

The School of Mathematics



THE UNIVERSITY
of EDINBURGH

Analysis on Leptospirosis Risk Factors Based on General Linear Mixed Models

by

Yile Shi

Dissertation Presented for the Degree of
MSc in Statistics with Data Science

August 2021

Supervised by
Dr. Gail Robertson and Dr. Amy Wilson

Executive Summary

Background: Dementia, a major international public health concern, seriously affects people's lives. It is reported that 50 million people around the world had been diagnosed with dementia in 2019, and this number is estimated to reach 152 million by 2050 [1]. Currently, the cure of dementia is still absent, hence effective prevention strategies become critical to reduce the risk of dementia and lessen its burden on society.

Research question: This research aims to shed light on the potential risk factors that affects individual's cognitive function, and further generate knowledge to inform the development of lifestyle interventions for dementia risk reduction. We consider 13 factors: age, gender, country, education, drinking behaviour, smoking, obesity, physical activity, chronic disease, working status, household finance, social connection and depression.

Data: We use a subset of easySHARE dataset from Survey of Health, Ageing and Retirement in Europe, which contains 57310 subjects recorded in 2013.

Methods: We construct Bayesian networks where the response of interest is individual's cognitive score. Modified network based on Hill-climbing algorithm is selected due to best performance through 10-fold cross-validation. Further, we use the network to estimate model parameters and plot line graphs to explore the trends of conditional probabilities of cognitive score in different groups.

Results: Age, gender, country, education and depression are shown to have significant effects on individual's cognitive impairment. Generally, age and depression have negative influence where individual's cognitive impairment gets more severe as age or depression level increases. Particularly, cognitive impairment seems to aggregate for individuals over 65 years or experiencing severe depression (figure ??). A positive effect of education level on cognition is observed, where people with higher education background usually have lower dementia risk. For the binary factor country, individuals living in countries with low average gross domestic product are more likely to suffer from dementia. Besides, female groups have higher dementia risk than males.

Acknowledgments

I am sincerely grateful to the supervisors of this project, Dr. Sara Wade and Dr. Cecilia Balocchi, as well as PhD supporter Steven Soutar, for their advice and guidance. I would also like to thank SHARE for providing the data and background information used in the research.

This report uses data from SHARE Waves 1 [2], 2 [3], 3 [4], 4 [5], 5 [6], 6 [7], 7 [8], and 8 [9]. The SHARE data collection has been funded by the European Commission, DG RTD through FP5 (QLK6-CT-2001-00360), FP6 (SHARE-I3: RII-CT-2006-062193, COMPARE: CIT5-CT-2005-028857, SHARELIFE: CIT4-CT-2006-028812), FP7 (SHARE-PREP: GA N°211909, SHARE-LEAP: GA N°227822, SHARE M4: GA N°261982, DASISH: GA N°283646) and Horizon 2020 (SHARE-DEV3: GA N°676536, SHARE-COHESION: GA N°870628, SERISS: GA N°654221, SSHOC: GA N°823782, SHARE-COVID19: GA N°101015924) and by DG Employment, Social Affairs & Inclusion through VS 2015/0195, VS 2016/0135, VS 2018/0285, VS 2019/0332, and VS 2020/0313. Additional funding from the German Ministry of Education and Research, the Max Planck Society for the Advancement of Science, the U.S. National Institute on Aging (U01_AG09740-13S2, P01_AG005842, P01_AG08291, P30_AG12815, R21_AG025169, Y1-AG-4553-01, IAG_BSR06-11, OGHA_04-064, HHSN271201300071C, RAG052527A) and from various national funding sources is gratefully acknowledged (see www.share-project.org).

This report also uses data from the generated easySHARE data set [10]. The easySHARE release 8.0.0 is based on SHARE Waves 1, 2, 3, 4, 5, 6, 7 and 8.

University of Edinburgh – Own Work Declaration

This sheet must be filled in, signed and dated - your work will not be marked unless this is done.

Name: Yile Shi

Matriculation Number: s2168022

Title of work: Analysis on Dementia Risk Factors Based on Bayesian Networks

We confirm that all this work is my own except where indicated, and that We have:

- Clearly referenced/listed all sources as appropriate
- Referenced and put in inverted commas all quoted text (from books, web, etc)
- Given the sources of all pictures, data etc. that are not my own
- Not made any use of the report(s) or essay(s) of any other student(s) either past or present
- Not sought or used the help of any external professional academic agencies for the work
- Acknowledged in appropriate places any help that We have received from others (e.g. fellow students, technicians, statisticians, external sources)
- Complied with any other plagiarism criteria specified in the Course handbook

We understand that any false claim for this work will be penalised in accordance with the University regulations (<https://teaching.maths.ed.ac.uk/main/msc-students/msc-programmes/statistics/data-science/assessment/academic-misconduct>).

Signature:

A handwritten signature in black ink that reads "Yile Shi". The signature is written in a cursive, slightly slanted style.

Date: 2022/8/10

Contents

1	Introduction	2
2	Exploratory Data Analysis	3
2.1	Erroneous Data and Duplicated Data	3
2.2	Missing Data	3
2.2.1	livestk_home	3
2.2.2	disthosp	3
2.2.3	location	4
2.2.4	occupation	4
2.3	Feature Selection	6
2.3.1	Nominal Variables	7
2.3.2	Numerical Variables	8
2.3.3	Geographical Variables	10
3	Implementation	12
3.1	General Linear Mixed Model	12
3.2	Model Construction and Selection	12
3.2.1	Step 1: Initial Model	12
3.2.2	Step 2: Find the Optimal Fixed Structure	12
3.2.3	Step 3: Find the Optimal Random Structure	13
4	Results	15
4.1	Age	15
4.2	Education	15
4.3	Depression	16
4.4	Gender	16
4.5	Country	16
5	Conclusions	17
	Appendix	20
	Word Count	21

1 Introduction

Dementia, a clinical state characterized by loss of function in multiple cognitive domains, becomes a serious public health concern worldwide [11]. Someone develops dementia every 3 seconds and current annual cost of dementia is estimated at 1 million dollars, which is set double in 2030 [1]. Because of the absence of effective treatment, prevention strategy to reduce dementia risk becomes an active research topic. This report aims to contribute to this topic by identifying potential risk factors behind dementia diagnosis. We expect that a better understanding of the influence of factors on cognitive impairment can help to inform the development of lifestyle interventions for dementia risk reduction.

We apply exploratory data analysis on easySHARE dataset from Survey of Health, Ageing and Retirement in Europe (SHARE) and extract a sample of 57310 individuals recorded in 2013. Based on the subset, we build Bayesian networks and select the one with the best performance. Our goal is to find, among the following factors, the most likely to be relevant to cognitive decline: age; gender; country; education; drinking behaviour; smoking; obesity; physical activity; chronic disease; working status; household finance; social isolation; and depression.

2 Exploratory Data Analysis

2.1 Erroneous Data and Duplicated Data

We start with going through the dataset to have some initial insight of the data and detect that there exist some fault data against our general knowledge. For example, according to the column `relationshiphh`, some sampled people are the daughters of their corresponding household heads while their genders are recoded as "Male". We fix these problems. Besides, we observe and drop some duplicated data which are possibly recorded by accident.

Note that there also exist some pairs of observations that are different in `sampleid` and `parent` but have same values in other columns. We decide to keep both of them in the dataset as we cannot tell whether these observations are duplicated or not.

2.2 Missing Data

Table 1 displays the number and the proportion (4 decimal places) of missing values in variables, arranged in descending order.

variable	count	proportion
<code>occupation</code>	343	0.3649
<code>disthosp</code>	259	0.2755
<code>livestk_home</code>	236	0.2511
<code>location</code>	236	0.2511
<code>landuse</code>	4	0.0043
<code>gender</code>	1	0.0011
<code>age</code>	1	0.0011

Table 1: Number and proportion of missing values in variables

According to the table above, we find that `landuse`, `gender` and `age` contain missing values less than 1%. Thus, we drop the corresponding observations as it doesn't lead to much loss of information of the original dataset.

As for `occupation`, `disthosp`, `livestk_home` and `location`, which contains a large amount of missing data respectively, they require exploration and discussion in more detail to determine a proper way to deal with missing data.

2.2.1 `livestk_home`

According to the data description, `livestk_home` is a binary variable indicating whether or not livestock is kept in the household of sampled person. The work of Cook et al in 2016 [?] found that the exposure to livestock could be an important risk factor for Leptospirosis, hence we might also expect the significant contribution to Leptospirosis diagnosis of this variable. However, from figure 1, we observe a significant imbalance in this column where most families of sampled people have livestock at their homes. In this case, we will not include this variable in further analysis as the imbalance probably leads to insignificant results of this variable. It becomes unnecessary to think of the missing values in this column.

2.2.2 `disthosp`

`disthosp` is the Euclidean distance from the sampled person's household to local hospital. Here, we assume that there is only one hospital in each village and people from that village only go to that hospital. Then, we use the mean distance to the local hospital of each village to impute the missing values in this column, which works for most villages. We observe that data in `disthosp` in village 12, 13 and 23 are completely missing, so we cannot obtain the mean distance to the local hospital and further fail to impute the missing distance in these villages. As a result, we drop the corresponding rows of village 12, 13 and 23.

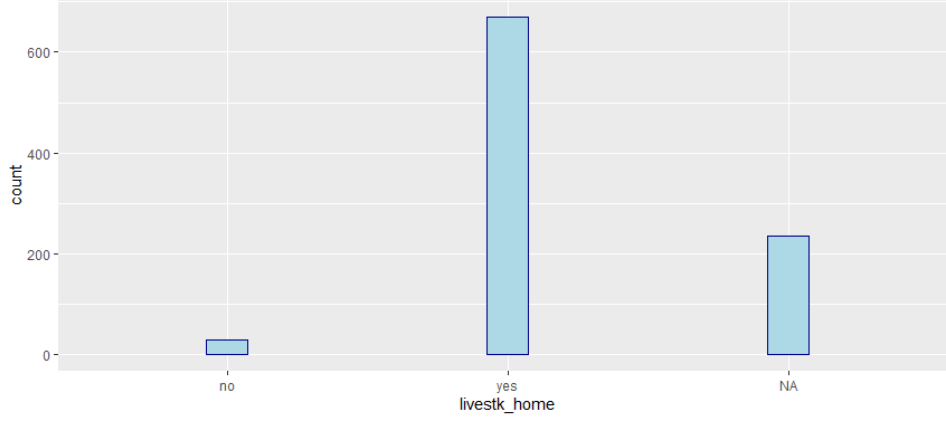


Figure 1: Distribution of `livestk_home`

2.2.3 location

`location` is the anonymised location of area where sampling was done, ranging from 1 to 19. This column has over 25% missing data. Again, we consider using the `village` variable to determine the corresponding location in the same row. However, two problems are detected:

- Similar to the case of `disthosp`, the locations of observations in some villages are completely missing, which makes it impossible to determine the correct locations.
- Some observations in the same village belong to different locations, e.g. some individuals from village 17 belong to location 9 while others belong to location 18. Without more information, we cannot decide which location they belong to.

As a result, we drop all missing values in this column instead of imputation due to these problems.

2.2.4 occupation

The column `occupation` represents the occupation of the person sampled, which could have influence on the prevalence of Leptospirosis. Specifically, [?] pointed out that people whose working places are closer to water or animals are more likely to suffer from Leptospirosis. Therefore, we may want to take this variable into consideration when modelling and deal with this column more carefully.

The occupation of a person could be associated with various aspects. We first consider the influence of individual's age. Figure 2 and 3 show the conditional distributions of `occupation` in juvenile (< 18 years old) and adult (≥ 18 years old) groups respectively.



Figure 2: Distribution of `occupation` in juvenile group

We observe very different distributions in two groups from the plots. For juveniles under 18 years old, they are most likely to be students, hence we impute the missing occupation for juvenile individuals with "Student", which is consistent with our general knowledge.

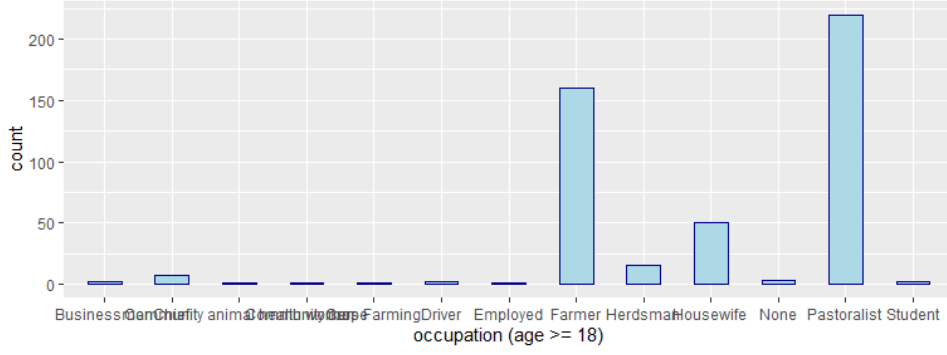


Figure 3: Distribution of `occupation` in adult group

The case for adults is more complex, as their occupations are also correlated with other variables. `landuse` denotes the characterization of the sampling site based on land use, which might be correlated with individual's occupation. Figure 4 displays the conditional distributions of adults' occupations in different `landuse` groups.

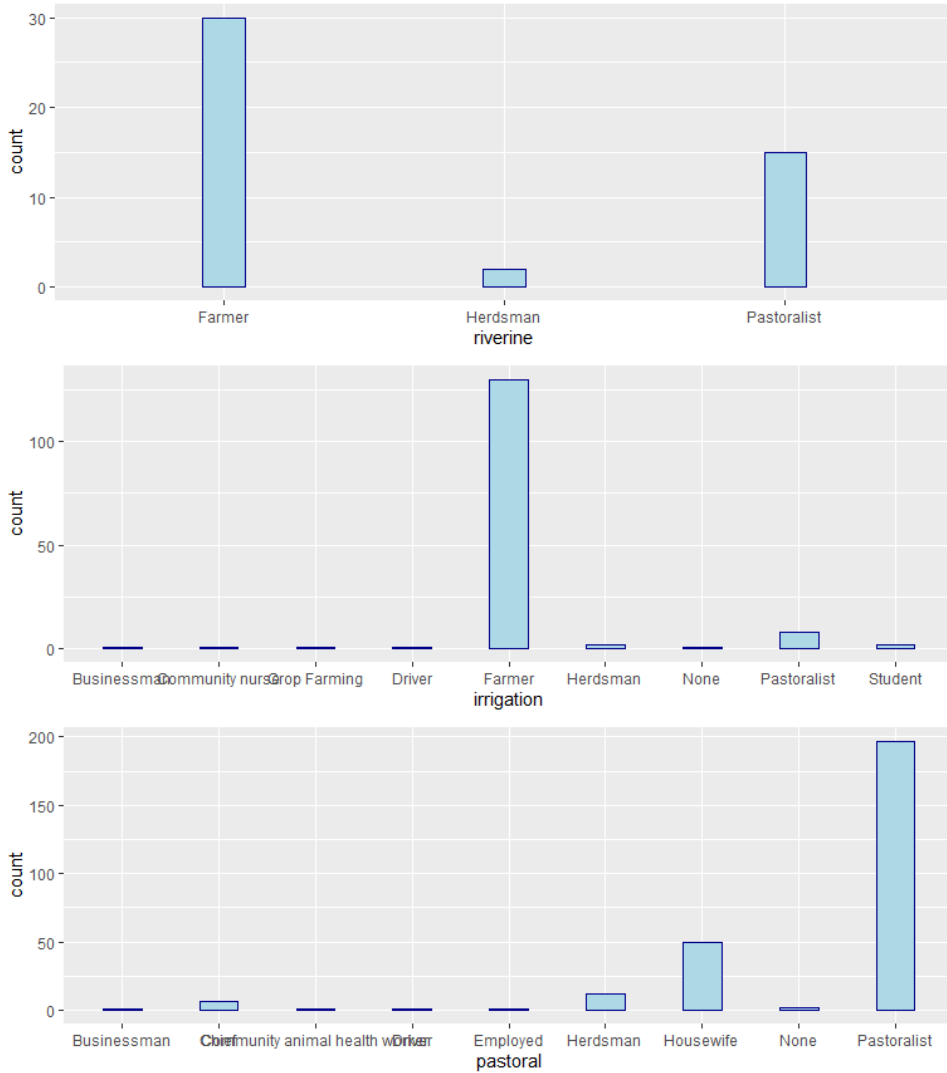


Figure 4: Conditional distributions of `occupation` in adult group

According to the plot, most sampled adults with land use as "riverine" or "irrigation" work as farmers, while individuals with land use as "pastoral" are most likely to be pastoralists. Moreover, holding the `landuse` as "pastoral", we observe significance difference among different gender and

constituency groups, as shown in Figure 5.

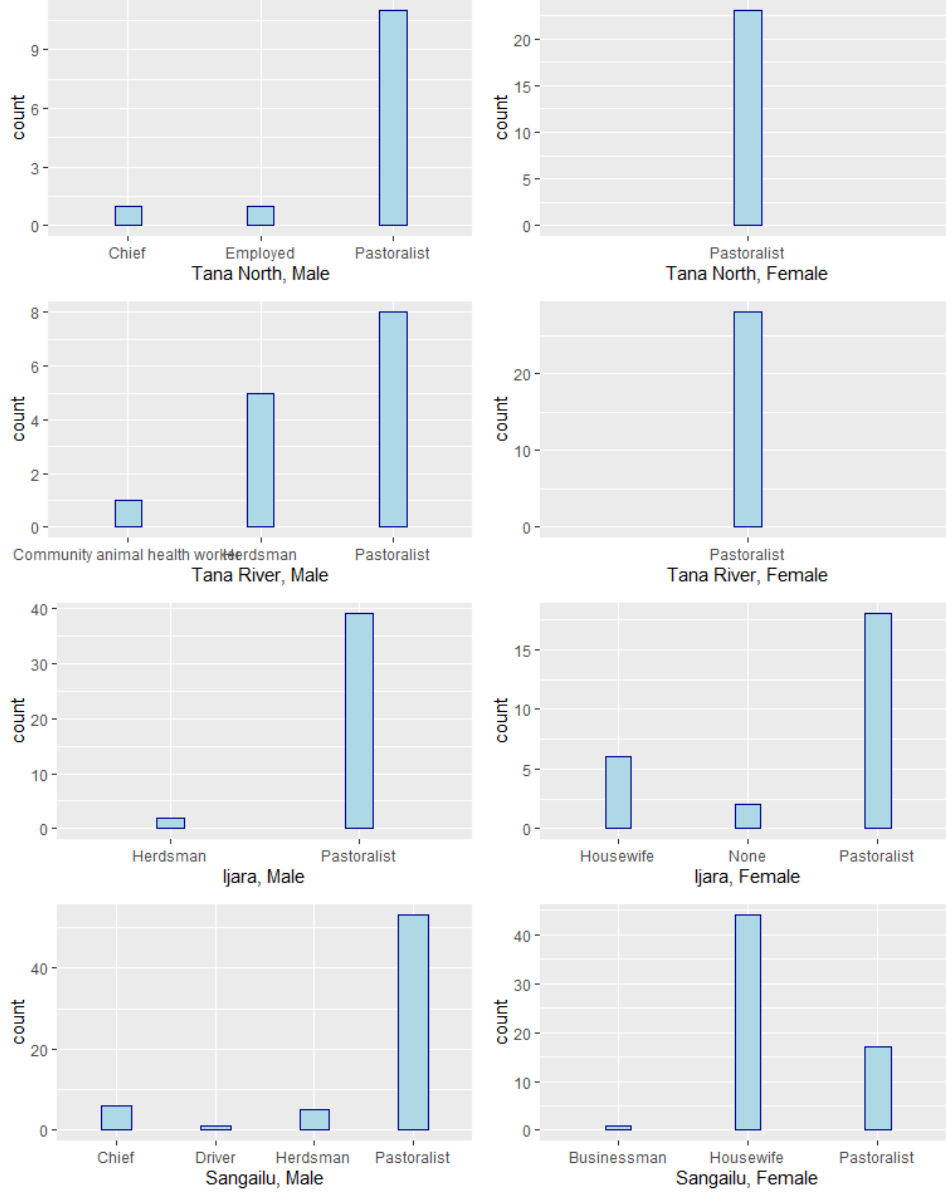


Figure 5: Conditional distributions of `occupation` in adult and "pastoral" group

In the "pastoral" stratification, male individuals across constituencies usually have their occupations as "Pastoralist", while female groups show various distributions. Females from "Tana North" and "Tana River" constituencies all work as pastoralists, but some of those from "Ijara" and "Sangailu" constituencies are housewives. Note that "Housewife" is the most popular occupation in "Sangailu" constituency, rather than "Pastoralist".

Therefore, we decide to drop observations with missing values in `occupation` instead of imputing them, as the cases are complicated and imputations probably lead to great model bias.

2.3 Feature Selection

After dealing with missing data, we now determine which columns could be used for modelling. The column `result`, which represents whether the sample is positive or negative on ELISA (enzyme-linked immunosorbent assay), is considered as the response variable in our model.

As for potential risk factors, according to the data description, identifier columns including `sampleid` and `parent` are excluded from our model, as well as `relationshiph` and `genhhid`. Besides, as mentioned before, we do not take the variable `livestk.home` into consideration. We select features from

other columns based on corresponding plots and the work of Goarant et al (2016) [?].

2.3.1 Nominal Variables

We start from factor variables `gender`, `occupation`, `landuse`, `hhoccup` and `hhgender`. The last two columns are nominal variables indicating the occupation and gender of the corresponding household head. Before any selection, we count the number of observations in each category of `occupation` and `hhoccup`. Moreover, we plot the stacked bars of `result`, grouped by these two columns respectively.

occupation		hhoccup	
category	count	category	count
Student	239	Pastoralist	345
Pastoralist	159	Farmer	247
Farmer	156	Chief	15
Housewife	36	Businessman	5
Herdsman	20	Civil Servant	4
Chief	6	Herdsman	3
Businessman	2	Crop Farming	2
Driver	2	None	2
Community animal health worker	1	Casual Labourer	1
Community nurse	1	Driver	1
Crop Farming	1		
Employed	1		
None	1		

Table 2: Number of observations in each category in `occupation` and `hhoccup`

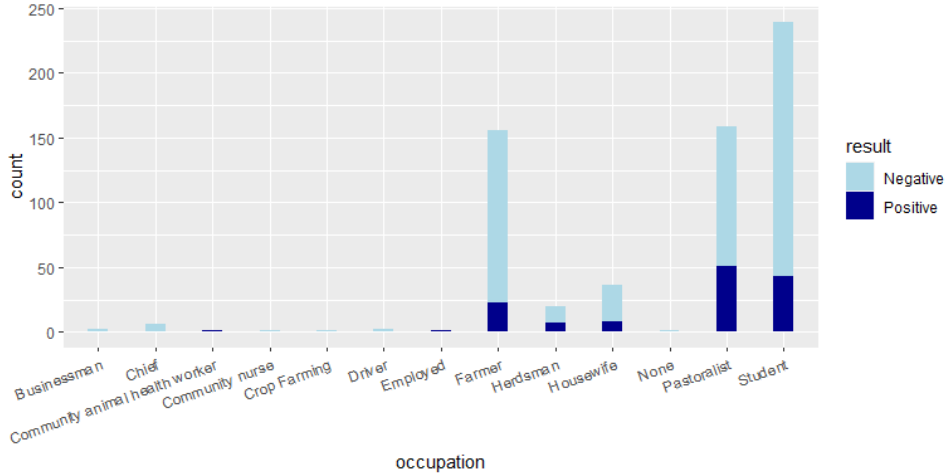


Figure 6: Stacked bar plot of `result` in `occupation`

According to Table 2 and Figure 6 and 7, most categories in two columns have observations less than 20. Small and insufficient group size would lead to unstable and biased results. On the other hand, numbers of occupation groups could result in larger degrees of freedom, which increases model complexity. As discussed before, we may use these variables in the model, as the occupation of a person or the corresponding household head could have influence on the infection of Leptospirosis [?]. Thus, it is necessary to deal with these problems properly.

We find that some occupation groups with similar meanings and risks of exposure could be merged. Thus, we merge "Crop Farming" with "Farmer" and merge "Herdsman" with "Pastoralist". However, most of the smaller occupation groups have very different risk profiles, e.g. drivers are unlikely to have the same contact with animal risk as farmers. In this case, we disregard such categories and discard corresponding observations.

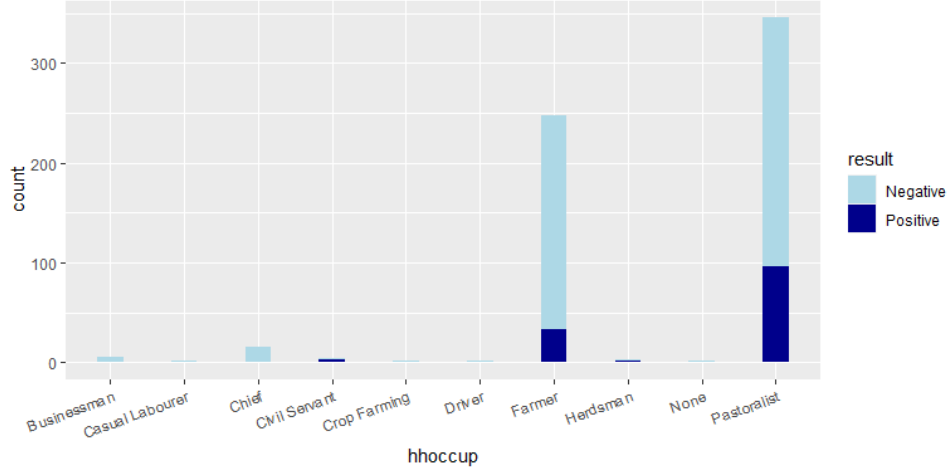


Figure 7: Stacked bar plot of **result** in **hhoccup**

Now, we plot stacked bar charts for each aforementioned factor and report the prevalence in each category of it.

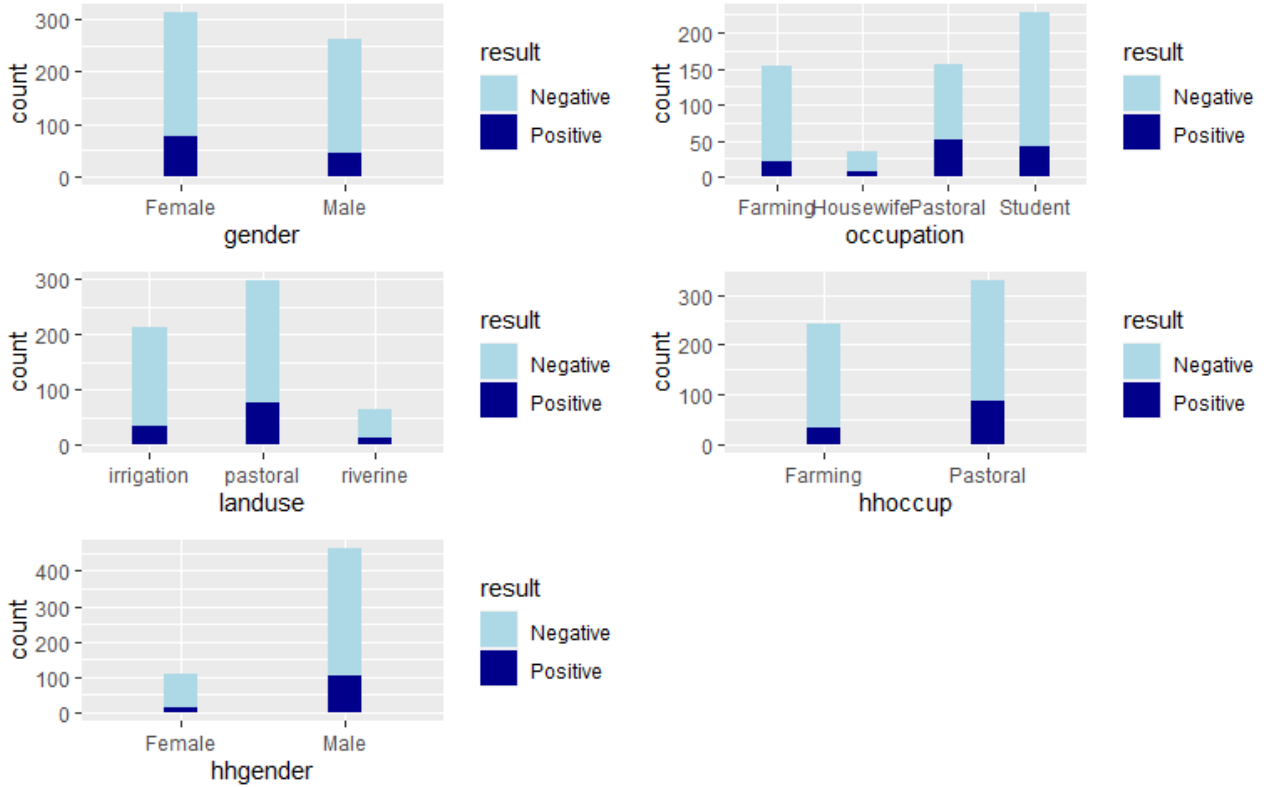


Figure 8: Stacked bar plot of **result** in factor variables

According to Figure 8 and Table 3, **gender**, **occupation**, **landuse** and **hhoccup** have significantly different prevalence among categories, hence we select them as feature candidates in modelling step. Although it is also shown that prevalence between **hhgender** groups are significantly different, we do not include this column in our models as it is imbalanced distributed where over 80% observations in this column are "Male".

2.3.2 Numerical Variables

We make box plots for numerical variables including **age**, **altitude**, **nmales**, **nfemales**, **famsize** and **disthosp**, grouped by **result**. Here, **altitude** denotes the altitude (recorded via GPS) of the village

factor	category	prevalence
gender	Female	0.2508
	Male	0.1821
occupation	Farmer	0.1465
	Student	0.1834
	Pastoralist	0.3257
	Housewife	0.2353
landuse	irrigation	0.1636
	pastoral	0.2597
	riverine	0.2090
hhoccup	Farmer	0.1325
	Pastoralist	0.2803
hhgender	Female	0.1481
	Male	0.2341

Table 3: Prevalence in each category of factor variables

where the sample was collected, `nmales` and `nfemales` represent the number of males and females in the household sampled respectively, and `famsize`, the sum of `nmales` and `nfemales`, is the number of people in the sampled household. Note that there exists an observation with `disthosp` over 500 kilometres, which is deleted as an outlier, because the values of this feature in other observations from the same village are much lower.

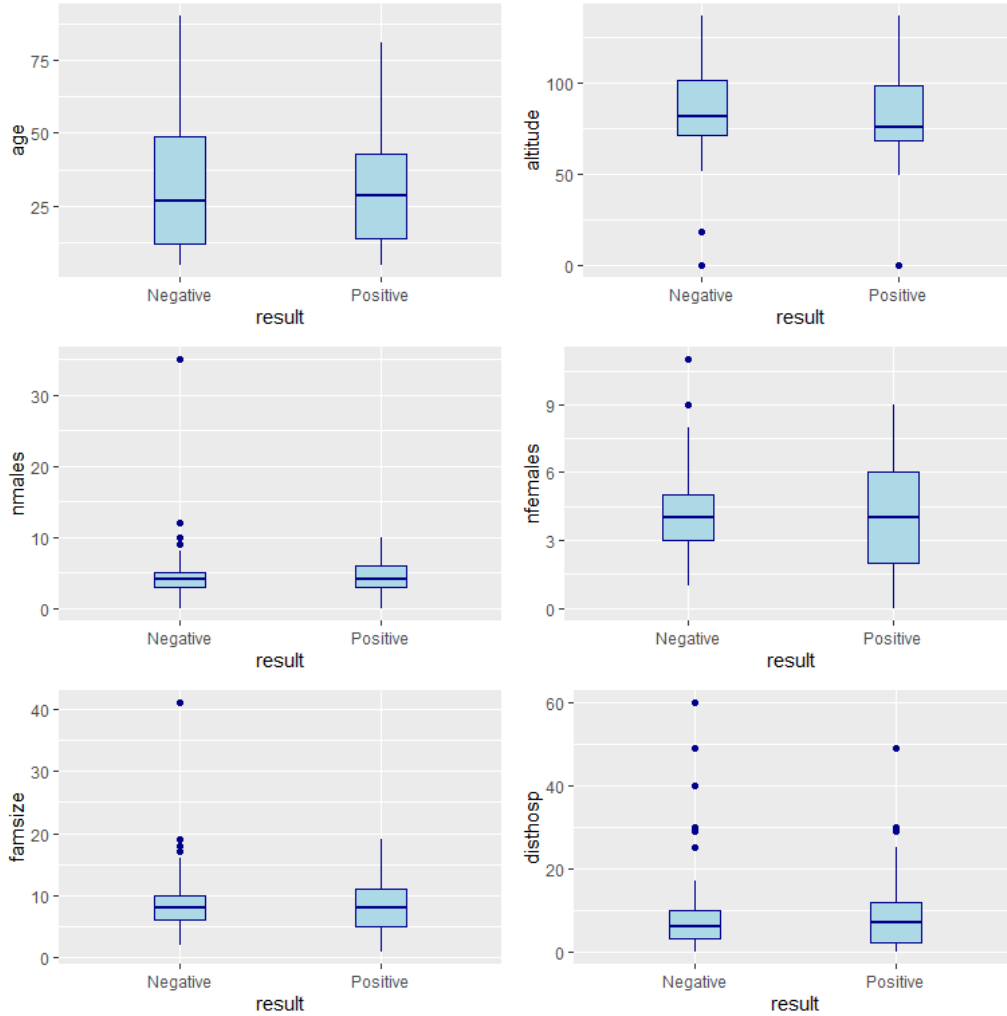


Figure 9: Box plots of numerical variables, grouped by `result`

Comparing with factors, the distributions of numerical variables do not show significant difference between "Negative" and "Positive" groups from Figure 9, which are probably insignificant in the models. Based on [?], we select **age**, **altitude**, **famsize** and **disthosp** as feature candidates for modelling.

2.3.3 Geographical Variables

village, **location** and **constituency** are nominal variables indicating the geographical information of sampled people, with a nested order **village** < **location** < **constituency**. As what we do for other nominal variables, we plot the stacked bar charts of the response variable and compute the prevalence, grouped by the levels in **village**, **location** and **constituency** respectively.

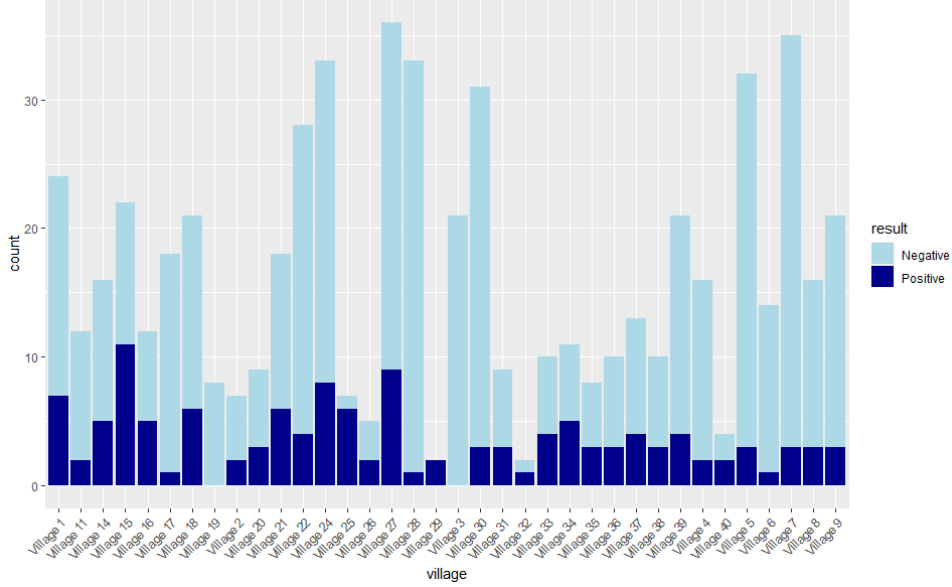


Figure 10: Stacked bar plots of **result** in **village**

category	prevalence	category	prevalence	category	prevalence
Village 1	0.2917	Village 16	0.4167	Village 29	1.0000
Village 2	0.2857	Village 17	0.0556	Village 30	0.0968
Village 3	0	Village 18	0.2857	Village 31	0.3333
Village 4	0.1250	Village 19	0.0000	Village 32	0.5000
Village 5	0.0938	Village 20	0.3333	Village 33	0.4000
Village 6	0.0714	Village 21	0.3333	Village 34	0.4545
Village 7	0.0857	Village 22	0.1429	Village 35	0.3750
Village 8	0.1875	Village 24	0.2424	Village 36	0.3000
Village 9	0.1429	Village 25	0.8571	Village 37	0.3077
Village 11	0.1667	Village 26	0.4000	Village 38	0.3000
Village 14	0.3125	Village 27	0.2500	Village 39	0.1905
Village 15	0.5000	Village 2	0.0303	Village 40	0.5000

Table 4: Prevalence in each category of **village**

According to the plots and tables, we may also consider the influence of these geographical variables on Leptospirosis diagnosis in our models due to different prevalence among groups. However, **village** and **location** contain a number of levels, which causes great model complexity with high degrees of freedom and further unstable estimates. Thus, we consider the models with the random effect of these geographical variables. The details of random effect will be explained in the next section.

Up to now, we have extracted a subset for modelling with 595 observations left and 11 feature candidates, including 4 nominal variables, 4 numerical variables and 3 geographical variables.

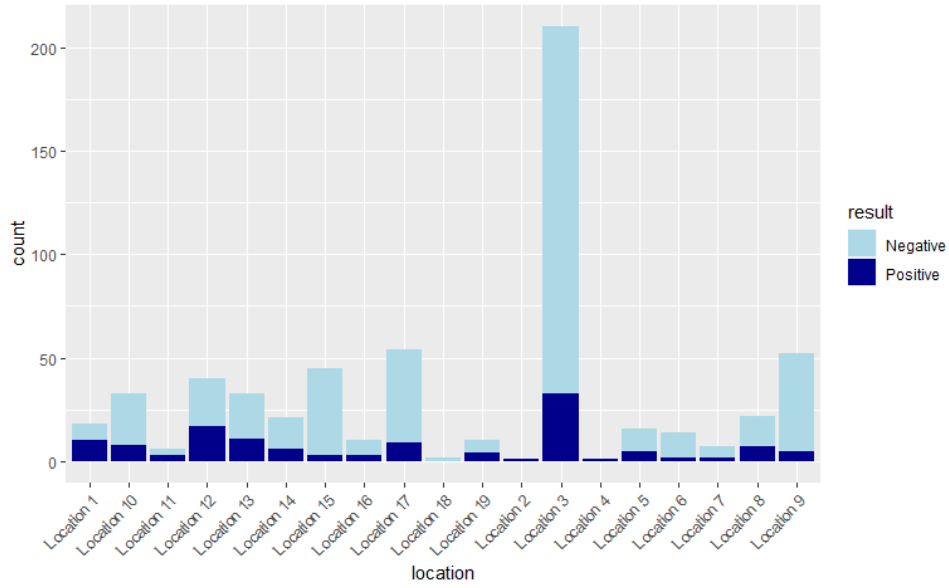


Figure 11: Stacked bar plots of **result** in **location**

category	prevalence	category	prevalence	category	prevalence
Location 1	0.5556	Location 7	0.2857	Location 13	0.3333
Location 2	1.0000	Location 8	0.3182	Location 14	0.2857
Location 3	0.1571	Location 9	0.0962	Location 15	0.0667
Location 4	1.0000	Location 10	0.2424	Location 16	0.3000
Location 5	0.3125	Location 11	0.5000	Location 17	0.1667
Location 6	0.1429	Location 12	0.4250	Location 18	0.0000
				Location 19	0.4000

Table 5: Prevalence in each category of **location**

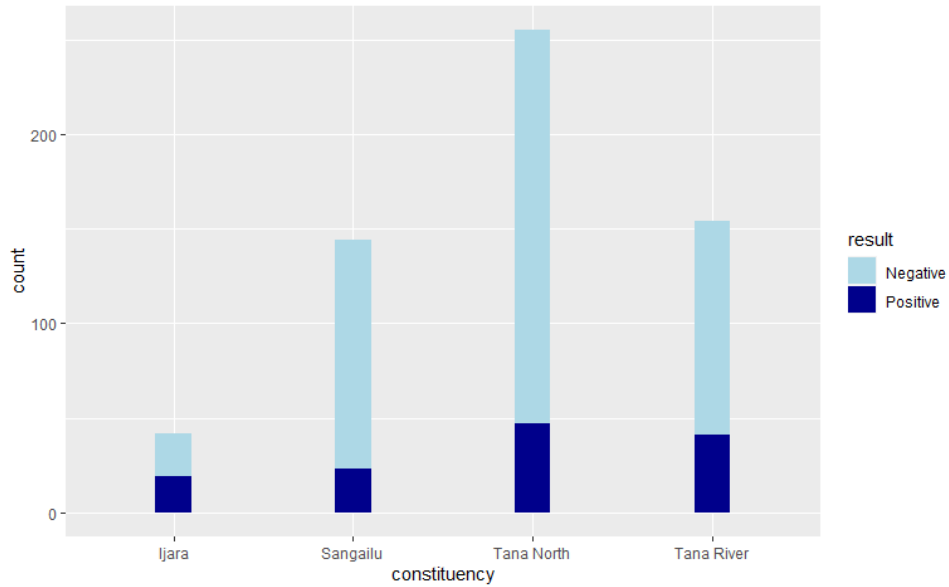


Figure 12: Stacked bar plots of **result** in **constituency**

category	prevalence	category	prevalence
Tana North	0.1843	Tana River	0.2662
Ijara	0.4524	Sangailu	0.1597

Table 6: Prevalence in each category of **constituency**

3 Implementation

In this section, we introduce the definition and purpose of General Linear Mixed Model (GLMM) with random effect used for modelling. Moreover, we explain the model construction, and further the comparison among models.

3.1 General Linear Mixed Model

The exposition in this subsection follows that in Zuur (2009) [?].

3.2 Model Construction and Selection

3.2.1 Step 1: Initial Model

Following Zuur’s logic, we start with a binomial GLM, i.e. logistic regression, containing all feature candidates excluding geographical variables. Here, we do not contain the interaction terms of variables in the model as we do not detect much interaction existing among features in EDA section. On the other hand, introducing interaction terms increases much model complexity due to a large number of potential interaction terms.

3.2.2 Step 2: Find the Optimal Fixed Structure

The `summary` function and `step` function are both used to determine variables with the optimal fixed effect among candidates. The former function obtains the feature significance based on the t test, while the latter one choose the optimal model and significant features by AIC and the χ^2 test. Table 7 show the significance of features in the initial GLM based on `summary` function, and Table 8 display the significant features based on `step` function.

	t value	$\Pr(> z)$
(Intercept)	-0.341	0.73317
genderMale	-2.428	0.01520 *
age	-1.798	0.07219 .
occupationHousewife	-0.757	0.44890
occupationPastoralist	0.754	0.45068
occupationStudent	-1.423	0.15463
landusepastoral	-2.347	0.01892 *
landuseriverine	-1.607	0.10810
altitude	-1.681	0.09283 .
famsize	0.346	0.72900
hhoccupPastoralist	3.181	0.00147 **
disthosp	0.301	0.76322
Significance codes: 0, ‘***’ 0.001, ‘**’ 0.01, ‘*’ 0.05, ‘.’ 0.1, ‘ ’ 1		

Table 7: Feature significance in initial GLM using `summary` function

Based on two tables, we select `gender` and `hhoccup` in further steps, which are both statistically significant at 0.05 significance level in two function results. Besides, `landuse` is also selected in fixed structure, as it is significant at 0.1 significance level in stepwise selection and its level is statistically significant at 0.05 level in model summary.

Although `occupation` shows its significance in stepwise algorithm, none of its levels is significant in model summary. Moreover, as mentioned in EDA section, we observe high correlation between `occupation` and `landuse`, which probably causes the problem of collinearity and poor model performance. Therefore, we do not include `occupation` in the fixed structure.

In conclusion, the optimal fixed structure consists of `gender`, `landuse` and `hhoccup`.

	Df	AIC	Pr(>Chi)
(none)		602.72	
altitude	1	603.45	0.0990659 .
age	1	603.86	0.0765869 .
landuse	2	604.59	0.0531339 .
occupation	3	605.55	0.0316715 *
gender	1	606.75	0.0141019 *
hhoccup	1	611.90	0.0008302 ***
Significance codes: 0, '***' 0.001, '**' 0.01, '*' 0.05, '.' 0.1, ' ' 1			

Table 8: Significant features in stepwise algorithm using `step` function

3.2.3 Step 3: Find the Optimal Random Structure

With the pre-determined fixed structure, we now introduce the random effect of geographical variables and determine GLMM with the best performance. As previously mentioned, we consider the random effect of `village` and `location` to investigate whether there is a village or location effect.

We start with fitting the random intercept model of `village`. To assess whether the model with mixed effect model is better than the ordinary binomial GLM, we refit the latter one over the selected features without random intercept and compare the models using `anova` function. The output is given in Table 9:

	npars	AIC	BIC	logLik	deviance	Chisq	Df	Pr(>Chisq)
GLM	5	603.63	625.57	-296.81	593.63			
GLMM(village)	6	601.87	628.20	-294.93	589.87	3.7611	1	0.05246 .
Significance codes: 0, '***' 0.001, '**' 0.01, '*' 0.05, '.' 0.1, ' ' 1								

Table 9: Comparison between GLM and GLMM of `village`

The lower AIC value in the table indicate that the mixed model including the random effect of village is preferred. Moreover, we extend the random part by building random intercept and slope models. Specifically, the `landuse` effect may be different per village, and the same may hold for the `hhoccup` effect. Table 10 show the output of model comparison using `anova` function:

	npars	AIC	BIC	logLik	deviance	Chisq	Df	Pr(>Chisq)
random intercept	6	601.87	628.20	-294.93	589.87			
random intercept + <code>landuse</code>	8	605.75	640.85	-294.87	589.75	0.122	2	0.9408
random intercept + <code>hhoccup</code>	11	611.48	659.76	-294.74	589.48	0.264	3	0.9667
Significance codes: 0, '***' 0.001, '**' 0.01, '*' 0.05, '.' 0.1, ' ' 1								

Table 10: Comparison among GLMMs of `village` with different random structures

Extending the model with random slopes gives no improvement as the AIC values of models with random slopes get higher. Now we look at the alternative random intercept model considering the location effect. Comparison between random intercept models of `village` and `location` using AIC command is displayed in Table 11:

	df	AIC
GLMM(village)	6	576.9177
GLMM(location)	6	577.3410

Table 11: Comparison between random intercept models of `village` and `location`

According to the table, the AIC value of the random intercept model with the village effect is slightly lower. Besides, since `village` contains more levels than `location`, the models considering the random effect on village aspect obtain more detailed results. Moreover, the random intercept and

slope models of **location** perform even worse. As a result, we prefer the random intercept model of **village**, with the following expression:

$$\text{logit}(p_{ij}) = \text{Intercept} + \beta_1 \times \text{gender} + \beta_2 \times \text{landuse} + \beta_3 \times \text{hhoccup} + \text{randomIntercept}_i + \epsilon_{ij}$$

Here, the index i refers to villages and j to the observation within a village. Moreover, randomintercept_i denotes the random intercept across villages and ϵ_{ij} is the within-village variation.

Note that here we do not consider the constituency effect in our model as **constituency** only contains 4 levels. It may be better to include this variable in model's fixed structure. However, after introducing the fixed effect of **constituency**, other variables in the fixed structure become insignificant. Moreover, the variance and standard deviation of the random effect of **village** or **location** turn to 0, indicating the insignificance of random effect. Further details are given in Appendix. Consequently, the influence of **constituency** is ruled out.

4 Results

The modified Bayesian network based on Hill-climbing algorithm is determined as the final network, shown in figure ?? . As the parents of cognitive score, **age**, **gender**, **country**, **education** and **depression** tend to have significant effects on respondents' cognitive function.

Next, we estimate the parameters of local distributions of these parent nodes. We use the classical Bayesian posterior estimator with imaginary sample size of 10000, which provides smoother and more robust estimates, comparing with the alternative method using maximum likelihood estimation. We plot line graphs to explore the trends of the conditional probabilities of cognitive score, conditioning on its parent nodes. Since model overfitting happens in stratifications with sparse counts when conditioning on too many nodes, plots of the basic model are wiggly. As a remedy, we apply ordered logistic regressions of cognitive score, grouped by gender, over other parent nodes. Figure ?? displays the improvement of line graphs after using ordered logistic regression.

With smoother line graphs, we discuss the influence of parent nodes on dementia risk:

4.1 Age

Age shows a negative influence on individual's cognitive function from the line graphs. Generally, as respondents age, the conditional probability of having severe cognitive impairment increases, while the probability of being cognitively normal drops. This is consistent with our common sense. As people get older, they suffer from worse health condition and more loneliness, which could lead to cognitive decline. The path "age \rightarrow chronic diseases \rightarrow depression \rightarrow cognitive score" in the network supports this view.

We condition on people's education and depression level. Given a specific gender group, we observe better cognitive function among respondents from countries with high aGDP in all age groups, as they have higher probability to be cognitively normal and lower probability to have severe cognitive impairment than those from countries with low aGDP. On the other hand, looking at individuals from the same country group, female respondents show higher chance to suffer from severe cognitive impairment than males, especially in higher age groups (e.g. over 85). Figure ?? displays the trends of cognitive score over age in different gender and country groups, among respondents with no education and low depression level.

Interestingly, we notice that most lines become steeper after 65 years old, indicating faster decline on people's cognition. Thus, we think that 65 could be an important time point where individual's cognitive decline aggregates. This is consistent with results from the study of Lee et al (2018) [17] in China. Furthermore, we particularly focus on the influence of other factors in the age groups around 65.

4.2 Education

Education level tends to have a positive effect on dementia risk reduction. According to line graphs, conditional probability of severe cognitive impairment decreases when respondents have higher education levels, and these people have higher chance to be cognitively normal. It has been proved that education, as a main contributor, stimulates people's cognitive function in early life according to the work of Black et al (2018) [18]. Besides, individual's education level might influence the work and further the income in midlife, which are also believed to have effects on cognition in other researches [14].

Similarly, we condition on age and depression level here. In each gender group, though the trends of lines are very similar, people from countries with low aGDP seems to have higher dementia risk, due to lower conditional probability of being cognitively normal and higher probability of having severe impairment on cognition. In each country group, again, female respondents could be more likely to suffer from cognitive impairment, particularly in low education levels. Figure ?? shows the trends of cognitive score over education levels in different gender and country groups, holding 65-69 years old and low depression level.

4.3 Depression

Depression is found to aggravate individual's cognitive impairment. From the line graphs, we observe that as depression level increases, the conditional probability of normal cognition drops and people are more likely to suffer from severe cognitive impairment. As a part of the prodrome and early stages of dementia, depression is associated with various possible psychological or physiological mechanisms. This can be supported by edges corresponding to depression across many domains in the network including "social \rightarrow depression", "household finance \rightarrow depression", "chronic diseases \rightarrow depression" and "obesity \rightarrow depression". Poor health condition, low income and social isolation all tends to cause the feeling of depression, and further influence individual's cognition.

Given age, gender and education level, the conditional probability of normal cognition of respondents from countries with high aGDP is significantly higher than those from countries with low aGDP at each depression level. Accordingly, the probability of suffering from severe cognitive impairment in countries with high aGDP is much lower than it in the other stratification. Conditioning on the country group, we find that women possibly have higher dementia risk than men from the plots. Figure ?? shows the trends of the conditional probabilities of different cognitive function levels over depression, holding 65-69 age group and no education level.

Moreover, we notice that segments from medium to high depression level become much steeper than those from low to medium level. According to the ranges of depression levels, respondents with at least 5 negative feelings are determined with high depression. Thus, we guess it could be a threshold that people experiencing over 4 negative feelings have much higher dementia risk than others.

4.4 Gender

Gender shows consistent influence on cognitive function across analysis on age, education and depression that generally females have higher dementia risk than males, holding other factors constant. In particular, this difference is observed to be more significant in elder age groups. These results are consistent with the conclusions in [13]. A possible explanation for females having higher probabilities of cognitive impairment may be that women generally survive to longer ages than men, making the proportion of women become larger in elder groups. On the other hand, due to their longevity, female respondents are more likely to live alone. This could result in more severe social isolation, increases their depression level and finally affects their cognition, which is supported by the path "female \rightarrow social isolation \rightarrow depression". Besides, we believe that gender could affect individual's cognitive function in other ways including the education level and obesity, from the network.

4.5 Country

Recall that we divide countries into two groups, based on their average gross domestic product in 2013. According to previous analysis, we conclude that respondents from countries with high aGDP tend to have lower dementia risk than those from countries with low aGDP. The influence of country can be various. From the network, we observe that country are significantly associated with nodes including education, obesity, drinking behaviour and household finance. Generally, countries with high aGDP probably have better welfare, which means they could provide higher-standard education and more advanced health care for their citizens than other countries, positively affecting the cognitive function. Moreover, provided with better welfare, people in these countries may live with more joy, hence they have lower depression level, which also reduces the probability of cognitive decline.

5 Conclusions

We construct a modified Bayesian network on a subset of easySHARE dataset to identify risk factors on dementia. As a result, age, gender, country, education and depression tend to have significant influence on individual’s cognitive function. Specifically, the probability of severe cognitive impairment rises as age or depression level increases. In particular, individuals over 65 years old or experiencing high-level depression have much higher dementia risk. Education shows opposite influence on cognitive function that people with higher education levels have lower risk. Gender and country affect individual’s cognition in a variety of ways. Generally, females have higher dementia risk than males; people from countries with low aGDP have higher probability to suffer from dementia.

There are some limitations of this research. First, although hearing loss, diabetes, and hypertension are believed to have important effect on cognitive impairment in [14], their influence could not be identified due to the lack of relevant columns in easySHARE data. We could only define a general factor ”chronic diseases” to count the number of chronic diseases the respondents had and explore its possible influence on dementia risk instead. To determine the specific effects of these diseases, we need to access data with relevant variables. Second, our research belongs to cross-sectional analysis, which focuses on data at a specific time point among different individuals and obtains conclusions for the general population. In this case, our research is lack of the exploration on longitudinal effects of some risk factors on dementia risk among particular respondents over the time.

References

- [1] Alzheimer’s Disease International. World alzheimer report 2019: attitudes to dementia. *Alzheimer’s Disease International: London*, 2019.
- [2] A. Börsch-Supan. *Survey of Health, Ageing and Retirement in Europe (SHARE) Wave 1*. Release version: 8.0.0. SHARE-ERIC. Data set, 2022. doi: [10.6103/SHARE.w1.800](https://doi.org/10.6103/SHARE.w1.800).
- [3] A. Börsch-Supan. *Survey of Health, Ageing and Retirement in Europe (SHARE) Wave 2*. Release version: 8.0.0. SHARE-ERIC. Data set, 2022. doi: [10.6103/SHARE.w2.800](https://doi.org/10.6103/SHARE.w2.800).
- [4] A. Börsch-Supan. *Survey of Health, Ageing and Retirement in Europe (SHARE) Wave 3*. Release version: 8.0.0. SHARE-ERIC. Data set, 2022. doi: [10.6103/SHARE.w3.800](https://doi.org/10.6103/SHARE.w3.800).
- [5] A. Börsch-Supan. *Survey of Health, Ageing and Retirement in Europe (SHARE) Wave 4*. Release version: 8.0.0. SHARE-ERIC. Data set, 2022. doi: [10.6103/SHARE.w4.800](https://doi.org/10.6103/SHARE.w4.800).
- [6] A. Börsch-Supan. *Survey of Health, Ageing and Retirement in Europe (SHARE) Wave 5*. Release version: 8.0.0. SHARE-ERIC. Data set, 2022. doi: [10.6103/SHARE.w5.800](https://doi.org/10.6103/SHARE.w5.800).
- [7] A. Börsch-Supan. *Survey of Health, Ageing and Retirement in Europe (SHARE) Wave 6*. Release version: 8.0.0. SHARE-ERIC. Data set, 2022. doi: [10.6103/SHARE.w6.800](https://doi.org/10.6103/SHARE.w6.800).
- [8] A. Börsch-Supan. *Survey of Health, Ageing and Retirement in Europe (SHARE) Wave 7*. Release version: 8.0.0. SHARE-ERIC. Data set, 2022. doi: [10.6103/SHARE.w7.800](https://doi.org/10.6103/SHARE.w7.800).
- [9] A. Börsch-Supan. *Survey of Health, Ageing and Retirement in Europe (SHARE) Wave 8*. Release version: 8.0.0. SHARE-ERIC. Data set, 2022. doi: [10.6103/SHARE.w8.800](https://doi.org/10.6103/SHARE.w8.800).
- [10] A. Börsch-Supan and S. Gruber. *easySHARE*. Release version: 8.0.0. SHARE-ERIC. Data set, 2022. doi: [10.6103/SHARE.easy.800](https://doi.org/10.6103/SHARE.easy.800).
- [11] Flávio Luiz Seixas, Bianca Zadrozny, Jerson Laks, Aura Conci, and Débora Christina Muchaluat Saade. A bayesian network decision model for supporting the diagnosis of dementia, alzheimer’s disease and mild cognitive impairment. *Computers in biology and medicine*, 51:140–158, 2014.
- [12] Eileen M Crimmins, Jung Ki Kim, Kenneth M Langa, and David R Weir. Assessment of cognition using surveys and neuropsychological assessment: the health and retirement study and the aging, demographics, and memory study. *Journals of Gerontology Series B: Psychological Sciences and Social Sciences*, 66(suppl_1):i162–i171, 2011.
- [13] Christopher R Beam, Cody Kaneshiro, Jung Yun Jang, Chandra A Reynolds, Nancy L Pedersen, and Margaret Gatz. Differences between women and men in incidence rates of dementia and alzheimer’s disease. *Journal of Alzheimer’s Disease*, 64(4):1077–1083, 2018.
- [14] Gill Livingston, Andrew Sommerlad, Vasiliki Orgeta, Sergi G Costafreda, Jonathan Huntley, David Ames, Clive Ballard, Sube Banerjee, Alistair Burns, Jiska Cohen-Mansfield, et al. Dementia prevention, intervention, and care. *The Lancet*, 390(10113):2673–2734, 2017.
- [15] Stefan Gruber, Christian Hunkler, and Stephanie Stuck. *Generating easySHARE: guidelines, structure, content and programming*. SHARE Working Paper Series (17-2014). Munich: MEA, Max Planck Institute for Social Law and Social Policy, 2014.
- [16] M Scutari. Learning bayesian networks with the bnlearn r package. *Journal of Statistical Software*, 35(3), 2010.
- [17] Allen TC Lee, Marcus Richards, Wai C Chan, Helen FK Chiu, Ruby SY Lee, and Linda CW Lam. Association of daily intellectual activities with lower risk of incident dementia among older chinese adults. *JAMA psychiatry*, 75(7):697–703, 2018.

- [18] Deborah Blacker and Jennifer Weuve. Brain exercise and brain outcomes: does cognitive activity really work to maintain your brain? *JAMA psychiatry*, 75(7):703–704, 2018.

Appendix

Programming part in this research is completed with R.

R packages:

For EDA section, we use the `tidyverse` and `gridExtra` packages. We construct and estimate the networks with `bnlearn` package and display the network with `Rgraphviz` package. `MASS` package is used to fit the ordered logistic regression on cognitive score. All histograms and line graphs are generated with `ggplot2`.

R code:

The complete code is available via this [Github repository](#).

Word Count

This report contains 4828 words, including executive summary, main text, references and appendix. The screenshot using **Analyse Text** function in TeXstudio is provided.