

# BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email [info.bmjopen@bmj.com](mailto:info.bmjopen@bmj.com)

# BMJ Open

## Compress the Curve: An Observational Study of Variations in COVID-19 Infections Across California Nursing Homes

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2020-042804
Article Type:	Original research
Date Submitted by the Author:	17-Jul-2020
Complete List of Authors:	Gopal, Ram; University of Warwick Han, Xu; Fordham University Yaraghi, Niam; Brookings Institution, Governance Studies
Keywords:	COVID-19, GERIATRIC MEDICINE, Health policy < HEALTH SERVICES ADMINISTRATION & MANAGEMENT

SCHOLARONE™  
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

## Original Investigation

**Title:** Compress the Curve: An Observational Study of Variations in COVID-19 Infections Across California Nursing Homes

Ram D. Gopal, PhD<sup>1</sup>, Xu Han, PhD<sup>2</sup>, Niam Yaraghi, PhD<sup>3,4,\*</sup>

<sup>1</sup>: Professor, Warwick Business School, University of Warwick

<sup>2</sup>: Assistant Professor, Gabelli School of Business, Fordham University

<sup>3</sup>: Assistant Professor, School of Business, University of Connecticut

<sup>4</sup>: Non-resident Fellow, Governance Studies, The Brookings Institution

*\*: All authors contributed equally. Authors are listed in alphabetical order of their last name.*

Corresponding author:

Niam Yaraghi

[Nyaraghi@brookings.edu](mailto:Nyaraghi@brookings.edu)

1 University Place, Stamford, CT 06901 Phone: (203) 251-9583

Word count: 2668

## Abstract

*Objective:* Nursing homes' residents and staff constitute the largest proportion of the fatalities associated with COVID-19 epidemic. Although there is a significant variation in COVID-19 outbreaks among the US nursing homes, we still do not know why such outbreaks are larger and more likely in some nursing homes than others. This research aims to understand why some nursing homes are more susceptible to larger COVID-19 outbreaks.

*Design:* Observational study of all nursing homes in the state of California until May 1st, 2020.

*Setting:* The state of California.

*Participants:* 713 long term care facilities in the State of California that participate in public reporting of COVID-19 infections as of May 1<sup>st</sup>, 2020 and their infections data could be matched with CMS database on ratings and governance features.

*Main Outcome Measure:* The number of reported COVID-19 infections among staff and residents.

*Results:* Study sample included 713 nursing homes. The size of outbreaks among residents in for-profit nursing homes is 13 times larger than their non-profit counterparts (log count = 2.57; 95% CI, 1.99 to 3.15; P<.001). Higher ratings in CMS-reported health inspections are associated with lower number of infections among both staff (log count = -0.20; 95% CI, -0.38 to -0.01; P = 0.04) and residents (log count = -0.20; 95% CI, -0.26 to -0.14; P<.001). Nursing homes with higher discrepancy between their CMS- and self-reported ratings have higher number of infections among their staff (log count = 0.42; 95% CI, 0.32 to 0.52; P<.001) and residents (log count = 0.13; 95% CI, 0.07 to 0.18; P<.001).

1  
2  
3 *Conclusions:* The size of COVID-19 outbreaks in nursing homes is associated with  
4  
5 their ratings and governance features. To prepare for the possible next waves of  
6  
7 COVID-19 epidemic, policy makers should use these insights to identify the nursing  
8  
9 homes who are more likely to experience large outbreaks.  
10  
11  
12  
13  
14  
15  
16

17 **Key words:** *COVID-19, Nursing Homes, Long-Term Care*  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Article Summary

### Strengths and limitations of this study

- examines the association between nursing home features and the likelihood and size of COVID-19 outbreaks amongst their staff and residents.
- develops and evaluates predictive models that can identify nursing homes with the highest chance of experiencing COVID-19 outbreaks.
- The findings are limited to nursing homes in the state of California.

## Introduction

Nursing homes have been most severely impacted by the COVID-19 pandemic owing to the advanced age and high number of comorbidities of their residents.<sup>1,2</sup> In Europe, as much as 57% of all deaths related to COVID-19 were at such facilities.<sup>3</sup> In the United States, nursing homes' residents and staff account for 34% of all COVID-19 fatalities.<sup>4</sup> Infection prevention and control at nursing homes and long-term facilities has therefore become a priority in managing the epidemic.<sup>5,6</sup>

Given the considerable variation in the prevalence and size of the COVID-19 outbreaks at nursing homes, the objective of this research is (1) to understand why some nursing homes are more susceptible to COVID-19 outbreaks, and (2) to develop predictive models that can identify such nursing homes so that they could be prioritized in efforts to prevent and contain next waves of the epidemic.<sup>7,8</sup>

## Methods

### *Patient and public involvement*

Patients had no influence on the research questions or outcomes of this research. No patients were involved in the design of this study. We used blind patient files; therefore, no patient recruitment took place. We only used data on the aggregated number of COVID patients and staff in the nursing homes as reported by the State of California and therefore no personal information of patients was used in this study. Given the nature of removing all personal information, there is no requirement to disseminate the information to patients.

### *Data Sources and Study Variables*

We collect data from various publicly available sources. The New York Times aggregates and provides data on COVID-19 cases per county.<sup>9</sup> California Department of Public Health (CDPH) provides data on the number of confirmed COVID-19



1  
2  
3 infections among staff and residents of nursing homes in the state.<sup>10</sup> CMS provides data  
4 on nursing home characteristics, including their self-reported ratings and CMS health  
5 inspections.<sup>11</sup> Applying the methods suggested by Han et. al,<sup>12</sup> we identify the nursing  
6 homes with significant discrepancies between their self-reported measures and  
7 independent CMS inspections. These methods rely on data that are only available for  
8 nursing homes in California and therefore, the scope of this study is also limited to  
9 nursing homes in California. After cleaning and merging the above-mentioned data  
10 sources, we analyse a final dataset consisting of 713 nursing homes in California.  
11 Details of the data cleaning and merging process is presented in Supplementary  
12 Appendix.

13  
14  
15 We examine the following outcomes in this study: whether a nursing home has at least  
16 one COVID-19 infection amongst its residents or staff, the number of confirmed  
17 COVID-19 infections among its residents, and the number of confirmed infections  
18 among its staff. We also calculate a fourth outcome that indicates the large outbreaks  
19 as the ones in which more than 10 members of staff or residents were infected with  
20 COVID-19. This threshold translates to approximately 95<sup>th</sup> percentile of the number of  
21 infected staff. Given that more residents are infected than staff, this threshold translates  
22 to 75<sup>th</sup> percentile of the number of residents.

23  
24  
25 The independent variables describe the severity of the COVID-19 outbreak in the  
26 surrounding area of a nursing home, its governance characteristics, as well as its ratings  
27 on quality, staffing and CMS inspections. Table 1 provides detailed description of the  
28 study variables.

### 29 *Statistical Analysis*

30  
31  
32 To answer the first research question and understand why some nursing homes are more  
33 susceptible to COVID-19 outbreaks, we apply Zero Inflated Bivariate Poisson (ZIBP)

1  
2  
3 regression. The model allows us to examine the effects of nursing homes' ratings,  
4  
5 governance features, and their surroundings on the likelihood and size of their COVID-  
6  
7 19 outbreaks. Econometric details of the model are provided by Walhin, 2001.<sup>13</sup>  
8  
9 Intuitively, our approach assumes that the number of zero's in the count of infected  
10  
11 staff and residents are generated either because the nursing home was in an area that  
12  
13 was less infected by the COVID-19 or because it implemented successful prevention  
14  
15 procedures to protect its staff and residents. Moreover, the model assumes that in a  
16  
17 nursing home, the number of infected staff covaries with the number of infected  
18  
19 residents since they can infect each other and since common infection prevention and  
20  
21 control policies apply to both groups. Taking this interdependency into account also  
22  
23 alleviates the concerns over the possible impact of omitted variables in our model. In  
24  
25 this particular context, because of the close proximity of residents and staff, the same  
26  
27 variables that could affect the number of infections among one group, would most likely  
28  
29 also impact the number of infections among the other group. The covariance coefficient  
30  
31 captures this interdependency in outcomes. As a sensitivity analysis, we also report the  
32  
33 results of zero-inflated double Poisson regression. In this model, the counts of  
34  
35 infections among staff and residents are assumed to be independent from each other.  
36  
37  
38 To answer the second research question and identify the nursing homes with the highest  
39  
40 risk of COVID-19 outbreaks, we use our models to predict the probability of  
41  
42 experiencing an infection and compare their performance with common machine  
43  
44 learning techniques, namely Neural Networks (NN) and Support Vector Machine with  
45  
46 Radial Basis Function kernel (SVM-RBF). Further details about these machine learning  
47  
48 techniques are provided in the Supplementary Appendix. We also measure the  
49  
50 performance of our models in predicting the nursing homes with highest risks of  
51  
52 experiencing large outbreaks with more than 10 infections.  
53  
54  
55  
56  
57  
58  
59  
60

## Results

### *Study Sample*

During the data cleaning and merging process, 493 nursing homes were eliminated from our final sample, either because their names were not matching across different datasets, or their ratings information is not available from CMS, or because their COVID-19 infections are not reported by CDPH. To ensure that the final sample is random and our results are not biased, we compared the eliminated nursing homes with the ones in the study sample. The results of two sample t-tests and logistic regression are presented in Supplementary Appendix. None of the observed governance factors affect the chance of being included in the sample. Amongst the remaining variables, while the difference with regards to quality ratings and county infections per 100K is statistically significant between the two groups, their magnitude is small and serve to make our estimates more conservative.

Study sample included 713 nursing homes in California. As reported in Table 1, as of May 1<sup>st</sup>, 2020, 23% of the study sample reported at least one COVID-19 infection among either their staff or residents. Of those, 31% experienced large outbreaks with more than 10 infections among either their staff or residents. The geographic spread of COVID-19 infections in California nursing homes is graphically presented in the Supplementary Appendix.

### *Preventing COVID-19 Infections*

As reported in the first panel of Table 2, the only variables with statistically significant impact on the chance of COVID-19 outbreaks at nursing homes are their size and the rate of infections per 100 thousand residents at the county in which they are located. For both of these variables, a one-unit of increase is associated with a 1% increase in the odds of experiencing at least one COVID-19 infection.

### *Controlling COVID-19 Outbreaks*

As reported in the second and third panel of Table 2, while the number of infections amongst both staff and residents increase with the size of the nursing home, they are not associated with the rate of infections per 100 thousand residents at the county in which the nursing home is located. This indicates that although the severity of COVID-19 epidemic in the surrounding area increases the chance of experiencing at least one infection at the nursing homes, it may not necessarily translate to larger outbreaks.

While the expected number of infected residents is 13 times higher in for-profit nursing homes, the number of infected staff in for-profit nursing homes is not statistically different from non-profit ones. Prior empirical research has repeatedly shown that for-profit nursing homes are inferior in many aspects of care quality.<sup>14-17</sup>

Occupancy rate is associated with a lower number of infections among staff such that a one percent increase in occupancy rate decreases the expected count of infections among staff by 2.4%.

Among the three different ratings, the CMS-reported health inspection rating is associated with a sizable decrease in the number of infections among both staff and residents. One unit of increase in CMS-reported health inspection ratings is associated with a 18% decrease in the expected number of infections in both staff and residents. A one-unit improvement in staffing rating is associated with a 23% decrease in the number of infections among residents. Note that better staff rating is highly dependent on higher ratio of staff to residents and the higher number of staff per resident would allow nursing homes to control infections more efficiency among their residents. While the observed association between ratings on health inspections and staffing with the number of infected staff and residents were expected, the association between self-reported quality ratings and the number of infections is the opposite of our expectations.

1  
2  
3 One unit of increase in self-reported quality ratings is associated with, respectively,  
4  
5 51% and 14% increase in infections among staff and residents. This finding is aligned  
6  
7 with the emerging stream of research that shows nursing homes embellish their self-  
8  
9 reported quality ratings and therefore these ratings may not always indicate better  
10  
11 quality of care for residents.<sup>12,18-21</sup> Our final variable, inflation score, quantifies the  
12  
13 discrepancy between the self- and CMS-reported ratings. The higher the discrepancy,  
14  
15 the more likely it is that the nursing home is overstating their quality measures. With a  
16  
17 one-unit increase in such discrepancy, the expected number of infections among staff  
18  
19 and residents increases by, 52% and 14%, respectively.  
20  
21  
22

### 23 *Improving the Quality Reporting System*

24  
25 CMS could solve these discrepancies and improve the reporting process by  
26  
27 implementing better inspection and auditing strategies.<sup>22</sup> Figure 1 shows how the  
28  
29 number of infections among staff and residents could be compressed had the self-  
30  
31 reported quality measures by nursing homes were truly reflecting their quality of care.  
32  
33 Given the importance of ratings for nursing homes,<sup>23</sup> with a reliable rating system with  
34  
35 no discrepancy between self- and CMS-reported measures, nursing homes would strive  
36  
37 to elevate their ratings through actual improvements in their quality of care. As shown  
38  
39 in the upper panel of Figure 1, compared to the current system, lower number of  
40  
41 predicted infections among staff would have been more frequent under an improved  
42  
43 rating system such that predicted average number of infections among staff would have  
44  
45 decreased from 1.85 to 1.52, which is equal to 17.6% fewer total infections across the  
46  
47 staff of all nursing homes. As shown in the lower panel of Figure 1, the same effect is  
48  
49 observed for nursing home residents. Had self-reported quality ratings were truly  
50  
51 reflecting the quality of care, the expected number of infections among residents of  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 nursing homes would have reduced from 8.67 to 8.15 which is equal to 5.8% fewer  
4  
5 total infections across the residents of all nursing homes.  
6

7  
8 Finally, the sizable covariance estimate (0.68; 95% CI 0.54 to 0.87; P=0.1) indicates  
9  
10 that the number of infected staff is not independent from the number of infected  
11  
12 residents. This observation empirically confirms our expectation of dependency  
13  
14 between the count of infections in staff and residents such that nursing homes with high  
15  
16 number of infected staff also have high number of infected residents. This finding was  
17  
18 expected as residents and staff are in close contact with each other and once infections  
19  
20 occur among the members of one group, it would be very difficult to prevent them in  
21  
22 the other group. More importantly, common infection control procedures implemented  
23  
24 by nursing homes would apply to both groups and prevent infections among both  
25  
26 groups.  
27  
28  
29

### 30 *Identifying Nursing Homes with Highest Chance of COVID-19 Infections & Outbreaks*

31  
32 Figure 2 compares the lift of the ZIBP model with those of NN and SVM-RBF. The  
33  
34 first 50 nursing homes are zoomed in at the top right corner of the figure. The ZIBP  
35  
36 model's performance is comparable with the common NN and SVM-RBF methods. For  
37  
38 the first 50 nursing homes, the rate of true positives of ZIBP model is between 2.45 and  
39  
40 2.73 times higher than that of a random selection model. The Area Under the Curve  
41  
42 (AUC) for ZIBP, NN and SVM-RBF models are respectively 0.68, 0.73, and 0.62.  
43  
44  
45

46  
47 Figure 3 presents the lifts of the ZIBP model in identifying the nursing homes with  
48  
49 large COVID-19 outbreaks among those that have confirmed at least ten infections. For  
50  
51 the first 50 nursing homes, ZIBP correctly identifies nursing homes with large  
52  
53 outbreaks among staff between 1.3 to 3.9 times better than a random selection model.  
54  
55 The model's performance for predicting large outbreaks among residents for the first  
56  
57 50 nursing homes is 1.5 to 2.1 times better than a random selection model.  
58  
59  
60

## Discussion

Staff and residents of nursing homes constitute the largest demographic of COVID-19 fatalities in the US. However, nursing homes have not been uniformly impacted by the epidemic; some have not experienced even a single infection while some others have been devastated by COVID-19 fatalities. To prepare for the possible next waves of the epidemic, it is critical to uncover the underlying reason of such variation and to explore the nursing homes' features that are associated with higher chance and size of outbreaks.

The aim of this research was to understand how publicly available data on nursing homes can explain the significant variation in the chance and size of COVID-19 infections at nursing homes, and to also develop predictive models that can identify the nursing homes with the highest chance and size of outbreaks.

Our results indicate that COVID-19 outbreaks are more likely to happen at larger nursing homes and those with higher rate of COVID-19 infections in the surrounding area.

Those with better staffing and health inspection ratings are more successful in controlling the outbreaks. Interestingly, self-reported quality ratings are associated with larger size of outbreaks. This counter-intuitive result could be further evidence that nursing homes exaggerate their self-reported quality measures. Higher discrepancy between self-reported measures and CMS-reported health inspections was associated with larger COVID-19 outbreaks.

The size of the outbreaks among residents is significantly higher in for-profit nursing homes which have been previously shown to also be of poorer quality in various aspects of care.<sup>14-17</sup>

1  
2  
3 The model developed in this research can correctly identify the nursing homes that are  
4 more likely to experience an infection or are at the highest risk of an outbreak.  
5  
6

7  
8 The insights of this research help policy makers to identify the nursing homes with the  
9 highest probability and size of COVID-19 outbreaks. This will allow them to prioritize  
10 such nursing homes in their efforts to control the epidemic. Such efforts could entail  
11 devoting more resources towards nursing homes with significantly higher risk or when  
12 feasible, temporarily transferring patients to different nursing homes to control the  
13 spread of the virus.  
14  
15

16  
17 This work leaves several areas for future research. First, given the variation in testing  
18 at different nursing homes, the number of confirmed infections may be undercounting  
19 the actual number of infections and therefore a more reliable measure would be the  
20 number of fatalities associated with COVID-19. Second, should temporal data become  
21 available, researchers can study growth curves of infections or deaths among staff and  
22 residents and examine their interlinked effects on each other. Third, should national  
23 data become available, we can test our contentions using a much larger sample at the  
24 national level. This would increase the external validity and generalizability of our  
25 findings.  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41

#### 42 **Author Contributions**

43  
44 RG and NY, designed the study. RG, XH, and NY had full access to all the data in the  
45 study and take responsibility for the integrity of the data and the accuracy of the data  
46 analyses. RG, XH, and NY analysed the data. RG and NY interpreted the data. NY  
47 drafted the manuscript. RG critically revised the manuscript.  
48  
49  
50  
51  
52

#### 53 **Funding Statement**

54  
55 This research received no specific grant from any funding agency in the public,  
56 commercial or not-for-profit sectors.  
57  
58  
59  
60



## Competing Interests

There are no competing interests for any of the authors.

## Data sharing statement

All data in this research are publicly available and their sources have been cited in the manuscript. Data on the discrepancy between self-reported and CMS-reported measures of nursing homes are available by request from the corresponding author.

## References

1. McMichael TM, Currie DW, Clark S, et al. Epidemiology of Covid-19 in a Long-Term Care Facility in King County, Washington. *N Engl J Med*. Published online March 27, 2020. doi:10.1056/NEJMoa2005412
2. Arentz M, Yim E, Klaff L, et al. Characteristics and Outcomes of 21 Critically Ill Patients With COVID-19 in Washington State. *JAMA*. 2020;323(16):1612-1614. doi:10.1001/jama.2020.4326
3. Comas-Herrera A, Zalakain J, Litwin C, Hsu AT, Lane N, Fernández J-L. Mortality associated with COVID-19 outbreaks in care homes: early international evidence. *Int Long-Term Care Policy Netw CPEC-LSE Internet*. Published online 2020.
4. Yourish K, Lai KKR, Ivory D, Smith M. One-Third of All U.S. Coronavirus Deaths Are Nursing Home Residents or Workers. *The New York Times*. <https://www.nytimes.com/interactive/2020/05/09/us/coronavirus-cases-nursing-homes-us.html>. Published May 11, 2020. Accessed May 12, 2020.
5. Bedford J, Enria D, Giesecke J, et al. COVID-19: towards controlling of a pandemic. *The Lancet*. 2020;395(10229):1015–1018.
6. Adalja AA, Toner E, Inglesby TV. Priorities for the US Health Community Responding to COVID-19. *JAMA*. 2020;323(14):1343-1344. doi:10.1001/jama.2020.3413
7. Xu S, Li Y. Beware of the second wave of COVID-19. *The Lancet*. 2020;395(10233):1321-1322. doi:10.1016/S0140-6736(20)30845-X
8. Leung K, Wu JT, Liu D, Leung GM. First-wave COVID-19 transmissibility and severity in China outside Hubei after control measures, and second-wave scenario planning: a modelling impact assessment. *The Lancet*. 2020;395(10233):1382-1393. doi:10.1016/S0140-6736(20)30746-7

- 1  
2  
3 9. The New York Times. California Coronavirus Map and Case Count. *The New*  
4 *York Times*. [https://www.nytimes.com/interactive/2020/us/california-](https://www.nytimes.com/interactive/2020/us/california-coronavirus-cases.html)  
5 [coronavirus-cases.html](https://www.nytimes.com/interactive/2020/us/california-coronavirus-cases.html). Accessed May 21, 2020.  
6  
7
- 8 10. California Department of Public Health. Skilled Nursing Facilities: COVID-19.  
9 Accessed May 21, 2020.  
10 [https://www.cdph.ca.gov/Programs/CID/DCDC/Pages/COVID-](https://www.cdph.ca.gov/Programs/CID/DCDC/Pages/COVID-19/SNFsCOVID_19.aspx)  
11 [19/SNFsCOVID\\_19.aspx](https://www.cdph.ca.gov/Programs/CID/DCDC/Pages/COVID-19/SNFsCOVID_19.aspx)  
12
- 13 11. Archived Datasets | Data.Medicare.gov. Data.Medicare.Gov. Accessed May 23,  
14 2020. <https://data.medicare.gov/data/archives/nursing-home-compare>  
15
- 16 12. Han X, Yaraghi N, Gopal R. Winning at all costs: Analysis of inflation in  
17 nursing homes' rating system. *Prod Oper Manag*. 2018;27(2):215–233.  
18
- 19 13. Walhin JF. Bivariate ZIP models. *Biom J*. 2001;43(2):147–160.  
20
- 21 14. Hillmer MP, Wodchis WP, Gill SS, Anderson GM, Rochon PA. Nursing Home  
22 Profit Status and Quality of Care: Is There Any Evidence of an Association?  
23 *Med Care Res Rev*. 2005;62(2):139-166. doi:10.1177/1077558704273769  
24
- 25 15. Comondore VR, Devereaux PJ, Zhou Q, et al. Quality of care in for-profit and  
26 not-for-profit nursing homes: systematic review and meta-analysis. *Bmj*.  
27 2009;339:b2732.  
28
- 29 16. Harrington C, Woolhandler S, Mullan J, Carrillo H, Himmelstein DU. Does  
30 investor ownership of nursing homes compromise the quality of care? *Am J*  
31 *Public Health*. 2001;91(9):1452–1455.  
32
- 33 17. Amirkhanyan AA, Kim HJ, Lambright KT. Does the public sector outperform  
34 the nonprofit and for-profit sectors? Evidence from a national panel study on  
35 nursing home quality and access. *J Policy Anal Manage*. 2008;27(2):326-353.  
36 doi:10.1002/pam.20327  
37
- 38 18. Johari K, Kellogg C, Vazquez K, Irvine K, Rahman A, Enguidanos S. Ratings  
39 game: an analysis of Nursing Home Compare and Yelp ratings. *BMJ Qual Saf*.  
40 2018;27(8):619-624. doi:10.1136/bmjqs-2017-007301  
41
- 42 19. Neuman MD, Wirtalla C, Werner RM. Association between skilled nursing  
43 facility quality indicators and hospital readmissions. *JAMA*. 2014;312(15):1542-  
44 1551. doi:10.1001/jama.2014.13513  
45
- 46 20. Sanghavi P, Pan S, Caudry D. Assessment of nursing home reporting of major  
47 injury falls for quality measurement on nursing home compare. *Health Serv Res*.  
48 2020;55(2):201-210. doi:10.1111/1475-6773.13247  
49
- 50 21. Fuller RL, Goldfield NI, Hughes JS, McCullough EC. Nursing Home Compare  
51 Star Rankings and the Variation in Potentially Preventable Emergency  
52 Department Visits and Hospital Admissions. *Popul Health Manag*.  
53 2019;22(2):144-152. doi:10.1089/pop.2018.0065  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57
22. Han X, Yaraghi N, Gopal R. Catching Them Red-Handed: Optimizing the Nursing Homes' Rating System. *ACM Trans Manag Inf Syst TMIS*. 2019;10(2):1–26.
  23. Werner RM, Konetzka RT, Polsky D. Changes in consumer demand following public reporting of summary quality ratings: an evaluation in nursing homes. *Health Serv Res*. 2016;51:1291–1309.

For peer review only

58 **Figures**  
59  
60

1  
2  
3 **Figure 1.** Impact of Improved Rating System on Infection Density Curves  
4

5 Note: The blue (solid) curve represents the density of predicted number of  
6 infections under current rating system while the red (dashed) curve shows the  
7 density of counterfactual number of infections had there been no discrepancy  
8 between self- and CMS-reported ratings. The vertical blue and red lines show  
9 the average number of predicted infections with and without discrepancy in  
10 ratings.  
11  
12  
13  
14  
15  
16

17 **Figure 2.** Comparison of Performance of ZIBP, NN, and SVM-RBF Models in  
18 Predicting at Least One Infection  
19  
20  
21  
22

23 **Figure 3.** Performance of ZIBP Model for Predicting Large Outbreaks (More than 10  
24 Infections) Among Staff and Residents  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57

58 **Tables**  
59  
60

**Table 1.** Sources and Descriptions of the Study Variables

Variable	Description	Source	Mean	Std. Dev.	Min	Max
<b>Outcomes</b>						
<b>Nursing home infected</b>	Indicates if the nursing home has at least one confirmed case of COVID-19 infection among its staff or residents	CDPH	0.23	0.42	0	1
<b>Confirmed residents</b>	The number of COVID-19 infections among the residents of nursing homes	CDPH	1.91	7.88	0	81
<b>Confirmed staff</b>	The number of COVID-19 infections among the staff of nursing homes	CDPH	0.41	2.19	0	26
<b>Large outbreak</b>	Among those nursing homes with at least 1 infection, indicates if the number of infected staff or residents is greater than 11.	Authors' calculation	0.31	0.46	0	1
<b>Severity of COVID-19 epidemic in the surrounding area</b>						
<b>County infections per 100K</b>	The rate of COVID-19 infections per 100,000 residents in the county in which the nursing home is located as of May 1 <sup>st</sup> , 2020.	New York Times	143.42	80.07	0	259.8
<b>Governance features</b>						
<b>For profit</b>	Indicates if the nursing home has a for-profit status	CMS	0.86	0.35	0	1
<b>Family council</b>	Indicates if a family council for the residents exists in the nursing home	CMS	0.2	0.4	0	1
<b>Certified beds</b>	The number of Medicaid? Certified beds in the nursing home	CMS	98.89	54.77	14	769
<b>Occupancy rate</b>	The ratio of residents to the total number of certified beds	Authors' calculation	0.87	0.12	0.14	1
<b>Inflation score</b>	Counts the number of years in which a significant discrepancy was observed between the self-reported quality measures and CMS-reported health inspections.	Authors' calculation	0.32	0.81	0	5
<b>Ratings</b>						
<b>Quality rating</b>	Self-reported indicator of quality of services as of 2017	CMS	4.59	0.87	0	5

<b>Staffing rating</b>	Self-reported indicator of staffing features as of 2017	CMS	3.41	1.13	0	5
<b>Health inspection rating</b>	CMS-reported indicator of health inspections ratings as of 2017	CMS	2.88	1.29	1	5

**Table 2.** Effects of study variables on the likelihood and the size of COVID-19 outbreaks

Parameter	Zero Inflated Bivariate Poisson Model			Zero Inflated Double Poisson Model		
	Estimate	(95% CI)	P Value	Estimate	(95% CI)	P Value
	e			e		
<b>Nursing Home (Likelihood of nursing home getting at least one COVID-19 infection)</b>						
Intercept	-2.41	(-4.48 to -0.34)	0.02	-1.76	(-3.75 to 0.24)	0.08
County infections per 100K	0.01	(0.01 to 0.02)	<.001	0.01	(0.01 to 0.02)	<.001
For profit	-0.3	(-0.88 to 0.28)	0.31	-0.27	(-0.85 to 0.31)	0.36
Family council	0.15	(-0.32 to 0.61)	0.53	0.21	(-0.26 to 0.67)	0.38
Certified beds	0.01	(0.01 to 0.02)	0.003	0.01	(0.01 to 0.02)	0.01
Occupancy rate	-0.18	(-1.97 to 1.62)	0.85	-0.98	(-2.69 to 0.74)	0.26
Inspection rating	-0.02	(-0.19 to 0.17)	0.91	-0.02	(-0.19 to 0.17)	0.90
Quality rating	-0.14	(-0.36 to 0.1)	0.25	-0.13	(-0.35 to 0.1)	0.27
Staffing rating	0.01	(-0.17 to 0.18)	0.96	-0.01	(-0.18 to 0.17)	0.96
Inflation score	0.05	(-0.18 to 0.28)	0.67	0.06	(-0.17 to 0.29)	0.61
<b>Infected Staff (number of staff with confirmed COVID-19 infections)</b>						
Intercept	0.29	(-2.02 to 2.59)	0.81	-0.43	(-2.1 to 1.25)	0.63
County infections per 100K	-0.01	(-0.01 to 0.01)	0.24	-0.01	(-0.01 to 0.01)	0.11
For profit	-0.27	(-0.84 to 0.3)	0.35	-0.16	(-0.55 to 0.24)	0.44
Family council	-0.06	(-0.56 to 0.45)	0.82	0.19	(-0.12 to 0.49)	0.24
Certified beds	0.01	(0.01 to 0.01)	<.001	0.01	(0.01 to 0.01)	0.02
Occupancy rate	-2.42	(-4.34 to -0.51)	0.01	-1.11	(-2.53 to 0.32)	0.13

Inspection rating	-0.2	(-0.38 to -0.01)	0.04	-0.16	(-0.28 to -0.03)	0.02
Quality rating	0.41	(0.13 to 0.68)	0.004	0.33	(0.15 to 0.52)	<.001
Staffing rating	0.11	(-0.07 to 0.28)	0.23	0.25	(0.12 to 0.37)	<.001
Inflation score	0.42	(0.32 to 0.52)	<.001	0.27	(0.19 to 0.35)	<.001
<b>Infected Residents (number of residents with confirmed COVID-19 infections)</b>						
Intercept	1.33	(0.33 to 2.33)	0.01	1.69	(0.84 to 2.55)	<.001
County infections per 100K	-0.01	(-0.01 to -0.01)	<.001	-0.01	(-0.01 to -0.01)	<.001
For profit	2.57	(1.99 to 3.15)	<.001	1.88	(1.51 to 2.26)	<.001
Family council	0.07	(-0.09 to 0.21)	0.40	0.1	(-0.04 to 0.24)	0.15
Certified beds	0.01	(0.01 to 0.01)	0.03	0.01	(-0.01 to 0.01)	0.13
Occupancy rate	-0.25	(-1.02 to 0.53)	0.53	-0.15	(-0.88 to 0.6)	0.71
Inspection rating	-0.2	(-0.26 to -0.14)	<.001	-0.2	(-0.26 to -0.14)	<.001
Quality rating	0.13	(0.05 to 0.21)	0.002	0.15	(0.08 to 0.23)	<.001
Staffing rating	-0.26	(-0.31 to -0.2)	<.001	-0.2	(-0.25 to -0.15)	<.001
Inflation score	0.13	(0.07 to 0.18)	<.001	0.11	(0.06 to 0.16)	<.001
Covariance	0.68	(0.54 to 0.87)	0.1			
<b>Fit Statistics</b>						
-2 log likelihood		4422.6			4561.7	
AIC		4484.6			4621.7	
BIC		4626.2			4758.8	

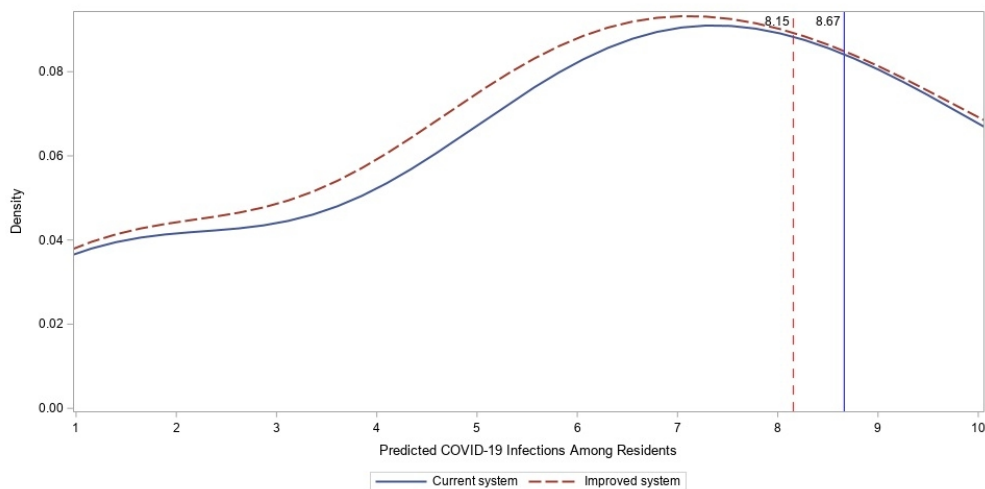
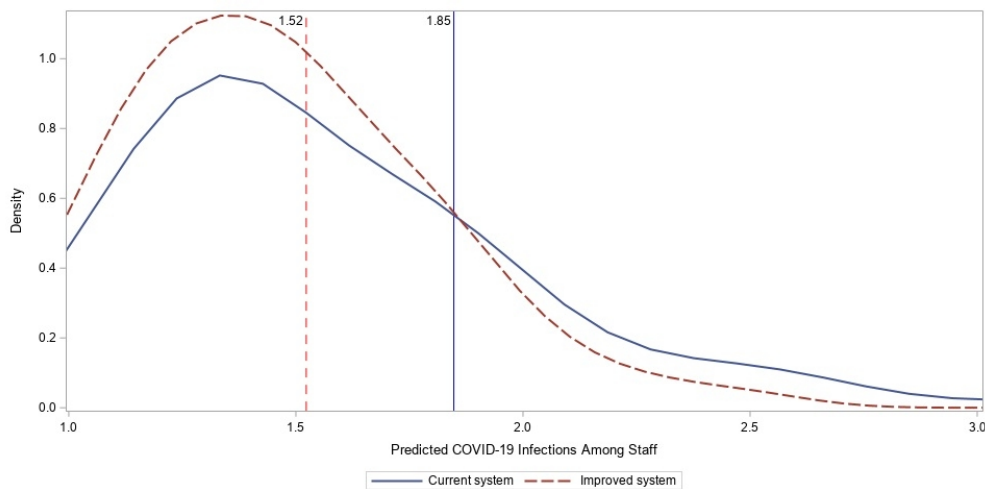


Figure 1

127x131mm (192 x 192 DPI)



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

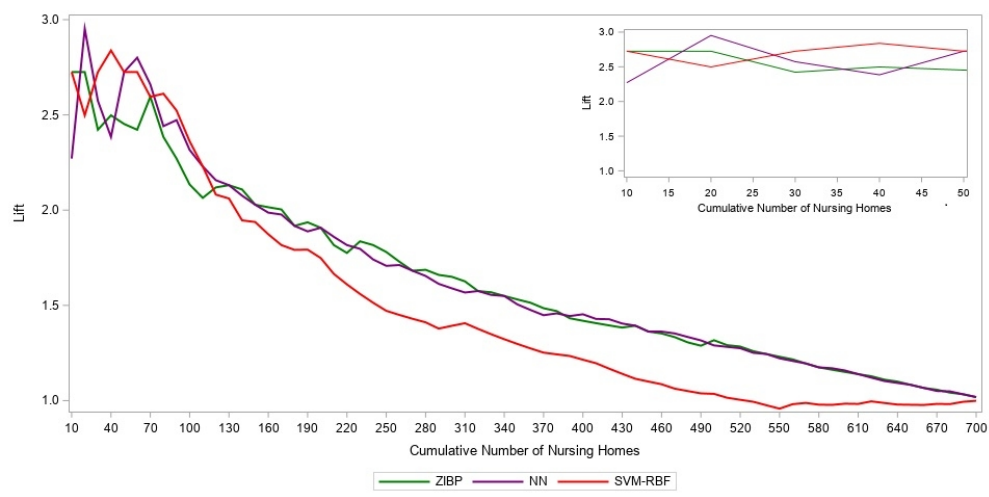


Figure 2

127x63mm (192 x 192 DPI)

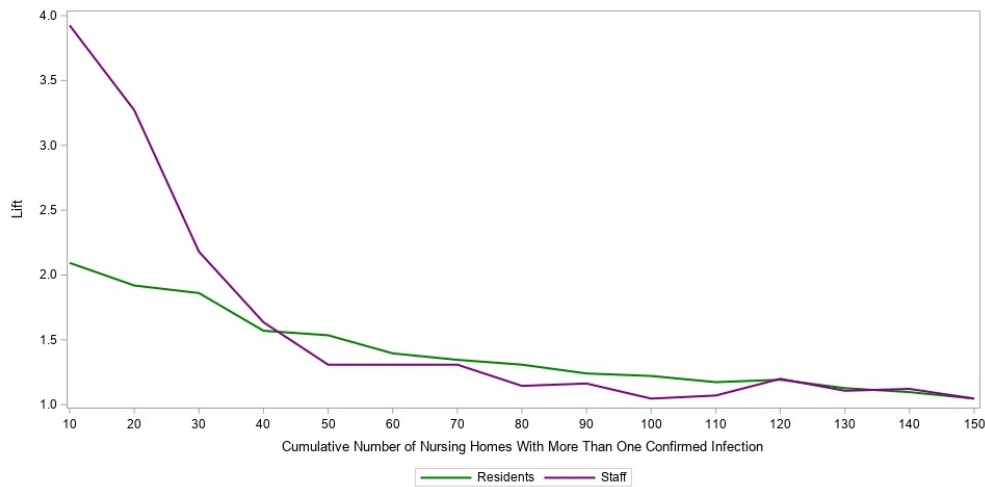


Figure 3

254x127mm (96 x 96 DPI)

# Supplementary Appendix

## Table of Contents

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Missing Observations .....	2
Machine learning Techniques.....	2
Figures .....	3
Figure S1. Study population and analysis sample.....	3
Figure S2: Spread of COVID-19 Infection Among California Nursing Homes.....	4
Figure S3: Receiver Operator Characteristic (ROC) Curves for Predicting at Least One Infection in Nursing Homes .....	5
Tables.....	6
Table S1: Logistic Regression Results for Estimating the Effects of Nursing Homes’ Features on Odds of Being Included in the Study Sample .....	6
Table S2: Results of Two-Sample t-Test for Equality of the Means of the Excluded and Included Observations.....	7

## Missing Observations

Data cleaning process is presented in Figure S1. 493 nursing homes were excluded from the study sample either due to the mismatch between their names across multiple datasets or because their COVID-19 infection data were not available in CDPH reports. To examine if the excluded nursing homes are similar to those included in the study sample, we conducted two logistic regression with the dependent variables set to be 1 to indicate if a record is included in the study sample and 0 otherwise. In the first logistic regression we only include governance features as independent variables, while in the second logistic regression we include all the features.

As reported in Table S1, both regression results show that none of the governance features are statistically significant, which indicates that the included records have no selection bias on governance features. Amongst the remaining variables, quality rating and county infections per 100k are significant are statistically significant yet the difference between the two groups is not substantial, as reported in Table S2. Further, the differences in these two variables across the two groups make our estimates more conservative.

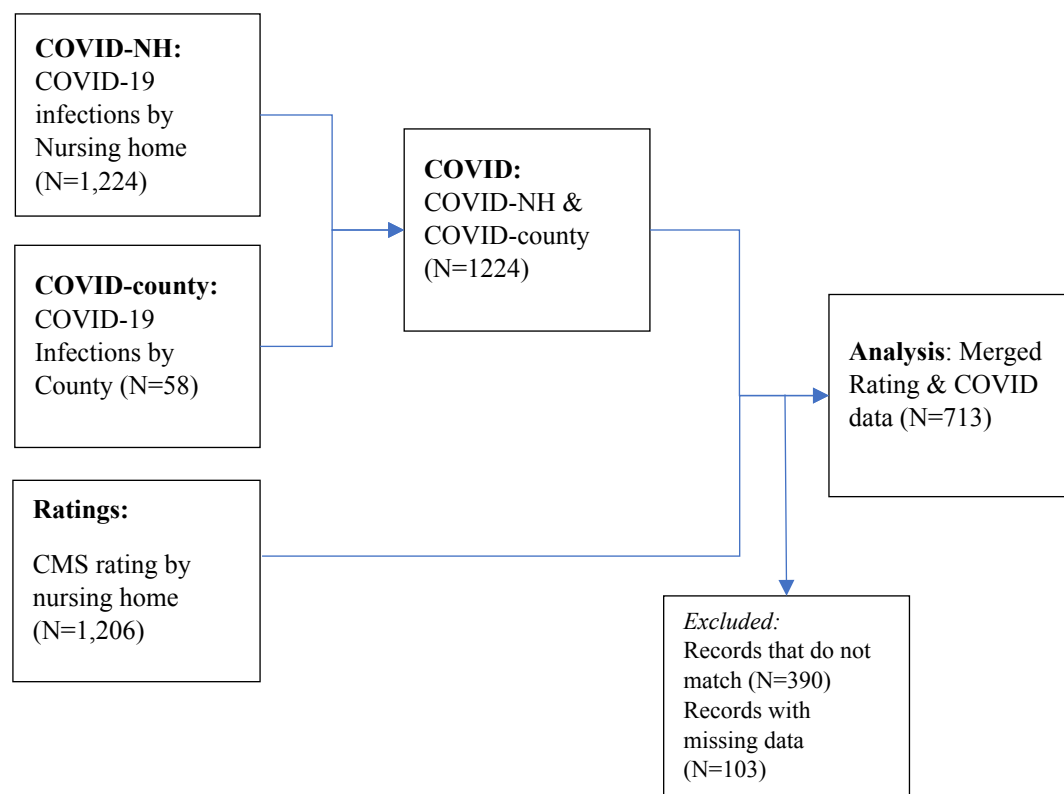
## Machine learning Techniques

We then apply machine learning techniques to predict the COVID-19 infection in nursing homes and compare the results with our model. In view that our problem has a highly nonlinear structure, advanced machine learning models that do not rely on data structure assumptions may provide a flexible and desired solution. We predict the nursing home level COVID-19 infection situation by using Neural Networks (NN) and Support Vector Machines (SVM) with RBF kernel function. Variable *NH* is used as the target variable in each model, and is equal to 1 if at least one patient or staff reported to be infected. The prediction features include nursing home governance features such as occupancy rate, number of certified beds, whether a family council presents, whether the nursing home is for profit or not, and inflation score evaluated from past years. The nursing homes' health inspection rating, staffing rating and quality rating are also included

in our prediction model. To capture the severity of COVID-19 epidemic in the surrounding area, we also incorporate county level COVID-19 infections per 100K population.

## Figures

Figure S1. Study population and analysis sample

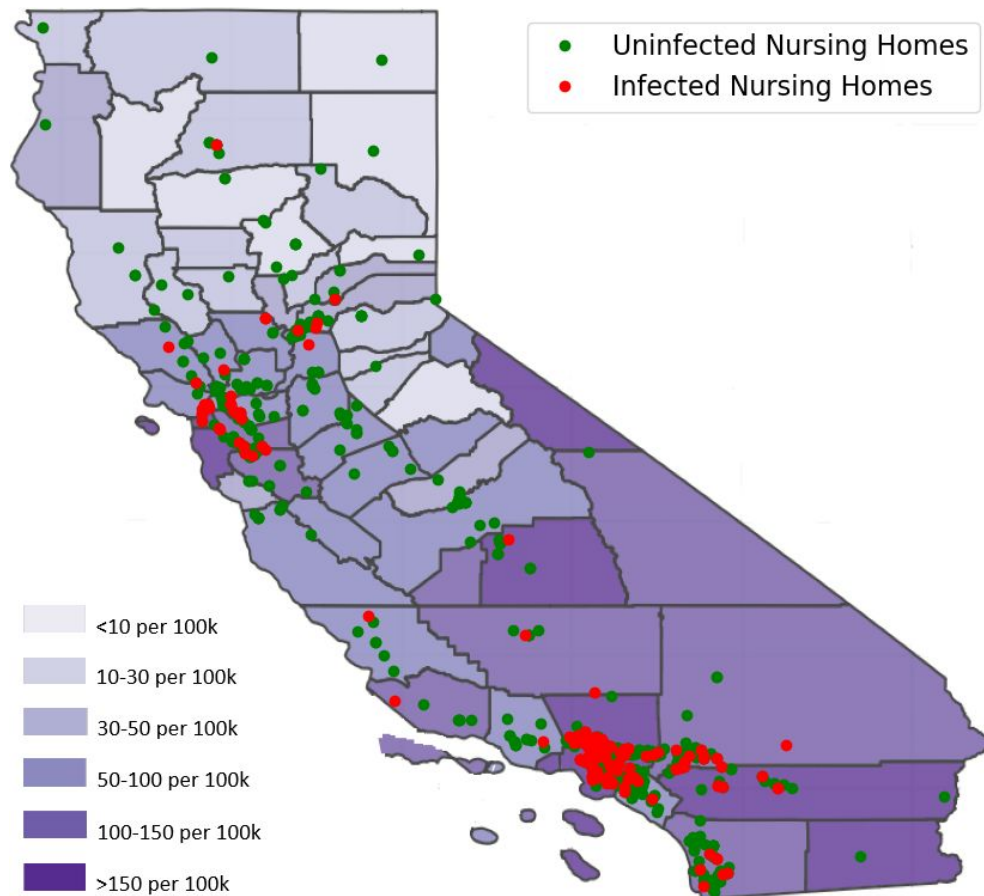


Note: Original CMS Rating for year 2017 data (*ratings*) include 1206 nursing homes. Original CA COVID-19 Infection by county (*COVID-county*) data as of April 30<sup>th</sup>, 2020 include on 58 counties. Original COVID-19 CA Infections by nursing homes (*COVID-NH*) data as of April 30<sup>th</sup>, 2020 include 1224 nursing homes.

We first merged *COVID-NH* and *COVID-county* data for all 1224 rows (0 record lost). We then merged the resulting data (*COVID*) with *ratings* data which resulted in 713 rows. 390 records were

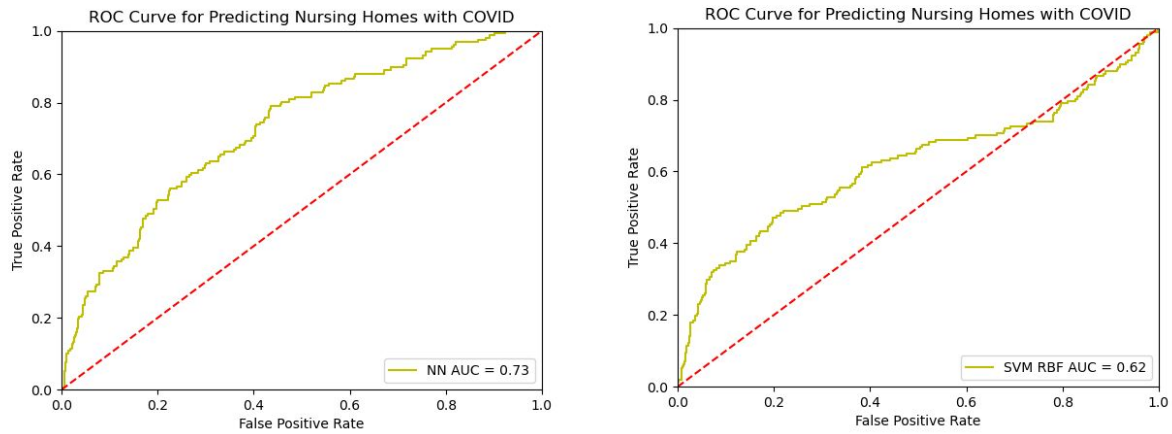
1  
2  
3 lost due to mismatch between the names of the facilities in the two datasets, and 103 records were  
4 lost for those nursing homes that did not report COVID 19 infection data or their ratings information  
5 is missing.  
6  
7  
8  
9

10  
11 **Figure S2: Spread of COVID-19 Infection Among California Nursing Homes**  
12  
13



44 Note: The figure presents the spread of COVID-19 infection among California nursing homes as  
45 of May 1<sup>st</sup>, 2020  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Figure S3: Receiver Operator Characteristic (ROC) Curves for Predicting at Least One Infection in Nursing Homes



Note: ROC for Nursing Home (NH) COVID-19 prediction using Neural Networks (NN), SVM with RBF kernel. The AUC is reported for each model: NN=0.73, SVM-RBF (default)=0.62

## Tables

Table S1: Logistic Regression Results for Estimating the Effects of Nursing Homes' Features on Odds of Being Included in the Study Sample

Parameter	Validation with Governance Features Only (Included vs. Excluded Records)			Validation with All Features (Included vs. Excluded Records)		
	Estimate	(95% CI)	P Value	Estimate	(95% CI)	P Value
Constant	0.1	(-0.72 to 0.92)	0.81	-0.66	(-2.09 to 0.76)	0.36
For profit	0.25	(-0.08 to 0.58)	0.14	0.29	(-0.1 to 0.68)	0.14
Family council	-0.19	(-0.49 to 0.12)	0.23	-0.07	(-0.4 to 0.26)	0.68
Certified beds	-0.0004	(-0.003 to 0.002)	0.71	-0.0008	(-0.003 to 0.002)	0.52
Occupancy rate	0.61	(-0.3 to 1.52)	0.19	0.56	(-0.62 to 1.74)	0.35
Inflation score	-0.04	(-0.2 to 0.12)	0.6	-0.03	(-0.2 to 0.14)	0.75
Quality rating				0.21	(0.07 to 0.36)	0.004
Staffing rating				0.002	(-0.14 to 0.14)	0.97
Health inspection rating				0.08	(-0.04 to 0.19)	0.21
County infections per 100K				-0.002	(-0.004 to -0.0007)	0.004



Table S2: Results of Two-Sample t-Test for Equality of the Means of the Excluded and Included Observations

Features	Excluded Records*	Included Records*	P Value**
For profit	0.82	0.86	0.11
Family council	0.21	0.18	0.21
Certified beds	99.6	98.0	0.65
Occupancy rate	0.85	0.86	0.14
Inflation score	0.32	0.31	0.83
Quality rating	4.43	4.57	0.01
Staffing rating	3.49	3.49	0.93
Health inspection rating	2.66	2.86	0.01
County infections per 100K	159.36	143.88	0.003

Note: \*: Reports the average value of features.

1  
2  
3 \*\*:P values are for two-tailed t-tests of the equality of the two means.  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For peer review only

**STROBE Statement**

Checklist of items that should be included in reports of observational studies

Section/Topic	Item No	Recommendation	Reported on Page No
<b>Title and abstract</b>	1	(a) Indicate the study's design with a commonly used term in the title or the abstract	1,2
		(b) Provide in the abstract an informative and balanced summary of what was done and what was found	2
<b>Introduction</b>			
Background/rationale	2	Explain the scientific background and rationale for the investigation being reported	5
Objectives	3	State specific objectives, including any prespecified hypotheses	5
<b>Methods</b>			
Study design	4	Present key elements of study design early in the paper	6
Setting	5	Describe the setting, locations, and relevant dates, including periods of recruitment, exposure, follow-up, and data collection	6
Participants	6	(a) <i>Cohort study</i> —Give the eligibility criteria, and the sources and methods of selection of participants. Describe methods of follow-up	6
		<i>Case-control study</i> —Give the eligibility criteria, and the sources and methods of case ascertainment and control selection. Give the rationale for the choice of cases and controls	
		<i>Cross-sectional study</i> —Give the eligibility criteria, and the sources and methods of selection of participants	
Variables	7	(b) <i>Cohort study</i> —For matched studies, give matching criteria and number of exposed and unexposed	6
		<i>Case-control study</i> —For matched studies, give matching criteria and the number of controls per case	
Data sources/measurement	8*	Clearly define all outcomes, exposures, predictors, potential confounders, and effect modifiers. Give diagnostic criteria, if applicable	6
Bias	9	For each variable of interest, give sources of data and details of methods of assessment (measurement). Describe comparability of assessment methods if there is more than one group	7 & Appendix
Study size	10	Describe any efforts to address potential sources of bias	6 & Appendix
Quantitative variables	11	Explain how the study size was arrived at	7
		Explain how quantitative variables were handled in the analyses. If applicable, describe which groupings were chosen and why	7
		(a) Describe all statistical methods, including those used to control for confounding	7
		(b) Describe any methods used to examine subgroups and interactions	7
		(c) Explain how missing data were addressed	7
Statistical methods	12	(d) <i>Cohort study</i> —If applicable, explain how loss to follow-up was addressed	7
		<i>Case-control study</i> —If applicable, explain how matching of cases and controls was addressed	7

*Cross-sectional study*—If applicable, describe analytical methods taking account of sampling strategy

(e) Describe any sensitivity analyses

7

Section/Topic	Item No	Recommendation	Reported on Page No
<b>Results</b>			
Participants	13*	(a) Report numbers of individuals at each stage of study—eg numbers potentially eligible, examined for eligibility, confirmed eligible, included in the study, completing follow-up, and analysed	8
		(b) Give reasons for non-participation at each stage	Appendix
		(c) Consider use of a flow diagram	Appendix
Descriptive data	14*	(a) Give characteristics of study participants (eg demographic, clinical, social) and information on exposures and potential confounders	8
		(b) Indicate number of participants with missing data for each variable of interest	8
		(c) <i>Cohort study</i> —Summarise follow-up time (eg, average and total amount)	
Outcome data	15*	<i>Cohort study</i> —Report numbers of outcome events or summary measures over time	
		<i>Case-control study</i> —Report numbers in each exposure category, or summary measures of exposure	
		<i>Cross-sectional study</i> —Report numbers of outcome events or summary measures	8,9
Main results	16	(a) Give unadjusted estimates and, if applicable, confounder-adjusted estimates and their precision (eg, 95% confidence interval). Make clear which confounders were adjusted for and why they were included	9,10
		(b) Report category boundaries when continuous variables were categorized	
		(c) If relevant, consider translating estimates of relative risk into absolute risk for a meaningful time period	
Other analyses	17	Report other analyses done—eg analyses of subgroups and interactions, and sensitivity analyses	10
<b>Discussion</b>			
Key results	18	Summarise key results with reference to study objectives	12
Limitations	19	Discuss limitations of the study, taking into account sources of potential bias or imprecision. Discuss both direction and magnitude of any potential bias	13
Interpretation	20	Give a cautious overall interpretation of results considering objectives, limitations, multiplicity of analyses, results from similar studies, and other relevant evidence	12
Generalisability	21	Discuss the generalisability (external validity) of the study results	13
<b>Other Information</b>			
Funding	22	Give the source of funding and the role of the funders for the present study and, if applicable, for the original study on which the present article is based	13

\*Give information separately for cases and controls in case-control studies and, if applicable, for exposed and unexposed groups in cohort and cross-sectional studies.

**Note:** An Explanation and Elaboration article discusses each checklist item and gives methodological background and published examples of transparent reporting. The STROBE checklist is best used in conjunction with this article (freely available on the Web sites of PLoS Medicine at <http://www.plosmedicine.org/>, Annals of Internal Medicine at <http://www.annals.org/>, and Epidemiology at <http://www.epidem.com/>). Information on the STROBE Initiative is available at [www.strobe-statement.org](http://www.strobe-statement.org).

3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

For peer review only

# BMJ Open

## Compress the Curve: Variations in COVID-19 Infections Across California Nursing Homes

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2020-042804.R1
Article Type:	Original research
Date Submitted by the Author:	05-Nov-2020
Complete List of Authors:	Gopal, Ram; University of Warwick Han, Xu; Fordham University Yaraghi, Niam; Brookings Institution, Governance Studies
<b>Primary Subject Heading</b>:	Geriatric medicine
Secondary Subject Heading:	Infectious diseases
Keywords:	COVID-19, GERIATRIC MEDICINE, Health policy < HEALTH SERVICES ADMINISTRATION & MANAGEMENT

SCHOLARONE™  
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

## Original Investigation

**Title:** Compress the Curve: Variations in COVID-19 Infections Across California  
Nursing Homes

Ram D. Gopal, PhD<sup>1</sup>, Xu Han, PhD<sup>2</sup>, Niam Yaraghi, PhD<sup>3,4,\*</sup>

<sup>1</sup>: Professor, Warwick Business School, University of Warwick

<sup>2</sup>: Assistant Professor, Gabelli School of Business, Fordham University

<sup>3</sup>: Assistant Professor, Miami Herbert Business School, University of Miami

<sup>4</sup>: Non-resident Fellow, Governance Studies, The Brookings Institution

*\*: All authors contributed equally. Authors are listed in alphabetical order of their last name.*

Corresponding author:

Niam Yaraghi

[Nyaraghi@brookings.edu](mailto:Nyaraghi@brookings.edu)

5250 University Drive Coral Gables, FL 33146

Phone: (305) 284-3314

Word count: 2668



## Abstract

*Objective:* Nursing homes' residents and staff constitute the largest proportion of the fatalities associated with COVID-19 epidemic. Although there is a significant variation in COVID-19 outbreaks among the US nursing homes, we still do not know why such outbreaks are larger and more likely in some nursing homes than others. This research aims to understand why some nursing homes are more susceptible to larger COVID-19 outbreaks.

*Design:* Observational study of all nursing homes in the state of California until May 1st, 2020.

*Setting:* The state of California.

*Participants:* 713 long term care facilities in the State of California that participate in public reporting of COVID-19 infections as of May 1<sup>st</sup>, 2020 and their infections data could be matched with data on ratings and governance features of nursing homes provided by CMS.

*Main Outcome Measure:* The number of reported COVID-19 infections among staff and residents.

*Results:* Study sample included 713 nursing homes. The size of outbreaks among residents in for-profit nursing homes is 12.7 times larger than their non-profit counterparts (log count = 2.54; 95% CI, 1.97 to 3.11; P<.001). Higher ratings in CMS-reported health inspections are associated with lower number of infections among both staff (log count = -0.19; 95% CI, -0.37 to -0.01; P = 0.05) and residents (log count = -0.20; 95% CI, -0.27 to -0.14; P<.001). Nursing homes with higher discrepancy between their CMS- and self-reported ratings have higher number of infections among their staff (log count = 0.41; 95% CI, 0.31 to 0.51; P<.001) and residents (log count = 0.13; 95% CI, 0.08 to 0.18; P<.001).

1  
2  
3 *Conclusions:* The size of COVID-19 outbreaks in nursing homes is associated with  
4  
5 their ratings and governance features. To prepare for the possible next waves of  
6  
7 COVID-19 epidemic, policy makers should use these insights to identify the nursing  
8  
9 homes who are more likely to experience large outbreaks.  
10  
11  
12  
13  
14  
15  
16

17 **Key words:** *COVID-19, Nursing Homes, Long-Term Care*  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Article Summary

### Strengths and limitations of this study

- Examines the association between nursing home features and the likelihood and size of COVID-19 outbreaks amongst their staff and residents.
- Develops and evaluates predictive models that can identify nursing homes with the highest chance of experiencing COVID-19 outbreaks.
- The findings of the study are limited by the fact that the study was conducted with data only from California.
- The number of COVID-19 cases reported by nursing homes may be subject to under-reporting.
- The dataset on nursing homes' features is based on the year 2017 which is two years prior to the outbreak.

## Introduction

Nursing homes have been most severely impacted by the COVID-19 pandemic owing to the advanced age and high number of comorbidities of their residents.<sup>1,2</sup> In Europe, as much as 57% of all deaths related to COVID-19 were at such facilities.<sup>3</sup> In the United States, nursing homes' residents and staff account for 34% of all COVID-19 fatalities.<sup>4</sup> Infection prevention and control at nursing homes and long-term facilities has therefore become a priority in managing the epidemic.<sup>5,6</sup>

Given the considerable variation in the prevalence and size of the COVID-19 outbreaks at nursing homes, the objective of this research is (1) to understand why some nursing homes are more susceptible to COVID-19 outbreaks, and (2) to develop predictive models that can identify such nursing homes so that they could be prioritized in efforts to prevent and contain next waves of the epidemic.<sup>7,8</sup>

## Methods

### *Patient and public involvement*

Patients had no influence on the research questions or outcomes of this research. No patients were involved in the design of this study. We used blind patient files; therefore, no patient recruitment took place. We only used data on the aggregated number of COVID-19 patients and staff in the nursing homes as reported by the State of California and therefore no personal information of patients was used in this study. Given the nature of removing all personal information, there is no requirement to disseminate the information to patients.

### *Data Sources and Study Variables*

We collected data from various publicly available sources. The New York Times aggregates and provides data on COVID-19 cases per county.<sup>9</sup> California Department of Public Health (CDPH) provides data on the number of confirmed COVID-19

1  
2  
3 infections among staff and residents of nursing homes in the state.<sup>10</sup> CMS provides data  
4 on nursing home characteristics, including their self-reported ratings and CMS health  
5 inspections.<sup>11</sup> A description of this data is provided in the next section. Applying the  
6 methods suggested by Han et. al,<sup>12</sup> we identified the nursing homes with significant  
7 discrepancies between their self-reported measures and independent CMS inspections  
8 for a consecutive 5-year period. We aggregated the results and used the number of  
9 years a nursing home is predicted to be a likely inflator as the overall inflation score for  
10 a nursing home. Therefore, an honest nursing home will have an inflation score of 0  
11 while an inflating nursing home can have an inflation score between 1 to 5, with 5 being  
12 the most severe. In our dataset, 19.25% of nursing homes were inflating their scores  
13 and some of these had a score of 5 indicating that they inflated their scores in all 5  
14 years.

15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31 These methods rely on data that are only available for nursing homes in California and  
32 therefore, the scope of this study is also limited to nursing homes in California. After  
33 cleaning and merging the above-mentioned data sources, we analysed a final dataset  
34 consisting of 713 nursing homes in California. Details of the data cleaning and merging  
35 process is presented in Supplementary Appendix.

36  
37  
38  
39  
40  
41  
42 We examined the following outcomes in this study: whether a nursing home has at least  
43 one COVID-19 infection amongst its residents or staff, the number of confirmed  
44 COVID-19 infections among its residents, and the number of confirmed infections  
45 among its staff. We also calculated a fourth outcome that indicates the large outbreaks  
46 as the ones in which more than 10 members of staff or residents were infected with  
47 COVID-19. This threshold translates to approximately 95<sup>th</sup> percentile of the number of  
48 infected staff. Given that more residents are infected than staff, this threshold translates  
49 to 75<sup>th</sup> percentile of the number of residents.

1  
2  
3 The independent variables describe the severity of the COVID-19 outbreak in the  
4 surrounding area of a nursing home, its governance characteristics, as well as its ratings  
5 on quality, staffing and CMS inspections. Table 1 provides detailed description of the  
6 study variables. Note that while almost all nursing homes have resident councils,  
7 only 20 percent of nursing homes have existing family councils. We included  
8 the existence of family council as a binary variable in our analysis with the  
9 contention that it may imply closer coordination and higher engagement with  
10 the families of the residents.  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

### 23 *Description of CMS' Nursing Home Compare System*

24 The CMS nursing home rating data consists of basic information about nursing facilities  
25 such as name, address, phone number, etc., as well as some key features used in our  
26 analysis, such the number of certified beds, whether the nursing home is for-profit or  
27 non-profit, whether the nursing home has a family council, etc.  
28  
29  
30  
31  
32  
33

34 The CMS nursing home rating data serves the CMS Nursing Home Compare System,  
35 in which nursing home ratings are generated based on three domains: Inspection,  
36 Staffing, and Quality measures. The Inspection is conducted and reported by CMS-  
37 certified inspectors annually. The other two domains are self-reported by nursing  
38 homes. The annual inspection investigates areas such as medication management,  
39 nursing home administration, environment, food service, and residents' rights and  
40 quality of life. The Staffing domain is evaluated based on the self-reported CMS  
41 Certification and Survey Provider Enhanced Reports (CASPER) staffing data. The two  
42 measures used are the total nursing hours and Registered Nursing (RN) hours and are  
43 adjusted for case-mix based on the Resource Utility Group (RUG-III) case-mix system  
44 derived from the Minimum Data Set (MDS). The staffing star rating is then updated by  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 the end of the quarter when raw data is collected. Note that with more recent changes,  
4 the Staffing data reported by nursing homes is subject to validation with nursing homes'  
5 payroll data reported through Payroll-Based Journal (PBJ). The Quality Measure rating  
6 uses quality measurement criteria, which covers both long-stay terms and short-stay  
7 terms. The quality measure star rating is updated by the end of each quarter by using  
8 the results from three most recent quarters.  
9

10 To calculate the star ratings, CMS first assigns an initial star rating to all nursing homes  
11 based on their annual inspection results. Nursing homes are then assigned star ratings  
12 for the Staffing and Quality Measures domains. The overall star rating is then  
13 calculated by considering the inspection rating as the baseline, increasing or decreasing  
14 by 1 star if any self-reported domain satisfies the conditions stated as follows. Both 4  
15 and 5 stars in staffing rating are qualified for obtaining additional overall star rating,  
16 while only 5 stars in quality measure is qualified. Additional conditions apply to nursing  
17 homes whose inspection ratings are only 1 star, and for nursing homes which are in the  
18 CMS's Special Focus Facility (SFF) program. The overall star rating is lowered by one  
19 star if any self-reported domain is 1 star. The overall star rating cannot be more than 5  
20 stars or less than 1 star. Detailed data from CMS on nursing homes is available online.<sup>13</sup>  
21

### 22 *Statistical Analysis*

23 To answer the first research question and understand why some nursing homes are more  
24 susceptible to COVID-19 outbreaks, we applied Zero Inflated Bivariate Poisson (ZIBP)  
25 regression. The model allows us to examine the effects of nursing homes' ratings,  
26 governance features, and their surroundings on the likelihood and size of their COVID-  
27 19 outbreaks. Econometric details of the model are provided by Walhin, 2001.<sup>14</sup>  
28 Conventional Poisson models are suitable for modelling count data, while the zero  
29 inflated variation of Poisson model is more suitable for modelling count data with  
30

1  
2  
3 excess zeros, especially when excess zeros are generated by a separate processes that  
4 could be modelled separately. This leads to a framework that consists of a logit model  
5 for estimating the excess zeros in addition to a Poisson count model. ZIBP model is an  
6 extension of zero inflated Poisson model and is best suited for situations in which the  
7 count data with excess zeros are generated for two outcomes that may be correlated. In  
8 cases were the outcome variables are independent, the model reduces to the product of  
9 two independent zero inflated Poisson regression models, referred to as Zero Inflated  
10 Double Poisson model. in our setting, the two count variables are the number of  
11 COVID-19 infections among staff, and residents. These counts include excess zeros  
12 since many nursing homes reported no COVID-19 cases, primarily because they are  
13 located in areas where at the time of the data collection, had not yet experienced  
14 significant surges in COVID-19 cases. These two counts are also correlated since they  
15 both happen at the same nursing home and the factors that give rise to them are common  
16 at the nursing home level.

17  
18 Intuitively, we assume that the number of zero's in the count of infected staff and  
19 residents are generated either because the nursing home was in an area that was less  
20 infected by the COVID-19 or because it implemented successful prevention procedures  
21 to protect its staff and residents. Moreover, we assume that in a nursing home, the  
22 number of infected staff covaries with the number of infected residents since they can  
23 infect each other and since common infection prevention and control policies apply to  
24 both groups. Taking this interdependency into account also alleviates the concerns over  
25 the possible impact of omitted variables in our model. In this context, because of the  
26 close proximity of residents and staff, the same variables that could affect the number  
27 of infections among one group, would most likely also impact the number of infections  
28 among the other group. The covariance coefficient captures this interdependency in  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 outcomes. As a sensitivity analysis, we also report the results of zero-inflated double  
4  
5 Poisson regression. In this model, the counts of infections among staff and residents are  
6  
7 assumed to be independent from each other. We use NLMIXED procedure in SAS  
8  
9 software to estimate our models.<sup>15,16</sup> Note that we have provided access to both the data  
10  
11 and the SAS code for this analysis.<sup>17,18</sup>  
12  
13

14 To answer the second research question and identify the nursing homes with the  
15  
16 highest risk of COVID-19 outbreaks, we used our models to predict the probability of  
17  
18 experiencing an infection and compared their performance with common machine  
19  
20 learning techniques, namely Neural Networks (NN) and Support Vector Machine with  
21  
22 Radial Basis Function kernel (SVM-RBF). Since our problem has a highly nonlinear  
23  
24 structure, advanced machine learning models such as NN and SVM that do not rely  
25  
26 on data structure assumptions may provide a flexible and desired solution. Variable  
27  
28 NH is used as the target variable in each model, and NH is equal to 1 if at least one  
29  
30 patient or staff reported to be infected. The prediction features include nursing home  
31  
32 governance features such as occupancy rate, number of certified beds, whether a  
33  
34 family council presents, whether the nursing home is for profit or not, and inflation  
35  
36 score evaluated from past years. The nursing homes' health inspection rating, staffing  
37  
38 rating and quality rating are also included. The machine learning models are  
39  
40 implemented in Python 3.7 with 70% data training and 30% data testing. The entire  
41  
42 dataset is used to plot the lift chart. We also measured the performance of our models  
43  
44 in predicting the nursing homes with highest risks of experiencing large outbreaks  
45  
46 with more than 10 infections.  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Results

### *Study Sample*

During the data cleaning and merging process, 493 nursing homes were eliminated from our final sample, either because their names were not matching across different datasets, or their ratings information is not available from CMS, or because their COVID-19 infections are not reported by CDPH. To ensure that the final sample is random and our results are not biased, we compared the eliminated nursing homes with the ones in the study sample. The results of two sample t-tests and logistic regression are presented in Supplementary Appendix. None of the observed governance factors affect the chance of being included in the sample. Amongst the remaining variables, while the difference with regards to quality ratings and county infections per 100K is statistically significant between the two groups, their magnitude is small and serve to make our estimates more conservative.

Study sample included 713 nursing homes in California. As reported in Table 1, as of May 1<sup>st</sup>, 2020, 23% of the study sample reported at least one COVID-19 infection among either their staff or residents. Of those, 31% experienced large outbreaks with more than 10 infections among either their staff or residents. The geographic spread of COVID-19 infections in California nursing homes is graphically presented in the Supplementary Appendix.

### *Preventing COVID-19 Infections*

According to the model selection criteria reported in Table 2, the ZIBP model provides a better fit as its AIC, BIC and -2Log Likelihood are all smaller than those of Zero Inflated Double Poisson model. We therefore report the estimates of the ZIBP model in the text. The coefficients in the first panel of Table 2 represent how the log odds of experiencing an infection changes with one unit of increase in the corresponding

1  
2  
3 predictor. As reported in the first panel of Table 2, the only variables with statistically  
4  
5 significant impact on the chance of COVID-19 outbreaks at nursing homes are their  
6  
7 size and the rate of infections per 100 thousand residents at the county in which they  
8  
9 are located. For both variables, a one-unit of increase is associated with a 1% increase  
10  
11 in the odds of experiencing at least one COVID-19 infection.  
12  
13

#### 14 *Controlling COVID-19 Outbreaks*

15  
16 The coefficients in the second and third panel of Table 2 represent how the expected  
17  
18 log count of the infections changes for each unit increase in the corresponding predictor.  
19  
20 As reported in the second and third panel of Table 2, the expected rate of infections  
21  
22 amongst both staff and residents increase with the size of the nursing home. This  
23  
24 indicates that although the severity of COVID-19 epidemic in the surrounding area  
25  
26 increases the chance of experiencing at least one infection at the nursing homes.  
27  
28

29  
30 While the size of outbreaks among residents is about 12.7 times higher in for-profit  
31  
32 nursing homes, the size of outbreak among staff in for-profit nursing homes is not  
33  
34 statistically different from non-profit ones. This is in line with prior empirical research  
35  
36 that has repeatedly shown that for-profit nursing homes are inferior in many aspects of  
37  
38 care quality.<sup>19-22</sup>  
39  
40

41  
42 Occupancy rate, which represents the ratio of the number of enrolled patients to the  
43  
44 number of certified beds of a nursing home, is associated with a lower rate of infections  
45  
46 among staff such that a one percent increase in occupancy rate decreases the expected  
47  
48 count of infections among staff by 2.4%.  
49  
50

51  
52 Among the three different ratings, the CMS-reported health inspection rating is  
53  
54 associated with a sizable decrease in the number of infections among both staff and  
55  
56 residents. One unit of increase in CMS-reported health inspection ratings is associated  
57  
58 with a 17% and 18% decrease in the expected number of infections in staff and  
59  
60

1  
2  
3 residents, respectively. A one-unit improvement in staffing rating is associated with a  
4  
5 23% decrease in the number of infections among residents. Note that better staff rating  
6  
7 is highly dependent on higher ratio of staff to residents and the higher number of staff  
8  
9 per resident would allow nursing homes to control infections more efficiency among  
10  
11 their residents. While the observed association between ratings on health inspections  
12  
13 and staffing with the number of infected staff and residents were expected, the  
14  
15 association between self-reported quality ratings and the number of infections is the  
16  
17 opposite of our expectations. One unit of increase in self-reported quality ratings is  
18  
19 associated with, respectively, 49% and 14% increase in infections among staff and  
20  
21 residents. This finding is aligned with the emerging stream of research that shows  
22  
23 nursing homes embellish their self-reported quality ratings and therefore these ratings  
24  
25 may not always indicate better quality of care for residents.<sup>12,23-26</sup> Our final variable,  
26  
27 inflation score, quantifies the discrepancy between the self- and CMS-reported ratings.  
28  
29 The higher the discrepancy, the more likely it is that the nursing home is overstating  
30  
31 their quality measures. With a one-unit increase in such discrepancy, the expected  
32  
33 number of infections among staff and residents increases by, 51% and 14%,  
34  
35 respectively.

#### 41 42 *Improving the Quality Reporting System*

43  
44 CMS could solve these discrepancies and improve the reporting process by  
45  
46 implementing better inspection and auditing strategies.<sup>27</sup> Figure 1 shows how the  
47  
48 number of infections among staff and residents could be compressed had the self-  
49  
50 reported quality measures by nursing homes were truly reflecting their quality of care.  
51  
52 Given the importance of ratings for nursing homes,<sup>28</sup> with a reliable rating system with  
53  
54 no discrepancy between self- and CMS-reported measures, nursing homes would strive  
55  
56 to elevate their ratings through actual improvements in their quality of care. As shown  
57  
58  
59  
60

1  
2  
3 in the upper panel of Figure 1, compared to the current system, lower number of  
4  
5 predicted infections among staff would have been more frequent under an improved  
6  
7 rating system such that predicted average number of infections among staff would have  
8  
9 decreased from 1.85 to 1.52, which is equal to 17.6% fewer total infections across the  
10  
11 staff of all nursing homes. As shown in the lower panel of Figure 1, the same effect is  
12  
13 observed for nursing home residents. Had self-reported quality ratings were truly  
14  
15 reflecting the quality of care, the expected number of infections among residents of  
16  
17 nursing homes would have reduced from 8.67 to 8.15 which is equal to 5.8% fewer  
18  
19 total infections across the residents of all nursing homes.  
20  
21  
22

23  
24 Finally, the sizable covariance estimate (0.68; 95% CI 0.54 to 0.87; P=0.1) indicates  
25  
26 that the number of infected staff is not independent from the number of infected  
27  
28 residents. This observation empirically confirms our expectation of dependency  
29  
30 between the count of infections in staff and residents such that nursing homes with high  
31  
32 number of infected staff also have high number of infected residents. This finding was  
33  
34 expected as residents and staff are in close contact with each other and once infections  
35  
36 occur among the members of one group, it would be very difficult to prevent them in  
37  
38 the other group. More importantly, common infection control procedures implemented  
39  
40 by nursing homes would apply to both groups and prevent infections among both  
41  
42 groups. Note that as discussed earlier, according to all the model selection criteria, the  
43  
44 ZIBP performs better than its competitors. This is not surprising since it has the  
45  
46 advantage of modelling and adjusting for the correlation between the count of infections  
47  
48 among staff and residents. In the Appendix, we provide further empirical details on the  
49  
50 correlation between the number of infections among residents and staff.  
51  
52  
53  
54

55  
56 *Identifying Nursing Homes with Highest Chance of COVID-19 Infections & Outbreaks*  
57  
58  
59  
60

1  
2  
3 Figure 2 compares the lift of the ZIBP model with those of NN and SVM-RBF. We use  
4 lift as a measure for the ability of the model at predicting or classifying cases with  
5 respect to random selection. Lift shows how much better our model works compared to  
6 a random selection model. The first 50 nursing homes are zoomed in at the top right  
7 corner of the figure. The ZIBP model's performance is comparable with the common  
8 NN and SVM-RBF methods. For the first 50 nursing homes, the rate of true positives  
9 of ZIBP model is between 2.45 and 2.73 times higher than that of a random selection  
10 model. The Area Under the Curve (AUC) for ZIBP, NN and SVM-RBF models are  
11 respectively 0.68, 0.73, and 0.62.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

24 Figure 3 presents the lifts of the ZIBP model in identifying the nursing homes with  
25 large COVID-19 outbreaks among those that have confirmed at least ten infections. For  
26 the first 50 nursing homes, ZIBP correctly identifies nursing homes with large  
27 outbreaks among staff between 1.3 to 3.9 times better than a random selection model.  
28 The model's performance for predicting large outbreaks among residents for the first  
29 50 nursing homes is 1.5 to 2.1 times better than a random selection model.  
30  
31  
32  
33  
34  
35  
36  
37

### 38 **Discussion**

39  
40 Staff and residents of nursing homes constitute the largest demographic of COVID-19  
41 fatalities in the US. However, nursing homes have not been uniformly impacted by the  
42 epidemic; some have not experienced even a single infection while some others have  
43 been devastated by COVID-19 fatalities. To prepare for the possible next waves of the  
44 epidemic, it is critical to uncover the underlying reason of such variation and to explore  
45 the nursing homes' features that are associated with higher chance and size of  
46 outbreaks.  
47  
48  
49  
50  
51  
52  
53  
54

55  
56 The aim of this research was to understand how publicly available data on nursing  
57 homes can explain the significant variation in the chance and size of COVID-19  
58  
59  
60

1  
2  
3 infections at nursing homes, and to also develop predictive models that can identify the  
4 nursing homes with the highest chance and size of outbreaks.  
5

6  
7  
8 Our results indicate that COVID-19 outbreaks are more likely to happen at larger  
9 nursing homes and those with higher rate of COVID-19 infections in the surrounding  
10 area. These factors have been shown to be associated with higher probability of  
11 experiencing infections by other researchers as well.<sup>29</sup>  
12  
13

14  
15 Those with better staffing and health inspection ratings are more successful in  
16 controlling the outbreaks. The association between staffing levels and likelihood of  
17 having COVID-19 infections among both staff and residents has been reported by other  
18 researchers as well.<sup>30</sup> Interestingly, higher self-reported quality ratings are associated  
19 with larger size of outbreaks. This counter-intuitive result could be further evidence  
20 that nursing homes exaggerate their self-reported quality measures. Higher discrepancy  
21 between self-reported measures and CMS-reported health inspections was associated  
22 with larger COVID-19 outbreaks.  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

34  
35 The size of the outbreaks among residents is significantly higher in for-profit nursing  
36 homes which have been previously shown to also be of poorer quality in various aspects  
37 of care.<sup>19-22</sup>  
38  
39  
40

41  
42 There is a complex relationship between the main variables in our models. For-profit  
43 NHs generally have lower nurse staffing, more deficiencies, are larger in size, and have  
44 a greater likelihood of inflating their ratings.<sup>31,32</sup> It is therefore not surprising that they  
45 were found to be more likely to have larger numbers of COVID infected residents and  
46 staff.  
47  
48  
49  
50  
51

52  
53 The model developed in this research can correctly identify the nursing homes that are  
54 more likely to experience an infection or are at the highest risk of an outbreak.  
55  
56  
57  
58  
59  
60

1  
2  
3 The insights of this research help policy makers to identify the nursing homes with the  
4 highest probability and size of COVID-19 outbreaks. This will allow them to prioritize  
5 such nursing homes in their efforts to control the epidemic. Such efforts could entail  
6 devoting more resources towards nursing homes with significantly higher risk or when  
7 feasible, temporarily transferring patients to different nursing homes to control the  
8 spread of the virus.  
9

10  
11 Our results show that our ZIBP model outperforms SVM and that the predictive ability  
12 of the NN is only modestly better than ZIBP model. That is, the application and  
13 comparison of these machine learning models with the results of the ZIBP model  
14 confirms that not only the ZIBP model can explain the relationship between various  
15 independent variables and COVID-19 infections at nursing homes, but it also offers  
16 competitive predictive performance.  
17

18  
19 An important takeaway from this research is the importance of data collection and  
20 transparency. Our research was made possible because of the availability of key  
21 information on COVID-19 infections in nursing homes in the US and publicly available  
22 data such as ownership, size, staffing, and key performance measures. Access to such  
23 data is invaluable in both understanding and taking preventive action to curb the  
24 COVID-19 infections in nursing homes. As such we hope that other industrialized  
25 nations take necessary steps to collect and disseminate such information to protect and  
26 safeguard the vulnerable residents in long-term care facilities.  
27

28  
29 This work leaves several areas for future research. First, given the variation in testing  
30 at different nursing homes, the number of confirmed infections may be undercounting  
31 the actual number of infections and therefore a more reliable measure would be the  
32 number of fatalities associated with COVID-19. Second, should temporal data become  
33 available, researchers can study growth curves of infections or deaths among staff and  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 residents and examine their interlinked effects on each other. Third, should national  
4 data become available, we can test our contentions using a much larger sample at the  
5 national level. This would increase the external validity and generalizability of our  
6 findings. Finally, when data from other states and other time becomes available, we can  
7 include a spatial random effect in the model to account for spatial dependencies  
8 between the infections at different nursing homes.  
9

10 One of the limitations of the study is that its data on nursing homes' features is collected  
11 in 2017 which is over two years prior to the outbreak. Although more recent data were  
12 available on the time of the study, the variable "inflation score" had to be adopted from  
13 the 2017 data. We should also note that 86 percent of CA nursing homes are for-profit  
14 and these nursing homes were probably more likely to under-report their infection rates  
15 and deaths than other nursing homes for fear of losing residents and revenue.<sup>33</sup>  
16

### 17 **Author Contributions**

18 RG and NY, designed the study. RG, XH, and NY had full access to all the data in the  
19 study and take responsibility for the integrity of the data and the accuracy of the data  
20 analyses. RG, XH, and NY analysed the data. RG and NY interpreted the data. NY  
21 drafted the manuscript. NY, and RG critically revised the manuscript.  
22

### 23 **Funding Statement**

24 This research received no specific grant from any funding agency in the public,  
25 commercial or not-for-profit sectors.  
26

### 27 **Competing Interests**

28 There are no competing interests for any of the authors.  
29

### 30 **Data sharing statement**

All data in this research are publicly available and their sources have been cited in the manuscript. Data on the discrepancy between self-reported and CMS-reported measures of nursing homes are available by request from the corresponding author.

## References

1. McMichael TM, Currie DW, Clark S, et al. Epidemiology of Covid-19 in a Long-Term Care Facility in King County, Washington. *N Engl J Med*. Published online March 27, 2020. doi:10.1056/NEJMoa2005412
2. Arentz M, Yim E, Klaff L, et al. Characteristics and Outcomes of 21 Critically Ill Patients With COVID-19 in Washington State. *JAMA*. 2020;323(16):1612-1614. doi:10.1001/jama.2020.4326
3. Comas-Herrera A, Zalakain J, Litwin C, Hsu AT, Lane N, Fernández J-L. Mortality associated with COVID-19 outbreaks in care homes: early international evidence. *Int Long-Term Care Policy Netw CPEC-LSE Internet*. Published online 2020.
4. Yourish K, Lai KKR, Ivory D, Smith M. One-Third of All U.S. Coronavirus Deaths Are Nursing Home Residents or Workers. *The New York Times*. <https://www.nytimes.com/interactive/2020/05/09/us/coronavirus-cases-nursing-homes-us.html>. Published May 11, 2020. Accessed May 12, 2020.
5. Bedford J, Enria D, Giesecke J, et al. COVID-19: towards controlling of a pandemic. *The Lancet*. 2020;395(10229):1015–1018.
6. Adalja AA, Toner E, Inglesby TV. Priorities for the US Health Community Responding to COVID-19. *JAMA*. 2020;323(14):1343-1344. doi:10.1001/jama.2020.3413
7. Xu S, Li Y. Beware of the second wave of COVID-19. *The Lancet*. 2020;395(10233):1321-1322. doi:10.1016/S0140-6736(20)30845-X
8. Leung K, Wu JT, Liu D, Leung GM. First-wave COVID-19 transmissibility and severity in China outside Hubei after control measures, and second-wave scenario planning: a modelling impact assessment. *The Lancet*. 2020;395(10233):1382-1393. doi:10.1016/S0140-6736(20)30746-7
9. The New York Times. California Coronavirus Map and Case Count. *The New York Times*. <https://www.nytimes.com/interactive/2020/us/california-coronavirus-cases.html>. Accessed May 21, 2020.
10. California Department of Public Health. Skilled Nursing Facilities: COVID-19. Accessed May 21, 2020. [https://www.cdph.ca.gov/Programs/CID/DCDC/Pages/COVID-19/SNFsCOVID\\_19.aspx](https://www.cdph.ca.gov/Programs/CID/DCDC/Pages/COVID-19/SNFsCOVID_19.aspx)

11. Archived Datasets | Data.Medicare.gov. Data.Medicare.Gov. Accessed May 23, 2020. <https://data.medicare.gov/data/archives/nursing-home-compare>
12. Han X, Yaraghi N, Gopal R. Winning at all costs: Analysis of inflation in nursing homes' rating system. *Prod Oper Manag.* 2018;27(2):215–233.
13. Calgary O. Archived Datasets | Data.Medicare.gov. Data.Medicare.Gov. Accessed October 27, 2020. <https://data.medicare.gov/data/archives/nursing-home-compare>
14. Walhin JF. Bivariate ZIP models. *Biom J.* 2001;43(2):147–160.
15. AlMuhayfith FE, Alzaid AA, Omair MA. On bivariate Poisson regression models. *J King Saud Univ-Sci.* 2016;28(2):178–189.
16. Guide SU. Version 9.SAS Inst. *Inc Cary NC.* Published online 2004.
17. SAS code for BMJ. figshare. doi:10.6084/m9.figshare.13179875.v1
18. Data for COVID-19 in California nursing homes. Published online October 30, 2020. doi:10.6084/m9.figshare.13148813.v3
19. Hillmer MP, Wodchis WP, Gill SS, Anderson GM, Rochon PA. Nursing Home Profit Status and Quality of Care: Is There Any Evidence of an Association? *Med Care Res Rev.* 2005;62(2):139-166. doi:10.1177/1077558704273769
20. Comondore VR, Devereaux PJ, Zhou Q, et al. Quality of care in for-profit and not-for-profit nursing homes: systematic review and meta-analysis. *Bmj.* 2009;339:b2732.
21. Harrington C, Woolhandler S, Mullan J, Carrillo H, Himmelstein DU. Does investor ownership of nursing homes compromise the quality of care? *Am J Public Health.* 2001;91(9):1452–1455.
22. Amirkhanyan AA, Kim HJ, Lambright KT. Does the public sector outperform the nonprofit and for-profit sectors? Evidence from a national panel study on nursing home quality and access. *J Policy Anal Manage.* 2008;27(2):326-353. doi:10.1002/pam.20327
23. Johari K, Kellogg C, Vazquez K, Irvine K, Rahman A, Enguidanos S. Ratings game: an analysis of Nursing Home Compare and Yelp ratings. *BMJ Qual Saf.* 2018;27(8):619-624. doi:10.1136/bmjqs-2017-007301
24. Neuman MD, Wirtalla C, Werner RM. Association between skilled nursing facility quality indicators and hospital readmissions. *JAMA.* 2014;312(15):1542-1551. doi:10.1001/jama.2014.13513
25. Sanghavi P, Pan S, Caudry D. Assessment of nursing home reporting of major injury falls for quality measurement on nursing home compare. *Health Serv Res.* 2020;55(2):201-210. doi:10.1111/1475-6773.13247

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60
26. Fuller RL, Goldfield NI, Hughes JS, McCullough EC. Nursing Home Compare Star Rankings and the Variation in Potentially Preventable Emergency Department Visits and Hospital Admissions. *Popul Health Manag.* 2019;22(2):144-152. doi:10.1089/pop.2018.0065
  27. Han X, Yaraghi N, Gopal R. Catching Them Red-Handed: Optimizing the Nursing Homes' Rating System. *ACM Trans Manag Inf Syst TMIS.* 2019;10(2):1-26.
  28. Werner RM, Konetzka RT, Polsky D. Changes in consumer demand following public reporting of summary quality ratings: an evaluation in nursing homes. *Health Serv Res.* 2016;51:1291-1309.
  29. Abrams HR, Loomer L, Gandhi A, Grabowski DC. Characteristics of U.S. Nursing Homes with COVID-19 Cases. *J Am Geriatr Soc.* 2020;68(8):1653-1656. doi:10.1111/jgs.16661
  30. Harrington C, Ross L, Chapman S, Halifax E, Spurlock B, Bakerjian D. Nurse Staffing and Coronavirus Infections in California Nursing Homes. *Policy Polit Nurs Pract.* 2020;21(3):174-186. doi:10.1177/1527154420938707
  31. McGregor MJ, Harrington C. COVID-19 and long-term care facilities: Does ownership matter? *CMAJ.* 2020;192(33):E961-E962. doi:10.1503/cmaj.201714
  32. Harrington C, Olney B, Carrillo H, Kang T. Nurse Staffing and Deficiencies in the Largest For-Profit Nursing Home Chains and Chains Owned by Private Equity Companies. *Health Serv Res.* 2012;47(1pt1):106-128. doi:10.1111/j.1475-6773.2011.01311.x
  33. Harrington C, Pollock AM, Sutaria S. Privatization of Nursing Homes in the United Kingdom and the United States. In: *The Privatization of Care, The Case of Nursing Homes.* Routledge; 2019:51-67.

## Figures

### Figure 1. Impact of Improved Rating System on Infection Density Curves

Note: The blue (solid) curve represents the density of predicted number of infections under current rating system while the red (dashed) curve shows the density of counterfactual number of infections had there been no discrepancy between self- and CMS-reported ratings. The vertical blue and red lines show the average number of predicted infections with and without discrepancy in ratings.

### Figure 2. Comparison of Performance of ZIBP, NN, and SVM-RBF Models in Predicting at Least One Infection

Note: The first 50 nursing homes are zoomed in at the top right corner of the figure. The lift of ZIBP model is presented in green, while the lifts of NN and SVM-RBF are presented with purple and red lines respectively.

### Figure 3. Performance of ZIBP Model for Predicting Large Outbreaks (More than 10 Infections) Among Staff and Residents

Note: The lifts of the ZIBP model for identifying large outbreaks among residents and staff are presented, respectively, by the green and purple line.

## Tables

**Table 1.** Sources and Descriptions of the Study Variables

Variable	Description	Source	Mean	Std. Dev.	Min	Max
<b>Outcomes</b>						
<b>Nursing home infected</b>	Indicates if the nursing home has at least one confirmed case of COVID-19 infection among its staff or residents	CDPH	0.23	0.42	0	1
<b>Confirmed residents</b>	The number of COVID-19 infections among the residents of nursing homes	CDPH	1.91	7.88	0	81
<b>Confirmed staff</b>	The number of COVID-19 infections among the staff of nursing homes	CDPH	0.41	2.19	0	26
<b>Large outbreak</b>	Among those nursing homes with at least 1 infection, indicates if the number of infected staff or residents is more than 10 infections.	Authors' calculation	0.31	0.46	0	1
<b>Severity of COVID-19 epidemic in the surrounding area</b>						
<b>County infections per 100K</b>	The rate of COVID-19 infections per 100,000 residents in the county in which the nursing home is located as of May 1 <sup>st</sup> , 2020.	New York Times	143.42	80.07	0	259.8
<b>Governance features</b>						
<b>For profit</b>	Indicates if the nursing home has a for-profit status	CMS	0.86	0.35	0	1
<b>Family council</b>	Indicates if a family council for the residents exists in the nursing home	CMS	0.2	0.4	0	1
<b>Certified beds</b>	The number of beds certified to provide care to Medicare and Medicaid beneficiaries	CMS	98.89	54.77	14	769
<b>Occupancy rate</b>	The ratio of residents to the total number of certified beds	Authors' calculation	0.87	0.12	0.14	1

<b>Inflation score</b>	Counts the number of years in which a significant discrepancy was observed between the self-reported quality measures and CMS-reported health inspections.	Authors' calculation	0.32	0.81	0	5
<b>Ratings</b>						
<b>Quality rating</b>	Self-reported indicator of quality of services as of 2017	CMS	4.59	0.87	0	5
<b>Staffing rating</b>	Self-reported measure of staffing hours as of 2017. This is based on a combination of registered nurse hours per resident day and the total nursing hours per resident day.	CMS	3.41	1.13	0	5
<b>Health inspection rating</b>	CMS-reported indicator of health inspections ratings as of 2017	CMS	2.88	1.29	1	5

**Table 2.** Effects of study variables on the likelihood and the size of COVID-19 outbreaks

	Zero Inflated Bivariate Poisson Model			Zero Inflated Double Poisson Model			
	Parameter	Estimate	(95% CI)	P Value	Estimate	(95% CI)	P Value
<b>Nursing Home (Likelihood of nursing home getting at least one COVID-19 infection)</b>							
Intercept	-2.34	(-4.41 to -0.28)	0.03	-1.76	(-3.75 to 0.24)	0.08	
County infections per 100K	0.01	(0.01 to 0.02)	<.001	0.01	(0.01 to 0.02)	<.001	
For profit	-0.36	(-0.94 to 0.22)	0.22	-0.27	(-0.85 to 0.31)	0.36	
Family council	0.19	(-0.28 to 0.64)	0.44	0.21	(-0.26 to 0.67)	0.38	
Certified beds	0.01	(0.01 to 0.02)	0.01	0.01	(0.01 to 0.02)	0.01	
Occupancy rate	-0.2	(-1.99 to 1.59)	0.83	-0.98	(-2.69 to 0.74)	0.26	
Inspection rating	-0.02	(-0.19 to 0.17)	0.9	-0.02	(-0.19 to 0.17)	0.90	
Quality rating	-0.14	(-0.36 to 0.1)	0.26	-0.13	(-0.35 to 0.1)	0.27	
Staffing rating	0.01	(-0.17 to 0.18)	0.97	-0.01	(-0.18 to 0.17)	0.96	

Inflation score	0.06	(-0.18 to 0.28)	0.67	0.06	(-0.17 to 0.29)	0.61
<b>Infected Staff (number of staff with confirmed COVID-19 infections)</b>						
Intercept	0.21	(-2.11 to 2.52)	0.87	-0.43	(-2.1 to 1.25)	0.63
County infections per				-0.01	(-0.01 to 0.01)	0.11
100K	-0.01	(-0.01 to 0.01)	0.23			
For profit	-0.21	(-0.78 to 0.37)	0.49	-0.16	(-0.55 to 0.24)	0.44
Family council	-0.04	(-0.54 to 0.46)	0.89	0.19	(-0.12 to 0.49)	0.24
Certified beds	0.01	(0.01 to 0.01)	<.001	0.01	(0.01 to 0.01)	0.02
Occupancy rate	-2.39	(-4.3 to -0.47)	0.02	-1.11	(-2.53 to 0.32)	0.13
Inspection rating				-0.16	(-0.28 to -	0.02
	-0.19	(-0.37 to -0.01)	0.05		0.03)	
Quality rating	0.4	(0.13 to 0.67)	0.01	0.33	(0.15 to 0.52)	<.001
Staffing rating	0.11	(-0.07 to 0.28)	0.23	0.25	(0.12 to 0.37)	<.001
Inflation score	0.41	(0.31 to 0.51)	<.001	0.27	(0.19 to 0.35)	<.001
<b>Infected Residents (number of residents with confirmed COVID-19 infections)</b>						
Intercept	1.36	(0.36 to 2.35)	0.01	1.69	(0.84 to 2.55)	<.001
County infections per				-0.01	(-0.01 to -	<.001
100K	-0.01	(-0.01 to -0.01)	<.001		0.01)	
For profit	2.54	(1.97 to 3.11)	<.001	1.88	(1.51 to 2.26)	<.001
Family council	0.07	(-0.09 to 0.21)	0.4	0.1	(-0.04 to 0.24)	0.15
Certified beds	0.01	(0.01 to 0.01)	0.04	0.01	(-0.01 to 0.01)	0.13
Occupancy rate	-0.24	(-1.01 to 0.54)	0.55	-0.15	(-0.88 to 0.6)	0.71
Inspection rating				-0.2	(-0.26 to -	<.001
	-0.2	(-0.27 to -0.14)	<.001		0.14)	
Quality rating	0.13	(0.05 to 0.21)	0.01	0.15	(0.08 to 0.23)	<.001
Staffing rating				-0.2	(-0.25 to -	<.001
	-0.26	(-0.31 to -0.2)	<.001		0.15)	
Inflation score	0.13	(0.08 to 0.18)	<.001	0.11	(0.06 to 0.16)	<.001
Covariance	0.69	(0.54 to 0.87)	0.01			
<b>Fit Statistics</b>						
-2 log likelihood			4422.7			4561.7
AIC			4484.7			4621.7
BIC			4626.4			4758.8



1  
2  
3 Note: The coefficients in the first panel represent how the log odds of experiencing an infection  
4 changes with one unit of increase in the corresponding predictor. The coefficients in the second  
5 and third panels represent how the expected log count of the infections changes for each unit  
6 increase in the corresponding predictor.  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For peer review only

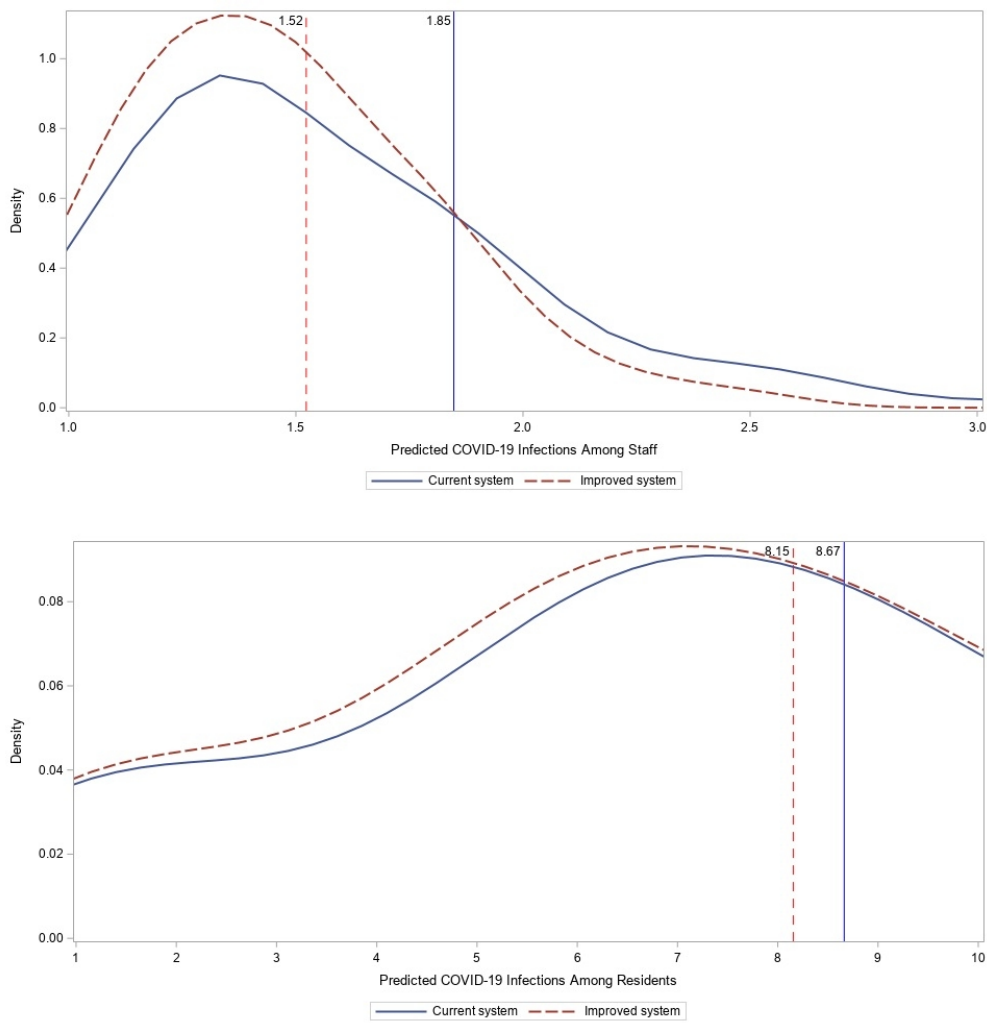


Figure 1

127x131mm (192 x 192 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

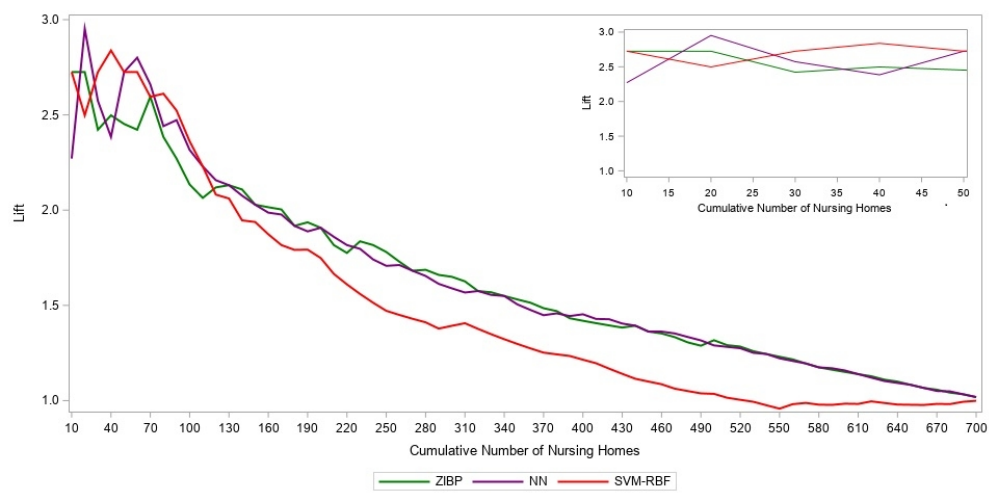


Figure 2

127x63mm (192 x 192 DPI)

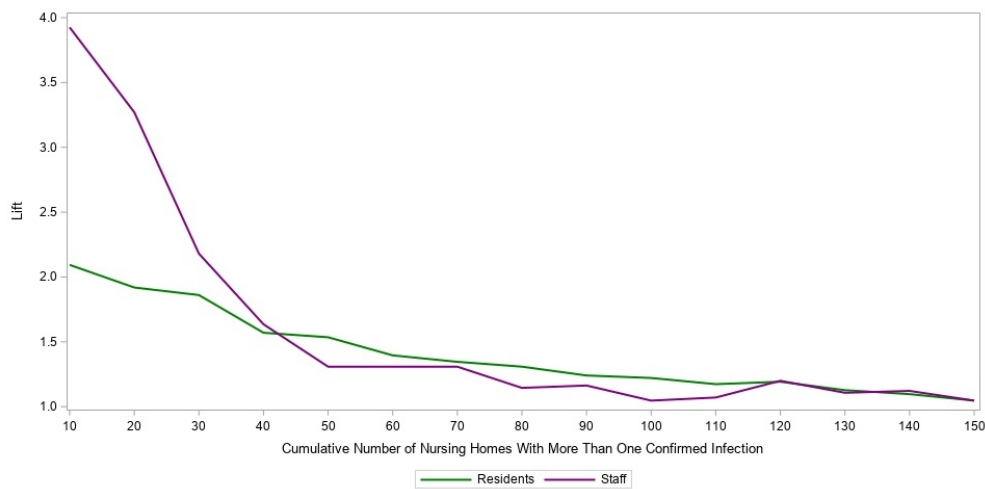


Figure 3

254x127mm (96 x 96 DPI)

## Supplementary Appendix for

### Compress the Curve: Variations in COVID-19 Infections Across California Nursing Homes

#### Table of Contents

Missing Observations.....	2
Machine learning Techniques.....	2
Bivariate and Double Poisson Estimates .....	3
Correlation Between Infections Among Staff and Residents .....	3
Figures .....	4
Figure S1. Study population and analysis sample.....	4
Figure S2: Spread of COVID-19 Infection Among California Nursing Homes.....	5
Figure S3: Receiver Operator Characteristic (ROC) Curves for Predicting at Least One Infection in Nursing Homes .....	6
Figure S4: Scatter plot of number of infections among staff and residents for those nursing homes that have experienced large outbreaks amongst both their staff and resident populations .....	6
Tables.....	7
Table S1: Logistic Regression Results for Estimating the Effects of Nursing Homes' Features on Odds of Being Included in the Study Sample .....	7
Table S2: Results of Two-Sample t-Test for Equality of the Means of the Excluded and Included Observations .....	8
Table S3: Replication of the main analysis results using Bivariate and Poisson Regression Models .....	9
Table S4: Confusion Matrix for SVM-RBF .....	10
Table S5: Confusion Matrix for NN .....	10
Table S6: Distribution of Infections Among Staff and Residents.....	10

## Missing Observations

Data cleaning process is presented in Figure S1. 493 nursing homes were excluded from the study sample either due to the mismatch between their names across multiple datasets or because their COVID-19 infection data were not available in CDPH reports. To examine if the excluded nursing homes are similar to those included in the study sample, we conducted two logistic regression with the dependent variables set to be 1 to indicate if a record is included in the study sample and 0 otherwise. In the first logistic regression we only include governance features as independent variables, while in the second logistic regression we include all the features.

As reported in Table S1, both regression results show that none of the governance features are statistically significant, which indicates that the included records have no selection bias on governance features. Amongst the remaining variables, quality rating and county infections per 100k are significant are statistically significant yet the difference between the two groups is not substantial, as reported in Table S2. Further, the differences in these two variables across the two groups make our estimates more conservative.

## Machine learning Techniques

We then apply machine learning techniques to predict the COVID-19 infection in nursing homes and compare the results with our model. In view that our problem has a highly nonlinear structure, advanced machine learning models that do not rely on data structure assumptions may provide a flexible and desired solution. We predict the nursing home level COVID-19 infection situation by using Neural Networks (NN) and Support Vector Machines (SVM) with RBF kernel function. Variable *NH* is used as the target variable in each model, and is equal to 1 if at least one patient or staff reported to be infected. The prediction features include nursing home governance features such as occupancy rate, number of certified beds, whether a family council presents, whether the nursing home is for profit or not, and inflation score evaluated from past years. The nursing homes' health inspection rating, staffing rating and quality rating are also included in our prediction model. To capture the severity of COVID-19 epidemic in the surrounding area, we also incorporate county level COVID-19 infections per 100K population.

## Bivariate and Double Poisson Estimates

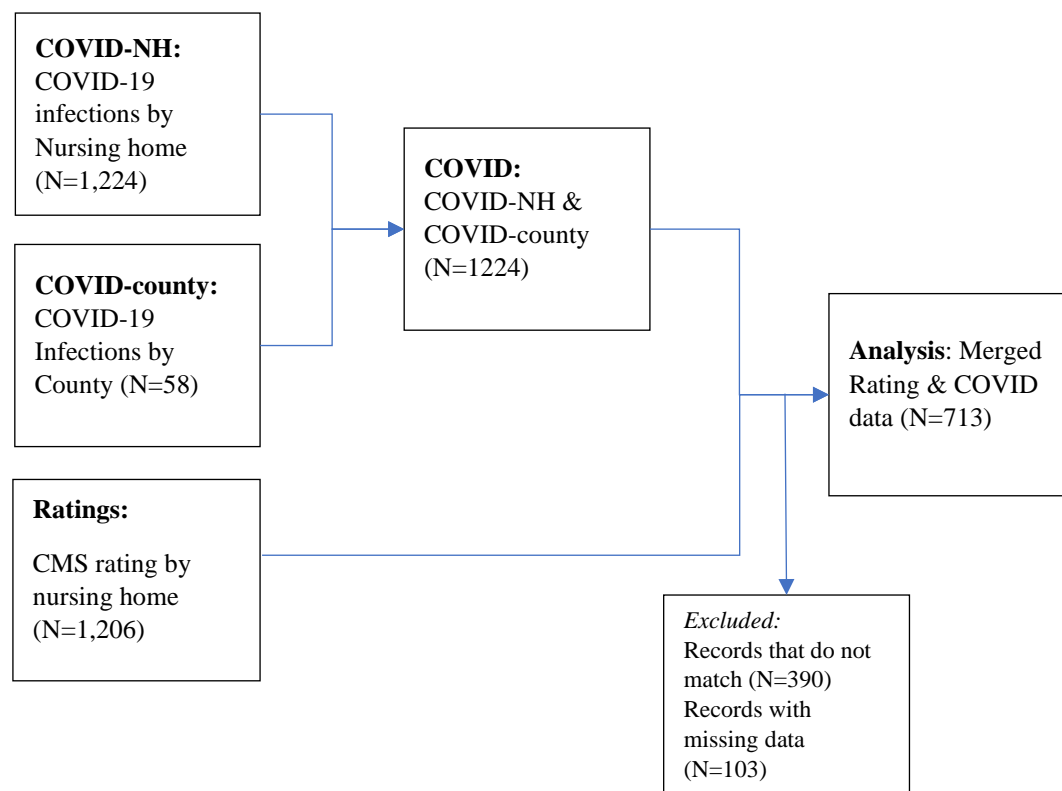
To test the robustness of our results and as a means of sensitivity analysis, we have replicated our main analysis using Bivariate and Double Poisson methods. The difference between these two methods and those reported in Table 2 of the main manuscript is these models do not assume an excess zero generating process and consider the outcome as a result of only two Poisson processes. In the Bivariate Poisson analysis, we assume that there is a correlation between the processes that give rise to the count of infections among staff and residents, while in the Double Poisson Regression, we assume independence between these two processes. The results are presented in Table S3. In comparison with the main results presented in the main table, the coefficients with larger sizes remain significant and close to their original estimates, while the smaller coefficients are not consistent with their original estimates. This is due to the fact that our dataset has significant excess zeros since most nursing homes had not reported infections many infections among either their staff or residents at the time of the study and therefore a zero inflated version of the Poisson models will be more appropriate for this setting.

## Correlation Between Infections Among Staff and Residents

To better examine the correlation between infections among staff and residents, we report the number and percentage of nursing homes with and without infections among their staff and residents in Table S6. We can observe that 91.75% of nursing homes with no infections among their residents also experienced no infections among their staff. Similarly, 54.21% of nursing homes that had at least one infection among their residents, also had at least one infection among their staff. In Figure S4, we show the scatter plot of number of infections among staff and residents for only those nursing homes that experienced a large outbreak among both their staff and residents. There is a clear correlation between the number of infections among staff and residents.

## Figures

Figure S1. Study population and analysis sample

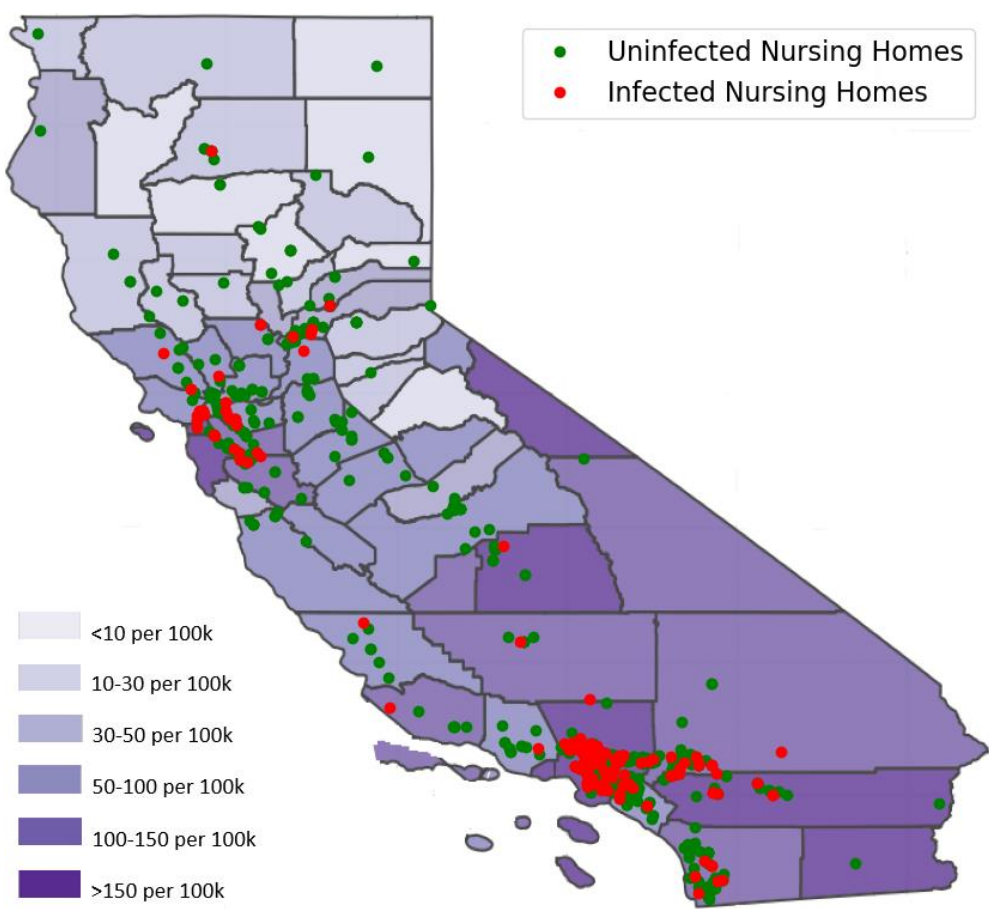


Note: Original CMS Rating for year 2017 data (*ratings*) include 1206 nursing homes. Original CA COVID-19 Infection by county (*COVID-county*) data as of April 30<sup>th</sup>, 2020 include on 58 counties. Original COVID-19 CA Infections by nursing homes (*COVID-NH*) data as of April 30<sup>th</sup>, 2020 include 1224 nursing homes.

We first merged *COVID-NH* and *COVID-county* data for all 1224 rows (0 record lost). We then merged the resulting data (*COVID*) with *ratings* data which resulted in 713 rows. 390 records were lost due to mismatch between the names of the facilities in the two datasets, and 103 records were lost for those nursing homes that did not report COVID 19 infection data or their ratings information is missing.

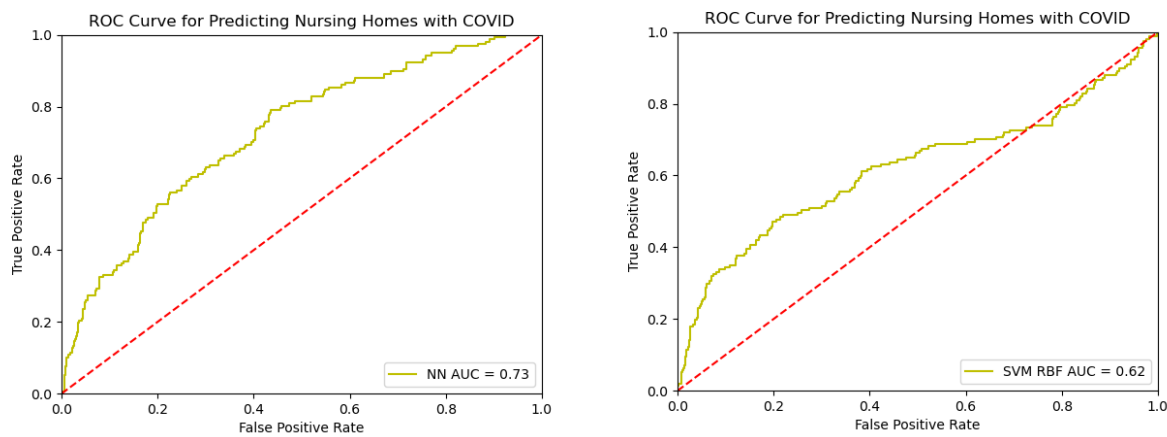


Figure S2: Spread of COVID-19 Infection Among California Nursing Homes



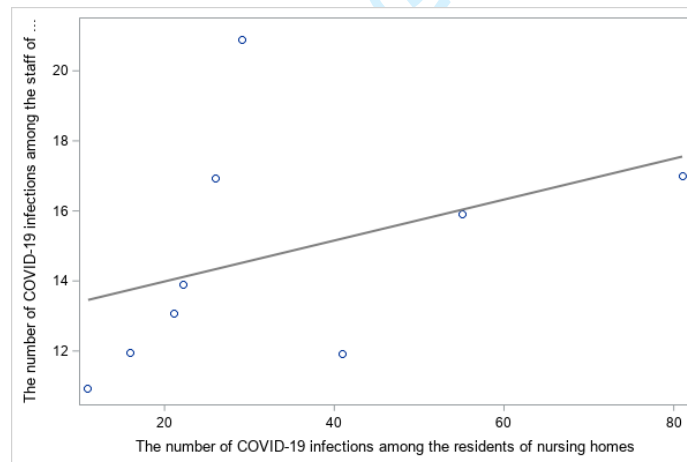
Note: The figure presents the spread of COVID-19 infection among California nursing homes as of May 1<sup>st</sup>, 2020

Figure S3: Receiver Operator Characteristic (ROC) Curves for Predicting at Least One Infection in Nursing Homes



Note: ROC for Nursing Home (NH) COVID-19 prediction using Neural Networks (NN), SVM with RBF kernel. The AUC is reported for each model: NN=0.73, SVM-RBF (default)=0.62

Figure S4: Scatter plot of number of infections among staff and residents for those nursing homes that have experienced large outbreaks amongst both their staff and resident populations



## Tables

Table S1: Logistic Regression Results for Estimating the Effects of Nursing Homes' Features on Odds of Being Included in the Study Sample

Parameter	Validation with Governance Features Only (Included vs. Excluded Records)			Validation with All Features (Included vs. Excluded Records)		
	Estimate	(95% CI)	P Value	Estimate	(95% CI)	P Value
Constant	0.1	(-0.72 to 0.92)	0.81	-0.66	(-2.09 to 0.76)	0.36
For profit	0.25	(-0.08 to 0.58)	0.14	0.29	(-0.1 to 0.68)	0.14
Family council	-0.19	(-0.49 to 0.12)	0.23	-0.07	(-0.4 to 0.26)	0.68
Certified beds	-0.0004	(-0.003 to 0.002)	0.71	-0.0008	(-0.003 to 0.002)	0.52
Occupancy rate	0.61	(-0.3 to 1.52)	0.19	0.56	(-0.62 to 1.74)	0.35
Inflation score	-0.04	(-0.2 to 0.12)	0.6	-0.03	(-0.2 to 0.14)	0.75
Quality rating				0.21	(0.07 to 0.36)	0.004
Staffing rating				0.002	(-0.14 to 0.14)	0.97
Health inspection rating				0.08	(-0.04 to 0.19)	0.21
County infections per 100K				-0.002	(-0.004 to -0.0007)	0.004

Note: Coefficients represent how the log odds of the dependent variable changes with one unit increase in the corresponding predictor

Table S2: Results of Two-Sample t-Test for Equality of the Means of the Excluded and Included Observations

Features	Excluded Records*	Included Records*	P Value**
For profit	0.82	0.86	0.11
Family council	0.21	0.18	0.21
Certified beds	99.6	98.0	0.65
Occupancy rate	0.85	0.86	0.14
Inflation score	0.32	0.31	0.83
Quality rating	4.43	4.57	0.01
Staffing rating	3.49	3.49	0.93
Health inspection rating	2.66	2.86	0.01
County infections per 100K	159.36	143.88	0.003

Note: \*: Reports the average value of features.

\*\* :P values are for two-tailed t-tests of the equality of the two means.

Table S3: Replication of the main analysis results using Bivariate and Poisson Regression Models

Parameter	Bivariate Poisson Model			Double Poisson Model		
	Estimate	(95% CI)	P Value	Estimate	(95% CI)	P Value
<b>Infected Staff (number of staff with confirmed COVID-19 infections)</b>						
Intercept	-3.9	(-5.97 to -1.83)	0.01	-3.29	(-4.7 to -1.88)	<.001
County infections per 100K	0.01	(0.01 to 0.01)	<.001	0.01	(0.01 to 0.01)	<.001
For profit	0.33	(-0.28 to 0.93)	0.3	0.01	(-0.37 to 0.39)	0.97
Family council	-0.08	(-0.59 to 0.43)	0.77	0.18	(-0.1 to 0.46)	0.21
Certified beds	0.01	(0.01 to 0.01)	<.001	0.01	(0.01 to 0.01)	<.001
Occupancy rate	-2.5	(-4.05 to -0.95)	0.01	-0.89	(-2.02 to 0.24)	0.13
Inspection rating	0.1	(-0.1 to 0.28)	0.35	-0.12	(-0.23 to -0.01)	0.05
Quality rating	0.25	(-0.05 to 0.54)	0.11	0.21	(0.03 to 0.39)	0.03
Staffing rating	0.12	(-0.06 to 0.29)	0.19	0.26	(0.14 to 0.38)	<.001
Inflation score	0.49	(0.39 to 0.59)	<.001	0.31	(0.23 to 0.39)	<.001
<b>Infected Residents (number of residents with confirmed COVID-19 infections)</b>						
Intercept	-2.1	(-3.01 to -1.19)	<.001	-1.46	(-2.2 to -0.71)	0.01
County infections per 100K	0.01	(0.01 to 0.01)	<.001	0.01	(0.01 to 0.01)	<.001
For profit	2.71	(2.12 to 3.31)	<.001	1.89	(1.5 to 2.28)	<.001
Family council	0.16	(0.02 to 0.3)	0.03	0.19	(0.06 to 0.31)	0.01
Certified beds	0.01	(0.01 to 0.01)	<.001	0.01	(0.01 to 0.01)	<.001
Occupancy rate	-0.08	(-0.66 to 0.51)	0.82	0.02	(-0.54 to 0.57)	0.96
Inspection rating	-0.2	(-0.25 to -0.14)	<.001	-0.21	(-0.26 to -0.16)	<.001
Quality rating	0.05	(-0.03 to 0.13)	0.2	0.08	(-0.01 to 0.15)	0.06
Staffing rating	-0.22	(-0.27 to -0.17)	<.001	-0.15	(-0.2 to -0.11)	<.001
Inflation score	0.13	(0.08 to 0.18)	<.001	0.13	(0.08 to 0.17)	<.001
Covariance	0.21	(0.18 to 0.25)	<.001			
<b>Fit Statistics</b>						
-2 log likelihood		8011.7			8468.6	
AIC		8053.7			8508.6	
BIC		8149.7			8600.0	

Table S4: Confusion Matrix for SVM-RBF

		ACTUAL CLASS	
		0	1
PREDICTED CLASS	0	142	2
	1	47	7

Table S5: Confusion Matrix for NN

		ACTUAL CLASS	
		0	1
PREDICTED CLASS	0	137	7
	1	37	17

Table S6: Distribution of Infections Among Staff and Residents

		INFECTIONS AMONG STAFF (%)	
		0	>=1
INFECTIONS AMONG RESIDENTS (%)	0	556 (91.75%)	50 (8.25%)
	>=1	49 (45.79%)	58 (54.21%)

## STROBE Statement

Checklist of items that should be included in reports of observational studies

Section/Topic	Item No	Recommendation	Reported on Page No
<b>Title and abstract</b>	1	(a) Indicate the study's design with a commonly used term in the title or the abstract	1,2
		(b) Provide in the abstract an informative and balanced summary of what was done and what was found	2
<b>Introduction</b>			
Background/rationale	2	Explain the scientific background and rationale for the investigation being reported	5
Objectives	3	State specific objectives, including any prespecified hypotheses	5
<b>Methods</b>			
Study design	4	Present key elements of study design early in the paper	6
Setting	5	Describe the setting, locations, and relevant dates, including periods of recruitment, exposure, follow-up, and data collection	6
Participants	6	(a) <i>Cohort study</i> —Give the eligibility criteria, and the sources and methods of selection of participants. Describe methods of follow-up	6
		<i>Case-control study</i> —Give the eligibility criteria, and the sources and methods of case ascertainment and control selection. Give the rationale for the choice of cases and controls	
		<i>Cross-sectional study</i> —Give the eligibility criteria, and the sources and methods of selection of participants	
Variables	7	(b) <i>Cohort study</i> —For matched studies, give matching criteria and number of exposed and unexposed	6
		<i>Case-control study</i> —For matched studies, give matching criteria and the number of controls per case	
Data sources/measurement	8*	Clearly define all outcomes, exposures, predictors, potential confounders, and effect modifiers. Give diagnostic criteria, if applicable	6
Bias	9	For each variable of interest, give sources of data and details of methods of assessment (measurement). Describe comparability of assessment methods if there is more than one group	7 & Appendix
Study size	10	Describe any efforts to address potential sources of bias	6 & Appendix
Quantitative variables	11	Explain how the study size was arrived at	7
		Explain how quantitative variables were handled in the analyses. If applicable, describe which groupings were chosen and why	7
		(a) Describe all statistical methods, including those used to control for confounding	7
		(b) Describe any methods used to examine subgroups and interactions	7
		(c) Explain how missing data were addressed	7
Statistical methods	12	(d) <i>Cohort study</i> —If applicable, explain how loss to follow-up was addressed	7
		<i>Case-control study</i> —If applicable, explain how matching of cases and controls was addressed	7

*Cross-sectional study*—If applicable, describe analytical methods taking account of sampling strategy

(e) Describe any sensitivity analyses

7

Section/Topic	Item No	Recommendation	Reported on Page No
<b>Results</b>			
Participants	13*	(a) Report numbers of individuals at each stage of study—eg numbers potentially eligible, examined for eligibility, confirmed eligible, included in the study, completing follow-up, and analysed	8
		(b) Give reasons for non-participation at each stage	Appendix
		(c) Consider use of a flow diagram	Appendix
Descriptive data	14*	(a) Give characteristics of study participants (eg demographic, clinical, social) and information on exposures and potential confounders	8
		(b) Indicate number of participants with missing data for each variable of interest	8
		(c) <i>Cohort study</i> —Summarise follow-up time (eg, average and total amount)	
Outcome data	15*	<i>Cohort study</i> —Report numbers of outcome events or summary measures over time	
		<i>Case-control study</i> —Report numbers in each exposure category, or summary measures of exposure	
		<i>Cross-sectional study</i> —Report numbers of outcome events or summary measures	8,9
Main results	16	(a) Give unadjusted estimates and, if applicable, confounder-adjusted estimates and their precision (eg, 95% confidence interval). Make clear which confounders were adjusted for and why they were included	9,10
		(b) Report category boundaries when continuous variables were categorized	
		(c) If relevant, consider translating estimates of relative risk into absolute risk for a meaningful time period	
Other analyses	17	Report other analyses done—eg analyses of subgroups and interactions, and sensitivity analyses	10
<b>Discussion</b>			
Key results	18	Summarise key results with reference to study objectives	12
Limitations	19	Discuss limitations of the study, taking into account sources of potential bias or imprecision. Discuss both direction and magnitude of any potential bias	13
Interpretation	20	Give a cautious overall interpretation of results considering objectives, limitations, multiplicity of analyses, results from similar studies, and other relevant evidence	12
Generalisability	21	Discuss the generalisability (external validity) of the study results	13
<b>Other Information</b>			
Funding	22	Give the source of funding and the role of the funders for the present study and, if applicable, for the original study on which the present article is based	13

\*Give information separately for cases and controls in case-control studies and, if applicable, for exposed and unexposed groups in cohort and cross-sectional studies.



**Note:** An Explanation and Elaboration article discusses each checklist item and gives methodological background and published examples of transparent reporting. The STROBE checklist is best used in conjunction with this article (freely available on the Web sites of PLoS Medicine at <http://www.plosmedicine.org/>, Annals of Internal Medicine at <http://www.annals.org/>, and Epidemiology at <http://www.epidem.com/>). Information on the STROBE Initiative is available at [www.strobe-statement.org](http://www.strobe-statement.org).

3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

For peer review only

# BMJ Open

## Compress the Curve: A Cross Sectional Study of Variations in COVID-19 Infections Across California Nursing Homes

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2020-042804.R2
Article Type:	Original research
Date Submitted by the Author:	05-Dec-2020
Complete List of Authors:	Gopal, Ram; University of Warwick Han, Xu; Fordham University Yaraghi, Niam; University of Miami, Miami Herbert Business School; Brookings Institution, Governance Studies
<b>Primary Subject Heading</b>:	Geriatric medicine
Secondary Subject Heading:	Infectious diseases
Keywords:	COVID-19, GERIATRIC MEDICINE, Health policy < HEALTH SERVICES ADMINISTRATION & MANAGEMENT

SCHOLARONE™  
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in this journal and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

## Original Investigation

**Title:** Compress the Curve: A Cross Sectional Study of Variations in COVID-19  
Infections Across California Nursing Homes

Ram D. Gopal, PhD<sup>1</sup>, Xu Han, PhD<sup>2</sup>, Niam Yaraghi, PhD<sup>3,4,\*</sup>

<sup>1</sup>: Professor, Warwick Business School, University of Warwick

<sup>2</sup>: Assistant Professor, Gabelli School of Business, Fordham University

<sup>3</sup>: Assistant Professor, Miami Herbert Business School, University of Miami

<sup>4</sup>: Non-resident Fellow, Governance Studies, The Brookings Institution

*\*: All authors contributed equally. Authors are listed in alphabetical order of their last name.*

Corresponding author:

Niam Yaraghi

[niamyaraghi@miami.edu](mailto:niamyaraghi@miami.edu)

5250 University Drive, Coral Gables, FL 33146

Phone: (305) 284-3314

Word count: 2668

## Abstract

*Objective:* Nursing homes' residents and staff constitute the largest proportion of the fatalities associated with COVID-19 epidemic. Although there is a significant variation in COVID-19 outbreaks among the US nursing homes, we still do not know why such outbreaks are larger and more likely in some nursing homes than others. This research aims to understand why some nursing homes are more susceptible to larger COVID-19 outbreaks.

*Design:* Observational study of all nursing homes in the state of California until May 1st, 2020.

*Setting:* The state of California.

*Participants:* 713 long term care facilities in the State of California that participate in public reporting of COVID-19 infections as of May 1<sup>st</sup>, 2020 and their infections data could be matched with data on ratings and governance features of nursing homes provided by CMS.

*Main Outcome Measure:* The number of reported COVID-19 infections among staff and residents.

*Results:* Study sample included 713 nursing homes. The size of outbreaks among residents in for-profit nursing homes is 12.7 times larger than their non-profit counterparts (log count = 2.54; 95% CI, 1.97 to 3.11; P<.001). Higher ratings in CMS-reported health inspections are associated with lower number of infections among both staff (log count = -0.19; 95% CI, -0.37 to -0.01; P = 0.05) and residents (log count = -0.20; 95% CI, -0.27 to -0.14; P<.001). Nursing homes with higher discrepancy between their CMS- and self-reported ratings have higher number of infections among their staff (log count = 0.41; 95% CI, 0.31 to 0.51; P<.001) and residents (log count = 0.13; 95% CI, 0.08 to 0.18; P<.001).

1  
2  
3 *Conclusions:* The size of COVID-19 outbreaks in nursing homes is associated with  
4 their ratings and governance features. To prepare for the possible next waves of  
5 COVID-19 epidemic, policy makers should use these insights to identify the nursing  
6 homes who are more likely to experience large outbreaks.  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16

17 **Key words:** *COVID-19, Nursing Homes, Long-Term Care*  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Article Summary

### Strengths and limitations of this study

- A bivariate Poisson model is employed to better capture the interdependencies of COVID-19 cases between staff and residents.
- Predictive models are developed to identify nursing homes with the highest chance of experiencing COVID-19 outbreaks.
- Data analyzed are only from California.
- The dataset on nursing homes' features is based on the year 2017.
- The number of COVID-19 cases reported by nursing homes may be subject to under-reporting.

## Introduction

Nursing homes have been most severely impacted by the COVID-19 pandemic owing to the advanced age and high number of comorbidities of their residents.<sup>1,2</sup> In Europe, as much as 57% of all deaths related to COVID-19 were at such facilities.<sup>3</sup> In the United States, nursing homes' residents and staff account for 34% of all COVID-19 fatalities.<sup>4</sup> Infection prevention and control at nursing homes and long-term facilities has therefore become a priority in managing the epidemic.<sup>5,6</sup>

Given the considerable variation in the prevalence and size of the COVID-19 outbreaks at nursing homes, the objective of this research is (1) to understand why some nursing homes are more susceptible to COVID-19 outbreaks, and (2) to develop predictive models that can identify such nursing homes so that they could be prioritized in efforts to prevent and contain next waves of the epidemic.<sup>7,8</sup>

## Methods

### *Patient and public involvement*

Patients had no influence on the research questions or outcomes of this research. No patients were involved in the design of this study. We used blind patient files; therefore, no patient recruitment took place. We only used data on the aggregated number of COVID-19 patients and staff in the nursing homes as reported by the State of California and therefore no personal information of patients was used in this study. Given the nature of removing all personal information, there is no requirement to disseminate the information to patients.

### *Data Sources and Study Variables*

We collected data from various publicly available sources. The New York Times aggregates and provides data on COVID-19 cases per county.<sup>9</sup> California Department of Public Health (CDPH) provides data on the number of confirmed COVID-19



1  
2  
3 infections among staff and residents of nursing homes in the state.<sup>10</sup> CMS provides data  
4 on nursing home characteristics, including their self-reported ratings and CMS health  
5 inspections.<sup>11</sup> A description of this data is provided in the next section. Applying the  
6 methods suggested by Han et. al,<sup>12</sup> we identified the nursing homes with significant  
7 discrepancies between their self-reported measures and independent CMS inspections  
8 for a consecutive 5-year period. We aggregated the results and used the number of  
9 years a nursing home is predicted to be a likely inflator as the overall inflation score for  
10 a nursing home. Therefore, an honest nursing home will have an inflation score of 0  
11 while an inflating nursing home can have an inflation score between 1 to 5, with 5 being  
12 the most severe. In our dataset, 19.25% of nursing homes were inflating their scores  
13 and some of these had a score of 5 indicating that they inflated their scores in all 5  
14 years.

15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31 These methods rely on data that are only available for nursing homes in California and  
32 therefore, the scope of this study is also limited to nursing homes in California. After  
33 cleaning and merging the above-mentioned data sources, we analysed a final dataset  
34 consisting of 713 nursing homes in California. Details of the data cleaning and merging  
35 process is presented in Supplementary Appendix.

36  
37  
38  
39  
40  
41  
42 We examined the following outcomes in this study: whether a nursing home has at least  
43 one COVID-19 infection amongst its residents or staff, the number of confirmed  
44 COVID-19 infections among its residents, and the number of confirmed infections  
45 among its staff. We also calculated a fourth outcome that indicates the large outbreaks  
46 as the ones in which more than 10 members of staff or residents were infected with  
47 COVID-19. This threshold translates to approximately 95<sup>th</sup> percentile of the number of  
48 infected staff. Given that more residents are infected than staff, this threshold translates  
49 to 75<sup>th</sup> percentile of the number of residents.

1  
2  
3 The independent variables describe the severity of the COVID-19 outbreak in the  
4 surrounding area of a nursing home, its governance characteristics, as well as its ratings  
5 on quality, staffing and CMS inspections. Table 1 provides detailed description of the  
6 study variables. Note that while almost all nursing homes have resident councils,  
7 only 20 percent of nursing homes have existing family councils. We included  
8 the existence of family council as a binary variable in our analysis with the  
9 contention that it may imply closer coordination and higher engagement with  
10 the families of the residents.  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

### 23 *Description of CMS' Nursing Home Compare System*

24 The CMS nursing home rating data consists of basic information about nursing facilities  
25 such as name, address, phone number, etc., as well as some key features used in our  
26 analysis, such the number of certified beds, whether the nursing home is for-profit or  
27 non-profit, whether the nursing home has a family council, etc.  
28  
29  
30  
31  
32  
33

34 The CMS nursing home rating data serves the CMS Nursing Home Compare System,  
35 in which nursing home ratings are generated based on three domains: Inspection,  
36 Staffing, and Quality measures. The Inspection is conducted and reported by CMS-  
37 certified inspectors annually. The other two domains are self-reported by nursing  
38 homes. The annual inspection investigates areas such as medication management,  
39 nursing home administration, environment, food service, and residents' rights and  
40 quality of life. The Staffing domain is evaluated based on the self-reported CMS  
41 Certification and Survey Provider Enhanced Reports (CASPER) staffing data. The two  
42 measures used are the total nursing hours and Registered Nursing (RN) hours and are  
43 adjusted for case-mix based on the Resource Utility Group (RUG-III) case-mix system  
44 derived from the Minimum Data Set (MDS). The staffing star rating is then updated by  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 the end of the quarter when raw data is collected. Note that with more recent changes,  
4 the Staffing data reported by nursing homes is subject to validation with nursing homes'  
5 payroll data reported through Payroll-Based Journal (PBJ). The Quality Measure rating  
6 uses quality measurement criteria, which covers both long-stay terms and short-stay  
7 terms. The quality measure star rating is updated by the end of each quarter by using  
8 the results from three most recent quarters.  
9

10 To calculate the star ratings, CMS first assigns an initial star rating to all nursing homes  
11 based on their annual inspection results. Nursing homes are then assigned star ratings  
12 for the Staffing and Quality Measures domains. The overall star rating is then  
13 calculated by considering the inspection rating as the baseline, increasing or decreasing  
14 by 1 star if any self-reported domain satisfies the conditions stated as follows. Both 4  
15 and 5 stars in staffing rating are qualified for obtaining additional overall star rating,  
16 while only 5 stars in quality measure is qualified. Additional conditions apply to nursing  
17 homes whose inspection ratings are only 1 star, and for nursing homes which are in the  
18 CMS's Special Focus Facility (SFF) program. The overall star rating is lowered by one  
19 star if any self-reported domain is 1 star. The overall star rating cannot be more than 5  
20 stars or less than 1 star. Detailed data from CMS on nursing homes is available online.<sup>13</sup>  
21

### 22 *Statistical Analysis*

23 To answer the first research question and understand why some nursing homes are more  
24 susceptible to COVID-19 outbreaks, we applied Zero Inflated Bivariate Poisson (ZIBP)  
25 regression. The model allows us to examine the effects of nursing homes' ratings,  
26 governance features, and their surroundings on the likelihood and size of their COVID-  
27 19 outbreaks. Econometric details of the model are provided by Walhin, 2001.<sup>14</sup>  
28 Conventional Poisson models are suitable for modelling count data, while the zero  
29 inflated variation of Poisson model is more suitable for modelling count data with  
30

1  
2  
3 excess zeros, especially when excess zeros are generated by a separate processes that  
4 could be modelled separately. This leads to a framework that consists of a logit model  
5 for estimating the excess zeros in addition to a Poisson count model. ZIBP model is an  
6 extension of zero inflated Poisson model and is best suited for situations in which the  
7 count data with excess zeros are generated for two outcomes that may be correlated. In  
8 cases were the outcome variables are independent, the model reduces to the product of  
9 two independent zero inflated Poisson regression models, referred to as Zero Inflated  
10 Double Poisson model. in our setting, the two count variables are the number of  
11 COVID-19 infections among staff, and residents. These counts include excess zeros  
12 since many nursing homes reported no COVID-19 cases, primarily because they are  
13 located in areas where at the time of the data collection, had not yet experienced  
14 significant surges in COVID-19 cases. These two counts are also correlated since they  
15 both happen at the same nursing home and the factors that give rise to them are common  
16 at the nursing home level.

17  
18 Intuitively, we assume that the number of zero's in the count of infected staff and  
19 residents are generated either because the nursing home was in an area that was less  
20 infected by the COVID-19 or because it implemented successful prevention procedures  
21 to protect its staff and residents. Moreover, we assume that in a nursing home, the  
22 number of infected staff covaries with the number of infected residents since they can  
23 infect each other and since common infection prevention and control policies apply to  
24 both groups. Taking this interdependency into account also alleviates the concerns over  
25 the possible impact of omitted variables in our model. In this context, because of the  
26 close proximity of residents and staff, the same variables that could affect the number  
27 of infections among one group, would most likely also impact the number of infections  
28 among the other group. The covariance coefficient captures this interdependency in  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 outcomes. As a sensitivity analysis, we also report the results of zero-inflated double  
4  
5 Poisson regression. In this model, the counts of infections among staff and residents are  
6  
7 assumed to be independent from each other. We use NLMIXED procedure in SAS  
8  
9 software to estimate our models.<sup>15,16</sup> Note that we have provided access to both the data  
10  
11 and the SAS code for this analysis.<sup>17,18</sup>  
12  
13

14 To answer the second research question and identify the nursing homes with the  
15  
16 highest risk of COVID-19 outbreaks, we used our models to predict the probability of  
17  
18 experiencing an infection and compared their performance with common machine  
19  
20 learning techniques, namely Neural Networks (NN) and Support Vector Machine with  
21  
22 Radial Basis Function kernel (SVM-RBF). Since our problem has a highly nonlinear  
23  
24 structure, advanced machine learning models such as NN and SVM that do not rely  
25  
26 on data structure assumptions may provide a flexible and desired solution. Variable  
27  
28 NH is used as the target variable in each model, and NH is equal to 1 if at least one  
29  
30 patient or staff reported to be infected. The prediction features include nursing home  
31  
32 governance features such as occupancy rate, number of certified beds, whether a  
33  
34 family council presents, whether the nursing home is for profit or not, and inflation  
35  
36 score evaluated from past years. The nursing homes' health inspection rating, staffing  
37  
38 rating and quality rating are also included. The machine learning models are  
39  
40 implemented in Python 3.7 with 70% data training and 30% data testing. The entire  
41  
42 dataset is used to plot the lift chart. We also measured the performance of our models  
43  
44 in predicting the nursing homes with highest risks of experiencing large outbreaks  
45  
46 with more than 10 infections.  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Results

### *Study Sample*

During the data cleaning and merging process, 493 nursing homes were eliminated from our final sample, either because their names were not matching across different datasets, or their ratings information is not available from CMS, or because their COVID-19 infections are not reported by CDPH. To ensure that the final sample is random and our results are not biased, we compared the eliminated nursing homes with the ones in the study sample. The results of two sample t-tests and logistic regression are presented in Supplementary Appendix. None of the observed governance factors affect the chance of being included in the sample. Amongst the remaining variables, while the difference with regards to quality ratings and county infections per 100K is statistically significant between the two groups, their magnitude is small and serve to make our estimates more conservative.

Study sample included 713 nursing homes in California. As reported in Table 1, as of May 1<sup>st</sup>, 2020, 23% of the study sample reported at least one COVID-19 infection among either their staff or residents. Of those, 31% experienced large outbreaks with more than 10 infections among either their staff or residents. The geographic spread of COVID-19 infections in California nursing homes is graphically presented in the Supplementary Appendix.

### *Preventing COVID-19 Infections*

According to the model selection criteria reported in Table 2, the ZIBP model provides a better fit as its AIC, BIC and -2Log Likelihood are all smaller than those of Zero Inflated Double Poisson model. We therefore report the estimates of the ZIBP model in the text. The coefficients in the first panel of Table 2 represent how the log odds of experiencing an infection changes with one unit of increase in the corresponding

1  
2  
3 predictor. As reported in the first panel of Table 2, the only variables with statistically  
4 significant impact on the chance of COVID-19 outbreaks at nursing homes are their  
5 size and the rate of infections per 100 thousand residents at the county in which they  
6 are located. For both variables, a one-unit of increase is associated with a 1% increase  
7 in the odds of experiencing at least one COVID-19 infection.  
8  
9  
10  
11  
12  
13

#### 14 *Controlling COVID-19 Outbreaks*

15  
16 The coefficients in the second and third panel of Table 2 represent how the expected  
17 log count of the infections changes for each unit increase in the corresponding predictor.  
18 As reported in the second and third panel of Table 2, the expected rate of infections  
19 amongst both staff and residents increase with the size of the nursing home. This  
20 indicates that although the severity of COVID-19 epidemic in the surrounding area  
21 increases the chance of experiencing at least one infection at the nursing homes.  
22  
23  
24  
25  
26  
27  
28  
29

30 While the size of outbreaks among residents is about 12.7 times higher in for-profit  
31 nursing homes, the size of outbreak among staff in for-profit nursing homes is not  
32 statistically different from non-profit ones. This is in line with prior empirical research  
33 that has repeatedly shown that for-profit nursing homes are inferior in many aspects of  
34 care quality.<sup>19–22</sup>  
35  
36  
37  
38  
39  
40  
41

42 Occupancy rate, which represents the ratio of the number of patients to the number of  
43 certified beds of a nursing home, is associated with a lower rate of infections among  
44 staff such that a one percent increase in occupancy rate decreases the expected count of  
45 infections among staff by 2.4%.  
46  
47  
48  
49  
50

51 Among the three different ratings, the CMS-reported health inspection rating is  
52 associated with a sizable decrease in the number of infections among both staff and  
53 residents. One unit of increase in CMS-reported health inspection ratings is associated  
54 with a 17% and 18% decrease in the expected number of infections in staff and  
55  
56  
57  
58  
59  
60

1  
2  
3 residents, respectively. A one-unit improvement in staffing rating is associated with a  
4  
5 23% decrease in the number of infections among residents. Note that better staff rating  
6  
7 is highly dependent on higher ratio of staff to residents and the higher number of staff  
8  
9 per resident would allow nursing homes to control infections more efficiency among  
10  
11 their residents. While the observed association between ratings on health inspections  
12  
13 and staffing with the number of infected staff and residents were expected, the  
14  
15 association between self-reported quality ratings and the number of infections is the  
16  
17 opposite of our expectations. One unit of increase in self-reported quality ratings is  
18  
19 associated with, respectively, 49% and 14% increase in infections among staff and  
20  
21 residents. This finding is aligned with the emerging stream of research that shows  
22  
23 nursing homes embellish their self-reported quality ratings and therefore these ratings  
24  
25 may not always indicate better quality of care for residents.<sup>12,23-26</sup> Our final variable,  
26  
27 inflation score, quantifies the discrepancy between the self- and CMS-reported ratings.  
28  
29 The higher the discrepancy, the more likely it is that the nursing home is overstating  
30  
31 their quality measures. With a one-unit increase in such discrepancy, the expected  
32  
33 number of infections among staff and residents increases by, 51% and 14%,  
34  
35 respectively.

#### 41 42 *Improving the Quality Reporting System*

43  
44 CMS could solve these discrepancies and improve the reporting process by  
45  
46 implementing better inspection and auditing strategies.<sup>27</sup> Figure 1 shows how the  
47  
48 number of infections among staff and residents could be compressed had the self-  
49  
50 reported quality measures by nursing homes were truly reflecting their quality of care.  
51  
52 Given the importance of ratings for nursing homes,<sup>28</sup> with a reliable rating system with  
53  
54 no discrepancy between self- and CMS-reported measures, nursing homes would strive  
55  
56 to elevate their ratings through actual improvements in their quality of care. As shown  
57  
58  
59  
60



1  
2  
3 in the upper panel of Figure 1, compared to the current system, lower number of  
4 predicted infections among staff would have been more frequent under an improved  
5 rating system such that predicted average number of infections among staff would have  
6 decreased from 1.85 to 1.52, which is equal to 17.6% fewer total infections across the  
7 staff of all nursing homes. As shown in the lower panel of Figure 1, the same effect is  
8 observed for nursing home residents. Had self-reported quality ratings were truly  
9 reflecting the quality of care, the expected number of infections among residents of  
10 nursing homes would have reduced from 8.67 to 8.15 which is equal to 5.8% fewer  
11 total infections across the residents of all nursing homes.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

23 Finally, the sizable covariance estimate (0.68; 95% CI 0.54 to 0.87; P=0.1) indicates  
24 that the number of infected staff is not independent from the number of infected  
25 residents. This observation empirically confirms our expectation of dependency  
26 between the count of infections in staff and residents such that nursing homes with high  
27 number of infected staff also have high number of infected residents. This finding was  
28 expected as residents and staff are in close contact with each other and once infections  
29 occur among the members of one group, it would be very difficult to prevent them in  
30 the other group. More importantly, common infection control procedures implemented  
31 by nursing homes would apply to both groups and prevent infections among both  
32 groups. Note that as discussed earlier, according to all the model selection criteria, the  
33 ZIBP performs better than its competitors. This is not surprising since it has the  
34 advantage of modelling and adjusting for the correlation between the count of infections  
35 among staff and residents. In the Appendix, we provide further empirical details on the  
36 correlation between the number of infections among residents and staff.  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54

55  
56 *Identifying Nursing Homes with Highest Chance of COVID-19 Infections & Outbreaks*  
57  
58  
59  
60

1  
2  
3 Figure 2 compares the lift of the ZIBP model with those of NN and SVM-RBF. We use  
4 lift as a measure for the ability of the model at predicting or classifying cases with  
5 respect to random selection. Lift shows how much better our model works compared to  
6 a random selection model. The first 50 nursing homes are zoomed in at the top right  
7 corner of the figure. The ZIBP model's performance is comparable with the common  
8 NN and SVM-RBF methods. For the first 50 nursing homes, the rate of true positives  
9 of ZIBP model is between 2.45 and 2.73 times higher than that of a random selection  
10 model. The Area Under the Curve (AUC) for ZIBP, NN and SVM-RBF models are  
11 respectively 0.68, 0.73, and 0.62.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

24 Figure 3 presents the lifts of the ZIBP model in identifying the nursing homes with  
25 large COVID-19 outbreaks among those that have confirmed at least ten infections. For  
26 the first 50 nursing homes, ZIBP correctly identifies nursing homes with large  
27 outbreaks among staff between 1.3 to 3.9 times better than a random selection model.  
28 The model's performance for predicting large outbreaks among residents for the first  
29 50 nursing homes is 1.5 to 2.1 times better than a random selection model.  
30  
31  
32  
33  
34  
35  
36  
37

### 38 Discussion

39 Staff and residents of nursing homes constitute the largest demographic of COVID-19  
40 fatalities in the US. However, nursing homes have not been uniformly impacted by the  
41 epidemic; some have not experienced even a single infection while some others have  
42 been devastated by COVID-19 fatalities. To prepare for the possible next waves of the  
43 epidemic, it is critical to uncover the underlying reason of such variation and to explore  
44 the nursing homes' features that are associated with higher chance and size of  
45 outbreaks.  
46  
47  
48  
49  
50  
51  
52  
53  
54

55 The aim of this research was to understand how publicly available data on nursing  
56 homes can explain the significant variation in the chance and size of COVID-19  
57  
58  
59  
60

1  
2  
3 infections at nursing homes, and to also develop predictive models that can identify the  
4 nursing homes with the highest chance and size of outbreaks.  
5

6  
7  
8 Our results indicate that COVID-19 outbreaks are more likely to happen at larger  
9 nursing homes and those with higher rate of COVID-19 infections in the surrounding  
10 area. These factors have been shown to be associated with higher probability of  
11 experiencing infections by other researchers as well.<sup>29</sup>  
12

13  
14  
15 Those with better staffing and health inspection ratings are more successful in  
16 controlling the outbreaks. The association between staffing levels and likelihood of  
17 having COVID-19 infections among both staff and residents has been reported by other  
18 researchers as well.<sup>30</sup> Interestingly, higher self-reported quality ratings are associated  
19 with larger size of outbreaks. This counter-intuitive result could be further evidence  
20 that nursing homes exaggerate their self-reported quality measures. Higher discrepancy  
21 between self-reported measures and CMS-reported health inspections was associated  
22 with larger COVID-19 outbreaks.  
23

24  
25  
26 The size of the outbreaks among residents is significantly higher in for-profit nursing  
27 homes which have been previously shown to also be of poorer quality in various aspects  
28 of care.<sup>19-22</sup>  
29

30  
31  
32 There is a complex relationship between the main variables in our models. For-profit  
33 NHs generally have lower nurse staffing, more deficiencies, are larger in size, and have  
34 a greater likelihood of inflating their ratings.<sup>31,32</sup> It is therefore not surprising that they  
35 were found to be more likely to have larger numbers of COVID infected residents and  
36 staff.  
37

38  
39  
40 The model developed in this research can correctly identify the nursing homes that are  
41 more likely to experience an infection or are at the highest risk of an outbreak.  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 The insights of this research help policy makers to identify the nursing homes with the  
4 highest probability and size of COVID-19 outbreaks. This will allow them to prioritize  
5 such nursing homes in their efforts to control the epidemic. Such efforts could entail  
6 devoting more resources towards nursing homes with significantly higher risk or when  
7 feasible, temporarily transferring patients to different nursing homes to control the  
8 spread of the virus.  
9

10  
11 Our results show that our ZIBP model outperforms SVM and that the predictive ability  
12 of the NN is only modestly better than ZIBP model. That is, the application and  
13 comparison of these machine learning models with the results of the ZIBP model  
14 confirms that not only the ZIBP model can explain the relationship between various  
15 independent variables and COVID-19 infections at nursing homes, but it also offers  
16 competitive predictive performance.  
17

18  
19 An important takeaway from this research is the importance of data collection and  
20 transparency. Our research was made possible because of the availability of key  
21 information on COVID-19 infections in nursing homes in the US and publicly available  
22 data such as ownership, size, staffing, and key performance measures. Access to such  
23 data is invaluable in both understanding and taking preventive action to curb the  
24 COVID-19 infections in nursing homes. As such we hope that other industrialized  
25 nations take necessary steps to collect and disseminate such information to protect and  
26 safeguard the vulnerable residents in long-term care facilities.  
27

28  
29 This work leaves several areas for future research. First, given the variation in testing  
30 at different nursing homes, the number of confirmed infections may be undercounting  
31 the actual number of infections and therefore a more reliable measure would be the  
32 number of fatalities associated with COVID-19. Second, should temporal data become  
33 available, researchers can study growth curves of infections or deaths among staff and  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 residents and examine their interlinked effects on each other. Third, should national  
4 data become available, we can test our contentions using a much larger sample at the  
5 national level. This would increase the external validity and generalizability of our  
6 findings. Finally, when data from other states and other time becomes available, we can  
7 include a spatial random effect in the model to account for spatial dependencies  
8 between the infections at different nursing homes.  
9

10 One of the limitations of the study is that its data on nursing homes' features is collected  
11 in 2017 which is over two years prior to the outbreak. Although more recent data were  
12 available on the time of the study, the variable "inflation score" had to be adopted from  
13 the 2017 data. We should also note that 86 percent of CA nursing homes are for-profit  
14 and these nursing homes were probably more likely to under-report their infection rates  
15 and deaths than other nursing homes for fear of losing residents and revenue.<sup>33</sup>  
16

### 17 **Author Contributions**

18 RG and NY, designed the study. RG, XH, and NY had full access to all the data in the  
19 study and take responsibility for the integrity of the data and the accuracy of the data  
20 analyses. RG, XH, and NY analysed the data. RG and NY interpreted the data. NY  
21 drafted the manuscript. NY, and RG critically revised the manuscript.  
22

### 23 **Funding Statement**

24 This research received no specific grant from any funding agency in the public,  
25 commercial or not-for-profit sectors.  
26

### 27 **Competing Interests**

28 There are no competing interests for any of the authors.  
29

### 30 **Data sharing statement**

All data in this research are publicly available and their sources have been cited in the manuscript.

## References

1. McMichael TM, Currie DW, Clark S, et al. Epidemiology of Covid-19 in a Long-Term Care Facility in King County, Washington. *N Engl J Med*. Published online March 27, 2020. doi:10.1056/NEJMoa2005412
2. Arentz M, Yim E, Klaff L, et al. Characteristics and Outcomes of 21 Critically Ill Patients With COVID-19 in Washington State. *JAMA*. 2020;323(16):1612-1614. doi:10.1001/jama.2020.4326
3. Comas-Herrera A, Zalakain J, Litwin C, Hsu AT, Lane N, Fernández J-L. Mortality associated with COVID-19 outbreaks in care homes: early international evidence. *Int Long-Term Care Policy Netw CPEC-LSE Internet*. Published online 2020.
4. Yourish K, Lai KKR, Ivory D, Smith M. One-Third of All U.S. Coronavirus Deaths Are Nursing Home Residents or Workers. *The New York Times*. <https://www.nytimes.com/interactive/2020/05/09/us/coronavirus-cases-nursing-homes-us.html>. Published May 11, 2020. Accessed May 12, 2020.
5. Bedford J, Enria D, Giesecke J, et al. COVID-19: towards controlling of a pandemic. *The Lancet*. 2020;395(10229):1015–1018.
6. Adalja AA, Toner E, Inglesby TV. Priorities for the US Health Community Responding to COVID-19. *JAMA*. 2020;323(14):1343-1344. doi:10.1001/jama.2020.3413
7. Xu S, Li Y. Beware of the second wave of COVID-19. *The Lancet*. 2020;395(10233):1321-1322. doi:10.1016/S0140-6736(20)30845-X
8. Leung K, Wu JT, Liu D, Leung GM. First-wave COVID-19 transmissibility and severity in China outside Hubei after control measures, and second-wave scenario planning: a modelling impact assessment. *The Lancet*. 2020;395(10233):1382-1393. doi:10.1016/S0140-6736(20)30746-7
9. The New York Times. California Coronavirus Map and Case Count. *The New York Times*. <https://www.nytimes.com/interactive/2020/us/california-coronavirus-cases.html>. Accessed May 21, 2020.
10. California Department of Public Health. Skilled Nursing Facilities: COVID-19. Accessed May 21, 2020. [https://www.cdph.ca.gov/Programs/CID/DCDC/Pages/COVID-19/SNFsCOVID\\_19.aspx](https://www.cdph.ca.gov/Programs/CID/DCDC/Pages/COVID-19/SNFsCOVID_19.aspx)
11. Archived Datasets | Data.Medicare.gov. Data.Medicare.Gov. Accessed May 23, 2020. <https://data.medicare.gov/data/archives/nursing-home-compare>

12. Han X, Yaraghi N, Gopal R. Winning at all costs: Analysis of inflation in nursing homes' rating system. *Prod Oper Manag.* 2018;27(2):215–233.
13. Calgary O. Archived Datasets | Data.Medicare.gov. Data.Medicare.Gov. Accessed October 27, 2020. <https://data.medicare.gov/data/archives/nursing-home-compare>
14. Walhin JF. Bivariate ZIP models. *Biom J.* 2001;43(2):147–160.
15. AlMuhayfith FE, Alzaid AA, Omair MA. On bivariate Poisson regression models. *J King Saud Univ-Sci.* 2016;28(2):178–189.
16. Guide SU. Version 9.SAS Inst. *Inc Cary NC.* Published online 2004.
17. SAS code for BMJ. figshare. doi:10.6084/m9.figshare.13179875.v1
18. Data for COVID-19 in California nursing homes. Published online October 30, 2020. doi:10.6084/m9.figshare.13148813.v3
19. Hillmer MP, Wodchis WP, Gill SS, Anderson GM, Rochon PA. Nursing Home Profit Status and Quality of Care: Is There Any Evidence of an Association? *Med Care Res Rev.* 2005;62(2):139-166. doi:10.1177/1077558704273769
20. Comondore VR, Devereaux PJ, Zhou Q, et al. Quality of care in for-profit and not-for-profit nursing homes: systematic review and meta-analysis. *Bmj.* 2009;339:b2732.
21. Harrington C, Woolhandler S, Mullan J, Carrillo H, Himmelstein DU. Does investor ownership of nursing homes compromise the quality of care? *Am J Public Health.* 2001;91(9):1452–1455.
22. Amirkhanyan AA, Kim HJ, Lambright KT. Does the public sector outperform the nonprofit and for-profit sectors? Evidence from a national panel study on nursing home quality and access. *J Policy Anal Manage.* 2008;27(2):326-353. doi:10.1002/pam.20327
23. Johari K, Kellogg C, Vazquez K, Irvine K, Rahman A, Enguidanos S. Ratings game: an analysis of Nursing Home Compare and Yelp ratings. *BMJ Qual Saf.* 2018;27(8):619-624. doi:10.1136/bmjqs-2017-007301
24. Neuman MD, Wirtalla C, Werner RM. Association between skilled nursing facility quality indicators and hospital readmissions. *JAMA.* 2014;312(15):1542-1551. doi:10.1001/jama.2014.13513
25. Sanghavi P, Pan S, Caudry D. Assessment of nursing home reporting of major injury falls for quality measurement on nursing home compare. *Health Serv Res.* 2020;55(2):201-210. doi:10.1111/1475-6773.13247
26. Fuller RL, Goldfield NI, Hughes JS, McCullough EC. Nursing Home Compare Star Rankings and the Variation in Potentially Preventable Emergency Department Visits and Hospital Admissions. *Popul Health Manag.* 2019;22(2):144-152. doi:10.1089/pop.2018.0065

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54
27. Han X, Yaraghi N, Gopal R. Catching Them Red-Handed: Optimizing the Nursing Homes' Rating System. *ACM Trans Manag Inf Syst TMIS*. 2019;10(2):1–26.
  28. Werner RM, Konetzka RT, Polsky D. Changes in consumer demand following public reporting of summary quality ratings: an evaluation in nursing homes. *Health Serv Res*. 2016;51:1291–1309.
  29. Abrams HR, Loomer L, Gandhi A, Grabowski DC. Characteristics of U.S. Nursing Homes with COVID-19 Cases. *J Am Geriatr Soc*. 2020;68(8):1653-1656. doi:10.1111/jgs.16661
  30. Harrington C, Ross L, Chapman S, Halifax E, Spurlock B, Bakerjian D. Nurse Staffing and Coronavirus Infections in California Nursing Homes. *Policy Polit Nurs Pract*. 2020;21(3):174-186. doi:10.1177/1527154420938707
  31. McGregor MJ, Harrington C. COVID-19 and long-term care facilities: Does ownership matter? *CMAJ*. 2020;192(33):E961-E962. doi:10.1503/cmaj.201714
  32. Harrington C, Olney B, Carrillo H, Kang T. Nurse Staffing and Deficiencies in the Largest For-Profit Nursing Home Chains and Chains Owned by Private Equity Companies. *Health Serv Res*. 2012;47(1pt1):106-128. doi:10.1111/j.1475-6773.2011.01311.x
  33. Harrington C, Pollock AM, Sutaria S. Privatization of Nursing Homes in the United Kingdom and the United States. In: *The Privatization of Care, The Case of Nursing Homes*. Routledge; 2019:51–67.

## 55 Figures

### 57 Figure 1. Impact of Improved Rating System on Infection Density Curves



1  
2  
3 Note: The blue (solid) curve represents the density of predicted number of  
4 infections under current rating system while the red (dashed) curve shows the  
5 density of counterfactual number of infections had there been no discrepancy  
6 between self- and CMS-reported ratings. The vertical blue and red lines show  
7 the average number of predicted infections with and without discrepancy in  
8 ratings.  
9  
10  
11  
12  
13

14  
15 **Figure 2.** Comparison of Performance of ZIBP, NN, and SVM-RBF Models in  
16 Predicting at Least One Infection  
17

18 Note: The first 50 nursing homes are zoomed in at the top right corner of the  
19 figure. The lift of ZIBP model is presented in green, while the lifts of NN and  
20 SVM-RBF are presented with purple and red lines respectively.  
21  
22  
23  
24  
25  
26  
27  
28

29  
30 **Figure 3.** Performance of ZIBP Model for Predicting Large Outbreaks (More than 10  
31 Infections) Among Staff and Residents  
32

33 Note: The lifts of the ZIBP model for identifying large outbreaks among  
34 residents and staff are presented, respectively, by the green and purple line.  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Tables

**Table 1.** Sources and Descriptions of the Study Variables

Variable	Description	Source	Mean	Std. Dev.	Min	Max
<b>Outcomes</b>						
<b>Nursing home infected</b>	Indicates if the nursing home has at least one confirmed case of COVID-19 infection among its staff or residents	CDPH	0.23	0.42	0	1
<b>Confirmed residents</b>	The number of COVID-19 infections among the residents of nursing homes	CDPH	1.91	7.88	0	81
<b>Confirmed staff</b>	The number of COVID-19 infections among the staff of nursing homes	CDPH	0.41	2.19	0	26
<b>Large outbreak</b>	Among those nursing homes with at least 1 infection, indicates if the number of infected staff or residents is more than 10 infections.	Authors' calculation	0.31	0.46	0	1
<b>Severity of COVID-19 epidemic in the surrounding area</b>						
<b>County infections per 100K</b>	The rate of COVID-19 infections per 100,000 residents in the county in which the nursing home is located as of May 1 <sup>st</sup> , 2020.	New York Times	143.42	80.07	0	259.8
<b>Governance features</b>						
<b>For profit</b>	Indicates if the nursing home has a for-profit status	CMS	0.86	0.35	0	1
<b>Family council</b>	Indicates if a family council for the residents exists in the nursing home	CMS	0.2	0.4	0	1
<b>Certified beds</b>	The number of beds certified to provide care to Medicare and Medicaid beneficiaries	CMS	98.89	54.77	14	769
<b>Occupancy rate</b>	The ratio of residents to the total number of certified beds	Authors' calculation	0.87	0.12	0.14	1
<b>Inflation score</b>	Counts the number of years in which a significant discrepancy was observed between the self-reported quality	Authors' calculation	0.32	0.81	0	5

	measures and CMS-reported health inspections.					
<b>Ratings</b>						
<b>Quality rating</b>	Self-reported indicator of quality of services as of 2017	CMS	4.59	0.87	0	5
<b>Staffing rating</b>	Self-reported measure of staffing hours as of 2017. This is based on a combination of registered nurse hours per resident day and the total nursing hours per resident day.	CMS	3.41	1.13	0	5
<b>Health inspection rating</b>	CMS-reported indicator of health inspections ratings as of 2017	CMS	2.88	1.29	1	5

**Table 2.** Effects of study variables on the likelihood and the size of COVID-19 outbreaks

	Zero Inflated Bivariate Poisson Model				Zero Inflated Double Poisson Model		
	Parameter	Estimate	(95% CI)	P Value	Estimate	(95% CI)	P Value
<b>Nursing Home (Likelihood of nursing home getting at least one COVID-19 infection)</b>							
Intercept	-2.34	(-4.41 to -0.28)	0.03	-1.76	(-3.75 to 0.24)	0.08	
County infections per 100K	0.01	(0.01 to 0.02)	<.001	0.01	(0.01 to 0.02)	<.001	
For profit	-0.36	(-0.94 to 0.22)	0.22	-0.27	(-0.85 to 0.31)	0.36	
Family council	0.19	(-0.28 to 0.64)	0.44	0.21	(-0.26 to 0.67)	0.38	
Certified beds	0.01	(0.01 to 0.02)	0.01	0.01	(0.01 to 0.02)	0.01	
Occupancy rate	-0.2	(-1.99 to 1.59)	0.83	-0.98	(-2.69 to 0.74)	0.26	
Inspection rating	-0.02	(-0.19 to 0.17)	0.9	-0.02	(-0.19 to 0.17)	0.90	
Quality rating	-0.14	(-0.36 to 0.1)	0.26	-0.13	(-0.35 to 0.1)	0.27	
Staffing rating	0.01	(-0.17 to 0.18)	0.97	-0.01	(-0.18 to 0.17)	0.96	
Inflation score	0.06	(-0.18 to 0.28)	0.67	0.06	(-0.17 to 0.29)	0.61	
<b>Infected Staff (number of staff with confirmed COVID-19 infections)</b>							

Intercept	0.21	(-2.11 to 2.52)	0.87	-0.43	(-2.1 to 1.25)	0.63
County infections per 100K	-0.01	(-0.01 to 0.01)	0.23	-0.01	(-0.01 to 0.01)	0.11
For profit	-0.21	(-0.78 to 0.37)	0.49	-0.16	(-0.55 to 0.24)	0.44
Family council	-0.04	(-0.54 to 0.46)	0.89	0.19	(-0.12 to 0.49)	0.24
Certified beds	0.01	(0.01 to 0.01)	<.001	0.01	(0.01 to 0.01)	0.02
Occupancy rate	-2.39	(-4.3 to -0.47)	0.02	-1.11	(-2.53 to 0.32)	0.13
Inspection rating	-0.19	(-0.37 to -0.01)	0.05	-0.16	(-0.28 to -0.03)	0.02
Quality rating	0.4	(0.13 to 0.67)	0.01	0.33	(0.15 to 0.52)	<.001
Staffing rating	0.11	(-0.07 to 0.28)	0.23	0.25	(0.12 to 0.37)	<.001
Inflation score	0.41	(0.31 to 0.51)	<.001	0.27	(0.19 to 0.35)	<.001
<b>Infected Residents (number of residents with confirmed COVID-19 infections)</b>						
Intercept	1.36	(0.36 to 2.35)	0.01	1.69	(0.84 to 2.55)	<.001
County infections per 100K	-0.01	(-0.01 to -0.01)	<.001	-0.01	(-0.01 to -0.01)	<.001
For profit	2.54	(1.97 to 3.11)	<.001	1.88	(1.51 to 2.26)	<.001
Family council	0.07	(-0.09 to 0.21)	0.4	0.1	(-0.04 to 0.24)	0.15
Certified beds	0.01	(0.01 to 0.01)	0.04	0.01	(-0.01 to 0.01)	0.13
Occupancy rate	-0.24	(-1.01 to 0.54)	0.55	-0.15	(-0.88 to 0.6)	0.71
Inspection rating	-0.2	(-0.27 to -0.14)	<.001	-0.2	(-0.26 to -0.14)	<.001
Quality rating	0.13	(0.05 to 0.21)	0.01	0.15	(0.08 to 0.23)	<.001
Staffing rating	-0.26	(-0.31 to -0.2)	<.001	-0.2	(-0.25 to -0.15)	<.001
Inflation score	0.13	(0.08 to 0.18)	<.001	0.11	(0.06 to 0.16)	<.001
Covariance	0.69	(0.54 to 0.87)	0.01			
<b>Fit Statistics</b>						
-2 log likelihood		4422.7			4561.7	
AIC		4484.7			4621.7	
BIC		4626.4			4758.8	

Note: The coefficients in the first panel represent how the log odds of experiencing an infection changes with one unit of increase in the corresponding predictor. The coefficients in the second and third panels represent how the expected log count of the infections changes for each unit increase in the corresponding predictor.

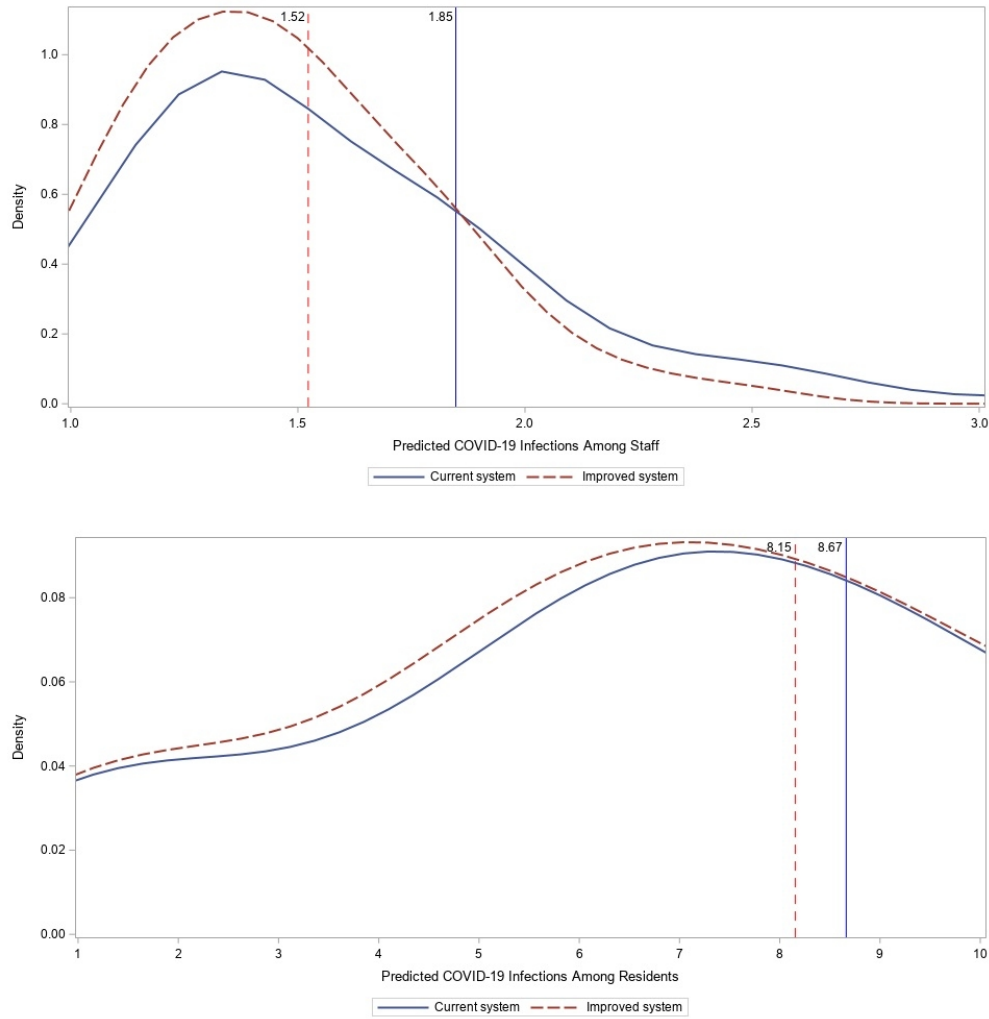


Figure 1

127x131mm (192 x 192 DPI)

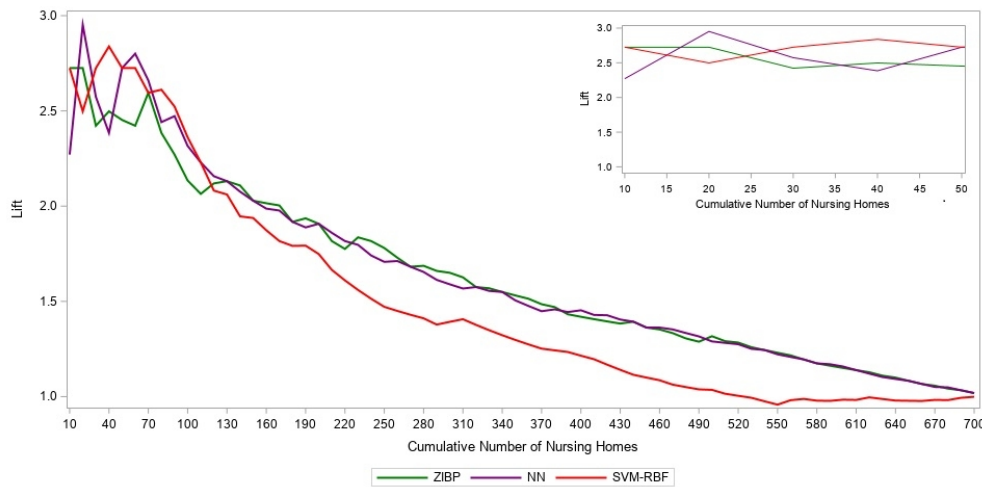


Figure 2

127x63mm (192 x 192 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

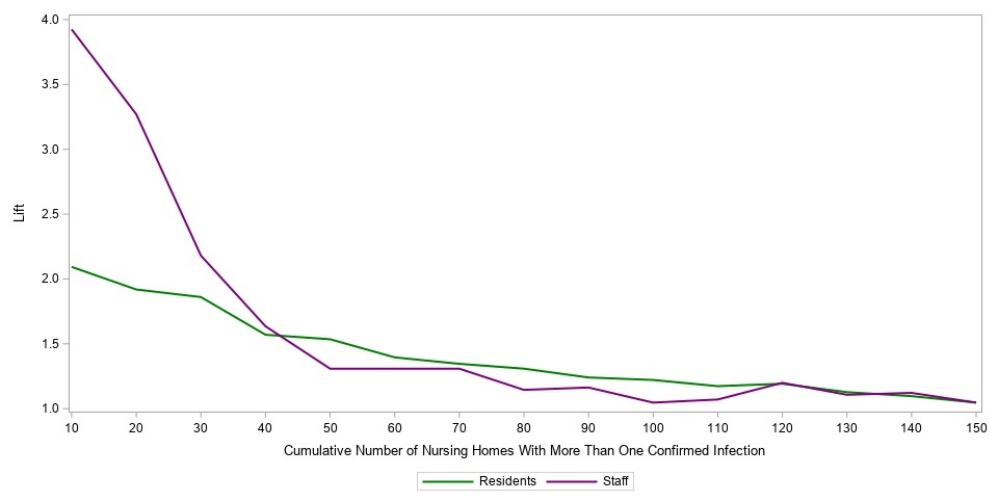


Figure 3

254x127mm (96 x 96 DPI)

## Supplementary Appendix for

# Compress the Curve: Variations in COVID-19 Infections Across California Nursing Homes

## Table of Contents

Missing Observations.....	2
Machine learning Techniques.....	2
Bivariate and Double Poisson Estimates .....	3
Correlation Between Infections Among Staff and Residents .....	3
Figures .....	4
Figure S1. Study population and analysis sample.....	4
Figure S2: Spread of COVID-19 Infection Among California Nursing Homes.....	5
Figure S3: Receiver Operator Characteristic (ROC) Curves for Predicting at Least One Infection in Nursing Homes .....	6
Figure S4: Scatter plot of number of infections among staff and residents for those nursing homes that have experienced large outbreaks amongst both their staff and resident populations .....	6
Tables.....	7
Table S1: Logistic Regression Results for Estimating the Effects of Nursing Homes' Features on Odds of Being Included in the Study Sample .....	7
Table S2: Results of Two-Sample t-Test for Equality of the Means of the Excluded and Included Observations .....	8
Table S3: Replication of the main analysis results using Bivariate and Poisson Regression Models .....	9
Table S4: Confusion Matrix for SVM-RBF .....	10
Table S5: Confusion Matrix for NN .....	10
Table S6: Distribution of Infections Among Staff and Residents.....	10



## Missing Observations

Data cleaning process is presented in Figure S1. 493 nursing homes were excluded from the study sample either due to the mismatch between their names across multiple datasets or because their COVID-19 infection data were not available in CDPH reports. To examine if the excluded nursing homes are similar to those included in the study sample, we conducted two logistic regression with the dependent variables set to be 1 to indicate if a record is included in the study sample and 0 otherwise. In the first logistic regression we only include governance features as independent variables, while in the second logistic regression we include all the features.

As reported in Table S1, both regression results show that none of the governance features are statistically significant, which indicates that the included records have no selection bias on governance features. Amongst the remaining variables, quality rating and county infections per 100k are significant are statistically significant yet the difference between the two groups is not substantial, as reported in Table S2. Further, the differences in these two variables across the two groups make our estimates more conservative.

## Machine learning Techniques

We then apply machine learning techniques to predict the COVID-19 infection in nursing homes and compare the results with our model. In view that our problem has a highly nonlinear structure, advanced machine learning models that do not rely on data structure assumptions may provide a flexible and desired solution. We predict the nursing home level COVID-19 infection situation by using Neural Networks (NN) and Support Vector Machines (SVM) with RBF kernel function. Variable *NH* is used as the target variable in each model, and is equal to 1 if at least one patient or staff reported to be infected. The prediction features include nursing home governance features such as occupancy rate, number of certified beds, whether a family council presents, whether the nursing home is for profit or not, and inflation score evaluated from past years. The nursing homes' health inspection rating, staffing rating and quality rating are also included in our prediction model. To capture the severity of COVID-19 epidemic in the surrounding area, we also incorporate county level COVID-19 infections per 100K population.

## Bivariate and Double Poisson Estimates

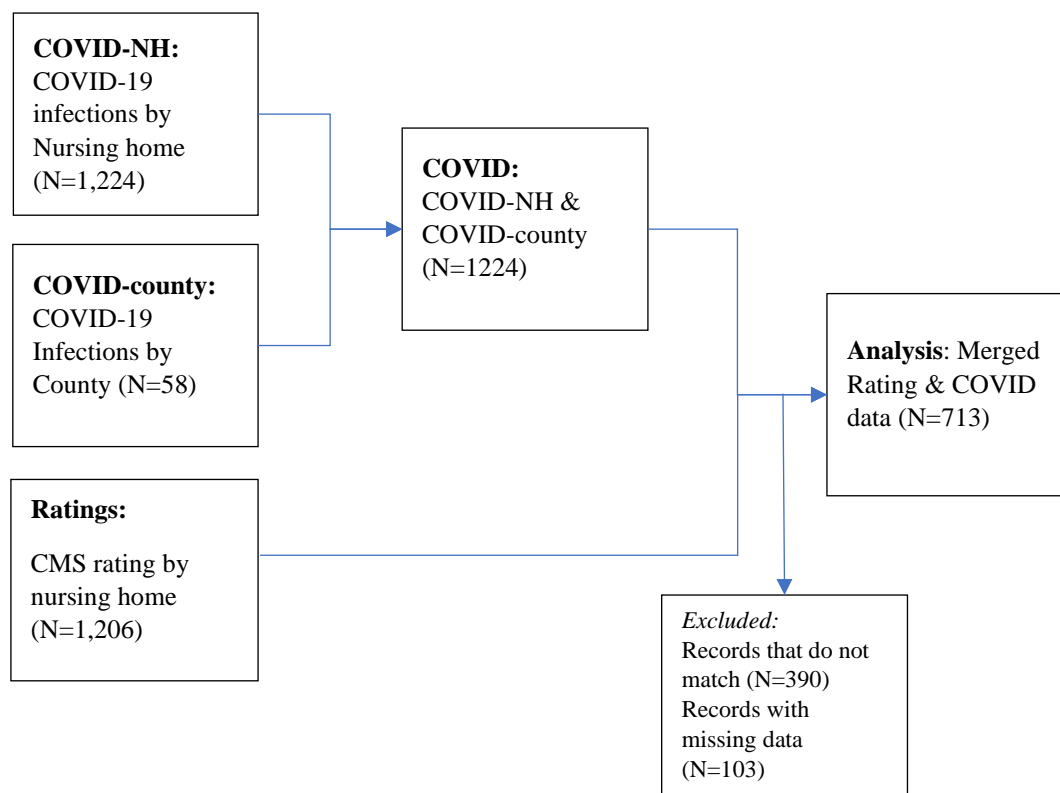
To test the robustness of our results and as a means of sensitivity analysis, we have replicated our main analysis using Bivariate and Double Poisson methods. The difference between these two methods and those reported in Table 2 of the main manuscript is these models do not assume an excess zero generating process and consider the outcome as a result of only two Poisson processes. In the Bivariate Poisson analysis, we assume that there is a correlation between the processes that give rise to the count of infections among staff and residents, while in the Double Poisson Regression, we assume independence between these two processes. The results are presented in Table S3. In comparison with the main results presented in the main table, the coefficients with larger sizes remain significant and close to their original estimates, while the smaller coefficients are not consistent with their original estimates. This is due to the fact that our dataset has significant excess zeros since most nursing homes had not reported infections many infections among either their staff or residents at the time of the study and therefore a zero inflated version of the Poisson models will be more appropriate for this setting.

## Correlation Between Infections Among Staff and Residents

To better examine the correlation between infections among staff and residents, we report the number and percentage of nursing homes with and without infections among their staff and residents in Table S6. We can observe that 91.75% of nursing homes with no infections among their residents also experienced no infections among their staff. Similarly, 54.21% of nursing homes that had at least one infection among their residents, also had at least one infection among their staff. In Figure S4, we show the scatter plot of number of infections among staff and residents for only those nursing homes that experienced a large outbreak among both their staff and residents. There is a clear correlation between the number of infections among staff and residents.

## Figures

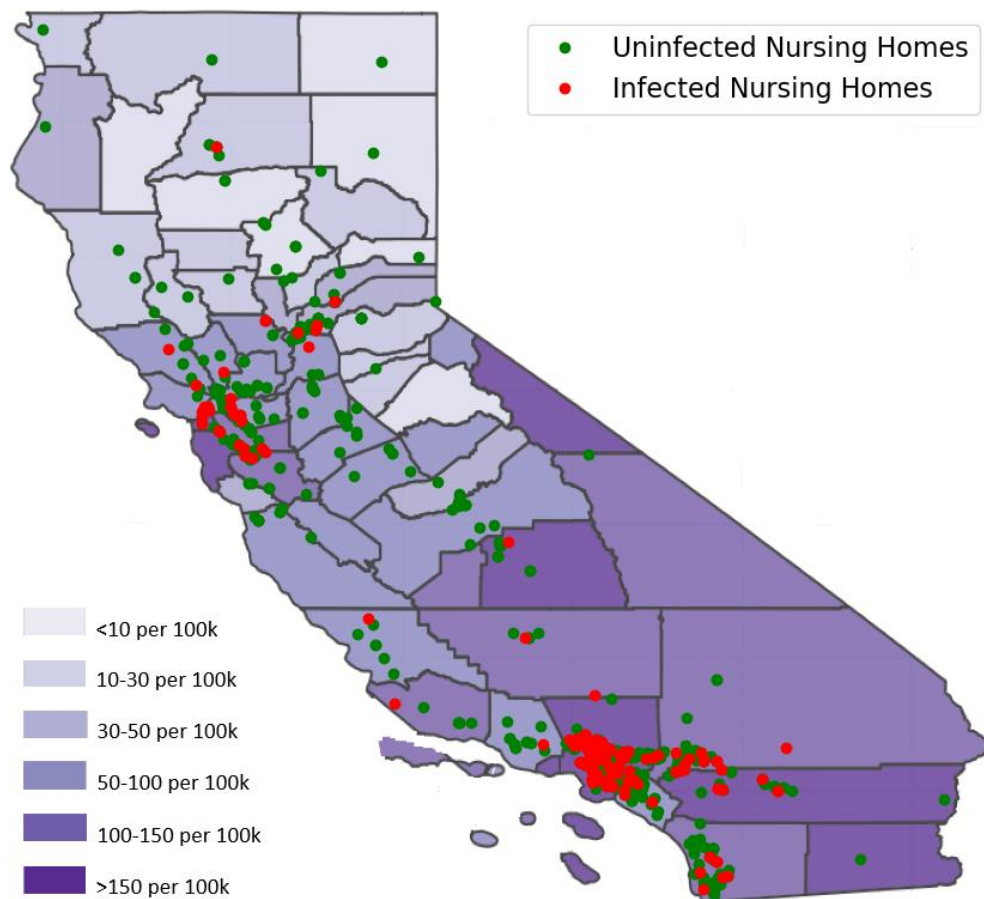
Figure S1. Study population and analysis sample



Note: Original CMS Rating for year 2017 data (*ratings*) include 1206 nursing homes. Original CA COVID-19 Infection by county (*COVID-county*) data as of April 30<sup>th</sup>, 2020 include on 58 counties. Original COVID-19 CA Infections by nursing homes (*COVID-NH*) data as of April 30<sup>th</sup>, 2020 include 1224 nursing homes.

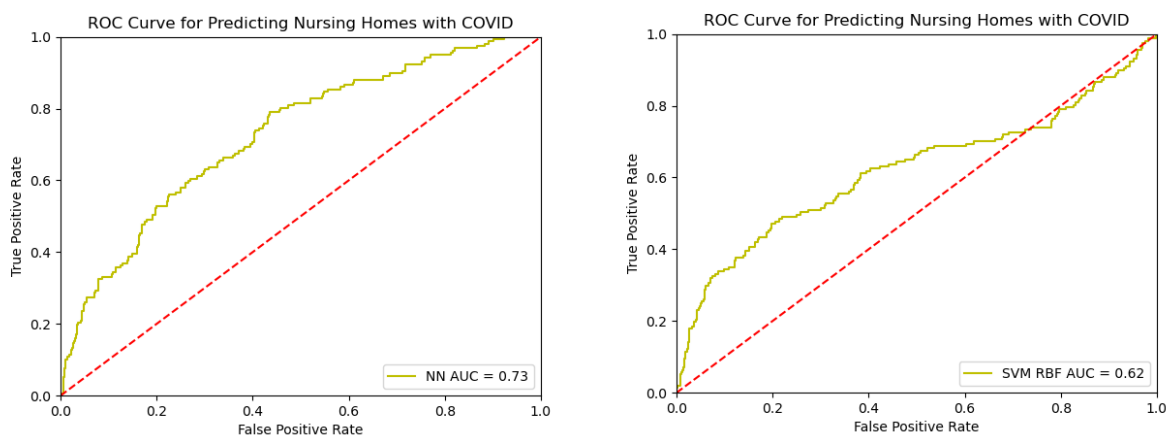
We first merged *COVID-NH* and *COVID-county* data for all 1224 rows (0 record lost). We then merged the resulting data (*COVID*) with *ratings* data which resulted in 713 rows. 390 records were lost due to mismatch between the names of the facilities in the two datasets, and 103 records were lost for those nursing homes that did not report COVID 19 infection data or their ratings information is missing.

Figure S2: Spread of COVID-19 Infection Among California Nursing Homes



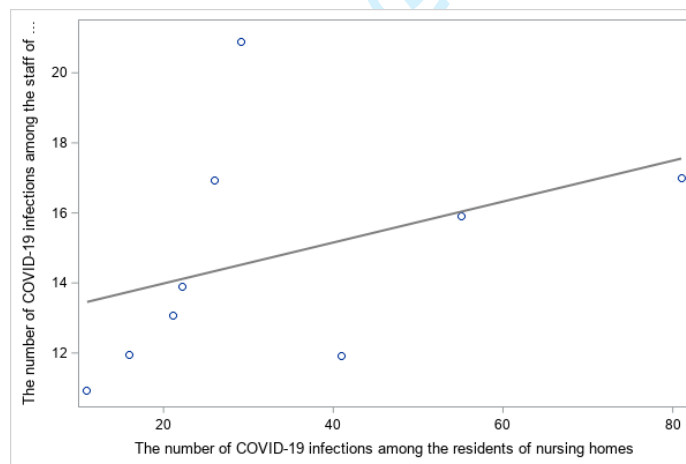
Note: The figure presents the spread of COVID-19 infection among California nursing homes as of May 1<sup>st</sup>, 2020

Figure S3: Receiver Operator Characteristic (ROC) Curves for Predicting at Least One Infection in Nursing Homes



Note: ROC for Nursing Home (NH) COVID-19 prediction using Neural Networks (NN), SVM with RBF kernel. The AUC is reported for each model: NN=0.73, SVM-RBF (default)=0.62

Figure S4: Scatter plot of number of infections among staff and residents for those nursing homes that have experienced large outbreaks amongst both their staff and resident populations



## Tables

Table S1: Logistic Regression Results for Estimating the Effects of Nursing Homes' Features on Odds of Being Included in the Study Sample

Parameter	Validation with Governance Features Only (Included vs. Excluded Records)			Validation with All Features (Included vs. Excluded Records)		
	Estimate	(95% CI)	P Value	Estimate	(95% CI)	P Value
Constant	0.1	(-0.72 to 0.92)	0.81	-0.66	(-2.09 to 0.76)	0.36
For profit	0.25	(-0.08 to 0.58)	0.14	0.29	(-0.1 to 0.68)	0.14
Family council	-0.19	(-0.49 to 0.12)	0.23	-0.07	(-0.4 to 0.26)	0.68
Certified beds	-0.0004	(-0.003 to 0.002)	0.71	-0.0008	(-0.003 to 0.002)	0.52
Occupancy rate	0.61	(-0.3 to 1.52)	0.19	0.56	(-0.62 to 1.74)	0.35
Inflation score	-0.04	(-0.2 to 0.12)	0.6	-0.03	(-0.2 to 0.14)	0.75
Quality rating				0.21	(0.07 to 0.36)	0.004
Staffing rating				0.002	(-0.14 to 0.14)	0.97
Health inspection rating				0.08	(-0.04 to 0.19)	0.21
County infections per 100K				-0.002	(-0.004 to -0.0007)	0.004

Note: Coefficients represent how the log odds of the dependent variable changes with one unit increase in the corresponding predictor

Table S2: Results of Two-Sample t-Test for Equality of the Means of the Excluded and Included Observations

Features	Excluded Records*	Included Records*	P Value**
For profit	0.82	0.86	0.11
Family council	0.21	0.18	0.21
Certified beds	99.6	98.0	0.65
Occupancy rate	0.85	0.86	0.14
Inflation score	0.32	0.31	0.83
Quality rating	4.43	4.57	0.01
Staffing rating	3.49	3.49	0.93
Health inspection rating	2.66	2.86	0.01
County infections per 100K	159.36	143.88	0.003

Note: \*: Reports the average value of features.

\*\* : P values are for two-tailed t-tests of the equality of the two means.

Table S3: Replication of the main analysis results using Bivariate and Poisson Regression Models

Parameter	Bivariate Poisson Model			Double Poisson Model		
	Estimate	(95% CI)	P Value	Estimate	(95% CI)	P Value
<b>Infected Staff (number of staff with confirmed COVID-19 infections)</b>						
Intercept	-3.9	(-5.97 to -1.83)	0.01	-3.29	(-4.7 to -1.88)	<.001
County infections per 100K	0.01	(0.01 to 0.01)	<.001	0.01	(0.01 to 0.01)	<.001
For profit	0.33	(-0.28 to 0.93)	0.3	0.01	(-0.37 to 0.39)	0.97
Family council	-0.08	(-0.59 to 0.43)	0.77	0.18	(-0.1 to 0.46)	0.21
Certified beds	0.01	(0.01 to 0.01)	<.001	0.01	(0.01 to 0.01)	<.001
Occupancy rate	-2.5	(-4.05 to -0.95)	0.01	-0.89	(-2.02 to 0.24)	0.13
Inspection rating	0.1	(-0.1 to 0.28)	0.35	-0.12	(-0.23 to -0.01)	0.05
Quality rating	0.25	(-0.05 to 0.54)	0.11	0.21	(0.03 to 0.39)	0.03
Staffing rating	0.12	(-0.06 to 0.29)	0.19	0.26	(0.14 to 0.38)	<.001
Inflation score	0.49	(0.39 to 0.59)	<.001	0.31	(0.23 to 0.39)	<.001
<b>Infected Residents (number of residents with confirmed COVID-19 infections)</b>						
Intercept	-2.1	(-3.01 to -1.19)	<.001	-1.46	(-2.2 to -0.71)	0.01
County infections per 100K	0.01	(0.01 to 0.01)	<.001	0.01	(0.01 to 0.01)	<.001
For profit	2.71	(2.12 to 3.31)	<.001	1.89	(1.5 to 2.28)	<.001
Family council	0.16	(0.02 to 0.3)	0.03	0.19	(0.06 to 0.31)	0.01
Certified beds	0.01	(0.01 to 0.01)	<.001	0.01	(0.01 to 0.01)	<.001
Occupancy rate	-0.08	(-0.66 to 0.51)	0.82	0.02	(-0.54 to 0.57)	0.96
Inspection rating	-0.2	(-0.25 to -0.14)	<.001	-0.21	(-0.26 to -0.16)	<.001
Quality rating	0.05	(-0.03 to 0.13)	0.2	0.08	(-0.01 to 0.15)	0.06
Staffing rating	-0.22	(-0.27 to -0.17)	<.001	-0.15	(-0.2 to -0.11)	<.001
Inflation score	0.13	(0.08 to 0.18)	<.001	0.13	(0.08 to 0.17)	<.001
Covariance	0.21	(0.18 to 0.25)	<.001			
<b>Fit Statistics</b>						
-2 log likelihood		8011.7			8468.6	
AIC		8053.7			8508.6	
BIC		8149.7			8600.0	



Table S4: Confusion Matrix for SVM-RBF

		ACTUAL CLASS	
		0	1
PREDICTED CLASS	0	142	2
	1	47	7

Table S5: Confusion Matrix for NN

		ACTUAL CLASS	
		0	1
PREDICTED CLASS	0	137	7
	1	37	17

Table S6: Distribution of Infections Among Staff and Residents

		INFECTIONS AMONG STAFF (%)	
		0	>=1
INFECTIONS AMONG RESIDENTS (%)	0	556 (91.75%)	50 (8.25%)
	>=1	49 (45.79%)	58 (54.21%)

**STROBE Statement**

Checklist of items that should be included in reports of observational studies

Section/Topic	Item No	Recommendation	Reported on Page No
<b>Title and abstract</b>	1	(a) Indicate the study's design with a commonly used term in the title or the abstract	1,2
		(b) Provide in the abstract an informative and balanced summary of what was done and what was found	2
<b>Introduction</b>			
Background/rationale	2	Explain the scientific background and rationale for the investigation being reported	5
Objectives	3	State specific objectives, including any prespecified hypotheses	5
<b>Methods</b>			
Study design	4	Present key elements of study design early in the paper	6
Setting	5	Describe the setting, locations, and relevant dates, including periods of recruitment, exposure, follow-up, and data collection	6
Participants	6	(a) <i>Cohort study</i> —Give the eligibility criteria, and the sources and methods of selection of participants. Describe methods of follow-up	6
		<i>Case-control study</i> —Give the eligibility criteria, and the sources and methods of case ascertainment and control selection. Give the rationale for the choice of cases and controls	
		<i>Cross-sectional study</i> —Give the eligibility criteria, and the sources and methods of selection of participants	
Variables	7	(b) <i>Cohort study</i> —For matched studies, give matching criteria and number of exposed and unexposed	6
		<i>Case-control study</i> —For matched studies, give matching criteria and the number of controls per case	
Data sources/measurement	8*	Clearly define all outcomes, exposures, predictors, potential confounders, and effect modifiers. Give diagnostic criteria, if applicable	6
Bias	9	For each variable of interest, give sources of data and details of methods of assessment (measurement). Describe comparability of assessment methods if there is more than one group	7 & Appendix
Study size	10	Describe any efforts to address potential sources of bias	6 & Appendix
Quantitative variables	11	Explain how the study size was arrived at	7
		Explain how quantitative variables were handled in the analyses. If applicable, describe which groupings were chosen and why	
		(a) Describe all statistical methods, including those used to control for confounding	
		(b) Describe any methods used to examine subgroups and interactions	
		(c) Explain how missing data were addressed	
Statistical methods	12	(d) <i>Cohort study</i> —If applicable, explain how loss to follow-up was addressed	7
		<i>Case-control study</i> —If applicable, explain how matching of cases and controls was addressed	7

*Cross-sectional study*—If applicable, describe analytical methods taking account of sampling strategy

(e) Describe any sensitivity analyses

7

Section/Topic	Item No	Recommendation	Reported on Page No
<b>Results</b>			
Participants	13*	(a) Report numbers of individuals at each stage of study—eg numbers potentially eligible, examined for eligibility, confirmed eligible, included in the study, completing follow-up, and analysed	8
		(b) Give reasons for non-participation at each stage	Appendix
		(c) Consider use of a flow diagram	Appendix
Descriptive data	14*	(a) Give characteristics of study participants (eg demographic, clinical, social) and information on exposures and potential confounders	8
		(b) Indicate number of participants with missing data for each variable of interest	8
		(c) <i>Cohort study</i> —Summarise follow-up time (eg, average and total amount)	
Outcome data	15*	<i>Cohort study</i> —Report numbers of outcome events or summary measures over time	
		<i>Case-control study</i> —Report numbers in each exposure category, or summary measures of exposure	
		<i>Cross-sectional study</i> —Report numbers of outcome events or summary measures	8,9
Main results	16	(a) Give unadjusted estimates and, if applicable, confounder-adjusted estimates and their precision (eg, 95% confidence interval). Make clear which confounders were adjusted for and why they were included	9,10
		(b) Report category boundaries when continuous variables were categorized	
		(c) If relevant, consider translating estimates of relative risk into absolute risk for a meaningful time period	
Other analyses	17	Report other analyses done—eg analyses of subgroups and interactions, and sensitivity analyses	10
<b>Discussion</b>			
Key results	18	Summarise key results with reference to study objectives	12
Limitations	19	Discuss limitations of the study, taking into account sources of potential bias or imprecision. Discuss both direction and magnitude of any potential bias	13
Interpretation	20	Give a cautious overall interpretation of results considering objectives, limitations, multiplicity of analyses, results from similar studies, and other relevant evidence	12
Generalisability	21	Discuss the generalisability (external validity) of the study results	13
<b>Other Information</b>			
Funding	22	Give the source of funding and the role of the funders for the present study and, if applicable, for the original study on which the present article is based	13

\*Give information separately for cases and controls in case-control studies and, if applicable, for exposed and unexposed groups in cohort and cross-sectional studies.

**Note:** An Explanation and Elaboration article discusses each checklist item and gives methodological background and published examples of transparent reporting. The STROBE checklist is best used in conjunction with this article (freely available on the Web sites of PLoS Medicine at <http://www.plosmedicine.org/>, Annals of Internal Medicine at <http://www.annals.org/>, and Epidemiology at <http://www.epidem.com/>). Information on the STROBE Initiative is available at [www.strobe-statement.org](http://www.strobe-statement.org).

For peer review only

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47