






OPEN ACCESS

Original research

External validation of the QCovid risk prediction algorithm for risk of COVID-19 hospitalisation and mortality in adults: national validation cohort study in Scotland

Colin R Simpson ^{1,2}, Chris Robertson,³ Steven Kerr ², Ting Shi ², Eleftheria Vasileiou,² Emily Moore,⁴ Colin McCowan,⁵ Utkarsh Agrawal,⁵ Annemarie Docherty,² Rachel Mulholland,² Josie Murray,⁶ Lewis Duthie Ritchie,⁷ Jim McMenamin,⁶ Julia Hippisley-Cox,⁸ Aziz Sheikh²

► Additional supplemental material is published online only. To view, please visit the journal online (<http://dx.doi.org/10.1136/thoraxjnl-2021-217580>).

For numbered affiliations see end of article.

Correspondence to

Professor Colin R Simpson, School of Health, Victoria University of Wellington, Wellington, New Zealand; colin.simpson@vuw.ac.nz

Received 4 May 2021

Accepted 11 October 2021

Published Online First

15 November 2021

ABSTRACT

Background The QCovid algorithm is a risk prediction tool that can be used to stratify individuals by risk of COVID-19 hospitalisation and mortality. Version 1 of the algorithm was trained using data covering 10.5 million patients in England in the period 24 January 2020 to 30 April 2020. We carried out an external validation of version 1 of the QCovid algorithm in Scotland.

Methods We established a national COVID-19 data platform using individual level data for the population of Scotland (5.4 million residents). Primary care data were linked to reverse-transcription PCR (RT-PCR) virology testing, hospitalisation and mortality data. We assessed the performance of the QCovid algorithm in predicting COVID-19 hospitalisations and deaths in our dataset for two time periods matching the original study: 1 March 2020 to 30 April 2020, and 1 May 2020 to 30 June 2020.

Results Our dataset comprised 5 384 819 individuals, representing 99% of the estimated population (5 463 300) resident in Scotland in 2020. The algorithm showed good calibration in the first period, but systematic overestimation of risk in the second period, prior to temporal recalibration. Harrell's C for deaths in females and males in the first period was 0.95 (95% CI 0.94 to 0.95) and 0.93 (95% CI 0.92 to 0.93), respectively. Harrell's C for hospitalisations in females and males in the first period was 0.81 (95% CI 0.80 to 0.82) and 0.82 (95% CI 0.81 to 0.82), respectively.

Conclusions Version 1 of the QCovid algorithm showed high levels of discrimination in predicting the risk of COVID-19 hospitalisations and deaths in adults resident in Scotland for the original two time periods studied, but is likely to need ongoing recalibration prospectively.

INTRODUCTION

In December 2019, a novel coronavirus (SARS-CoV-2) emerged in Wuhan, China.¹ WHO declared the outbreak a public health emergency of international concern on 30 January 2020, and then a pandemic on 11 March 2020. As of 15 September 2021, WHO has reported more than 225 million

Key messages

What is the key question?

- Does the QCovid algorithm accurately predict the risk of COVID-19 hospitalisation and death in Scotland?

What is the bottom line?

- The algorithm performed well according to a number of metrics we evaluated.

Why read on?

- It is important to validate the QCovid risk prediction algorithm because it is being used in the UK to inform shielding and vaccine prioritisation policies.

confirmed cases globally and over 4.6 million deaths.¹

Rapid, large-scale observational epidemiological studies have been used to identify the characteristics of people who are at greatest risk of COVID-19 hospitalisation and death, and to develop risk scoring systems.^{2–5} These studies have been used to guide policy for public health interventions, for example, lockdown measures, patient shielding and prioritisation for vaccination.⁶ The QCovid algorithm is one such risk scoring system that predicts the probability of COVID-19 hospitalisation and death. It was commissioned by the chief medical officer for England on behalf of the UK government. Version 1 of the algorithm was trained using data from 1205 general practices (n=10.5 million patients) in England using data drawn from the QResearch database for the period 24 January 2020 to 30 April 2020.² This period covers the 'first wave' of the pandemic in the UK, where testing and treatments for the disease were limited. Candidate predictor variables were selected on the basis of clinical plausibility, association with outcomes in other respiratory diseases and availability. A model selection process was followed that included removing variables with associated HRs close to 1, and variables whose predicted effect were clinically



► <http://dx.doi.org/10.1136/thoraxjnl-2021-218169>



© Author(s) (or their employer(s)) 2022. Re-use permitted under CC BY-NC. No commercial re-use. See rights and permissions. Published by BMJ.

To cite: Simpson CR, Robertson C, Kerr S, et al. *Thorax* 2022;**77**:497–504.

counterintuitive. The QCovid algorithm has been used by the UK government to inform policies on shielding and vaccine prioritisation for England.⁷

Following a request from the Scottish Government, we sought to externally validate the QCovid algorithm for the adult population resident in Scotland.

METHODS

Study design

Approximately 99% of the residents of Scotland were registered with primary care facilities that provide a comprehensive array of healthcare services. During the acute phase of the pandemic, community based COVID-19 hubs (a general practitioner (GP)-led service designed to segment patients and reduce the risk of nosocomial infections) were established. We developed retrospective cohorts drawn from patients registered with any primary care practice in Scotland from the period 1 March 2020 to 24 October 2020.

Datasets

We used data from all 940 Scottish primary care practices. Clinical data collected by primary care practitioners in Scotland have consistently been shown to be of high quality (90% completeness and accuracy⁸) and of great utility in epidemiological research.^{9–12} These were linked to the Electronic Communication of Surveillance in Scotland (national database for all virology testing including NHS (National Health Service) and UK Government test centre data), the Scottish Morbidity Record (record of hospitalisation data), and National Records Scotland (death certification) data as part of the Early Pandemic Evaluation and Enhanced Surveillance of COVID-19 (EAVE II) platform.³ A more detailed description of the data can be found in our cohort profile.¹³

Selection criteria

Any individual in the relevant linked dataset between the ages of 19 and 100 was included. Individuals who had an outcome event (COVID-19 hospitalisation or death) in the first period (1 March 2020–30 April 2020) were excluded from any analysis in the second period (1 May 2020–30 June 2020) (figure 1).

Exposures

Exposure variables were those used in the final selection of version 1 of the QCovid algorithm.² These are detailed in box 1.

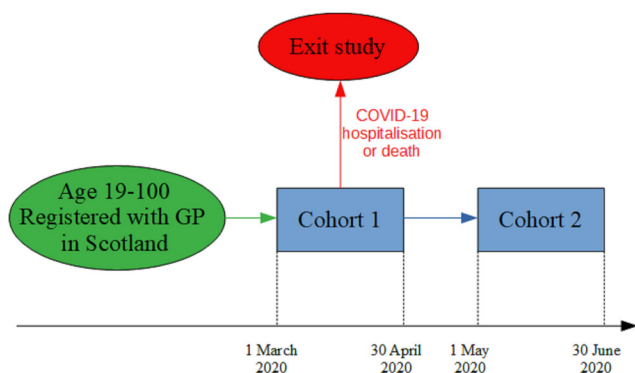


Figure 1 Study design. GP, general practitioner.

Box 1 Predictor variables in the QCovid algorithm

- ▶ Age in years (continuous).
- ▶ Townsend deprivation score (continuous).
- ▶ Accommodation (neither homeless nor care home, care home or nursing home).
- ▶ Ethnicity in 10 categories (Bangladeshi, Black African, Black Caribbean, Chinese, Indian, Mixed, Pakistani, White British, White Other, Other, Unknown).
- ▶ Body mass index (kg/m²).
- ▶ Chronic kidney disease (CKD)—(no CKD, CKD3, CKD4, CKD5, unknown).
- ▶ Learning disability (no learning disability, Down's syndrome, other learning disability).
- ▶ Chemotherapy in last 12 months (none, chemotherapy group A, B, C, unknown).
- ▶ Respiratory cancer.
- ▶ Radiotherapy in last 6 months.
- ▶ Solid organ transplant.
- ▶ Prescribed immunosuppressant medication by general practitioner.
- ▶ Prescribed leukotriene or long-acting beta blockers.
- ▶ Prescribed regular prednisolone.
- ▶ Sickle cell disease.
- ▶ Diabetes.
- ▶ Chronic obstructive pulmonary disease.
- ▶ Asthma.
- ▶ Rare pulmonary diseases.
- ▶ Pulmonary hypertension or pulmonary fibrosis.
- ▶ Coronary heart disease.
- ▶ Stroke.
- ▶ Atrial fibrillation.
- ▶ Congestive cardiac failure.
- ▶ Venous thromboembolism.
- ▶ Peripheral vascular disease.
- ▶ Congenital heart disease.
- ▶ Dementia.
- ▶ Parkinson's disease.
- ▶ Epilepsy.
- ▶ Rare neurological conditions.
- ▶ Cerebral palsy.
- ▶ Severe mental illness (bipolar disorder, schizophrenia, severe depression).
- ▶ Osteoporotic fracture.
- ▶ Rheumatoid arthritis or systemic lupus erythematosus.
- ▶ Cirrhosis of the liver.

All variables were taken as the most recent recorded value prior to the index date in the relevant dataset wherever available.

Outcomes

The primary outcomes were time to COVID-19 hospitalisation (hospitalisation with reverse-transcription PCR (RT-PCR) positive COVID-19 test within 28 days prior to admission and up to 2 days after admission, or admission with ICD-10 codes for COVID-19) and time to COVID-19 death (all-cause certified death 28 days postpositive RT-PCR test from National Records Scotland).

Missing data

Chemotherapy data were not available, so all individuals in the cohorts were assigned to the 'none' category for this

variable. We also did not have data available indicating whether the individual had a bone marrow or stem cell transplant in the last 6 months, whether they had received radiotherapy in the last 6 months, and whether they had received a solid organ transplant. The values of these variables (chemocat, p_marrow6, p_radio6 and p_solidtransplant) were set to 'none' in the cohorts. For all other comorbidities/treatments, a missing value was taken to indicate absence of that comorbidity/treatment.

Ethnicity data were not available, and all individuals in both cohorts were assigned to 'white British'. The most fine-grained residential location information available in our dataset was data zone, which is a geographical designation comprising of groups of UK Census output areas. Output areas typically consist of ~300 people, whereas data zones typically consist of 500–1000 people.¹⁴ Townsend Deprivation Scores (TDS)¹⁵ for each output area were obtained from the 2011 UK census.¹⁶ We took the median value of TDS for the output areas comprising each data zone in order to get a deprivation score for each data zone. Missing values for TDS were replaced with the mean value for the cohort. Missing values in the housing category variable were taken to indicate the individual was neither homeless, nor resident in a care home.

We used ordinary least squares regression with all other independent variables included as predictors to impute missing values for body mass index (BMI). There is some evidence of an association between higher BMI and lower levels of socioeconomic status in developed countries.¹⁷ Sex is known to be associated with BMI, as is coronary artery disease and diabetes.^{18 19}

There were no missing values for any of the other independent variables.

Model validation

We applied version 1 of the QCovid algorithm to males and females in the validation dataset and computed Harrell's Concordance,²⁰ the Brier scores, Royston's D,²¹ R^2 ²¹ and observed-expected ratio for the two time periods 1 March 2020 to 30 April 2020, and 1 May 2020 to 30 June 2020. Harrell's Concordance is a performance metric that characterises the tendency for people with higher risk scores to have earlier events. The Brier score is a measure of forecast accuracy

that is equal to the mean squared prediction error. Royston's D is a measure of 'separation' between survival curves for individuals with different characteristics. R^2 is a measure of the proportion of variation in survival time explained by the model. Observed-expected ratio is the number of observed events divided by the expected number of events predicted by the model. We made plots of observed vs expected risk by vigintiles of predicted risk. We recalibrated the algorithm in the second time period by scaling predicted risks by a multiplicative constant so that expected total number of events predicted was equal to observed total number of events.

Reporting

This study is reported in accordance with the Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis guidelines.²²

RESULTS

Characteristics of the study population

The total number of people in our dataset was 5 384 819, representing 99% of the entire population (5 463 300) estimated to be residing in Scotland in 2020.

After applying selection criteria, there were 4 392 014 individuals in the first time period cohort, and 4 382 281 individuals in the second time period cohort.

There were 5519 COVID-19 hospitalisations and 2693 COVID-19 deaths in the first time period. There were 5446 COVID-19 hospitalisations and 1300 COVID-19 deaths in the second time period. Hospitalisation and mortality tended to be positively associated with age and comorbidities (online supplemental tables 1 and 2).

Performance statistics

Table 1 shows Harrell's C, R^2 , Royston's D, Brier score and observed-expected ratio for the QCovid algorithm in predicting COVID-19 hospitalisations and deaths in our dataset for males and females in the first and second periods, respectively. Overall, the algorithm performed well according to these metrics. For predicting the risk of COVID-19 death in females in the first period, their values were: R^2 0.72 (95% CI 0.71 to 0.73); Royston's D 3.28 (95% CI 3.20 to 3.37); Harrell's C 0.95

Table 1 Performance metrics for COVID-19 hospitalisation by sex and time period

	COVID-19 death		COVID-19 hospitalisation	
	Females	Males	Females	Males
1 March 2020–30 April 2020				
R^2	0.72 (0.71–0.73)	0.69 (0.68–0.70)	0.47 (0.46–0.49)	0.48 (0.46–0.49)
Royston's D	3.28 (3.20–3.37)	3.06 (2.98–3.15)	1.93 (1.87–2.00)	1.96 (1.90–2.02)
Harrell's C	0.95 (0.94–0.95)	0.93 (0.92–0.93)	0.81 (0.80–0.82)	0.82 (0.81–0.82)
Brier Score	0.0022 (0.0008–0.0035)	0.0043 (0.0031–0.0056)	0.0011 (0.0009–0.0013)	0.0021 (0.0017–0.0025)
Observed-expected ratio	1.94	1.55	1.13	1.04
1 May 2020–30 June 2020				
R^2	0.75 (0.74–0.76)	0.75 (0.73–0.76)	0.47 (0.44–0.50)	0.54 (0.51–0.57)
Royston's D	3.56 (3.44–3.68)	3.50 (3.37–3.63)	1.92 (1.80–2.04)	0.83 (0.82–0.85)
Harrell's C	0.96 (0.95–0.96)	0.95 (0.94–0.96)	0.79 (0.78–0.81)	0.83 (0.82–0.85)
Brier Score	0.0006 (0.0004–0.0009)	0.0006 (0.0004–0.0008)	0.0003 (0.0002–0.00036)	0.0002 (0.0002–0.0002)
Observed-expected ratio	1.07	0.73	0.37	0.26

Observed-expected ratios were calculated prior to recalibration in the second period. Brier scores were calculated after recalibration.

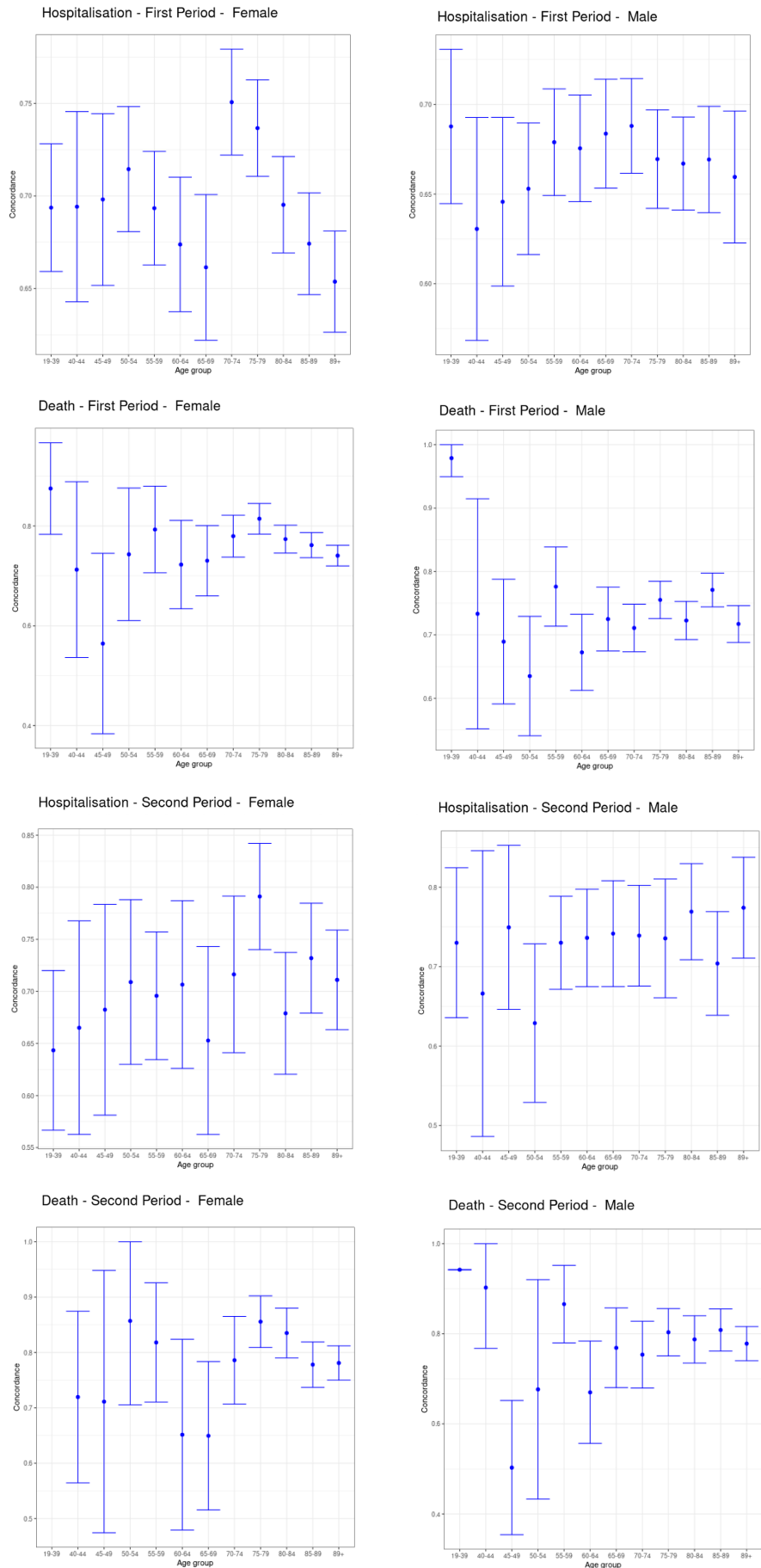
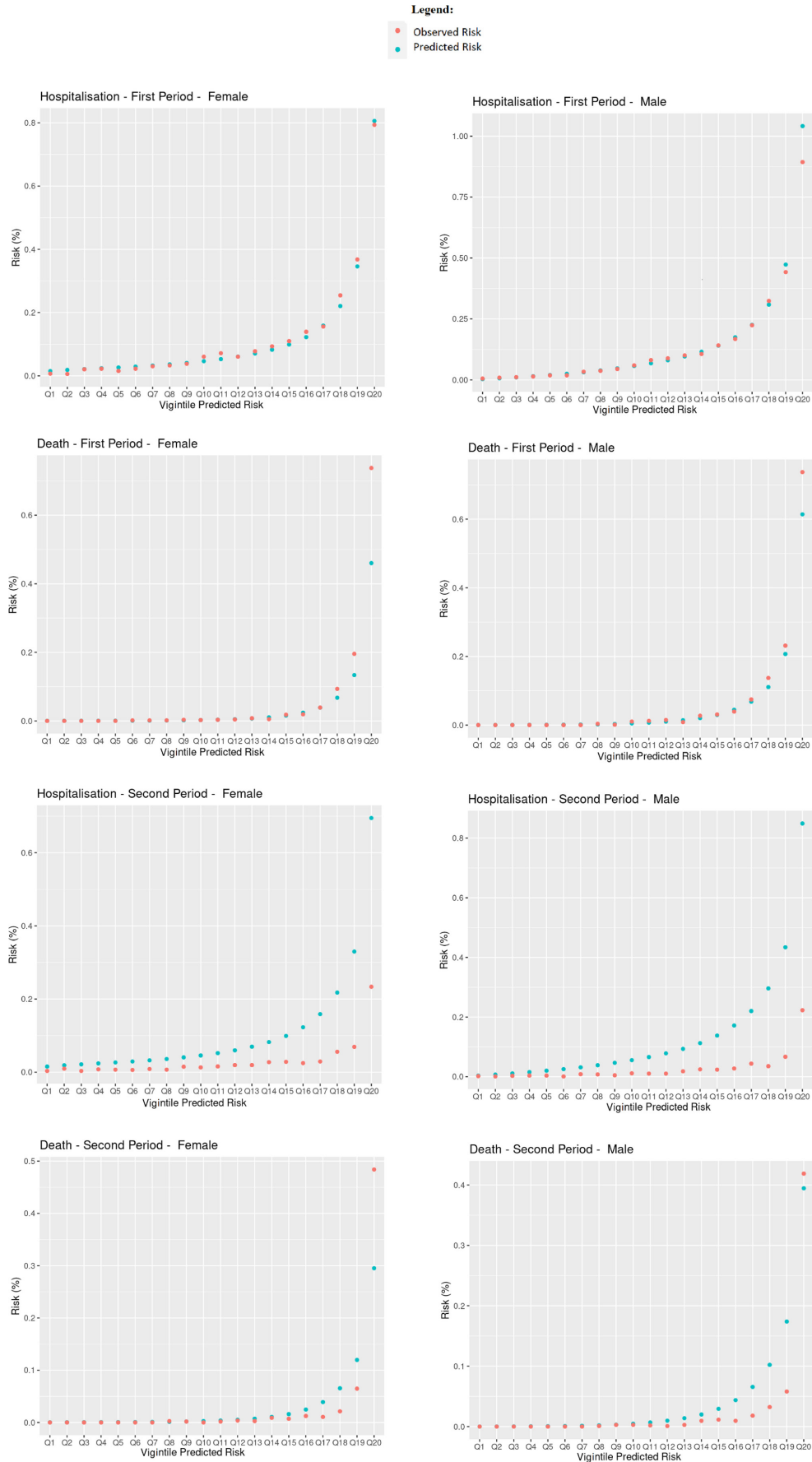


Figure 2 Harrell's C stratified by age, sex and period.



Observed risk is calculated as the proportion of individuals in the given risk vigintile who had the event

Figure 3 Observed and predicted risk.

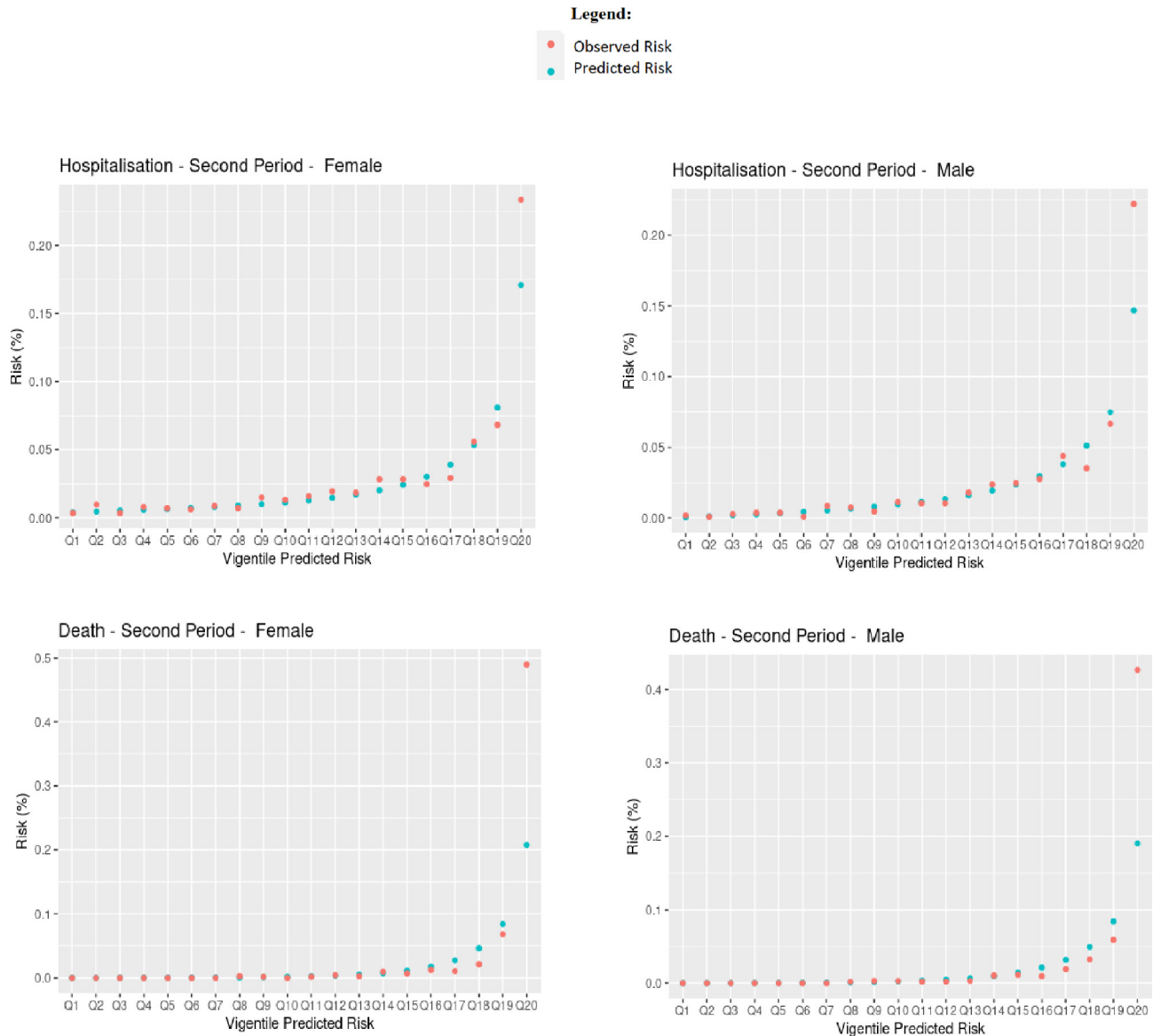


Figure 4 Observed and predicted risk in second period, recalibrated.

(95% CI 0.94 to 0.95); Brier score 0.0022 (95% CI 0.0008 to 0.0035); Observed-expected ratio 1.94. For predicting the risk of COVID-19 hospitalisation in females in the first time period, we found R^2 0.47 (95% CI 0.46 to 0.49); Royston's D 1.93 (95% CI 1.87 to 2.00); Harrell's C 0.81 (95% CI 0.80 to 0.82); Brier score 0.0011 (95% CI 0.0009 to 0.0013); Observed-expected ratio 1.13. The performance metrics for males were mostly of a similar magnitude. The high values for Harrell's C, and low values for the mean squared prediction error (Brier score) are particularly notable. Figure 2 shows Harrell's C stratified by sex, period and age group. Concordance was relatively high, with 95% CIs tending to get smaller as age increased, likely due to the larger number of events. The calibration plots in figure 3 overall showed good agreement between observed and predicted risks, particularly in the first period, but with a tendency to overpredict hospitalisation and death for those at higher predicted risk in the second period, as reflected in the low observed-expected ratios (table 1). Figure 4 shows the results after recalibrating in the second time period. Agreement between observed and predicted risk improved markedly, though risk in the highest vigintile was underpredicted.

DISCUSSION

This is the first national external validation of the QCovid algorithm for both COVID-19 hospitalisations and deaths. We found that the algorithm performed well against a range of performance metrics in males and females resident in Scotland for both the time periods under investigation. The algorithm showed good calibration in the first time period, and was improved in the second time period after recalibrating.

The QResearch database was used to train the QCovid algorithm.²³ It links together primary and secondary healthcare records, RT-PCR testing results and mortality records. As of April 2020, 1205 general practices in England were contributing to the QResearch, with coverage of ~10.5 million patients. The QResearch database has been used extensively to develop risk prediction algorithms across the NHS (National Health Service).

Overall, the cohort used in the derivation of the QCovid algorithm and the validation cohorts used in this paper were statistically quite similar with respect to the marginal distributions of patient characteristics used in the algorithm. A notable exception is BMI; ~62% of the derivation cohort had BMI in the 18.5–30

range, whereas only ~32% of the validation cohort used in this paper had a BMI in the same range. There are also likely to be significant differences in ethnic background; 64.51% of the derivation cohort had 'white' ethnicity, whereas ~96% of the population of Scotland have 'white' ethnicity according to the 2011 Scottish census.¹⁴

Our study had a number of important strengths. We developed a unique linked dataset covering 99% of the population resident in Scotland. The EAVE II database³ is one of the few national individual patient-level linked research databases in the world.²⁴ We evaluated the performance of the QCovid algorithm according to all metrics used in the original paper¹ and for identical time periods to facilitate comparison of results. We used binned plots for the Brier score and observed-expected ratio as our chosen measures of calibration because we believe they are pragmatic and informative for policy-makers.

However, our work has several limitations. We did not have data available for chemotherapy treatment, bone marrow or stem cell transplants, and solid organ transplants. The values of these variables were set to 'none' for all individuals. We believe this was reasonable for evaluating the overall performance of the algorithm because these treatments were extremely rare in the original derivation cohort.² These conditions are associated with slightly higher predicted risk of COVID-19 hospitalisation and death in the QCovid algorithm. Therefore, individuals with these conditions had slightly lower predicted risk than if this information was available. We also did not have access to ethnicity data, so all individuals were set to 'White British'. We believe modal substitution for ethnicity was reasonable because the most recent Scottish census indicated that 96% of the residents of Scotland identified their ethnicity as 'white'. We believe that minority ethnicities are more likely to have their ethnicity recorded by GPs in our dataset compared with white ethnicity and therefore missing values in the ethnicity field are likely to be disproportionately white compared with the population. *Ceteris paribus*, QCovid predicts slightly higher risk of hospitalisation and death for individuals of non-white ethnicity. Members of ethnic minority groups whose ethnicity was not recorded in our dataset will therefore have been assigned a slightly lower predicted risk of COVID-19 hospitalisation and death than if their ethnicity had been available. There was significant missingness in the BMI data, with 2 289 759 (52.1%) missing values in the 1 March 2020–30 April 2020 cohort, and 2 114 639 (48.3%) missing values in the 1 May 2020–30 June 2020 cohort. We used ordinary least squares regression to impute these missing values for BMI. We considered multiple imputation, but decided against it because it would not have been feasible given the high degree of missingness and the compute resources available to us. We believe that the average of multiple imputations would likely be similar to the mean predicted by OLS. We did not expect this to significantly affect the results of this validation exercise. The most fine-grained residential location information available in our dataset was data zone, which typically consists of multiple 2011 UK census output areas. We took the median value of the TDS for the output areas comprising each data zone in order to get a deprivation score for each data zone. Missing values of TDS were replaced with the average value for the cohort. Higher levels of deprivation as measured by TDS was associated with increased predicted risk of COVID-19 hospitalisation and death in the QCovid algorithm. The direction of the effect of having more finely

grained residential location data available would have on QCovid predicted risk of COVID-19 hospitalisation and death is dependent on whether the TDS for an individual's output area is higher or lower than the median TDS in the output areas comprising their data zone. The algorithm had a tendency to overpredict risk of hospitalisation and death for those at higher predicted risk in the second time period. However, after recalibration there was good agreement between observed and predicted risk. This suggests that while QCovid risk scores showed good levels of discrimination of outcome, risk predictions from the QCovid algorithm may require recalibration in clinical practice.

QCovid has been used by the UK government to identify the clinically extremely vulnerable for shielding advice and to inform vaccine prioritisation policies in determining priority groups, in particular for those in Joint Committee on Vaccination and Immunisation category 6.⁶ Other potential applications of the algorithm include deciding who should be prioritised for treatments or boosters. QCovid is a 'living' risk prediction algorithm, in the sense that it can be trained on new data as these become available, and should show responsiveness to changing circumstances. The time period we studied corresponds to the 'first wave' of the pandemic in the UK during which treatments were limited. Availability and use of treatments have since improved. Since the period of study in this paper, several new SARS-Cov-2 variants have emerged, vaccines have seen widespread roll out in the UK, and policies on non-pharmaceutical interventions have evolved over time. Versions 2 and 3 of the algorithm are currently under development. Vaccination status and evidence of prior infection are planned to be used as predictors in version 3. The findings from this validation work have been communicated to the Scottish government.

Author affiliations

¹School of Health, Victoria University of Wellington, Wellington, New Zealand

²Usher Institute, The University of Edinburgh, Edinburgh, UK

³Department of Mathematics and Statistics, University of Strathclyde, Glasgow, UK

⁴Information Services Division, Public Health Scotland, Edinburgh, UK

⁵School of Medicine, University of St Andrews, St Andrews, UK

⁶Health Protection Scotland, Public Health Scotland, Glasgow, UK

⁷Academic Primary Care, University of Aberdeen, Aberdeen, UK

⁸Primary Care Health Sciences, University of Oxford, Oxford, UK

Twitter Aziz Sheikh @DrAzizSheikh

Acknowledgements The authors would like to thank staff at Public Health Scotland, Albasoft Ltd, the general practices that contributed data to the study and the EAVE II Collaborators. AS, JM and CR serve on The Scottish Government's COVID-19 Chief Medical Officer's Advisory Group and the New and Emerging Respiratory Virus Threats Advisory Group (NERVTAG) Risk Stratification Subgroup.

Contributors AS and JH-C conceptualised the study. CR carried out the formal analysis. CRS wrote the initial draft of the manuscript. SK and EM assisted with the statistical analysis. SK wrote later versions of the manuscript. All authors assisted with review and editing. CR, EM and EV have verified the underlying data. CR is the guarantor for this work.

Funding Medical Research Council (MR/R008345/1), National Institute for Health Research Health Technology Assessment Programme, funded through the UK Research and Innovation Industrial Strategy Challenge Fund Health Data Research UK.

Competing interests JH-C reports grants from MRC, grants from Wellcome Trust, grants from NIHR, during the conduct of the study; other from ClinRisk, outside the submitted work. AS reports grants from NIHR, grants from MRC, grants from HDR UK, during the conduct of the study.

Patient consent for publication Not applicable.

Ethics approval Ethical permission for this study was granted from South East Scotland Research Ethics Committee 02 (12/SS/0201). The Public Benefit and Privacy Panel Committee of Public Health Scotland, approved the linkage and analysis of the deidentified datasets for this project (1920-0279).

Provenance and peer review Not commissioned; externally peer reviewed.

Data availability statement All code, metadata and documentation for this project is publicly available at <https://github.com/EAVE-II/Qcovid-validation>. A data dictionary is available at <https://github.com/EAVE-II/EAVE-II-data-dictionary>. Most of the data that were used in this study are highly sensitive and will not be made available publicly.

Open access This is an open access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited, appropriate credit is given, any changes made indicated, and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>.

ORCID iDs

Colin R Simpson <http://orcid.org/0000-0002-5194-8083>

Steven Kerr <http://orcid.org/0000-0002-3643-7859>

Ting Shi <http://orcid.org/0000-0002-4101-4535>

REFERENCES

- World Health Organization. Coronavirus disease (Covid-19) outbreak. Available: <https://www.who.int/emergencies/diseases/novel-coronavirus-2019>
- Clift AK, Coupland CAC, Keogh RH, *et al*. Living risk prediction algorithm (QCOVID) for risk of hospital admission and mortality from coronavirus 19 in adults: national derivation and validation cohort study. *BMJ* 2020;371:m3731.
- Simpson CR, Robertson C, Vasileiou E, *et al*. Early pandemic evaluation and enhanced surveillance of COVID-19 (EAVE II): protocol for an observational study using linked Scottish national data. *BMJ Open* 2020;10:e039097.
- The Scottish Government. Coronavirus (Covid-19): modelling the epidemic (issue no. 23), 2020. Available: <https://www.gov.scot/publications/coronavirus-Covid-19-modelling-epidemic-issue-no-23/> [Accessed 18 Dec 2020].
- Knight SR, Ho A, Pius R, *et al*. Risk stratification of patients admitted to hospital with covid-19 using the ISARIC WHO clinical characterisation protocol: development and validation of the 4C mortality score. *BMJ* 2020;370:m3339.
- GOV.UK: Joint Committee on vaccination and immunisation: advice on priority groups for COVID-19 vaccination, 2020. Available: <https://www.gov.uk/government/publications/priority-groups-for-coronavirus-covid-19-vaccination-advice-from-the-jcvi-30-december-2020/joint-committee-on-vaccination-and-immunisation-advice-on-priority-groups-for-covid-19-vaccination-30-december-2020>
- GOV.UK. New technology to help identify those at high risk from COVID-19, 2021. Available: <https://www.gov.uk/government/news/new-technology-to-help-identify-those-at-high-risk-from-Covid-19>
- Whitelaw FG, Nevin SL, Milne RM, *et al*. Completeness and accuracy of morbidity and repeat prescribing records held on general practice computers in Scotland. *Br J Gen Pract* 1996;46:181–6.
- Sullivan FM, Swan IRC, Donnan PT, *et al*. Early treatment with prednisolone or acyclovir in Bell's palsy. *N Engl J Med* 2007;357:1598–607.
- Dreischulte T, Donnan P, Grant A, *et al*. Safer prescribing — a trial of education, informatics, and financial incentives. *N Engl J Med* 2016;374:1053–64.
- SCOT-HEART Investigators, Newby DE, Adamson PD, *et al*. Coronary CT angiography and 5-year risk of myocardial infarction. *N Engl J Med* 2018;379:924–33.
- Ford I, Murray H, Packard CJ, *et al*. Long-term follow-up of the West of Scotland coronary prevention study. *N Engl J Med* 2007;357:1477–86.
- Mulholland RH, Vasileiou E, Simpson CR, *et al*. Cohort profile: early pandemic evaluation and enhanced surveillance of COVID-19 (EAVE II) database. *Int J Epidemiol* 2021;50:dyab028.
- Scotland's census. Available: <https://www.scotlandscensus.gov.uk/>
- Townsend P, Phillimore P, Beattie A. *Health and deprivation: inequality and the North*. London: Routledge, 1988.
- 2011 UK census Townsend deprivation scores. Available: <https://www.statistics.digitalresources.jisc.ac.uk/dataset/2011-uk-townsend-deprivation-scores>
- McLaren L. *Socioeconomic status and obesity, epidemiologic reviews.*, 2007: 29, 29–48.
- Romero-Corral A, Montori VM, Somers VK, *et al*. Association of bodyweight with total mortality and with cardiovascular events in coronary artery disease: a systematic review of cohort studies. *Lancet* 2006;368:666–78.
- Bays HE, Chapman RH, Grandy S, *et al*. The relationship of body mass index to diabetes mellitus, hypertension and dyslipidaemia: comparison of data from two national surveys. *Int J Clin Pract* 2007;61:737–47.
- Harrell FE, Califf RM, Pryor DB, *et al*. Evaluating the yield of medical tests. *JAMA* 1982;247:2543–6.
- Royston P, Sauerbrei W. A new measure of prognostic separation in survival data. *Stat Med* 2004;23:723–48.
- Collins GS, Reitsma JB, Altman DG, *et al*. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the TRIPOD statement. *BMJ* 2015;350:g7594.
- QResearch. Available: <https://www.qresearch.org>
- Simpson CR, Beever D, Challen K, *et al*. The UK's pandemic influenza research portfolio: a model for future research on emerging infections. *Lancet Infect Dis* 2019;19:e295–300.

Supplementary File**Supplemental Table 1: COVID-19 hospitalisation and mortality by demographics and medical conditions, 1 March – 30 April 2020**

Category	Cohort:	Hospitalisations:	Deaths:
Total	4,392,014 (100%)	5,519 (100%)	2,693 (100%)
Sex:			
Male	2,117,752 (48.22 %)	2,919 (52.89 %)	1,402 (52.06%)
Female	2,274,262 (51.78 %)	2,600 (47.11 %)	1,291 (47.94%)
Age band:			
19-29	779,120 (17.74 %)	109 (1.97 %)	1 (0.04%)
30-39	722,944 (16.46 %)	206 (3.73 %)	6 (0.22%)
40-49	676,503 (15.40 %)	402 (7.28 %)	33 (1.23 %)
50-59	791,298 (18.02 %)	875 (15.85 %)	111 (4.12%)
60-69	651,244 (14.83 %)	933 (16.91 %)	218 (8.10 %)
70-79	481,239 (10.96 %)	1,245 (22.56 %)	661 (24.55 %)
80-89	240,395 (5.47 %)	1,323 (23.97 %)	1,087(40.36%)
90+	49,271 (1.12 %)	426 (7.72 %)	576 (21.39 %)
BMI (integer):			
0-19	112,163 (2.55 %)	184 (3.33 %)	133 (4.94%)
20-24	631,075 (14.37 %)	717 (12.99 %)	446 (16.56 %)
25-29	788,663 (17.96 %)	1,157 (20.96 %)	474 (17.60 %)
30-34	447,812 (10.20 %)	859 (15.56 %)	283 (10.51 %)
35-39	19,106 (4.35 %)	432 (7.83 %)	124 (4.60 %)
40-51	103,436 (2.36 %)	280 (5.07 %)	55 (2.04 %)
Missing	2,289,759 (52.13%)	1,890 (34.35%)	1,535 (57.00)
Townsend Deprivation Score Quintile:			
1 (most affluent)	879,499 (20.02 %)	843 (15.27 %)	387 (14.37 %)
2	878,397 (20.00 %)	991 (17.96 %)	525 (19.49 %)
3	874,417 (19.91%)	1,095 (19.84 %)	589 (21.87 %)
4	869,324 (19.79 %)	1,222 (22.14 %)	495 (18.38 %)

5 (least affluent)	864,246 (19.68 %)	1,326 (24.03 %)	660 (24.51 %)
Missing	26,221 (0.60%)	42 (0.76%)	37 (1.37%)
Chronic Kidney Disease:			
No CKD	4,248,993 (96.74 %)	4,586 (83.09 %)	2,058 (76.42%)
CKD 3	133,142 (3.03 %)	823 (14.91 %)	565 (20.98 %)
CKD 4	5,683 (0.13%)	53 (0.96 %)	45 (1.67 %)
CKD 5	4,197 (0.1%)	57 (1.03 %)	25 (0.93 %)
Atrial Fibrillation	97,204 (2.21%)	577 (10.45%)	375 (13.92%)
Asthma	503,416 (11.46%)	797 (14.44%)	209 (7.76%)
Blood cancer	19,421 (0.44%)	105 (1.90%)	43 (1.60%)
Congestive Cardiac Failure	44,168(1.01%)	319 (5.78%)	184 (6.83%)
Cerebral Palsy	5,255 (0.12%)	13 (0.24%)	4 (0.15%)
Coronary heart disease	184,136 (4.19%)	846 (15.33%)	540 (20.05%)
Liver cirrhosis	21, 556 (0.49%)	77 (1.40%)	32 (1.19%)
Congenital heart disease	34,020 (0.77%)	104 (1.88%)	49 (1.82%)
COPD	122,428 (2.79%)	575 (10.42%)	341 (12.66%)
Dementia	37,181 (0.85%)	400 (7.25%)	749 (27.81%)
Diabetes Type 1	19,282 (0.44%)	49 (0.89%)	13 (0.48%)
Diabetes Type 2	233,889 (5.33%)	1,059 (19.19%)	474 (17.6%)
Epilepsy	57,233 (1.30%)	142 (2.57%)	70 (2.60%)
Hip, wrist, spine, humerus fracture	169,547 (3.86%)	557 (10.09%)	437 (16.23%)
Neurological conditions	16,456 (0.37%)	52 (0.94%)	24 (0.89%)
Parkinson's	9,106 (0.21%)	92(1.67%)	80 (2.97%)
Pulmonary hypertension	8,298 (0.19%)	68 (1.23%)	45 (1.67%)
Cystic fibrosis, bronchiectasis or alveolitis.	21,104 (0.48%)	103 (1.87%)	44 (1.63%)
Peripheral vascular disease	40,279 (0.92%)	210 (3.81%)	126 (4.68%)
SLE or rheumatoid arthritis	42,343 (0.96%)	143 (2.59%)	59 (2.19%)
Lung, oral cancer	10,212 (0.23%)	80 (1.45%)	44 (1.63%)
Severe mental illness	486,310 (11.07%)	902 (16.34%)	364 (13.52%)
Sickle cell disease or combined immune deficiency	2,706(0.6%)	6 (0.11%)	1 (0.04%)

syndrome			
Stroke, transient ischaemic attack	110,004 (2.50%)	648 (11.74%)	452 (16.78%)
Venous thromboembolism	68,763 (1.57%)	332 (6.02%)	181 (6.72%)

BMI - body mass index, COPD – chronic obstructive pulmonary disease, SLE - systemic lupus erythematosus, motor neurone disease, multiple sclerosis, myaesthesia, or Huntingtons's chorea.

Supplemental Table 2: Covid-19 hospitalisation and mortality by demographics and medical conditions, 1 May – 30 June 2020

Category:	Cohort:	Hospitalisations:	Deaths:
Total	4,382,821	5,446	1,300
Sex:			
Male	2,113,159 (48.21 %)	2,690 (49.39 %)	599 (46.08%)
Female	2,269,662 (51.79 %)	2,756 (50.61 %)	701 (53.92%)
Age band:			
19-29	778,219 (17.76 %)	149 (2.74 %)	0 (0.00%)
30-39	722,60 (16.48 %)	260 (4.77 %)	1 (0.08%)
40-49	676,027 (15.42 %)	473 (8.69 %)	17 (1.31 %)
50-59	790,64 (18.04 %)	977 (17.94 %)	38 (2.92 %)
60-69	650,315 (14.84 %)	971 (17.83 %)	94 (7.23 %)

70-79	479,512 (10.94 %)	1,095 (20.11 %)	268 (20.62%)
80-89	237,794 (5.43 %)	1,156 (21.23 %)	519 (39.92%)
90+	48,054 (1.10 %)	365 (6.70 %)	363 (27.92%)
BMI (integer):			
0-19	111,507 (2.54 %)	206 (3.78 %)	99 (7.62 %)
20-24	629,181 (14.36%)	722 (13.26 %)	224 (17.23%)
25-29	786,739 (17.95 %)	1,167 (21.43 %)	239 (18.38%)
30-34	446,822 (10.19 %)	873 (16.03 %)	149 (11.46%)
35-39	190,695 (4.35 %)	439 (8.06 %)	52 (4.00 %)
40-51	103,238 (2.36 %)	297 (5.45 %)	28 (2.15 %)
Missing	2,114,639 (48.25%)	3,704 (68.01)	791 (60.85%)
Townsend Deprivation Score Quintile:			
1 (most affluent)	878,046 (20.03 %)	861 (15.81 %)	226 (17.38%)
2	876,583 (20.00 %)	958 (17.59 %)	275 (21.15%)
3	872,528 (19.91%)	1,056 (19.39 %)	285 (21.92%)
4	867,483 (19.79 %)	1,231 (22.60 %)	241 (18.54%)
5 (least affluent)	862,162 (19.67 %)	1,312 (24.09 %)	256 (19.69%)
Missing	26,019 (0.59%)	28 (0.51%)	17 (1.30%)
Chronic Kidney Disease:			
No CKD	4,241,17 (96.78%)	4,563 (83.79%)	974 (74.92%)
CKD 3	131,656 (3.00%)	780 (14.32%)	293 (22.54%)
CKD 4	5,561 (0.13%)	52 (0.95%)	25 (1.92%)
CKD 5	4,128 (0.09%)	51 (0.94%)	8 (0.62%)
Atrial Fibrillation	96,079 (2.19%)	562 (10.32%)	230 (17.69%)
Asthma	502,372 (11.46%)	867 (15.92%)	121 (9.31%)
Blood cancer	19,241 (0.44%)	101 (1.85%)	25 (1.92%)
Congestive Cardiac Failure	43,523 (0.99%)	309 (5.67%)	121 (9.31%)
Cerebral Palsy	5,240 (0.12%)	13 (0.24%)	0 (0.00%)
Coronary heart disease	182,610 (4.17%)	834 (15.31%)	291 (22.38%)
Liver cirrhosis	21,399 (0.49%)	84 (1.54%)	29 (2.23%)
Congenital heart disease	33,861 (0.77%)	99 (1.82%)	25 (1.92%)
COPD	121,353 (2.77%)	555 (10.19%)	155 (11.92%)
Dementia	35,575 (0.81%)	374 (6.87%)	468 (36.00%)
Diabetes Type 1	19,236(0.44%)	53 (0.97%)	9 (0.69%)

Diabetes Type 2	232,416 (5.30%)	1,015 (18.64%)	256 (19.69%)
Epilepsy	56,987 (1.30%)	154 (2.83%)	37 (2.85%)
Hip, wrist, spine, humerus fracture	168,348 (3.84%)	539 (9.90%)	231 (17.77%)
Neurological conditions	16,386 (0.37%)	45 (0.83%)	12 (0.92%)
Parkinson's	8,910 (0.20%)	73 (1.34%)	38 (2.92%)
Pulmonary hypertension	8,145 (0.19%)	67 (1.23%)	17 (1.31%)
Cystic fibrosis, bronchiectasis or alveolitis.	20,870 (0.48%)	100 (1.84%)	21 (1.62%)
Peripheral vascular disease	39,877 (0.91%)	199 (3.65%)	57 (4.38%)
SLE or rheumatoid arthritis	42,156 (0.96%)	135 (2.48%)	36 (2.77%)
Lung, oral cancer	9,934 (0.23%)	61 (1.12%)	17 (1.31%)
Severe mental illness	484,985 (11.07%)	968 (17.77%)	194 (14.92%)
Sickle cell disease or combined immune deficiency syndrome	2,696 (0.06%)	7 (0.13%)	0 (0.00%)
Stroke, transient ischaemic attack	108,763 (2.48%)	596 (10.94%)	262 (20.15%)
Venous thromboembolism	68,193 (1.56%)	323 (5.93%)	111 (8.54%)

BMI - body mass index, COPD – chronic obstructive pulmonary disease,

SLE - systemic lupus erythematosus, motor neurone disease, multiple sclerosis, myaesthesia, or Huntingtons's chorea.