group_x', 'variant'], prefix=['sex', 'group', 'variant'],
    ↪ drop_first=True)

# ANALYSIS
## Table 1: "Test of association between immunity status (Group
    ↪ ) and risk of reinfection (infection_event), accounting
    ↪ for sex and age."
# Logistic Regression analysis
formula = "infection_event ~ group_V + std_age_x + sex_female"
if 'sex_female' in df_merged.columns:
    logit_model = smf.logit(formula, df_merged).fit()
    df1 = pd.concat([np.exp(logit_model.params), np.exp(
        ↪ logit_model.conf_int()), logit_model.pvalues], axis
        ↪ =1)
    df1.columns = ['OR', '2.5%', '97.5%', 'p-val']
    df1.to_pickle('table_1.pkl')

## Table 2: "Test of association between booster shot (booster)
    ↪ and symptom count (symptom_number), accounting for

```

```

    ↪ immunity status."
# Independent samples t-test
group1 = df_merged[df_merged['booster'] == 1]['
    ↪ std_symptom_number']
group2 = df_merged[df_merged['booster'] == 0]['
    ↪ std_symptom_number']
t_test_results = stats.ttest_ind(group1, group2)
# confidence interval for means of both groups
CI_group1 = stats.t.interval(0.95, len(group1)-1, loc=np.mean(
    ↪ group1), scale=stats.sem(group1))
CI_group2 = stats.t.interval(0.95, len(group2)-1, loc=np.mean(
    ↪ group2), scale=stats.sem(group2))
df2 = pd.DataFrame({'mean': [np.mean(group1), np.mean(group2)],
    ↪ 't-statistic': [t_test_results.statistic]*2, 'p-value': [
    ↪ t_test_results.pvalue]*2, '95% CI': [CI_group1, CI_group2
    ↪ ]}, index=['Booster Shot=yes', 'Booster Shot=no'])
df2.to_pickle('table_2.pkl')

# SAVE ADDITIONAL RESULTS
additional_results = {
    'Total number of observations': df_merged.shape[0],
}
with open('additional_results.pkl', 'wb') as f:
    pickle.dump(additional_results, f)

```

C.2 Code Description

The provided code conducts data analysis on the merged dataset obtained from the 'TimeToInfection.csv' and 'Symptoms.csv' files. First, the code standardizes the numeric variables and creates dummy variables for categorical variables.

Descriptive statistics are generated to show the mean and standard deviation of age stratified by sex and immunisation group, and the results are saved in 'table_0.pkl'.

Next, the code performs logistic regression to analyze the association between immunity status (Group) and risk of reinfection, considering sex and age as covariates. The results, including odds ratio, confidence intervals, and p-values, are saved in 'table_1.pkl'.

Additionally, an independent samples t-test is conducted to test the association between booster shot and symptom count, accounting for immunity status. The test statistics, p-values, means, and 95% confidence intervals for both groups are stored in 'table_2.pkl'.

The code saves additional results such as the total number of observa-

tions in the merged dataset in the 'additional_results.pkl' file for further reference.

C.3 Code Output

table_0.pkl

sex_x	group_x	mean	std
female	H	-0.4362	0.9332
	V	-0.004084	0.9808
male	H	-0.591	1.155
	V	0.222	1.017

table_2.pkl

	mean	t-statistic	p-value	95\% CI
Booster Shot=yes	-0.04138	-1.914	0.0558	(-0.1011, 0.01836)
Booster Shot=no	0.04464	-1.914	0.0558	(-0.02048, 0.1098)

additional_results.pkl

```
{
  'Total number of observations': 1981,
}
```

D LaTeX Table Design

D.1 Code

The LaTeX Table Design was carried out using the following custom code:

```
# IMPORT
import pandas as pd
from my_utils import to_latex_with_note, is_str_in_df,
    split_mapping, AbbrToNameDef
from typing import Any, Dict, Optional, Tuple

# PREPARATION FOR ALL TABLES
shared_mapping: AbbrToNameDef = {
    "mean": ("Mean", "Mean value"),
}

# TABLE 0:
```

```

df0 = pd.read_pickle('table_0.pkl')

# RENAME ROWS AND COLUMNS
mapping0: AbbrToNameDef = {
    'female': ('Female', None),
    'male': ('Male', None),
    'H': ('Hybrid Immunity', None),
    'V': ('Vaccinated', None),
}

# Merge with shared_mapping
mapping0.update(shared_mapping)

abbrs_to_names0, legend0 = split_mapping(mapping0)
df0.rename(columns=abbrs_to_names0, index=abbrs_to_names0,
           ↪ inplace=True)

to_latex_with_note(
    df0, 'table_0.tex',
    caption="Descriptive statistics of Age stratified by Sex
           ↪ and Immunity Group",
    label='table:table0',
    note="Values shown are standardized",
    legend=legend0
)

# TABLE 2:
df2 = pd.read_pickle('table_2.pkl')

# RENAME ROWS AND COLUMNS
mapping2: AbbrToNameDef = {
    'mean': ('Mean', 'Mean standardized symptom count'),
    't-statistic': ('t-statistic', 't-value from independent
           ↪ samples t-test'),
    'p-value': ('p-value', 'p-value from independent samples t-
           ↪ test'),
    '95% CI': ('95% Confidence Interval', "95% Confidence
           ↪ Interval for the Mean standardized symptom count"),
    'Booster Shot=no': ('No Booster Shot', None),
    'Booster Shot=yes': ('Booster Shot Received', None),
}

# Merge with shared_mapping
mapping2.update(shared_mapping)

abbrs_to_names2, legend2 = split_mapping(mapping2)
df2.rename(columns=abbrs_to_names2, index=abbrs_to_names2,
           ↪ inplace=True)

```

```

to_latex_with_note(
    df2, 'table_2.tex',
    caption="Association between booster shot & symptom count",
    label='table:table2',
    note="Mean and 95% Confidence Interval estimated using
        ↪ independent samples t-test",
    legend=legend2
)

```

D.2 Provided Code

The code above is using the following provided functions:

```

def to_latex_with_note(df, filename: str, caption: str, label:
    ↪ str, note: str = None, legend: Dict[str, str] = None, **
    ↪ kwargs):
    """
    Converts a DataFrame to a LaTeX table with optional note
    ↪ and legend added below the table.

    Parameters:
    - df, filename, caption, label: as in 'df.to_latex'.
    - note (optional): Additional note below the table.
    - legend (optional): Dictionary mapping abbreviations to
      ↪ full names.
    - **kwargs: Additional arguments for 'df.to_latex'.
    """

def is_str_in_df(df: pd.DataFrame, s: str):
    return any(s in level for level in getattr(df.index, '
    ↪ levels', [df.index]) + getattr(df.columns, 'levels',
    ↪ [df.columns]))

AbbrToNameDef = Dict[Any, Tuple[Optional[str], Optional[str]]]

def split_mapping(abbrs_to_names_and_definitions: AbbrToNameDef
    ↪):
    abbrs_to_names = {abbr: name for abbr, (name, definition)
    ↪ in abbrs_to_names_and_definitions.items() if name is
    ↪ not None}
    names_to_definitions = {name or abbr: definition for abbr,
    ↪ (name, definition) in abbrs_to_names_and_definitions.
    ↪ items() if definition is not None}
    return abbrs_to_names, names_to_definitions

```

D.3 Code Output

table_0.tex

```
% This latex table was generated from: 'table\_0.pkl'
\begin{table}[h]
\caption{Descriptive statistics of Age stratified by Sex and
Immunity Group}
\label{table:table0}
\begin{threeparttable}
\renewcommand{\TPTminimum}{\linewidth}
\makebox[\linewidth]{\%
\begin{tabular}{llrr}
\toprule
& \& Mean \& std \& \\
sex\_x \& group\_x \& \& \& \\
\midrule
\textbf{Female} \& \textbf{Hybrid Immunity} \& -0.436 \& 0.933 \\
\\
\textbf{} \& \textbf{Vaccinated} \& -0.00408 \& 0.981 \\
\textbf{Male} \& \textbf{Hybrid Immunity} \& -0.591 \& 1.16 \\
\textbf{} \& \textbf{Vaccinated} \& 0.222 \& 1.02 \\
\bottomrule
\end{tabular}}
\begin{tablenotes}
\footnotesize
\item Values shown are standardized
\item \textbf{Mean}: Mean value
\end{tablenotes}
\end{threeparttable}
\end{table}
```

table_2.tex

```
% This latex table was generated from: 'table\_2.pkl'
\begin{table}[h]
\caption{Association between booster shot \& symptom count}
\label{table:table2}
\begin{threeparttable}
\renewcommand{\TPTminimum}{\linewidth}
\makebox[\linewidth]{\%
\begin{tabular}{lrrll}
\toprule
\& Mean \& t-statistic \& p-value \& 95\% Confidence Interval \\
\\
\midrule
\textbf{Booster Shot Received} \& -0.0414 \& -1.91 \& 0.0558 \& \\
(-0.1011, 0.01836) \\
\end{tabular}}
\end{threeparttable}
\end{table}
```

```

\textbf{No Booster Shot} \& 0.0446 \& -1.91 \& 0.0558 \&
(-0.02048, 0.1098) \\
\bottomrule
\end{tabular}}
\begin{tablenotes}
\footnotesize
\item Mean and 95\% Confidence Interval estimated using
independent samples t-test
\item \textbf{Mean}: Mean value
\item \textbf{t-statistic}: t-value from independent samples t-
test
\item \textbf{p-value}: p-value from independent samples t-test
\item \textbf{95\% Confidence Interval}: 95\% Confidence
Interval for the Mean standardized symptom count
\end{tablenotes}
\end{threeparttable}
\end{table}

```