# Characteristics associated with poor health

## Introduction

In the 2021 census the population were asked 'How is your health in general?' and could rate their health on a scale of 'Very good', 'Good', 'Fair', 'Poor', or 'Very poor'. Using the census 2021 data, a model was developed which gives a probability of someone reporting being in poor or very poor health, based on their characteristics such as employment status, ethnicity, age etc[1].

## What does this analysis tell us?

Multiple logistic regression is a statistical modelling technique for quantifying the strength of association between the occurrence of an event (e.g. poor health), and a set of characteristics[2]. In this case the technique:

1. identified the best set of variables, from those available, that were most associated with self-reported poor health,
2. produced a model using these variables which *predicts* the probability of self-reporting poor health, and
3. shows which characteristics were associated with a higher probability of poor health.

It is important to note that just because variables are associated or correlated with self-reported poor health, this does not necessarily imply that one causes the other i.e. correlation does not imply causation.

## Headlines

When other variables are controlled for, the analysis produced the following highlights[3];

- For individuals above middle age, those with 'Portuguese or Madeiran' ethnicity were the most likely to report poor health compared to other ethnic backgrounds
- For people over the age of 65, those with 'White Other' ethnicity were the second most likely to report poor health
- Individuals with 'Black', 'Asian', 'Mixed', and 'White British' ethnicities showed no significant difference in reporting poor health to those with 'White Jersey' ethnicity
- Those reporting 'Other' or 'Bisexual' sexual orientation were around three to four times more likely to report poor health than those who reported being 'Straight' or 'Gay'
- Those living in social rental accommodation were around three times more likely to report poor health than those in owner-occupied accommodation
- Those living in qualified rental or other[4] types of accommodation (including non-qualified rental), were around twice as likely to report poor health than those in owner-occupied accommodation
- Someone aged 65 and retired was ten times more likely to report poor health than someone aged 40 working in a higher skilled non-office job or non-routine office job

---

[1] The 2021 Census also asked a question about whether respondents had long term health conditions and whether they affected their activities of daily living – this analysis focuses on the more general concept of self-reported health rather than disabilities

[2] See 'Appendix 3: Methodology' for detailed description of steps taken

[3] For each example, characteristics not specified are controlled for by using the reference characteristics listed in 'Appendix 2: Definitions', with the exception of employment status for over 65's where this was set as 'Retired'.

[4] See 'Appendix 2: Definitions' for definition of 'other' types of Tenure

## The *best set* of variables that were associated with poor health

A significant association was seen between whether someone reported being in poor health and a number of census variables – including sex, level of education, marital status, and year residency began – when each variable was looked at separately. For a list of variables that were associated with self-reported poor health, when analysed separately, see Table 1 in Appendix 1.

Logistic regression identifies the best <u>set</u> of variables that were associated with poor health. So even though, for example, level of education and occupation were each *individually* associated with poor health, the analysis showed that only occupation was needed to achieve the best possible prediction of the probability of poor health. Including both variables didn't improve the prediction, because the two variables were inter-related.

The logistic regression found that the following <u>set</u> of characteristics were most associated with someone being more likely to report poor health[5]:

**Older age / Retired**

**Unemployed, looking after the home, or off work due to sickness**

**Sexual orientation reported as 'Bisexual' or 'Other'**

**Tenures other than owner-occupied, including social, qualified and non-qualified rent**

**Commute methods of passengers in a vehicle or working from home**

**Households made up of a single adult, or a single parent, and those living in communal establishments**

**Manual, routine, or lower-skilled non-office jobs**

**'Portuguese or Madeiran' ethnicity and above middle age**

**Pensioners with 'White Other' ethnicity**

---

[5] *For categorical variables (all variables excluding age) this is in comparison to reference groups in the logistic regression. See 'Appendix 2: Definitions' for more details on each variable and groups.*

# Predicting the probability of self-reported poor health for different characteristics
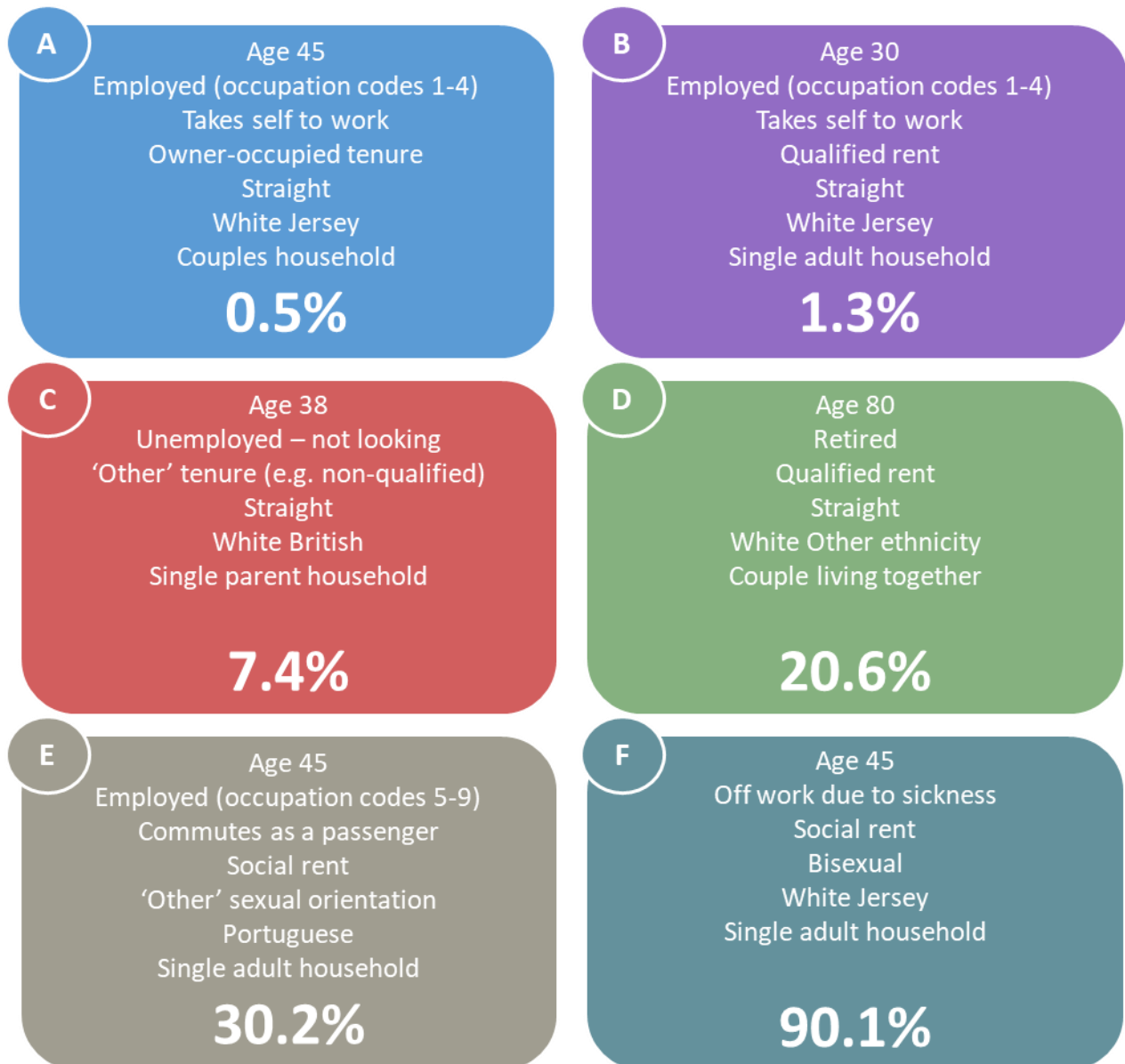
Across the whole adult population (aged 16 or over), around one in twenty (4.8%) reported having 'Poor' or 'Very poor' health in the 2021 Jersey census.

The logistic regression model enables a calculation of the probability of reporting poor health, and how it changes according to various characteristics.

Six case studies are shown below, along with the probability of someone with those characteristics reporting poor health, as predicted by the regression model.

Where this probability is greater than the average across Jersey (4.8%), individuals with this set of characteristics were *more* likely to report being in poor health than the general population. Where the probability is lower than 4.8%, individuals with that set of characteristics were *less* likely to report being in poor health than the general population.

*Figure 1: Probabilities of reporting poor health with given characteristics (see Appendix 2 for definitions)*

**A**
Age 45
Employed (occupation codes 1-4)
Takes self to work
Owner-occupied tenure
Straight
White Jersey
Couples household

**0.5%**

**B**
Age 30
Employed (occupation codes 1-4)
Takes self to work
Qualified rent
Straight
White Jersey
Single adult household

**1.3%**

**C**
Age 38
Unemployed – not looking
'Other' tenure (e.g. non-qualified)
Straight
White British
Single parent household

**7.4%**

**D**
Age 80
Retired
Qualified rent
Straight
White Other ethnicity
Couple living together

**20.6%**

**E**
Age 45
Employed (occupation codes 5-9)
Commutes as a passenger
Social rent
'Other' sexual orientation
Portuguese
Single adult household

**30.2%**

**F**
Age 45
Off work due to sickness
Social rent
Bisexual
White Jersey
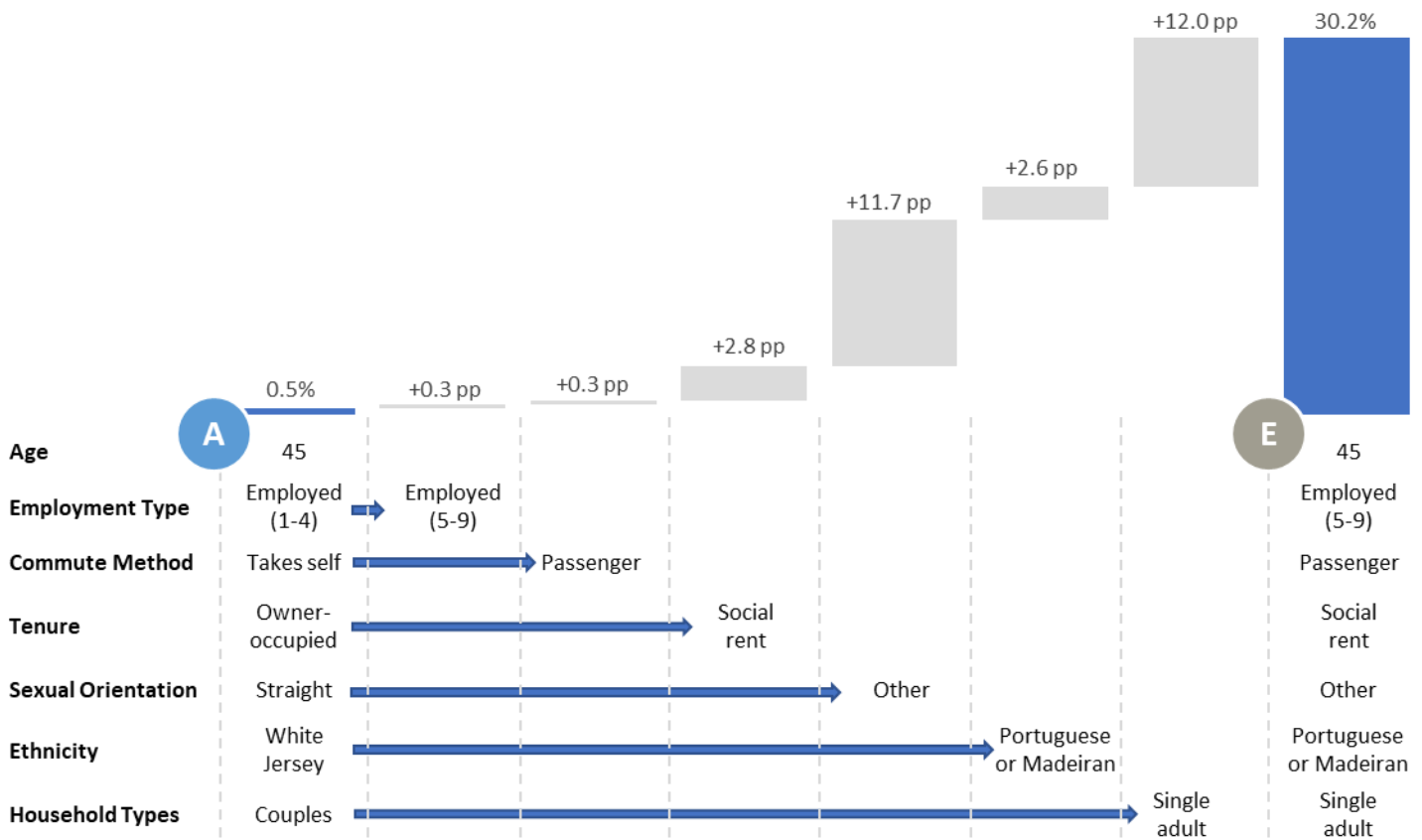Single adult household

**90.1%**

## Effect of changing characteristics

Figure 2 shows how the probability of reporting poor health can increase when characteristics are changed (to characteristics associated with a higher likelihood of reporting poor health). By making the change one by one it can be seen how much each is contributing to the difference in probability between examples A and E in Figure 1.

For example, Case study A is employed in an occupation code 1-4 job[6]; commutes by taking themselves (driving, cycling, walking etc.); their tenure is 'Owner-occupied'; sexual orientation is 'Straight'; ethnicity is 'White Jersey'; and household type is 'Couples'. If all characteristics are kept the same but employment type is changed to employed in a job with an occupation code 5-9, this increases their probability of reporting poor health by 0.3 percentage points (pp). If commute method is then changed to being a 'Passenger' (e.g. in a car, bus, taxi etc.), this increases their probability of reporting poor health by an additional 0.3 percentage points (pp), and so on until the characteristics are equivalent to Case study E.

*Figure 2: Waterfall chart showing <u>probability of reporting poor health</u> as individual characteristics are changed*
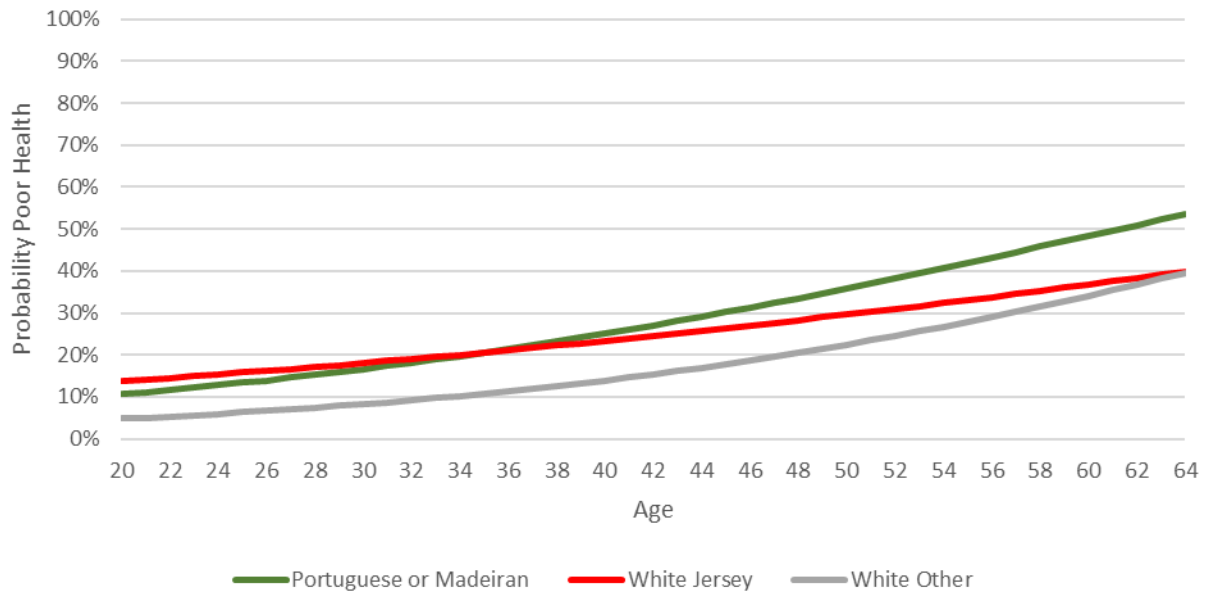


## Interaction effects

The model used an interaction effect between age and ethnicity. Interactions can be used in regression to test for the joint effect of two or more predictor variables (in this case, age and ethnicity) on an outcome variable (in this case, self-reported poor health) – for example ethnicity may have one particular association with the probability of reporting poor health for younger ages, and a different direction or strength of association for older ages.

Figure 3 shows how age interacted with ethnicity, keeping the other characteristics the same as Case study E: – employed in a job with an occupation code 5-9, commuting as a 'Passenger', 'Social rent' tenure, 'Other' sexual

---

[6] *Occupation codes 1-4 represent higher skilled non-office jobs or non-routine office jobs, see Appendix 2: Definitions for full description*
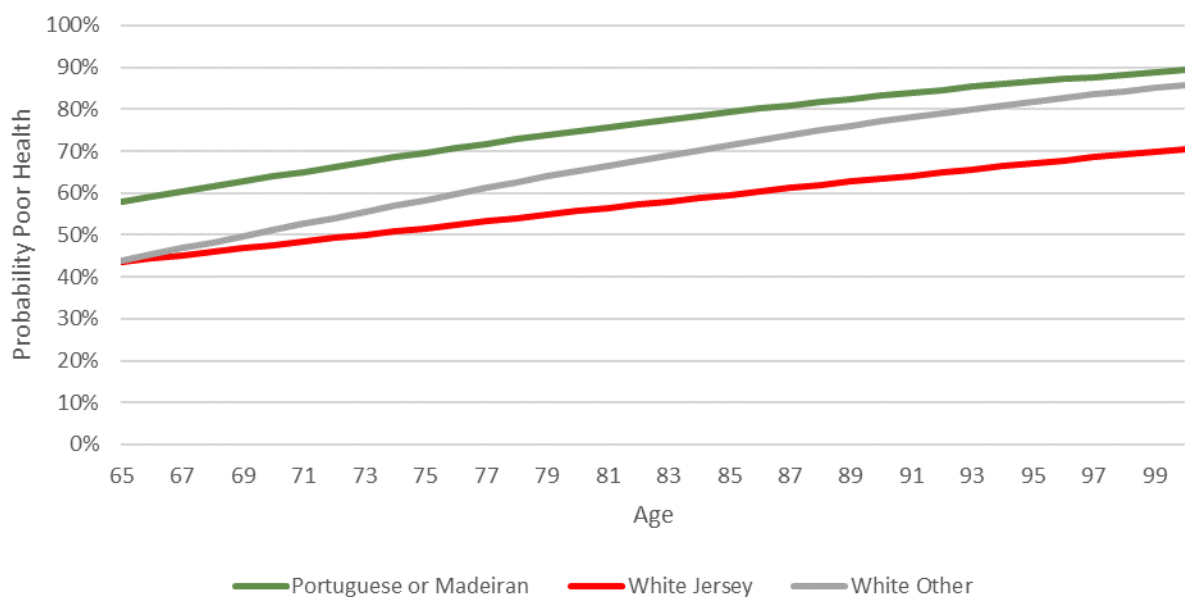
orientation, and 'Single adult' household. The probability of reporting poor health for those with 'Portuguese or Madeiran' ethnicity was higher than for those with 'White Jersey' ethnicity when aged over 36 years, and the difference increased up to age 65 years. For those with 'White Other' ethnicity, the probability of reporting poor health was lower than those with 'White Jersey' ethnicity throughout most of the working age group.

*Figure 3: Probability of poor health by age and ethnicity[7] with characteristics of 'Employed (5-9)', commuting as a 'Passenger', tenure of 'Social rent', 'Other' sexual orientation, and 'Single adult' household*



A different picture was seen for those above working age. Figure 4 shows the probability of poor health for an equivalent individual to Case study E[8]. For this age group, while those with 'Portuguese or Madeiran' ethnicity continued to have a higher probability of reporting poor health, those with 'White Other' ethnicities also had a higher probability of reporting poor health than those with 'White Jersey' ethnicity.

*Figure 4: Probability of poor health by age and ethnicity with characteristics of 'Retired', tenure of 'Social rent', 'Other' sexual orientation, and 'Single pensioner' household*



---

[7] *Ethnicity groups not shown were not significantly different to 'White Jersey' in the model*

[8] *Someone aged 65+ cannot be in the 'Single adult' group and so this needs to be changed to 'Single pensioner', it is also more likely that someone over 65+ will be retired*

# Characteristics which were associated with a higher probability of poor health

Figures 5a to 5e illustrate, for each variable, which specific characteristics were associated with a higher probability of poor health, and which were associated with the lowest probability of poor health. To enable comparison, a control set of characteristics was specified and one variable adjusted each time to show how the probability changed.

The control set of characteristics was: someone living as a couple in owner-occupied accommodation, working in an occupation 1-4 role, who takes themselves to work, and is of 'White Jersey' ethnicity. An arbitrary control age of 50 years-old was chosen.

Figure 5a shows that probability of reporting poor health was lowest for those employed in occupation codes 1-4[9] who take themselves to work, and highest for those who were unemployed and not looking for work[10].

*Figure 5a: The association between employment type and commute method, and probability of reporting poor health[11]*
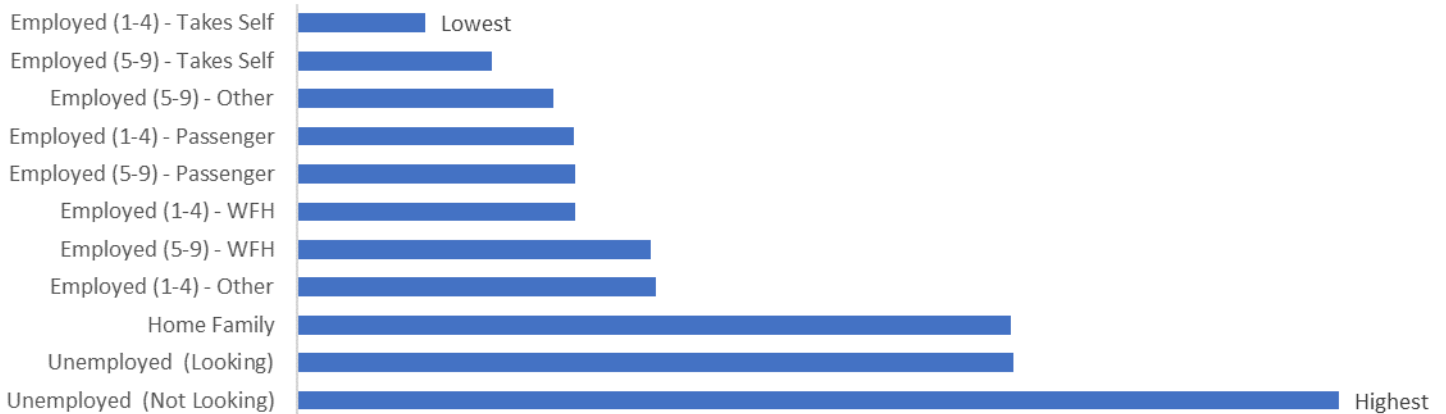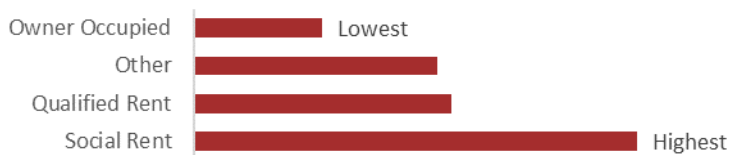


Figure 5b shows that the probability of reporting poor health was lowest for those in owner-occupied accommodation, and highest for those in social rental accommodation.

Figure 5b: The association between tenure and probability of reporting poor health[11]



---

[9] *Occupation codes 1-4 represent higher skilled non-office jobs or non-routine office jobs, see 'Appendix 2: Definitions' for full description*

[10] *Employment type of off work due to sickness gave the highest probability of poor health, as would be expected, but has been excluded from Figure 5a to analyse other characteristics more closely (as including would change the scale notably)*

[11] *See 'Appendix 2: Definitions' for more details on each variable and groups*

Figure 5c shows that the probability of reporting poor health was lowest for those with 'Black' ethnicity[12], and highest for 'Portuguese or Madeiran' ethnicity. This order holds true when age is set at 50 years old, however, when age is changed to 70 (and employment type to 'Retired') those with 'White Other' ethnicity moved to having the second highest probability of reporting poor health, as opposed to fifth in the control example.

Figure 5c: The association between ethnicity and probability of reporting poor health[13]



Figure 5d shows that the probability of reporting poor health was lowest for those reporting sexual orientation as 'Straight', and highest for those who reported 'Other'.

Figure 5d: The association between sexual orientation and probability of reporting poor health[13]
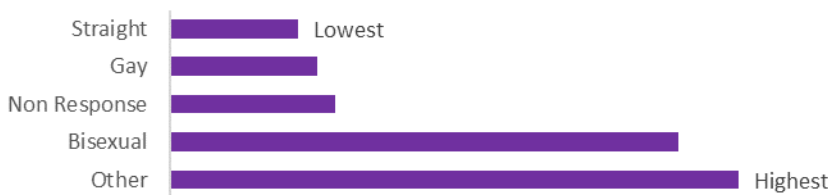


Figure 5e shows that the probability of reporting poor health was lowest for those living in a 'Couples' household, and highest for those living as a 'Single Adult'[14].

Figure 5e: The association between household type and probability of reporting poor health[13]



---

[12] Although the lowest probability, 'Black' ethnicty along with 'Asian', 'Mixed', and 'White British', were not considered significantly different to 'White Jersey' in the model

[13] See 'Appendix 2: Definitions' for more details on each variable and group

[14] Note that household type of 'Single Pensioner' has been excluded from Figure 5e as this comparison is controlling for age being 50 years-old

# Appendix 1: Census variables associated with self-reported poor health

Table 1 shows the variables significantly associated with self-reported poor health when looked at separately[15].

Many of the variables are interrelated, for example age is likely to impact on employment type, occupation, education, marital status and so on.

It is also important to note that just because variables are associated or correlated with self-reported poor health, this does not necessarily imply that one causes the other i.e. correlation does not imply causation.

*Table 1: Variables significantly associated with self-reported poor health when analysed separately*

| Variable[16] | Significant Association[17] | Strength of Association[18] |
|---|---|---|
| Employment Type / Commute Method[19] | Yes | Moderate |
| Age | Yes | Low |
| Tenure | Yes | Low |
| Education | Yes | Low |
| Year residency began | Yes | Low |
| Household Type | Yes | Low |
| Marital Status | Yes | Low |
| Sexual Orientation | Yes | Little |
| Ethnicity | Yes | Little |
| Sex | Yes | Little |

---

[15] *Using Chi-squared and Cramér's V statistics to test the significance of the relationship between each variable independently with self-reported poor health*

[16] *For variables included in the final model full descriptions of these can be seen in 'Appendix 2: Definitions'*

[17] *P-values are all very small (<0.01) indicating a significant association at a 99% confidence level*

[18] *Using Cramér's V statistic with thresholds: >0.7 – 'Very strong', >0.5 – 'High', > 0.3 – 'Moderate', > 0.1 – 'Low', 'Little' otherwise*

[19] *Variables were grouped because commute method is only applicable when someone is employed*

# Appendix 2: Definitions

*Table 2: Variable and Group Definitions. **Reference characteristics** (the 'control') are highlighted in **bold***

| Variable | Group / Characteristic | Definition |
|---|---|---|
| Employment Type | **Employed (1-4)** | Higher-skilled non-office job or non-routine office job based on occupation (SOC2010) code major groups - 1: Managers, directors and senior officials / 2: Professional occupations / 3: Associate professional and technical occupations / 4: Administrative and secretarial occupations |
| | Employed (5-9) | Manual, routine, or lower-skilled non-office jobs based on occupation (SOC2010) code major groups – 5: Skilled trades occupations / 6: Caring, Leisure and Other Service Occupations / 7: Sales and Customer Service Operatives / 8: Process, Plant and Machine Operatives / 9: Elementary Occupations |
| | In Education / Other | In education in any form (full time / part time / unemployed also in education) |
| | Home Family | Looking after home and / or family |
| | Unemployed - Looking | Unemployed looking for a job |
| | Unemployed - Not Looking | Unemployed not looking for a job |
| | Retired | Retired from paid work |
| | Sickness | Unable to work because of long-term sickness or disability |
| Commute Method (Only applicable when employed) | **Takes Self** | Private car as driver (with or without passengers), motorcycle or moped, cycle or electric bicycle, walk |
| | Passenger | Private car as passenger, bus, taxi |
| | WFH | Work mainly from home |
| | Other | Other methods or not applicable, census question required respondent to write in |
| Tenure | **Owner-Occupied** | Owned by the occupiers |
| | Social Rent | Social housing rent ('Andium homes' previously States housing, housing trust and parish rent) |
| | Other | Staff or service accommodation, registered lodging house, lodger paying rent in private household, other non-qualified accommodation |
| | Qualified Rent | Qualified private rent |
| Age | N/A | Age of individual (continuous variable) |
| Sexual Orientation | **Straight** | Straight / Heterosexual |
| | Gay | Gay or Lesbian |
| | Non-Response | This question in the census is voluntary so some individuals chose not to answer |
| | Bisexual | Bisexual |
| | Other | Other sexual orientation, census question required respondent to write in |
| Ethnicity | **White Jersey** | White Jersey |
| | White British | White British |
| | Portuguese or Madeiran | White Portuguese or Madeiran |
| | White Other | Irish, French, Polish, Romanian, South African, Other White background |
| | Asian | Indian, Thai, Chinese, Other Asian background |
| | Mixed | Asian and Black, Black and White, White and Asian, Other Mixed background |
| | Black | African, Caribbean, Other Black background |
| Household Type | **Couples** | Couple with adult (not dependent) children, couple with dependent children, adult couple, couple with one pensioner, couple pensioners |
| | Single Parents | Single parent with dependent or adult (not dependent) children |
| | Single Adult | Single adult |
| | Single Pensioner | Single pensioner |
| | Communal | Communal establishments (care home, children's home, hostel, detention centre, hotel, staff accommodation) |
| | Other | Any other household type e.g. two or more unrelated persons |

# Appendix 3: Methodology

**Logistic regression**

To explore the association between poor health and individuals' characteristics, a logistic regression was used. Logistic regression is a statistical modelling technique for quantifying the strength of association between the occurrence of an event (e.g. poor health), and a set of characteristics. The model can be used to infer the independent relationship between the event and a particular characteristic of interest while "adjusting" or "controlling" for other characteristics, which may be related to both the event and the characteristic of interest.

Stepwise regression was used in conjunction with logistic regression, this is a step-by-step iterative construction of a regression model to select independent variables to be used in a final model. It involves adding or removing potential explanatory variables in succession and testing for statistical significance after each iteration. There was also a notable element of manual iterations to establish whether adding / removing variables or using different variable groupings impacted the accuracy of predicting poor health.

For initial iterations, Alteryx software was used for data preparation and the logistic regression. For the final version of the model the logistic regression was run in RStudio software using the 'glm' and 'step' functions, so interactions between variables could be included.

Steps taken for modelling
1. The dataset was randomly split into training (70%) and test (30%) datasets to train and test the model
2. As only around 4.8% of the population had poor health an under-sampling technique was used on the training dataset[20]
3. As data was not collected in the census for some questions relating to under 16s (e.g. sexual orientation) the regression was run on those aged 16 and over
4. Some variables were grouped further or combined to explore significant relationships with poor health or reduce multicollinearity (see next page)
5. The final model also excluded the census questions on health conditions[21] as these are too closely related to the dependent variable of poor health
6. Due to undersampling on the training dataset, probabilities of poor health were then calibrated back to the census population parameters, based around Bayes Minimum Risk Theory[22]

The final list of variables included in the logistic regression can be seen in 'Appendix 2: Definitions', with the addition of the interaction between Age and Ethnicity[23]. Other variables in the census 2021 dataset were explored, such as sex, level of education, marital status and year residency began, but either deemed not significant by the model (some before or after stepwise regression) or did not increase the accuracy of predicting poor health, when included in addition to the variables used in the final model.

---

[20] *Under sampling involves randomly removing records to make the distribution of poor vs not poor health equal, as this results in a more balanced accuracy of predicting poor / not poor health. If undersampling was not used the model would result in a high level of accuracy <u>overall</u> but a lower level of accuracy for predicting the outcome of poor health (recall).*
[21] *"Do you have any physical or mental health conditions or illnesses lasting or expecting to last 12 months or more?" and "Do any of your conditions or illnesses reduce your ability to carry out day-to-day activities?"*
[22] *(PDF) Calibrating Probability with Undersampling for Unbalanced Classification (researchgate.net)*
[23] *Interactions can be used to test for the joint effect of two or more predictor variables on an outcome variable. This allows us to explore how the relationships between dependent and independent variables differ by context. In this case the interaction between Age and Ethnicity improves the Akaike information criterion (AIC) in stepwise regression, therefore has been kept in the model as significant.*

**Limitations**

- The analysis was based on self-reported health status, which is a subjective rather than an objective measure of poor health. Self-reported measures are often subjective and can reflect differences between socio-demographic groups in terms of their likelihood to report having poor health even for the same objective health status.
- The exact outputs of a regression model can have slight variations according to which techniques and order steps are used. However, the insights would be expected to be broadly similar.
- For some specific characteristics where the population sizes are small, there may not be sufficient data to identify a *statistically* significant association of that characteristic with reporting poor health, but it is possible that an association exists.

**Multicollinearity**

Multicollinearity (also known as collinearity) is where explanatory (independent) variables in a regression model are highly correlated with each other. However, an important assumption of multivariate regression is that explanatory variables are *not* too highly correlated with one another, as this can affect the stability and interpretation of the regression estimates.

Several steps were taken to reduce collinearity, for example, by combining employment type and commute method, as if someone is not employed, commute method will not be applicable. Household types have also been grouped in a way that attempts not to distinguish between pensioners and non-pensioners, as this is too closely correlated with age and employment type. In the final model many of the variables are still correlated with one another, removing or grouping them further led to reduced accuracy and removed insightful results. However, when testing for multicollinearity using Variance Inflation Factors (VIF)[24] all values are below 5 without interactions[25], indicating only low or moderate correlation[26].

---

[24] *using the 'VIF' function from the 'car' package in RStudio software*
[25] *With interactions results in a high VIF for Ethnicity and the interaction term as it is contained within both, however this is to be expected and outputs between the models with and without interaction have been compared to ensure there are no unexpected swings.*
[26] *A VIF value of between 1 and 5 is considered as having moderate correlation, but not severe enough to warrant corrective measures -* *Multicollinearity in Regression Analysis: Problems, Detection, and Solutions*

# Appendix 4: Regression Outputs

*Table 3: Logistic regression ANOVA table*

| Variable | Degrees of freedom | Deviance | P-value | Significant[27] |
|---|---|---|---|---|
| Employment Type / Commute Method | 13 | 2129.1 | 0.0000 | Yes |
| Tenure | 3 | 221.0 | 0.0000 | Yes |
| Age | 1 | 160.4 | 0.0000 | Yes |
| Ethnicity | 6 | 50.0 | 0.0000 | Yes |
| Household Type | 5 | 48.3 | 0.0000 | Yes |
| Sexual Orientation | 4 | 33.4 | 0.0000 | Yes |
| Age*Ethnicity | 6 | 23.9 | 0.0005 | Yes |

*Table 4: Test data confusion matrix (1 = Poor health / 0 = Not poor health)*

| | | Actual | |
|---|---|---|---|
| | | 0 | 1 |
| Predicted | 0 | 19,260 (74.0%) | 269 (1.0%) |
| | 1 | 5,479 (21.0%) | 1,029 (4.0%) |

*Table 5: Test data model prediction accuracy metrics ( 1 = Poor health / 0 = Not poor health)*

| Accuracy | Accuracy 0 | Accuracy 1 (Recall) | Precision |
|---|---|---|---|
| 77.9% | 77.9% | 79.3% | 15.8% |

Metrics in Table 5 have been calculated using the counts in Table 4, and the formulas following Table 6.

*Table 6: Confusion matrix format*

| | | Actual | |
|---|---|---|---|
| | | Negative | Positive |
| Predicted | Negative | True Negative (TN) | False Negative (FN) |
| | Positive | False Positive (FP) | True Positive (TP) |

Accuracy is the total number of correct predictions divided by the total in the test dataset

$$Accuracy = \frac{TN + TP}{TN + FN + FP + TP}$$

Accuracy 0, or the accuracy of predicting not having reported poor health, also known as specificity and true negative rate, is the number of correct negative predictions divided by the total number of negatives

$$Accuracy\ 0\ = \frac{TN}{TN + FP}$$

Accuracy 1, or the accuracy of predicting reported poor health, also known as recall, sensitivity, or true positive rate, is the number of correct positive predictions divided by the total number of positives. This metric is what has been prioritised when creating the model as the accuracy of predicting poor health is of most interest.

$$Accuracy\ 1\ (Recall) = \frac{TP}{FN + TP}$$

Precision is the number of correct positive predictions divided by the total number of positive predictions

$$Precision\ = \frac{TP}{FP + TP}$$

---

[27] *P-values all very small (<0.01) therefore all are significant at 99% confidence level*

*Table 7: Regression Coefficient Output*

| Variable | Group / Characteristic | Estimate | Std. Error | z value | Pr(>\|z\|) | Significance |
|---|---|---|---|---|---|---|
| | (Intercept) | -3.70 | 0.21 | -17.60 | 0.00 | *** |
| Employment Status / Commute Method | Ref: Employed (1-4) - Takes Self | | | | | |
| | Employed (1-4) - Passenger | 0.77 | 0.24 | 3.21 | 0.00 | *** |
| | Employed (1-4) - WFH | 0.78 | 0.17 | 4.61 | 0.00 | *** |
| | Employed (1-4) - Other | 1.03 | 0.74 | 1.40 | 0.16 | |
| | Employed (5-9) - Takes Self | 0.42 | 0.13 | 3.22 | 0.00 | *** |
| | Employed (5-9) - Passenger | 0.78 | 0.22 | 3.59 | 0.00 | *** |
| | Employed (5-9) - WFH | 1.02 | 0.29 | 3.50 | 0.00 | *** |
| | Employed (5-9) - Other | 0.69 | 0.79 | 0.88 | 0.38 | |
| | In Education / Other | 0.44 | 0.32 | 1.40 | 0.16 | |
| | HomeFamily | 1.74 | 0.18 | 9.50 | 0.00 | *** |
| | Unemployed - looking | 1.74 | 0.22 | 7.83 | 0.00 | *** |
| | Unemployed - Not looking | 2.13 | 0.27 | 8.03 | 0.00 | *** |
| | Retired | 1.65 | 0.14 | 11.56 | 0.00 | *** |
| | Sickness | 4.12 | 0.18 | 22.68 | 0.00 | *** |
| Tenure | Ref: Owner Occupied | | | | | |
| | Social Rent | 1.25 | 0.10 | 12.35 | 0.00 | *** |
| | Other | 0.65 | 0.16 | 4.07 | 0.00 | *** |
| | QualifiedRent | 0.70 | 0.10 | 7.22 | 0.00 | *** |
| Household Type | Ref: Couples | | | | | |
| | Single Parents | 0.38 | 0.13 | 2.88 | 0.00 | *** |
| | Single Adult | 0.66 | 0.12 | 5.48 | 0.00 | *** |
| | Single Pensioner | -0.09 | 0.12 | -0.73 | 0.46 | |
| | Communal | 0.57 | 0.23 | 2.50 | 0.01 | ** |
| | Other | 0.19 | 0.11 | 1.79 | 0.07 | * |
| Sexual Orientation | Ref: Straight | | | | | |
| | Gay | 0.14 | 0.31 | 0.43 | 0.67 | |
| | Non Response | 0.25 | 0.10 | 2.51 | 0.01 | ** |
| | Bisexual | 1.39 | 0.33 | 4.24 | 0.00 | *** |
| | Other | 1.51 | 0.51 | 2.94 | 0.00 | *** |
| Age | Age | 0.03 | 0.00 | 9.31 | 0.00 | *** |
| Age* Ethnicity | Ref: White Jersey | | | | | |
| | White British | 0.00 | 0.00 | -0.64 | 0.52 | |
| | White PortMadeira | 0.02 | 0.01 | 2.49 | 0.01 | ** |
| | White Other | 0.03 | 0.01 | 3.66 | 0.00 | *** |
| | Asian | 0.02 | 0.03 | 0.59 | 0.56 | |
| | Mixed | -0.02 | 0.02 | -0.84 | 0.40 | |
| | Black | -0.01 | 0.06 | -0.18 | 0.86 | |
| Ethnicity | Ref: White Jersey | | | | | |
| | White British | -0.12 | 0.27 | -0.44 | 0.66 | |
| | White PortMadeira | -0.68 | 0.41 | -1.66 | 0.10 | * |
| | White Other | -1.67 | 0.41 | -4.06 | 0.00 | *** |
| | Asian | -1.66 | 1.34 | -1.24 | 0.21 | |
| | Mixed | 0.63 | 1.01 | 0.63 | 0.53 | |
| | Black | -1.36 | 2.89 | -0.47 | 0.64 | |

## Appendix 5: Notes

This analysis was produced by a Statistics Jersey project team funded by the Covid Recovery Fund. The Covid Recovery Insights Project team are using administrative datasets from across the Government of Jersey to produce insights on which socio-demographic groups were more affected by the covid pandemic, and therefore how best to support our community to recover from the pandemic.

This initial work, using Census 2021 data, provides context and knowledge which will feed into the wider project. Further outputs will be produced from the project team through 2023.

The 2021 Census was run during the Covid-19 pandemic; as such, a number of restrictions were in place. See www.gov.je/census for further information on the operations, methodology, and reports of the census. Each census bulletin, and the full report, include notes on quality assurance and methodology.

## Contacts

Ian Cope | Chief Statistician | chiefstatistician@gov.je
Sarah Davis | Head of Analytics and Statistics Enablement Team | s.davis2@gov.je