

## Rapid Epidemiological Analysis of Comorbidities and Treatments as risk factors for COVID-19 in Scotland (REACT-SCOT): a population-based case-control study

Paul M McKeigue<sup>1 3</sup>, Amanda Weir<sup>3</sup>, Jen Bishop<sup>3</sup>, Stuart J McGurnaghan<sup>2</sup>, Sharon Kennedy<sup>7</sup>, David McAllister<sup>4 3</sup>, Chris Robertson<sup>5 3</sup>, Rachael Wood<sup>7</sup>, Nazir Lone<sup>1</sup>, Janet Murray<sup>3</sup>, Thomas M Caparrotta<sup>2</sup>, Alison Smith-Palmer<sup>3</sup>, David Goldberg<sup>3</sup>, Jim McMenamin<sup>3</sup>, Colin Ramsay<sup>3</sup>, Sharon Hutchinson<sup>6 3</sup>, Helen M Colhoun<sup>2 3</sup>

- 1** Usher Institute, College of Medicine and Veterinary Medicine, University of Edinburgh, Teviot Place, Edinburgh EH8 9AG, Scotland. PM - Professor of Genetic Epidemiology and Statistical Genetics. NL - Clinical Senior Lecturer in Critical Care
- 2** Institute of Genetics and Molecular Medicine, College of Medicine and Veterinary Medicine, University of Edinburgh, Western General Hospital Campus, Crewe Road, Edinburgh EH4 2XUC, Scotland. HC - AXA Chair in Medical Informatics and Epidemiology. TC - Sir George Alberti Doctoral Fellow in Pharmacoepidemiology.
- 3** Public Health Scotland, Meridian Court, 5 Cadogan Street, Glasgow G2 6QE
- 4** Institute of Health and Wellbeing, University of Glasgow, 1 Lilybank Gardens, Glasgow G12 8RZ. DM - Wellcome Trust Intermediate Clinical Fellow and Beit Fellow
- 5** Department of Mathematics and Statistics, University of Strathclyde, 16 Richmond Street, Glasgow G1 1XQ. CR - Professor of Public Health Epidemiology
- 6** School of Health and Life Sciences, Glasgow Caledonian University. SH - Professor of Epidemiology and Population Health
- 7** NHS Information Services Division (Public Health Scotland), Gyle Square, 1 South Gyle Crescent, Edinburgh, EH12 9EB. RW - Consultant in Maternal and Child Health.

On behalf of Public Health Scotland COVID-19 Health Protection Study Group

1

## Abstract

*Background*– The objectives of this study were to identify risk factors for severe COVID-19 and to lay the basis for risk stratification based on demographic data and health records.

*Methods* – The design was a matched case-control study. Severe cases were all those with a positive nucleic acid test for SARS-CoV-2 in the national database who had entered a critical care unit or died within 28 days of the first positive test. Seven controls per case matched for sex, age and primary care practice were selected from the population register. All diagnostic codes from the past five years of hospitalisation records and all drug codes from prescriptions dispensed during the past nine months were extracted. Rate ratios for severe COVID-19 were estimated by conditional logistic regression.

*Findings* – There were 2755 severe cases. In a logistic regression using the age-sex distribution of the national population, the odds ratios for severe disease were 2.4 for a 10-year increase in age and 1.81 for male sex. In the case-control analysis, the strongest risk factor was residence in a care home, with rate ratio (95% CI) 16.2 (13.9, 18.8). Univariate rate ratios (95% CIs) for conditions listed by public health agencies as conferring high risk were 4.26 (2.90, 6.24) for Type 1 diabetes, 1.83 (1.65, 2.02) for Type 2 diabetes, 1.63 (1.47, 1.81) for ischemic heart disease, 2.51 (2.29, 2.75) for other heart disease, 2.03 (1.85, 2.22) for chronic lower respiratory tract disease, 6.0 (4.4, 8.3) for chronic kidney disease, 4.79 (4.28, 5.35) for neurological disease, 4.82 (3.23, 7.20) for chronic liver disease and 2.88 (1.94, 4.29) for immune deficiency or suppression.

74% of cases and 48% of controls had at least one listed condition (49% of cases and 8% of controls under age 40). Severe disease was associated with encashment of at least one prescription in the past nine months and with at least one hospital admission in the past five years [rate ratios 3.89 (3.15, 4.80)] and 3.10 (2.79, 3.43) respectively] even after adjusting for the listed conditions. In those without listed conditions significant associations with severe disease were seen across many hospital diagnoses and drug categories. Age and sex provided 2 bits of information for discrimination. A model based on demographic variables, listed conditions, hospital diagnoses and prescriptions provided an additional 1.07 bits (C-statistic 0.805).

*Conclusions* – Along with older age and male sex, severe COVID-19 is strongly associated with past medical history across all age groups. Many comorbidities beyond the risk conditions designated by public health agencies contribute to this. A risk classifier that uses all the information available in health records, rather than only a limited set of conditions, will more accurately discriminate between low-risk and high-risk individuals who may require shielding until the epidemic is over.

## Background

Case series from many countries have suggested that in those with severe COVID-19 the prevalence of diabetes and cardiovascular disease is higher than expected. For example in a large UK series the commonest co-morbidities were cardiac disease, diabetes, chronic pulmonary disease and asthma [1]. However there are also anecdotal reports of apparently healthy young persons succumbing to disease [2].

Quantification of the risk associated with characteristics and co-morbidities has been limited by the lack of comparisons with the background population [3–5]. Two recent studies in the UK have included population comparators and have reported associations of in hospital test positive persons and COVID-19 death in hospital with co-morbidities including diabetes, asthma and heart disease [6,7]. These studies have focused on conditions presumptively listed by public health agencies as increasing risk for COVID-19 based on case series data.

Here we examine the frequency of sociodemographic factors and these listed conditions in all people with severe COVID-19 disease in Scotland compared to matched controls from the general population. In those without listed conditions we report a systematic examination of the hospitalisation record and prescribing history in severe COVID-19 cases compared to controls. The objectives were to identify risk factors for severe COVID-19 and to lay the basis for risk stratification based on a predictive model.

## Methods

### Case definition

The Electronic Communication of Surveillance in Scotland (ECOSS) database captures all virology testing in all NHS laboratories nationally. All individuals testing positive for nucleic acid for SARS-CoV-2 up to 13 May 2020 in ECOSS were ascertained for this study. Using the Community Health Index (CHI) identifier contained in ECOSS (the CHI number is a unique identifier used in all care systems in Scotland) linkage to other datasets was carried out. Hospital admissions from the time of testing were obtained from the RAPID database a daily return of current hospitalisations each day. Admissions to critical care were obtained from the Scottish Intensive Care Society and Audit Group (SICSAG) database that covers admissions to critical care [comprising adult intensive care units (ICUs), high dependency units (HDUs) and combined ICU / HDU units] across Scotland and has returned a daily census of patients in critical care from the beginning of the COVID-19 epidemic. Death registrations up to 15 May 2020 were obtained from linkage to the National Register of Scotland.

Severe or fatal COVID-19 was defined by a record of entering critical care in the SICSAG database, or death within 28 days of a positive nucleic acid test, regardless of the cause of death given on the death certificate. By restricting the case definition to those cases that were fatal or received critical care, we ensured complete ascertainment of all test-positive cases that were severe enough to have been fatal without critical care, whatever selection policies may have determined admission to hospital or entry to critical care.

### Matched controls

For each test-positive case, we ascertained ten matched controls of the same sex, one-year age band and registered with the same primary care practice who were alive on the date of the first test in the case using the Community Health Index (CHI) database. After removing embarkations there were seven controls per case. As this is an incidence

density sampling design, it is possible and correct for an individual to appear in the dataset more than once, initially as a control and subsequently as a case.

## Demographic data

Residence in a care home was ascertained from the CHI database. Socioeconomic status was assigned as the Scottish Index of Multiple Deprivation (SIMD), an indicator based on postal code. Ethnicity was assigned based on applying a name classification algorithm (ONOMAP) [8] to the names in the CHI database. For 72% of controls and 85% of cases self-assigned ethnicity, based on the categories used in the Census, had been recorded in Scottish Morbidity Records (SMR). Cross-tabulation of 83020 records for which both name classification and SMR records of ethnicity were available showed that the ONOMAP algorithm had sensitivity of 90% and specificity of 99.42% for classifying South Asian ethnicity, but misclassified most of those who identified as African, Caribbean or Black.

## Morbidity and drug prescribing

For all cases and controls, ICD-10 diagnostic codes were extracted from the last five years of hospital discharge records in the Scottish Morbidity Record (SMR01), excluding records of discharges less than 25 days before testing positive for SARS-CoV-2 and using all codes on the discharge. Diagnostic coding under ICD chapters 5 (Mental, Behavioural and Neurodevelopmental) and 15 (Pregnancy) is incomplete as most psychiatric and maternity unit returns are not captured in SMR01. British National Formulary (BNF) drug codes were extracted from the last year of encashed prescriptions, excluding those encashed less than 25 days before testing positive for SARS-CoV-2. The BNF groups drugs by 2-digit chapter codes. For this analysis prescription codes from chapters 14 and above, mostly for dressings and appliances but also including vaccines were grouped as “Other”.

We began by scoring a specific list of conditions that have been designated as risk conditions for COVID-19 by public health agencies [9]. A separate list of conditions designates “clinically extremely vulnerable” individuals who have been advised to shield themselves completely since early in the epidemic: this list includes solid organ transplant recipients, people receiving chemotherapy for cancer, and people with cystic fibrosis or leukaemia. We did not separately tabulate these conditions as we expected these individuals to be underrepresented among cases if shielding was adequate.

The eight listed conditions were scored based on diagnostic codes in any hospital discharge record during the last five years, or encashed prescription of a drug for which the only indications are in that group of diagnostic codes. The R script included as supplementary material contains the derivations of these variables from ICD-10 codes and BNF drug codes. Diagnosed cases of diabetes were identified through linkage to the national diabetes register (SCI-Diabetes), with a clinical classification of diabetes type as Type 1, Type 2 or Other/Unknown.

## Statistical methods

To estimate the relation of cumulative incidence and mortality from COVID-19 to age and sex, logistic regression models were fitted to the proportions of cases and non-cases in the Scottish population, using the estimated population of Scotland in mid-year 2019 which were available by one-year age group up to age 90 years. To allow for possible non-linearity of the relationship of the logit of risk to age, we also fitted generalized additive models, implemented in the R function `gam::gam`, with default smoothing function.

For the case-control study, all estimates of associations with severe COVID-19 were based on conditional logistic regression, implemented as Cox regression in the R function `survival::clogit`. Among those cases and controls without any of the pre-defined conditions we then further examined associations of ICD-10 and BNF chapter with severe COVID-19. Restriction of cases and controls, for instance to exclude those with any listed condition, may generate strata that do not contain at least one case and at least one control, but these strata are ignored by the conditional logistic regression model as they do not contribute to the conditional likelihood. With incidence density sampling, the odds ratios in conditional logistic regression models are equivalent to rate ratios. Note that odds ratios in a matched case control study are based on the conditional likelihood and the unconditional odds ratios calculable from the frequencies of exposure in cases and controls will differ from these and should not be used [10]. Although matching on primary care practice will match to some extent for associated variables such as care home residence, socioeconomic disadvantage and prescribing practice, the effects of these variables are still estimated correctly by the conditional odds ratios but with less precision than in an unmatched study of the same size [10].

To construct risk prediction models, we used stepwise regression alternating between forward and backward steps to maximize the AIC, implemented in the R function `stats::step`. The performance of the risk prediction model in classifying cases versus non-cases of severe COVID-19 was examined by 4-fold cross-validation. We calculated the performance calculated over all test folds using the C-statistic but also using the “expected information for discrimination”  $\Lambda$  expressed in bits [11]. The use of bits (logarithms to base 2) to quantify information is standard in information theory: one bit can be defined as the quantity of information that halves the hypothesis space. Although readers may be unfamiliar with the expected information for discrimination  $\Lambda$ , it has several properties that make it more useful than the C-statistic for quantifying increments in the performance of a risk prediction model [11]. A key advantage of using  $\Lambda$  is that contributions of independent predictors can be added. Thus in this study we can add the predictive information from a logistic model of age and sex in the general population to the predictive information provided by other risk factors from the case-control study matched for age and sex.

## Results

### Incidence and mortality from severe COVID-19 in the Scottish population

Figure 1 shows the relationships of incidence and mortality rates to age for each sex separately. The relationship of mortality to age is almost exactly linear on a logit scale, and the lines for male and female mortality are almost parallel. In models that included age and sex as covariates, the odds ratio associated with a 10-year increase in age was 2.4 for all severe disease and 3.46 for fatal disease. The odds ratio associated with male sex was 1.81 for all severe disease and 1.82 for fatal disease. For severe cases as defined in this study, the sex differential is narrow up to about age 50 but widens between ages 50 and 70 years. Thus at younger ages the ratio of critical care admissions to total fatalities is higher in women than in men, but that at later ages the ratio of critical admissions to total fatalities is higher in men.

## Risk factors

### Sociodemographic factors

Table 1 shows univariate associations of demographic factors with severe disease. Residence in a care home was by far the strongest risk factor for severe disease. Higher risk of severe disease was also associated with socioeconomic deprivation. Associations with ethnicity are shown for the full dataset based on name classification and separately for the subset of cases and controls in whom ethnicity had been recorded in the Scottish Morbidity Record. With Whites as reference category, the rate ratio (95% CI) associated with South Asian ethnicity was 1.20 (0.84, 1.70) based on name classification and 1.24 (0.73, 2.11), based on the subset with SMR records. The numbers of cases in other non-White ethnic groups were too sparse to tabulate separately.

### Factors derived from hospitalisation and prescribing records

Prevalence of the listed conditions in cases and controls by age band is shown in Table 2. 32 (49%) of the cases aged under 40 years had at least one listed condition, compared with only 43 (8%) of the controls. In those aged 75+ years 1331 (82%) of the cases and 6591 (59%) of the controls had at least one listed condition. Among those aged under 40 years, 55 (85%) of the cases and 312 (60%) of the controls had either a hospital admission in the last five years or a dispensed prescription in the last year. Differences in prescription rates between cases and controls narrowed with increasing age.

Over all age groups, 2050 (74%) of severe cases and 9460 (48%) of controls had at least one of the listed conditions. As shown in Table 3, all the listed conditions were more frequent in cases than controls except for immune conditions in the 75+ age group. The rate ratio associated with type 1 diabetes was higher than that for type 2 diabetes. The rate ratio was 1.63 (1.47, 1.81) for ischemic heart disease compared to 2.51 (2.29, 2.75) for the broad category “other heart disease”. In multivariate analysis ischemic heart disease was not independently associated with severity whereas other heart disease remained strongly associated. In those without a listed condition 654 (93%) of the cases and 8335 (82%) of the controls had either a recent admission or a prescription. In those aged under 60 years without a listed condition, 174 (86%) of the cases and 1654 (65%) had either a recent admission or a prescription.

Supplementary Tables S1 to S3 examine these associations by age group, with the 0-39 and 40-59 year age bands combined. All listed conditions were associated with severe disease in each age band. In those aged under 60 years, the rate ratio was 6.0 (3.3, 11.0) for Type 1 diabetes and 3.48 (2.56, 4.74) for Type 2 diabetes. The multivariate analyses shown in Table 3 and S1 to S3 show that overall and in each age group any admission to hospital in the past five years were strongly and independently associated with severe disease even after adjusting for care home residence and listed conditions. Dispensing of any prescription in the past year was associated with severe disease in multivariate analyses in the two younger age bands. Table 4 shows that in each age group the proportion of fatal cases who had not had either a hospital admission in the last five years or a dispensed prescription in the last year was very low.

In a sensitivity analysis in which we also included deaths registered with mention of COVID-19 on the death certificate in the definition of severe cases (Table S4) the same pattern of differences in prior admissions and prescribing history and in listed conditions between cases and controls was found, with the difference in residential care home status being somewhat greater. Such deaths were not included in our primary outcome definition as misclassification rates for COVID assignment as cause are unknown.

## Systematic analysis of diagnoses associated with severe disease 223

The association of severe COVID-19 with prior hospital admission was examined further 224  
by testing for association of hospitalisations at each ICD-10 chapter level with severe 225  
COVID-19, among those without any of the listed conditions. These results are shown 226  
in Table 5. In univariate analyses, almost all ICD-10 chapters, with the exception of 227  
Chapters 7 (eye) Chapters 8 (ear) and Chapter 15 (pregnancy) were associated with 228  
increased risk of severe disease. Note that hospital diagnoses classified under the 229  
pregnancy chapter here are derived from admissions with pregnancy related medical 230  
conditions to non-obstetric units only, as obstetric returns are not in the SMR01 231  
dataset. In a multivariate analysis the most significant association was with diagnoses 232  
in ICD chapter 2 (neoplasms). Supplementary Table S5 extracts univariate associations 233  
with ICD-10 subchapters in those without any listed conditions. This table is filtered to 234  
show only subchapters for which the univariate p-value is <0.001 and where there are at 235  
least 50 cases and controls with a diagnosis in this subchapter. This shows that many 236  
subchapter diagnoses are associated with markedly higher risk of severe COVID-19. 237

## Associations of prescribed drugs with severe disease 238

As shown in Table 3 and supplementary tables S1 to S3, encashment of at least one 239  
prescription in the last year was associated with severe disease. The univariate rate 240  
ratio associated with this variable varies from 4.41 (3.21, 6.05) in those aged under 60 241  
years to 3.46 (2.26, 5.30) in those aged 75 years and over. In a multivariate analysis 242  
adjusting for care home residence, any hospital admission and listed conditions, these 243  
rate ratios were reduced to 2.57 (1.83, 3.61) and 1.32 (0.84, 2.08) respectively. 244

To investigate this further, we partitioned the “Any prescription” variable into 245  
indicator variables for each chapter of the British National Formulary, in which drugs 246  
are grouped by broad indication, and restricted the analysis to those without one of the 247  
listed conditions. Table 6 shows these associations. In univariate analyses, prescriptions 248  
in almost all BNF chapters were associated with severe disease. In a multivariate 249  
analysis of all chapters, most of these associations were weaker. The BNF chapters with 250  
the strongest independent associations with severe disease were chapters 1 251  
(gastrointestinal) and chapters 4 (central nervous system). Other chapters associated 252  
with severe disease were 2 (cardiovascular), 5 (infections), 9 (nutrition and blood) and 253  
14+ (other, mostly dressings and appliances). 254

## Construction of a multivariate risk prediction model 255

To evaluate the contribution of the listed conditions to risk prediction, and the 256  
incremental contribution of other information in hospitalisation and prescription records 257  
after assigning these conditions, predictive models were constructed from three sets of 258  
variables: a baseline set consisting only of demographic variables, a set that included 259  
indicator variables for each listed condition, and an extended set that included 260  
demographic, variables, indicator variables for listed conditions and indicator variables 261  
for hospital diagnoses in each ICD-10 chapter and prescriptions in each BNF chapter. 262

For each variable set, a stepwise regression procedure was carried out using 263  
alternating forward-backward selection. The variables retained with each variable set 264  
are shown in Table S6. Coefficients for specific conditions here should not be interpreted 265  
as effect estimates, as global variables for any hospital diagnosis and any listed 266  
condition have been included in the model. The predictive performance of the model 267  
chosen by stepwise regression was estimated by 4-fold cross-validation. Observed and 268  
predicted case status were compared within each stratum over all test folds. Table 7 269

shows that using the extended set increased the C-statistic from 0.777 to 0.805 and the expected information for discrimination  $\Lambda$  from 0.89 bits to 1.07 bits.

This estimate of 1.07 bits for the information conditional on age and sex obtained from the matched case-control study can be added to the information for discrimination 2 bits obtained from the logistic regression on age and sex in the population using age and sex to estimate the total information for discrimination of a risk classifier that would be obtained in the population as 3.07 bits.

Figure 2 shows the distribution of the weight of evidence favouring case over control status from the model based on the extended variable set with a footnote explaining how  $\Lambda$  is derived. This shows, as expected for a multifactorial classifier, that the distributions are approximately Gaussian: there is no clear divide between high-risk and low-risk individuals of the same age and sex. Figure 3 shows the receiver operating characteristic curve with a footnote explaining its derivation from the distribution of the weights of evidence.

## Discussion

### Sociodemographic factors

This analysis confirms that risk for severe COVID-19 is associated with increasing age, male sex and socioeconomic deprivation. The slope of the relationship of severe disease (on the scale of log odds) to age is less steep than the slope of the relationship of fatal disease to age. Residence in a care home was associated with a 16-fold increased rate of severe COVID-19 in this age matched analysis, reduced to 12-fold by adjustment for listed conditions. This excess risk is likely to reflect both the spread of the epidemic in care homes and residual confounding by frailty.

Although the numbers of cases and controls of non-White ethnicity are small and the assignment of ethnicity is incomplete, the results give some indication of the likely upper bound of the absolute numbers of severe cases in non-White ethnic groups up to now. The only non-White ethnic group with any sizeable numbers is the South Asian category and we found no clear evidence of any elevation in risk in this group compared to Whites. Reports from England [7] found elevation in risks for some non-White groups. In the OpenSAFELY study risk ratios for fatal COVID-19 of 1.7 in those recorded as Black and 1.6 in those recorded as Asian, in comparison with those recorded as White, persisted after adjustment for comorbidities and socioeconomic status. In a study of risk factors for hospitalized disease in the UK Biobank cohort, adjustment for health care worker status and other social variables attenuated but did not fully explain the elevated crude risk ratios associated with non-White ethnicity [6,12]. The relative socioeconomic position of ethnic groups in Scotland is different to that in England, so it is plausible that the relation of health status to ethnicity will also differ. For example in the 2011 Scottish Census 1.6% of the population reported South Asian ethnicity. Among the 1.0% who identified as Pakistani or Bangladeshi the proportion living in the most deprived neighbourhoods was not higher than the national average [13]. Future work may allow more complete assignment of ethnicity and disaggregation of broad categories based on continent of origin.

### Co-morbidities

We have confirmed that the moderate risk conditions designated by the NHS and other agencies [9] are associated with increased risk of severe COVID-19. However the rate ratios associated with these conditions vary with age - for example the rate ratio associated with diabetes is higher at younger ages. The rate ratios of 4.3 for Type 1 diabetes and 1.8 for Type 2 diabetes are broadly similar to those reported in UK



Biobank and in the OpenSAFELY studies. We confirm the higher risk with asthma and chronic lung disease and liver disease reported in these and earlier studies. Of note other heart disease is more strongly associated than ischaemic heart disease. This category includes conditions such as atrial fibrillation, cardiomyopathies and heart failure. Over all age groups, 74% of severe cases had at least one of these listed conditions. Among cases and controls without these conditions, not surprisingly, neoplasms were associated with severe COVID-19; we had omitted it from the pre-specified list as in the current dataset we cannot separately identify those who are currently receiving chemotherapy or radiotherapy for whom shielding is advised. We have not attempted to estimate the risk associated with these conditions for which shielding is recommended, as the observed risk will depend on the adequacy of shielding rather than on the risk to those exposed to the epidemic. In patients without any listed conditions, further systematic evaluation of past hospitalisation history did not reveal a sparse set of underlying conditions; instead many diagnoses were associated with severe COVID-19.

Media reports of apparently healthy young people succumbing to severe COVID-19 have disseminated the message that all are at risk of disease whatever their age or health status. However we found that half of cases under 40 years had at least one of the listed conditions and among those who did not have one of these conditions, the proportions who had at least one prior hospitalisation or dispensed prescription were higher in cases than in controls. In all age groups, very few of the fatal cases had not had either a hospital admission in the past five years or a dispensed prescription in the past year.

An important finding of this study was the strong association of severe COVID-19 with having encashed at least one prescription in the past year, only partly explained by higher rates of prescribing among those with listed conditions. Partitioning of this association between BNF chapters, which represent broad indication-based drug classes, showed that the strongest association was with prescription of Chapter 1 drugs, prescribed for gastrointestinal conditions, which are not generally listed as risk factors for severe COVID-19. Also associated were those in the nervous system, cardiovascular, and nutritional and blood chapters. Although it is likely that most associations of severe COVID-19 with drug prescribing are attributable to the indications for which these drugs were prescribed, or more diffuse frailty especially in older persons, causal effects of drugs or direct effects of polypharmacy on susceptibility cannot be ruled out. These associations are explored in an accompanying paper.

## Relevance to policy

As lockdown restrictions are eased, there is general agreement that vulnerable individuals will require shielding, even if the restart of the epidemic can be slowed or suppressed by mass testing, contact tracing and isolation of those who test positive. The “stratify and shield” policy option [14], in which high-risk individuals comprising up to 15% of the population are shielded for a defined period while the epidemic is allowed to run relatively quickly in low-risk individuals until population-level immunity is attained, depends critically on informative risk discrimination. So too does the similarly named “segment and shield” option [15] which has the opposite objective of keeping transmissions low.

As awareness grows of how risk varies between individuals, individuals will seek information about their own level of risk. A key implication of our results is that risk of severe or fatal disease is multifactorial and that the rate ratio of 5.8 associated with a 20-year increase in age is stronger than that associated with common diseases such as Type 2 diabetes and asthma that are listed as conditions associated with high risk. A corollary of this is that a crude classification based on assigning all persons with a listed condition to a group for whom shielding is recommended will have poor specificity, as one quarter of those aged 60-74 years in the population have at least one of the listed

conditions we examined. It will also exclude many people at high risk because they have multiple risk factors each of small effect. The only way to optimize risk classification so as to ensure equity with respect to risk is to construct a classifier that uses all available information to assign a risk score. Our results show that this is possible in principle, though for this preliminary study we have not used the full repertoire of machine learning methods available for this type of problem. In Scotland it is technically possible to use existing electronic health records to calculate a risk score for every individual in the population, though more work would be required to develop this as a basis for official advice and individual decisions.

## Methodological strengths and weaknesses

Most reports of disease associations with COVID-19 have been case series. There have been few reports based on evaluating these associations in the population through cohort or case-control studies. With this matched case control design using incidence density sampling, we have been able to estimate rate ratios conditional on age and sex. An unpublished analysis from England explored the association of similar set of risk conditions with in-hospital COVID-19 deaths, but did not systematically evaluate the rest of the medical record including prescription records. Although we have records of encashment of prescriptions, we do not at present have access to other primary care data, which would contain additional information on morbidity and measurements such as body mass index. A strength of our study however is that hospital discharge diagnoses are coded to ICD-10 by trained coders, in contrast to the coding systems used in primary care databases that do not map to recognized disease classifications. Associations with ethnicity and other sociodemographic factors are not necessarily generalizable from Scotland to other populations.

## Conclusion

This study confirms that risk of severe COVID-19 is associated with sociodemographic factors and with chronic conditions such as diabetes, asthma, circulatory disease and others. However the associations with pre-existing disease are not just with a small set of conditions that contribute to risk, but with many conditions as demonstrated by associations with past medical and prescribing history in relation to multiple physiological systems. As countries attempt to emerge from lockdown whilst protecting vulnerable individuals, multivariate classifiers rather than crude rule-based approaches will be needed to define those most at risk of developing severe disease.

## Declarations

### Information governance

This study was conducted under approvals from the Privacy Advisory Committee ref 44/13 and Public Benefit Privacy Protection amendment 1617-0147. Datasets were de-identified before analysis.

### Conflicts of interest

The authors declare no conflicts of interest.

## Acknowledgements

We thank all staff in critical care units who submitted data to the SICSAG database, the Scottish Morbidity Record Data Team, the staff of the National Register of Scotland, the Public Health Scotland Terminology Services, the HPS COVID-19 Laboratory & Testing cell and the NHS Scotland Diagnostic Virology Laboratories, and Nicola Rowan (HPS) for coordinating this collaboration.

## Public Health Scotland COVID-19 Health Protection Study Group

Alice Whettlock<sup>1</sup>, Allan McLeod<sup>1</sup>, Andrew Gasiorowski<sup>1</sup>, Andrew Merrick<sup>1</sup>, Andy McAuley<sup>1</sup>, April Went<sup>1</sup>, Calum Purdie<sup>1</sup>, Colin Fischbacher<sup>1</sup>, Colin Ramsay<sup>1</sup>, David Bailey<sup>1</sup>, David Henderson<sup>1</sup>, Diogo Marques<sup>1</sup>, Eisin McDonald<sup>1</sup>, Genna Drennan<sup>1</sup>, Graeme Gowans<sup>1</sup>, Graeme Reid<sup>1</sup>, Heather Murdoch<sup>1</sup>, Jade Carruthers<sup>1</sup>, Janet Fleming<sup>1</sup>, Jade Carruthers<sup>1</sup>, Joseph Jasperse<sup>1</sup>, Josie Murray<sup>1</sup>, Karen Heatlie<sup>1</sup>, Lindsay Mathie<sup>1</sup>, Lorraine Donaldson<sup>1</sup>, Martin Paton<sup>1</sup>, Martin Reid<sup>1</sup>, Melissa Llano<sup>1</sup>, Michelle Murphy-Hall<sup>1</sup>, Paul Smith<sup>1</sup>, Ros Hall<sup>1</sup>, Ross Cameron<sup>1</sup>, Susan Brownlie<sup>1</sup>, Adam Gaffney<sup>2</sup>, Aynsley Milne<sup>2</sup>, Christopher Sullivan<sup>2</sup>, Edward McArdle<sup>2</sup>, Elaine Glass<sup>2</sup>, Johanna Young<sup>2</sup>, William Malcolm<sup>2</sup>, Jodie McCoubrey<sup>2</sup>

<sup>1</sup> Health Protection Scotland (Public Health Scotland), Meridian Court, 5 Cadogan Street, Glasgow G2 6QE.

<sup>2</sup> NHS National Services Scotland, Meridian Court, 5 Cadogan Street, Glasgow G2 6QE.

## Supplementary material

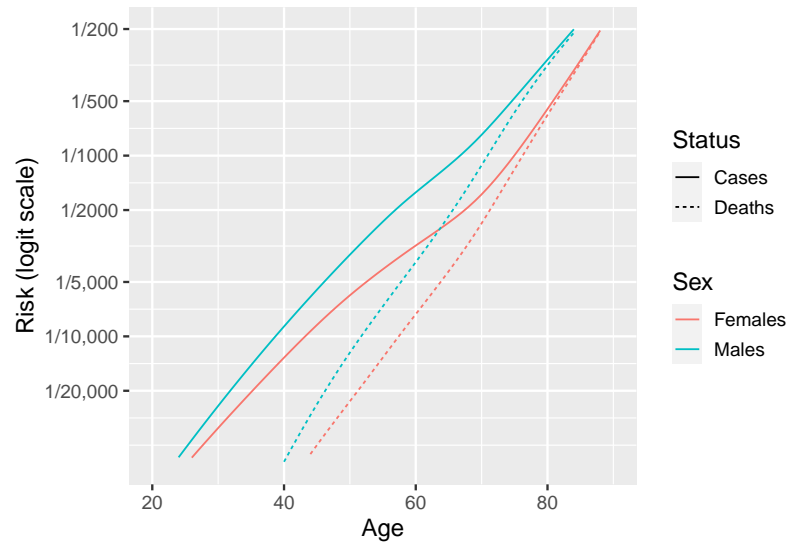
The R and Rmarkdown scripts used to generate this article, which include the code used to derive variables, will be made available with this manuscript.

## References

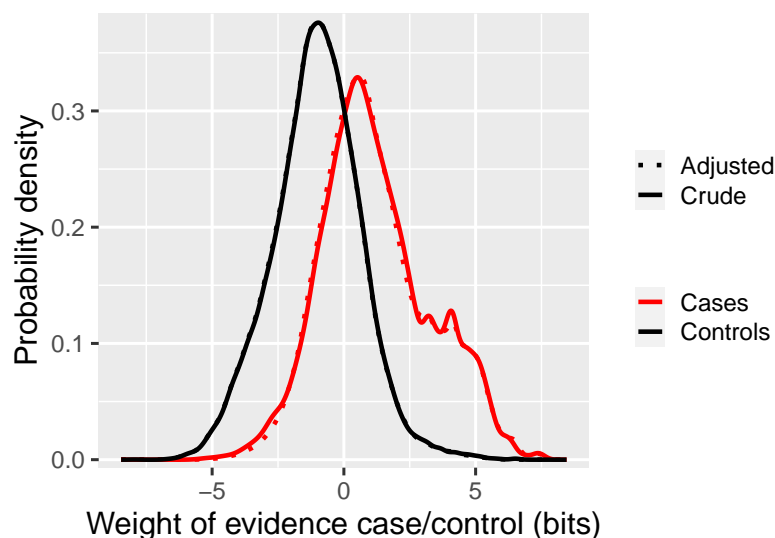
1. Docherty AB, Harrison EM, Green CA, Hardwick HE, Pius R, Norman L, et al. Features of 16,749 hospitalised UK patients with COVID-19 using the ISARIC WHO Clinical Characterisation Protocol. medRxiv. 2020; 2020.04.23.20076042. doi:10.1101/2020.04.23.20076042
2. McGoogan C, Steafel E. Why are young, healthy people dying of coronavirus? The symptoms to look out for. The Telegraph. <https://www.telegraph.co.uk/health-fitness/body/why-young-healthy-people-dying-coronavirus/>; 2020.
3. Zhou F, Yu T, Du R, Fan G, Liu Y, Liu Z, et al. Clinical course and risk factors for mortality of adult inpatients with COVID-19 in Wuhan, China: A retrospective cohort study. Lancet. 2020;395: 1054–1062. doi:10.1016/S0140-6736(20)30566-3
4. Grasselli G, Zangrillo A, Zanella A, Antonelli M, Cabrini L, Castelli A, et al. Baseline Characteristics and Outcomes of 1591 Patients Infected With SARS-CoV-2 Admitted to ICUs of the Lombardy Region, Italy. JAMA. 2020;323: 1574–1581. doi:10.1001/jama.2020.5394
5. Guan W-j, Ni Z-y, Hu Y, Liang W-h, Ou C-q, He J-x, et al. Clinical Characteristics of Coronavirus Disease 2019 in China. New England Journal of Medicine. 2020. doi:10.1056/NEJMoa2002032
6. Niedzwiedz CL, O'Donnell CA, Jani BD, Demou E, Ho FK, Celis-Morales C, et al. Ethnic and socioeconomic differences in SARS-CoV2 infection in the UK Biobank cohort study. medRxiv. 2020; 2020.04.22.20075663. doi:10.1101/2020.04.22.20075663

7. "The OpenSAFELY Collaborative", Williamson E, Walker AJ, Bhaskaran KJ, Bacon S, Bates C, et al. OpenSAFELY: Factors associated with COVID-19-related hospital death in the linked electronic health records of 17 million adult NHS patients. medRxiv. 2020; 2020.05.06.20092999. doi:10.1101/2020.05.06.20092999 455-458
8. Lakha F, Gorman DR, Mateos P. Name analysis to classify populations by ethnicity in public health: Validation of Onomap in Scotland. Public Health. 2011;125: 688–696. doi:10.1016/j.puhe.2011.05.003 459-461
9. NHS. Who's at higher risk from coronavirus - Coronavirus (COVID-19). nhs.uk. <https://www.nhs.uk/conditions/coronavirus-covid-19/people-at-higher-risk-from-coronavirus/whos-at-higher-risk-from-coronavirus/>; 462-463
10. Breslow NE, Day NE, Halvorsen KT, Prentice RL, Sabai C. Estimation of multiple relative risk functions in matched case-control studies. Am J Epidemiol. 1978;108: 299–307. doi:10.1093/oxfordjournals.aje.a112623 464-466
11. McKeigue P. Quantifying performance of a diagnostic test as the expected information for discrimination: Relation to the C-statistic. Stat Methods Med Res. 2019;28: 1841–1851. doi:10.1177/0962280218776989 467-470
12. Ho FK, Celis-Morales CA, Gray SR, Katikireddi SV, Niedzwiedz CL, Hastie C, et al. Modifiable and non-modifiable risk factors for COVID-19: Results from UK Biobank. medRxiv. 2020; 2020.04.28.20083295. doi:10.1101/2020.04.28.20083295 471-473
13. Walsh D, Buchanan D, Douglas A, Erdman J, Fischbacher C, McCartney G, et al. Increasingly Diverse: The Changing Ethnic Profiles of Scotland and Glasgow and the Implications for Population Health. Appl Spatial Analysis. 2019;12: 983–1009. doi:10.1007/s12061-018-9281-7 474-477
14. McKeigue PM, Colhoun HM. Evaluation of "stratify and shield" as a policy option for ending the COVID-19 lockdown in the UK. medRxiv. 2020; 2020.04.25.20079913. doi:10.1101/2020.04.25.20079913 478-480
15. Bunnik BAD van, Morgan ALK, Bessell P, Calder-Gerver G, Zhang F, Haynes S, et al. Segmentation and shielding of the most vulnerable members of the population as elements of an exit strategy from COVID-19 lockdown. medRxiv. 2020; 2020.05.04.20090597. doi:10.1101/2020.05.04.20090597 481-483

## Figures



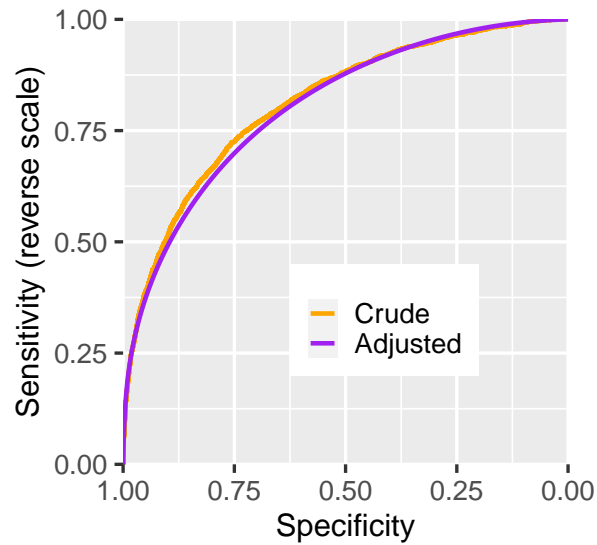
**Fig 1.** Incidence of severe and fatal COVID-19 in Scotland by age and sex: generalized additive models fitted to severe and fatal cases for males and females separately



**Fig 2.** Cross-validation of model chosen by stepwise regression using extended variable set: class-conditional distributions of weight of evidence

### Footnote for Figure 2

For each individual, the risk prediction model outputs the posterior probability of being a case, which can also be expressed as the posterior odds. Dividing the posterior odds by the prior odds gives the likelihood ratio favouring case over non-case status for an individual. The weight of evidence  $W$  is the logarithm of this ratio. The distributions of  $W$  in cases and controls in the test data are plotted in Figure 2. For a classifier, the further apart these curves are, the better the predictive performance. The expected information for discrimination  $\Lambda$  is the average of the mean of the distribution of  $W$  in cases and minus 1 times the mean of the distribution of  $W$  in controls. The distributions have been adjusted by taking a weighted average to make them mathematically consistent [11].



**Fig 3.** Cross-validation of model chosen by stepwise regression using extended variable set: receiver operating characteristic curve

### Footnote for Figure 3

The crude receiver operator characteristic (ROC) curve is computed by calculating at each value of the risk score the sensitivity and specificity of a classifier that uses this value as the threshold for classifying cases and non-cases. The C-statistic is the area under this curve, computed as the probability of correctly classifying a case/noncase pair using the score, evaluated over all possible such pairs in the dataset.

497  
498  
499  
500  
501  
502

## Tables

**Table 1.** Univariate associations of severe disease with demographic factors

	Controls (19670)	Cases (2755)	Rate ratio (95% CI)	<i>p</i> -value
Ethnicity based on name classification				
White	19274 (98%)	2694 (98%)		
South Asian	236 (1%)	41 (1%)	1.20 (0.84, 1.70)	0.3
Other	125 (1%)	12 (0%)	0.65 (0.35, 1.18)	0.2
SIMD quintile				
1 - most deprived	4555 (23%)	737 (27%)		
2	4291 (22%)	632 (23%)	0.89 (0.78, 1.00)	0.06
3	3494 (18%)	507 (18%)	0.85 (0.74, 0.98)	0.02
4	3518 (18%)	482 (18%)	0.78 (0.67, 0.89)	$4 \times 10^{-4}$
5 - least deprived	3670 (19%)	394 (14%)	0.54 (0.46, 0.64)	$5 \times 10^{-14}$
Care home	971 (5%)	836 (30%)	16.2 (13.9, 18.8)	$9 \times 10^{-289}$
Ethnicity based on Scottish Morbidity Record				
White	13905 (99%)	2302 (98%)		
South Asian	82 (1%)	22 (1%)	1.24 (0.73, 2.11)	0.4
Black	20 (0%)	4 (0%)	1.26 (0.41, 3.85)	0.7
Other	90 (1%)	15 (1%)	0.95 (0.53, 1.69)	0.9



**Table 2.** Frequencies of risk factors in cases and controls, by age group

	0-39 years		40-59 years		60-74 years		75+ years	
	Controls (521)	Cases (65)	Controls (2724)	Cases (379)	Controls (5273)	Cases (686)	Controls (11152)	Cases (1625)
Care home	0 (0%)	1 (2%)	4 (0%)	14 (4%)	50 (1%)	88 (13%)	917 (8%)	733 (45%)
Any prescription	276 (53%)	53 (82%)	1879 (69%)	341 (90%)	4530 (86%)	651 (95%)	10656 (96%)	1601 (99%)
Any admission	119 (23%)	36 (55%)	984 (36%)	236 (62%)	2628 (50%)	528 (77%)	7610 (68%)	1395 (86%)
Any listed condition	43 (8%)	32 (49%)	638 (23%)	209 (55%)	2188 (41%)	478 (70%)	6591 (59%)	1331 (82%)
Diagnosis or prescription	312 (60%)	55 (85%)	2015 (74%)	357 (94%)	4670 (89%)	670 (98%)	10778 (97%)	1618 (100%)
Type 1 diabetes	5 (1%)	3 (5%)	24 (1%)	17 (4%)	35 (1%)	9 (1%)	19 (0%)	13 (1%)
Type 2 diabetes	1 (0%)	2 (3%)	170 (6%)	68 (18%)	745 (14%)	180 (26%)	1851 (17%)	360 (22%)
Other/unknown type	1 (0%)	1 (2%)	5 (0%)	4 (1%)	14 (0%)	2 (0%)	37 (0%)	11 (1%)
Ischaemic heart disease	0 (0%)	1 (2%)	99 (4%)	32 (8%)	551 (10%)	134 (20%)	2101 (19%)	413 (25%)
Other heart disease	2 (0%)	7 (11%)	129 (5%)	66 (17%)	710 (13%)	206 (30%)	3340 (30%)	797 (49%)
Asthma or chronic airway disease	36 (7%)	25 (38%)	345 (13%)	102 (27%)	1019 (19%)	246 (36%)	2505 (22%)	537 (33%)
Chronic kidney disease or transplant recipient	0 (0%)	0 (0%)	3 (0%)	18 (5%)	12 (0%)	20 (3%)	72 (1%)	36 (2%)
Neurological (except epilepsy) or dementia	3 (1%)	8 (12%)	40 (1%)	38 (10%)	162 (3%)	109 (16%)	1216 (11%)	544 (33%)
Liver disease	1 (0%)	1 (2%)	8 (0%)	13 (3%)	32 (1%)	15 (2%)	24 (0%)	12 (1%)
Immune deficiency or suppression	2 (0%)	2 (3%)	16 (1%)	12 (3%)	40 (1%)	11 (2%)	31 (0%)	10 (1%)

**Table 3.** Associations of severe disease with listed conditions over all age groups

	Controls (19670)	Cases (2755)	Univariate		Multivariate	
			Rate ratio (95% CI)	p-value	Rate ratio (95% CI)	p-value
Care home	971 (5%)	836 (30%)	16.2 (13.9, 18.8)	$9 \times 10^{-289}$	11.7 (9.9, 13.8)	$9 \times 10^{-190}$
Any prescription	17341 (88%)	2646 (96%)	3.89 (3.15, 4.80)	$1 \times 10^{-36}$	2.10 (1.68, 2.62)	$6 \times 10^{-11}$
Any admission	11341 (58%)	2195 (80%)	3.10 (2.79, 3.43)	$5 \times 10^{-103}$	1.89 (1.68, 2.13)	$2 \times 10^{-25}$
Type 1 diabetes	83 (0%)	42 (2%)	4.26 (2.90, 6.24)	$1 \times 10^{-13}$	2.19 (1.41, 3.42)	$5 \times 10^{-4}$
Type 2 diabetes	2767 (14%)	610 (22%)	1.83 (1.65, 2.02)	$2 \times 10^{-31}$	1.62 (1.44, 1.81)	$1 \times 10^{-16}$
Other/unknown type	57 (0%)	18 (1%)	2.59 (1.52, 4.43)	$5 \times 10^{-4}$	1.70 (0.91, 3.19)	0.1
Ischaemic heart disease	2751 (14%)	580 (21%)	1.63 (1.47, 1.81)	$1 \times 10^{-19}$	1.10 (0.97, 1.24)	0.1
Other heart disease	4181 (21%)	1076 (39%)	2.51 (2.29, 2.75)	$2 \times 10^{-86}$	1.39 (1.24, 1.55)	$1 \times 10^{-8}$
Asthma or chronic airway disease	3905 (20%)	910 (33%)	2.03 (1.85, 2.22)	$3 \times 10^{-54}$	1.54 (1.39, 1.71)	$1 \times 10^{-16}$
Chronic kidney disease or transplant recipient	87 (0%)	74 (3%)	6.0 (4.4, 8.3)	$5 \times 10^{-28}$	4.27 (2.94, 6.21)	$3 \times 10^{-14}$
Neurological (except epilepsy) or dementia	1421 (7%)	699 (25%)	4.79 (4.28, 5.35)	$5 \times 10^{-167}$	1.98 (1.73, 2.27)	$3 \times 10^{-23}$
Liver disease	65 (0%)	41 (1%)	4.82 (3.23, 7.20)	$1 \times 10^{-14}$	2.17 (1.38, 3.40)	$7 \times 10^{-4}$
Immune deficiency or suppression	89 (0%)	35 (1%)	2.88 (1.94, 4.29)	$2 \times 10^{-7}$	1.30 (0.82, 2.08)	0.3

**Table 4.** Proportions of fatal cases and matched controls without and with a dispensed prescription or hospital diagnosis, by age group

	Controls	Fatal cases
<b>Age &lt;60</b>		
No scrip or diagnosis	946 (27%)	4 (3%)
Scrip or diagnosis	2616 (73%)	123 (97%)
<b>Age 60-74</b>		
No scrip or diagnosis	611 (11%)	8 (2%)
Scrip or diagnosis	4908 (89%)	432 (98%)
<b>Age 75+</b>		
No scrip or diagnosis	375 (3%)	6 (0%)
Scrip or diagnosis	10872 (97%)	1524 (100%)

**Table 5.** Associations of severe disease with hospital diagnoses in last 5 years, in those without any listed condition

	Univariate				Multivariate			
	Controls (10210)	Cases (705)	Rate ratio (95% CI)	p-value	Rate ratio (95% CI)	p-value	Rate ratio (95% CI)	p-value
Ch.1 infectious	290 (3%)	59 (8%)	3.49 (2.35, 5.17)	$5 \times 10^{-10}$	1.56 (0.99, 2.45)		1.56 (0.99, 2.45)	0.05
Ch.2 Neoplasms	657 (6%)	80 (11%)	2.29 (1.69, 3.10)	$7 \times 10^{-8}$	1.81 (1.28, 2.55)		1.81 (1.28, 2.55)	$8 \times 10^{-4}$
Ch.3 blood	155 (2%)	28 (4%)	3.54 (2.07, 6.06)	$4 \times 10^{-6}$	1.58 (0.86, 2.90)		1.58 (0.86, 2.90)	0.1
Ch.4 Endocrine	216 (2%)	35 (5%)	2.65 (1.66, 4.22)	$4 \times 10^{-5}$	1.09 (0.64, 1.87)		1.09 (0.64, 1.87)	0.7
Ch.5 Mental	264 (3%)	69 (10%)	3.68 (2.54, 5.35)	$7 \times 10^{-12}$	1.99 (1.28, 3.09)		1.99 (1.28, 3.09)	0.002
Ch.6 nervous	144 (1%)	19 (3%)	2.09 (1.15, 3.77)	0.01	1.09 (0.56, 2.12)		1.09 (0.56, 2.12)	0.8
Ch.7 eye	911 (9%)	57 (8%)	1.04 (0.74, 1.46)	0.8	0.99 (0.68, 1.43)		0.99 (0.68, 1.43)	0.9
Ch.8 ear	43 (0%)	6 (1%)	2.33 (0.84, 6.43)	0.1	1.65 (0.55, 4.93)		1.65 (0.55, 4.93)	0.4
Ch.9 circulatory	291 (3%)	39 (6%)	2.00 (1.31, 3.06)	0.001	1.10 (0.67, 1.79)		1.10 (0.67, 1.79)	0.7
Ch.10 respiratory	226 (2%)	49 (7%)	3.82 (2.52, 5.79)	$3 \times 10^{-10}$	2.10 (1.32, 3.36)		2.10 (1.32, 3.36)	0.002
Ch.11 digestive	1177 (12%)	132 (19%)	1.99 (1.57, 2.53)	$2 \times 10^{-8}$	1.32 (1.00, 1.73)		1.32 (1.00, 1.73)	0.05
Ch.12 skin	208 (2%)	26 (4%)	1.80 (1.11, 2.91)	0.02	0.78 (0.44, 1.38)		0.78 (0.44, 1.38)	0.4
Ch.13 musculoskeletal	794 (8%)	92 (13%)	1.96 (1.49, 2.59)	$2 \times 10^{-6}$	1.34 (0.98, 1.84)		1.34 (0.98, 1.84)	0.07
Ch.14 genitourinary	690 (7%)	105 (15%)	2.57 (1.94, 3.41)	$6 \times 10^{-11}$	1.57 (1.13, 2.19)		1.57 (1.13, 2.19)	0.008
Ch.15 Pregnancy	19 (0%)	0 (0%)	0.00 (0.00, Inf)	1	0.00 (0.00, Inf)		0.00 (0.00, Inf)	1
Ch.17 Congenital	20 (0%)	4 (1%)	5.8 (1.4, 23.5)	0.01	3.13 (0.64, 15.34)		3.13 (0.64, 15.34)	0.2
Ch.18 Symptoms	1041 (10%)	144 (20%)	2.40 (1.88, 3.05)	$1 \times 10^{-12}$	1.30 (0.96, 1.75)		1.30 (0.96, 1.75)	0.09
Ch.19 Injury	573 (6%)	96 (14%)	2.48 (1.85, 3.33)	$1 \times 10^{-9}$	0.84 (0.32, 2.18)		0.84 (0.32, 2.18)	0.7
Ch.20 External	639 (6%)	109 (15%)	2.70 (2.03, 3.58)	$6 \times 10^{-12}$	1.93 (0.76, 4.86)		1.93 (0.76, 4.86)	0.2
Ch.21 Health factors	1111 (11%)	122 (17%)	1.69 (1.32, 2.17)	$4 \times 10^{-5}$	0.80 (0.58, 1.09)		0.80 (0.58, 1.09)	0.2

**Table 6.** Associations of severe disease with prescribed drugs in those without any listed condition

	Univariate			Multivariate		
	Controls (10210)	Cases (705)	Rate ratio (95% CI)	<i>p</i> -value	Rate ratio (95% CI)	<i>p</i> -value
BNF 1 Gastro	3569 (35%)	380 (54%)	2.39 (1.99, 2.87)	$8 \times 10^{-21}$	1.64 (1.33, 2.03)	$3 \times 10^{-6}$
BNF 2 Cardiovascular	4809 (47%)	378 (54%)	1.60 (1.32, 1.93)	$1 \times 10^{-6}$	1.32 (1.08, 1.61)	0.006
BNF 3 Respiratory	757 (7%)	75 (11%)	1.68 (1.26, 2.24)	$4 \times 10^{-4}$	1.20 (0.87, 1.64)	0.3
BNF 4 Nervous	3836 (38%)	398 (56%)	2.39 (2.00, 2.86)	$2 \times 10^{-21}$	1.62 (1.32, 2.00)	$5 \times 10^{-6}$
BNF 5 Infections	2028 (20%)	221 (31%)	1.90 (1.56, 2.32)	$3 \times 10^{-10}$	1.35 (1.08, 1.68)	0.008
BNF 6 Endocrine	1621 (16%)	147 (21%)	1.50 (1.19, 1.89)	$5 \times 10^{-4}$	1.12 (0.87, 1.45)	0.4
BNF 7 Obstetrics	1279 (13%)	88 (12%)	0.98 (0.74, 1.29)	0.9	0.70 (0.52, 0.95)	0.02
BNF 8 Malignant	171 (2%)	18 (3%)	1.54 (0.84, 2.83)	0.2	0.96 (0.50, 1.87)	0.9
BNF 9 Nutrition	1704 (17%)	206 (29%)	2.58 (2.06, 3.23)	$1 \times 10^{-16}$	1.69 (1.32, 2.17)	$3 \times 10^{-5}$
BNF 10 Musculoskeletal	2050 (20%)	193 (27%)	1.56 (1.28, 1.91)	$1 \times 10^{-5}$	0.96 (0.76, 1.20)	0.7
BNF 11 Eye	1107 (11%)	67 (10%)	0.93 (0.68, 1.26)	0.6	0.72 (0.52, 1.00)	0.05
BNF 12 Ear	895 (9%)	67 (10%)	1.21 (0.90, 1.64)	0.2	0.82 (0.60, 1.14)	0.2
BNF 13 Skin	1976 (19%)	201 (29%)	1.59 (1.30, 1.94)	$8 \times 10^{-6}$	1.03 (0.82, 1.30)	0.8
BNF 14 Other	2007 (20%)	235 (33%)	2.09 (1.70, 2.56)	$2 \times 10^{-12}$	1.58 (1.26, 2.00)	$1 \times 10^{-4}$

**Table 7.** Prediction of severe COVID-19: cross-validation of models chosen by stepwise regression

	Cases / controls	Crude C- statistic	Adjusted C- statistic	Crude $\Lambda$ (bits)	Adjusted $\Lambda$ (bits)	Test log- likelihood (nats)
Demographic only	2724 / 19509	0.738	0.716	0.66	0.58	0.0
Demographic + listed conditions	2724 / 19509	0.793	0.777	0.96	0.89	379.7
Extended variable set	2724 / 19509	0.813	0.805	1.12	1.07	605.0

## Supplementary tables

**Table S1.** Associations of severe disease with listed conditions in those aged less than 60

	Univariate			Multivariate		
	Controls (3245)	Cases (444)	Rate ratio (95% CI)	p-value	Rate ratio (95% CI)	p-value
Care home	4 (0%)	15 (3%)	88.9 (11.6, 679.9)	$2 \times 10^{-5}$	18.5 (2.3, 150.9)	0.007
Any prescription	2155 (66%)	394 (89%)	4.41 (3.21, 6.05)	$4 \times 10^{-20}$	2.57 (1.83, 3.61)	$5 \times 10^{-8}$
Any admission	1103 (34%)	272 (61%)	3.24 (2.62, 4.01)	$1 \times 10^{-27}$	1.71 (1.33, 2.19)	$2 \times 10^{-5}$
Type 1 diabetes	29 (1%)	20 (5%)	6.0 (3.3, 11.0)	$4 \times 10^{-9}$	3.20 (1.61, 6.35)	$9 \times 10^{-4}$
Type 2 diabetes	171 (5%)	70 (16%)	3.48 (2.56, 4.74)	$2 \times 10^{-15}$	2.53 (1.79, 3.57)	$1 \times 10^{-7}$
Other/unknown type	6 (0%)	5 (1%)	7.2 (2.2, 23.6)	0.001	3.35 (0.75, 14.99)	0.1
Ischaemic heart disease	99 (3%)	33 (7%)	2.64 (1.72, 4.05)	$8 \times 10^{-6}$	0.97 (0.57, 1.62)	0.9
Other heart disease	131 (4%)	73 (16%)	4.74 (3.44, 6.54)	$3 \times 10^{-21}$	1.81 (1.21, 2.72)	0.004
Asthma or chronic airway disease	381 (12%)	127 (29%)	3.04 (2.39, 3.86)	$1 \times 10^{-19}$	1.58 (1.20, 2.08)	0.001
Chronic kidney disease or transplant recipient	3 (0%)	18 (4%)	40.1 (11.7, 137.0)	$4 \times 10^{-9}$	17.7 (4.3, 73.6)	$8 \times 10^{-5}$
Neurological (except epilepsy) or dementia	43 (1%)	46 (10%)	8.6 (5.5, 13.5)	$5 \times 10^{-21}$	4.04 (2.41, 6.78)	$1 \times 10^{-7}$
Liver disease	9 (0%)	14 (3%)	13.1 (5.5, 31.5)	$9 \times 10^{-9}$	4.30 (1.60, 11.57)	0.004
Immune deficiency or suppression	18 (1%)	14 (3%)	5.6 (2.8, 11.4)	$2 \times 10^{-6}$	0.86 (0.30, 2.52)	0.8



**Table S2.** Associations of severe disease with listed conditions in those aged 60-74 years

	Univariate			Multivariate		
	Controls (5273)	Cases (686)	Rate ratio (95% CI)	p-value	Rate ratio (95% CI)	p-value
Care home	50 (1%)	88 (13%)	19.6 (13.0, 29.7)	$5 \times 10^{-45}$	11.1 (7.0, 17.5)	$5 \times 10^{-25}$
Any prescription	4530 (86%)	651 (95%)	3.54 (2.45, 5.11)	$2 \times 10^{-11}$	1.73 (1.17, 2.55)	0.006
Any admission	2628 (50%)	528 (77%)	3.49 (2.89, 4.22)	$2 \times 10^{-38}$	2.12 (1.71, 2.63)	$9 \times 10^{-12}$
Type 1 diabetes	35 (1%)	9 (1%)	2.24 (1.07, 4.72)	0.03	0.77 (0.30, 1.97)	0.6
Type 2 diabetes	745 (14%)	180 (26%)	2.20 (1.82, 2.66)	$5 \times 10^{-16}$	1.74 (1.41, 2.15)	$2 \times 10^{-7}$
Other/unknown type	14 (0%)	2 (0%)	1.26 (0.28, 5.56)	0.8	0.75 (0.13, 4.31)	0.7
Ischaemic heart disease	551 (10%)	134 (20%)	2.10 (1.70, 2.60)	$1 \times 10^{-11}$	1.17 (0.91, 1.49)	0.2
Other heart disease	710 (13%)	206 (30%)	2.99 (2.47, 3.62)	$1 \times 10^{-29}$	1.35 (1.07, 1.70)	0.01
Asthma or chronic airway disease	1019 (19%)	246 (36%)	2.43 (2.04, 2.90)	$3 \times 10^{-23}$	1.65 (1.35, 2.01)	$8 \times 10^{-7}$
Chronic kidney disease or transplant recipient	12 (0%)	20 (3%)	12.7 (6.2, 26.1)	$4 \times 10^{-12}$	8.3 (3.6, 19.1)	$8 \times 10^{-7}$
Neurological (except epilepsy) or dementia	162 (3%)	109 (16%)	6.3 (4.8, 8.3)	$9 \times 10^{-41}$	2.66 (1.93, 3.67)	$3 \times 10^{-9}$
Liver disease	32 (1%)	15 (2%)	3.63 (1.94, 6.78)	$6 \times 10^{-5}$	1.55 (0.78, 3.10)	0.2
Immune deficiency or suppression	40 (1%)	11 (2%)	2.11 (1.08, 4.11)	0.03	0.72 (0.31, 1.65)	0.4

**Table S3.** Associations of severe disease with listed conditions in those aged 75 years and over

	Univariate			Multivariate		
	Controls (11152)	Cases (1625)	Rate ratio (95% CI)	p-value	Rate ratio (95% CI)	p-value
Care home	917 (8%)	733 (45%)	15.3 (13.0, 18.0)	$3 \times 10^{-239}$	11.8 (9.9, 14.1)	$2 \times 10^{-165}$
Any prescription	10656 (96%)	1601 (99%)	3.46 (2.26, 5.30)	$1 \times 10^{-8}$	1.32 (0.84, 2.08)	0.2
Any admission	7610 (68%)	1395 (86%)	2.79 (2.40, 3.23)	$1 \times 10^{-41}$	1.73 (1.45, 2.08)	$2 \times 10^{-9}$
Type 1 diabetes	19 (0%)	13 (1%)	5.5 (2.6, 11.6)	$7 \times 10^{-6}$	3.21 (1.27, 8.13)	0.01
Type 2 diabetes	1851 (17%)	360 (22%)	1.52 (1.34, 1.74)	$2 \times 10^{-10}$	1.44 (1.25, 1.67)	$1 \times 10^{-6}$
Other/unknown type	37 (0%)	11 (1%)	2.33 (1.18, 4.61)	0.02	1.78 (0.79, 4.01)	0.2
Ischaemic heart disease	2101 (19%)	413 (25%)	1.45 (1.28, 1.64)	$6 \times 10^{-9}$	1.09 (0.94, 1.27)	0.2
Other heart disease	3340 (30%)	797 (49%)	2.22 (1.99, 2.47)	$1 \times 10^{-46}$	1.39 (1.21, 1.60)	$2 \times 10^{-6}$
Asthma or chronic airway disease	2505 (22%)	537 (33%)	1.71 (1.53, 1.92)	$4 \times 10^{-20}$	1.47 (1.28, 1.68)	$3 \times 10^{-8}$
Chronic kidney disease or transplant recipient	72 (1%)	36 (2%)	3.36 (2.22, 5.10)	$1 \times 10^{-8}$	2.91 (1.81, 4.69)	$1 \times 10^{-5}$
Neurological (except epilepsy) or dementia	1216 (11%)	544 (33%)	4.29 (3.78, 4.87)	$4 \times 10^{-112}$	1.75 (1.50, 2.05)	$1 \times 10^{-12}$
Liver disease	24 (0%)	12 (1%)	3.60 (1.78, 7.30)	$4 \times 10^{-4}$	1.98 (0.86, 4.54)	0.1
Immune deficiency or suppression	31 (0%)	10 (1%)	2.02 (0.95, 4.32)	0.07	1.75 (0.74, 4.14)	0.2

**Table S4.** Frequencies of risk factors in cases and controls when deaths with mention of COVID-19 on the death certificate were included as cases, by age group

	0-39 years		40-59 years		60-74 years		75+ years	
	Controls (544)	Cases (68)	Controls (3252)	Cases (448)	Controls (6715)	Cases (869)	Controls (18136)	Cases (2656)
Care home	0 (0%)	1 (1%)	4 (0%)	19 (4%)	64 (1%)	160 (18%)	1877 (10%)	1467 (55%)
Any prescription	287 (53%)	56 (82%)	2254 (69%)	397 (89%)	5790 (86%)	832 (96%)	17353 (96%)	2615 (98%)
Any admission	124 (23%)	39 (57%)	1170 (36%)	278 (62%)	3321 (49%)	666 (77%)	12523 (69%)	2270 (85%)
Any listed condition	44 (8%)	34 (50%)	779 (24%)	246 (55%)	2734 (41%)	625 (72%)	10841 (60%)	2188 (82%)
Diagnosis or prescription	325 (60%)	58 (85%)	2414 (74%)	415 (93%)	5961 (89%)	851 (98%)	17538 (97%)	2641 (99%)
Type 1 diabetes	5 (1%)	3 (4%)	28 (1%)	18 (4%)	42 (1%)	11 (1%)	34 (0%)	19 (1%)
Type 2 diabetes	1 (0%)	3 (4%)	208 (6%)	78 (17%)	928 (14%)	214 (25%)	2946 (16%)	561 (21%)
Other/unknown type	1 (0%)	1 (1%)	7 (0%)	5 (1%)	19 (0%)	2 (0%)	67 (0%)	14 (1%)
Ischaemic heart disease	0 (0%)	1 (1%)	118 (4%)	38 (8%)	689 (10%)	163 (19%)	3462 (19%)	634 (24%)
Other heart disease	2 (0%)	8 (12%)	160 (5%)	80 (18%)	895 (13%)	268 (31%)	5573 (31%)	1259 (47%)
Asthma or chronic airway disease	37 (7%)	27 (40%)	424 (13%)	125 (28%)	1288 (19%)	326 (38%)	4038 (22%)	870 (33%)
Chronic kidney disease or transplant recipient	0 (0%)	0 (0%)	4 (0%)	19 (4%)	13 (0%)	22 (3%)	112 (1%)	54 (2%)
Neurological (except epilepsy) or dementia	3 (1%)	8 (12%)	47 (1%)	48 (11%)	198 (3%)	162 (19%)	2062 (11%)	999 (38%)
Liver disease	1 (0%)	1 (1%)	9 (0%)	15 (3%)	39 (1%)	20 (2%)	46 (0%)	19 (1%)
Immune deficiency or suppression	2 (0%)	2 (3%)	19 (1%)	13 (3%)	48 (1%)	16 (2%)	51 (0%)	12 (0%)

**Table S5.** Univariate associations of severe disease with hospital diagnoses by ICD subchapters in those without any listed conditions: rows retained are those with  $p < 0.001$  and at least 50 cases and controls

	Controls (10210)	Cases (705)	Rate ratio (95% CI)	p-value
Other Bacterial Diseases	63 (1%)	17 (2%)	4.60 (2.13, 9.93)	$1 \times 10^{-4}$
Bacterial And Viral Infectious Agents	102 (1%)	20 (3%)	4.49 (2.21, 9.13)	$3 \times 10^{-5}$
Malignant Neoplasms Of Digestive Organs	41 (0%)	15 (2%)	12.5 (4.5, 35.1)	$2 \times 10^{-6}$
Malignant Neuroendocrine Tumors	48 (0%)	15 (2%)	12.6 (4.5, 35.4)	$2 \times 10^{-6}$
Metabolic Disorders	115 (1%)	25 (4%)	3.68 (2.07, 6.56)	$9 \times 10^{-6}$
Mental Disorders Due To Known Physiological Conditions	89 (1%)	40 (6%)	7.1 (3.6, 13.7)	$8 \times 10^{-9}$
Other Acute Lower Respiratory Infections	81 (1%)	26 (4%)	9.2 (4.3, 19.6)	$1 \times 10^{-8}$
Other Diseases Of Intestines	560 (5%)	68 (10%)	1.97 (1.43, 2.72)	$4 \times 10^{-5}$
Acute Kidney Failure And Chronic Kidney Disease	145 (1%)	40 (6%)	5.7 (3.2, 10.3)	$4 \times 10^{-9}$
Other Diseases Of The Urinary System	311 (3%)	58 (8%)	3.01 (2.04, 4.44)	$3 \times 10^{-8}$
Symptoms And Signs Involving The Nervous And Musculoskeletal Systems	141 (1%)	45 (6%)	8.2 (4.3, 15.4)	$9 \times 10^{-11}$
Symptoms And Signs Involving Cognition, Perception, Emotional State And Behavior	104 (1%)	26 (4%)	3.46 (1.92, 6.23)	$4 \times 10^{-5}$
General Symptoms And Signs	218 (2%)	42 (6%)	4.07 (2.52, 6.58)	$1 \times 10^{-8}$
Abnormal Findings On Diagnostic Imaging And In Function Studies, Without Diagnosis	79 (1%)	17 (2%)	5.2 (2.4, 11.7)	$5 \times 10^{-5}$
Injuries To The Hip And Thigh	112 (1%)	28 (4%)	4.27 (2.20, 8.30)	$2 \times 10^{-5}$
Complications Of Surgical And Medical Care, Not Elsewhere Classified	119 (1%)	26 (4%)	3.32 (1.86, 5.94)	$5 \times 10^{-5}$
Slipping, Tripping, Stumbling And Falls	335 (3%)	58 (8%)	2.45 (1.66, 3.61)	$6 \times 10^{-6}$
Misadventures To Patients During Surgical And Medical Care	52 (1%)	14 (2%)	5.5 (2.2, 13.5)	$2 \times 10^{-4}$
Surgical And Other Medical Procedures As The Cause Of Abnormal Reaction Of The Patient, Or Of Later Complication, Without Mention Of Misadventure At The Time Of The Procedure	143 (1%)	29 (4%)	3.31 (1.89, 5.79)	$3 \times 10^{-5}$

**Table S6.** Stepwise regression: variables retained in model for severe disease

	log rate ratio	<i>p</i> -value
Care/nursing home	2.29	$8 \times 10^{-160}$
SIMD - quintile 1 as reference		
SIMD.quintile 2	0.04	0.6
SIMD.quintile 3	-0.08	0.3
SIMD.quintile 4	-0.15	0.06
SIMD.quintile 5 - least deprived	-0.29	0.002
Diabetes - non-diabetic as reference		
Type 1 diabetes	0.53	0.02
Type 2 diabetes	0.34	$1 \times 10^{-7}$
Other/unknown type	0.43	0.2
Other heart disease	0.28	$1 \times 10^{-6}$
Asthma or chronic airway disease	0.30	$2 \times 10^{-8}$
Chronic kidney disease or transplant recipient	1.34	$2 \times 10^{-12}$
Neurological (except epilepsy) or dementia	0.56	$1 \times 10^{-15}$
Liver disease	0.61	0.008
Any admission	0.49	$7 \times 10^{-15}$
Any prescription	0.26	0.03
BNF 1 Gastro	0.22	$7 \times 10^{-5}$
BNF 4 Nervous	0.34	$8 \times 10^{-9}$
BNF 5 Infections	0.16	0.002
BNF 6 Endocrine	0.15	0.006
BNF 9 Nutrition	0.30	$2 \times 10^{-8}$
BNF 11 Eye	-0.19	0.005
BNF 14 Other	0.21	$7 \times 10^{-5}$