

Lifestyle Choice Matters: Exploring Population Health Disparities and the Effect of Tobacco and Alcohol on Life Expectancy*

Yufei Liu

November 5, 2023

Life expectancy is an important measurement of population health and there are many significant influencing factors such as genetic and lifestyle. This paper examines the relationship between the life expectancy, the prevalence of tobacco products, and the alcohol consumption using data from Global Health Observatory data repository from the WHO. Using regional comparisons, we find as prevalence of tobacco increases, the life expectancy at age 60 decreases while as alcohol consumption increases, the life expectancy increases. We also find regional, temporal, and sexual differences in life expectancy at age 60. Further work could consider additional influencing variables and look at more complex models to capture non-linear relationships.

1 Introduction

Life expectancy refers the average number of years a person could expect to live, starting from birth (for life expectancy at birth) or other age groups such as for age 60, if sex- and age-specific mortality rates holds constant for a specific year and living area (World Health Organization 2023). Since the life expectancy varies across years and regions, it is the key metric to assess population health (Roser, Ortiz-Ospina, and Ritchie 2013). Several socioeconomic, genetic, lifestyle, nutritional, and environmental factors influence the life expectancy (Rahman, Rana, and Khanam 2022).

In this paper, we investigate how lifestyle factors influence life expectancy at age 60, and focus on two specific aspects, prevalence of tobacco products and alcohol consumption. We use R (R Core Team 2022) to analyze the relationship between the life expectancy, the prevalence

*Code and data are available at: https://github.com/Florence-Liu/life_expectancy

of tobacco, and the alcohol consumption using data from Global Health Observatory (GHO) data repository from the World Health Organization (WHO). We construct multiple linear regression models in which life expectancy at age 60 for all WHO regions is explained by five variables, sex, region, year, prevalence of tobacco, and alcohol consumption. We find that females have a higher life expectancy at age 60 than males across years, and life expectancy has a large regional difference as well as tobacco use and alcohol consumption. We also find as the prevalence of tobacco increases, the expected life expectancy decreases while as the alcohol consumption increases, the expected life expectancy also increases. Further work could introduce more complex models to capture non-linear relationships and give region-specific studies as well as taking into account other socioeconomic and environmental factors.

The remainder of this paper is structured as follows: Section 2 discusses the data with Section 2.1 including information about data collection and data cleaning results, and Section 2.2 including graphs and tables representing relationships between variables and some discussions; Section 3 introduces two linear regression models we employ and specifies parameters; Section 4 shows the estimates of fitted models and criteria for comparing two models; Section 5 includes discussions about model results and implications with Section 5.1 making a brief summary of what the paper has done, Section 5.2 showing findings based on model results and data visualization, and Section 5.3 talking about the weaknesses and future improvement.

2 Data

We used R to do the analysis in this paper (R Core Team 2022). We used packages `tidyverse` (Wickham et al. 2019), `janitor` (Firke 2021), and `here` (Müller 2020) to clean and load the data as well as create figures, and `knitr` (Xie 2014) and `modelsummary` (Arel-Bundock 2022) to generate tables. The color style of the figures was created referring to a R colors cheet-sheet (Wei 2021).

2.1 Data description

The datasets used in this paper was obtained from World Health Organization (WHO) Global Health Observatory data repository, and is publicly available from the WHO website. We utilized 3 datasets: life expectancy and healthy life expectancy by WHO region (2020), SDG Target 3.5 Substance abuse by WHO region (2023), and SDG Target 3.a Tobacco control by WHO region (2022). The life expectancy dataset contains life expectancy and healthy life expectancy at birth and at age 60 for different sexes and WHO regions in year 2000, 2010, 2015, and 2019. The life expectancy values were estimated based on mortality data from civil registration. The substance abuse dataset contains total alcohol per capita (aged 15+) consumption (total APC) for different sexes and WHO regions in year 2000, 2005, 2010, 2015, and 2019. The tobacco control dataset contains the percentage of the population aged 15+

who currently use tobacco products based on population-based surveys for different sexes and WHO regions in year 2000, 2005, 2010, 2015, 2018, 2019, and 2020.

We specifically selected life expectancy at age 60 since it would be more associated with either tobacco or alcohol usage than life expectancy at birth, and data in year 2000, 2010, 2015, and 2019. Then we merged the three datasets into one by region, sex, and year. Our cleaned data for analysis contains variables:

- **region:** WHO regions
- **year:** Year data collected
- **sex:** Sex at birth
- **life_expectancy:** The average number of years that a person of age 60 could expect to live
- **alcohol_consumption:** Total alcohol per capita aged 15+ consumption in litre
- **prevalence_of_tobacco:** The percentage of the population aged 15+ who currently use any tobacco products.

2.2 Data Visualization

Figure 1 shows how the prevalence of tobacco use relates to life expectancy at age 60 for two sexes, male and female. We see an overall decreasing trend that when the prevalence of tobacco use increases, the mean life expectancy at age 60 decreases. For the two sex groups, we could discover that there are clear clusters of the points, indicating the mean prevalence of tobacco may be different in the two groups. The two dashed lines is the linear fitted line for each group and the green line represents the linear fitted line for both group as a whole. It shows that the two variables interact and we should consider an interaction term between sex and prevalence of tobacco in our model.

Figure 2 shows how total alcohol per capita consumption relates to life expectancy at age 60 for male and female separately. The two dashed lines is the linear fitted line for each group and the green line representing the linear fitted line for both group as a whole. We see that Simpson's Paradox occurs that the life expectancy increases when the alcohol consumption increases for each sex group, however, the trend is reversed when we combine the two groups. It also indicates that we should include an interaction term in our model.

Table 1: Summary of average life expectancy, prevalence of tobacco, and alcohol consumption across years

		Mean	Mean	Mean
	Year	life expectancy (year)	prevalence of tobacco (%)	alcohol consumption (litre)
	2000	18.3	31.2	5.2
	2010	19.5	25.2	5.6

		Mean	Mean	Mean
	Year	life expectancy (year)	prevalence of tobacco (%)	alcohol consumption (litre)
	2015	20.1	22.9	5.6
	2019	20.5	21.4	5.3

Table 2: Summary of average life expectancy, prevalence of tobacco, and alcohol consumption for different WHO regions

		Mean	Mean	Mean
Region	life expectancy (year)	prevalence of tobacco (%)	alcohol consumption (litre)	
Africa	16.9	13.4	4.9	
Americas	22.0	21.2	7.8	
Eastern Mediterranean	18.0	22.2	0.3	
Europe	21.2	29.2	10.2	
Global	20.2	26.6	5.5	
South-East Asia	18.2	37.6	3.2	
Western Pacific	20.8	26.0	5.9	

Figure 3 shows life expectancy at age 60 for different years and different sexes. In general, life expectancy at age 60 increased from 2000 to 2019 for both sexes. We could see that in each year, the mean life expectancy at age 60 for female is much higher than that for male along with larger range and variance. Also, the distribution of life expectancy for both sexes are skewed with one whisker longer than the other, indicating more observations centered at higher values. There is no outliers for all the boxplots.

Table 1 and Table 2 shows summary tables for the mean values of life expectancy at age 60, prevalence of tobacco, and alcohol consumption across years and regions. We could see that the mean life expectancy increased from 2000 to 2019, which is consistent with Figure 3, while the mean prevalence of tobacco decreased across years and the mean alcohol consumption is relatively stable with minor changes. The potential inverse relationship between prevalence of tobacco and life expectancy is consistent with what we found in Figure 1. For different WHO regions, the three values differ a lot. Africa has the lowest mean life expectancy at age 60 while Americas has the highest mean life expectancy at age 60. However, Africa also has the lowest prevalence of tobacco, which contradicts our previous findings. Also, the mean alcohol consumption has a large range with the lowest value at 0.3 for Eastern Mediterranean and highest value at 10.2 for Europe. These findings in region difference could be further investigated since there might be potential cultural and socio-economic reasons behind.

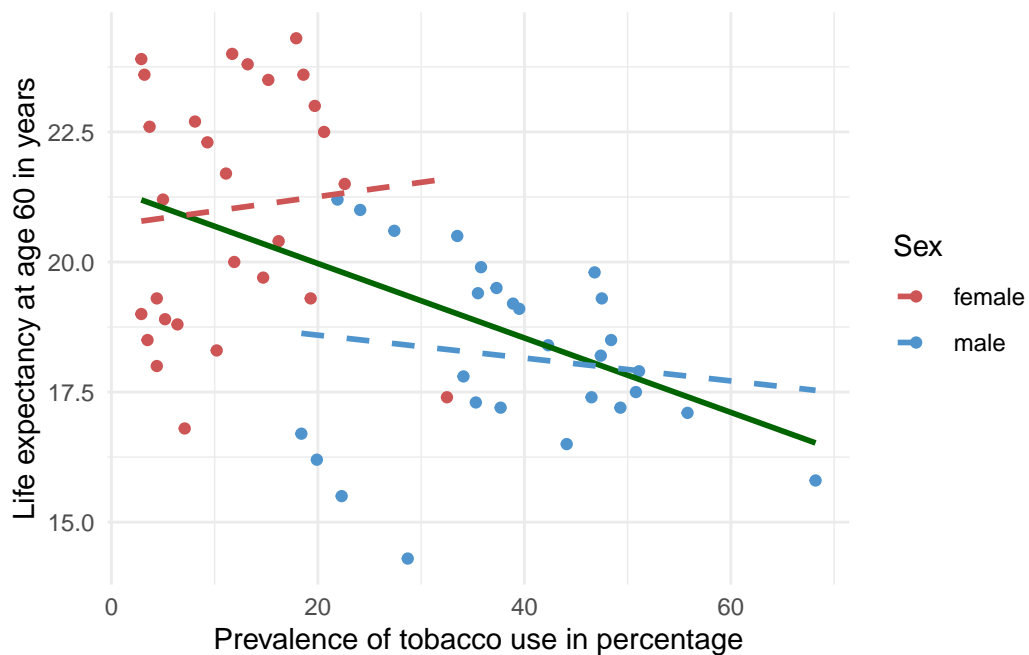


Figure 1: The effect of prevalence of tobacco on life expectancy at age 60 for different sex

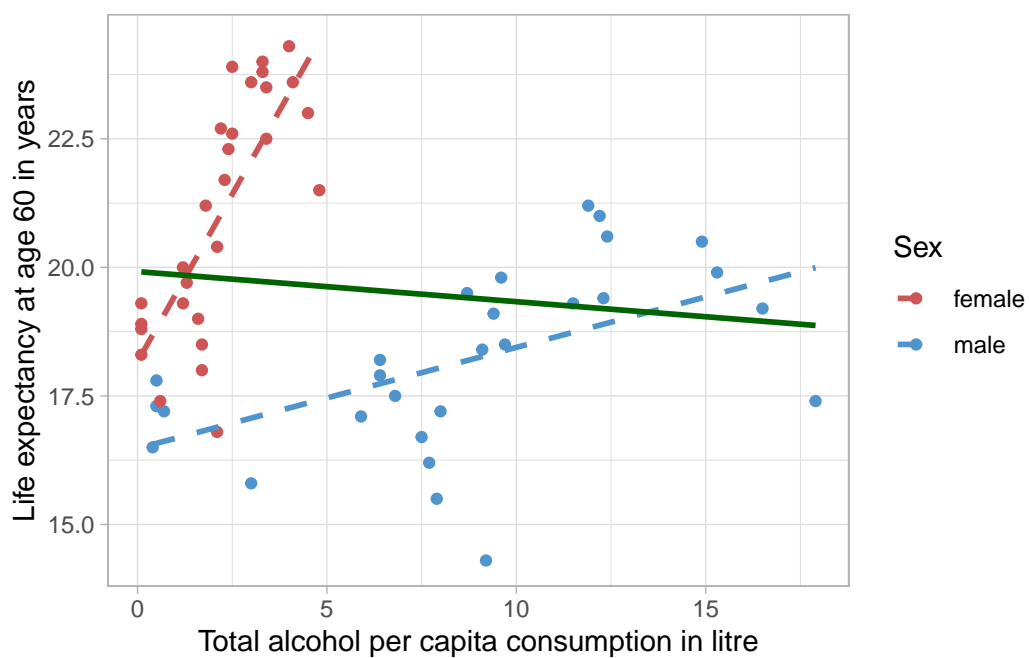


Figure 2: The effect of total alcohol per capita consumption on life expectancy at age 60 for different sex

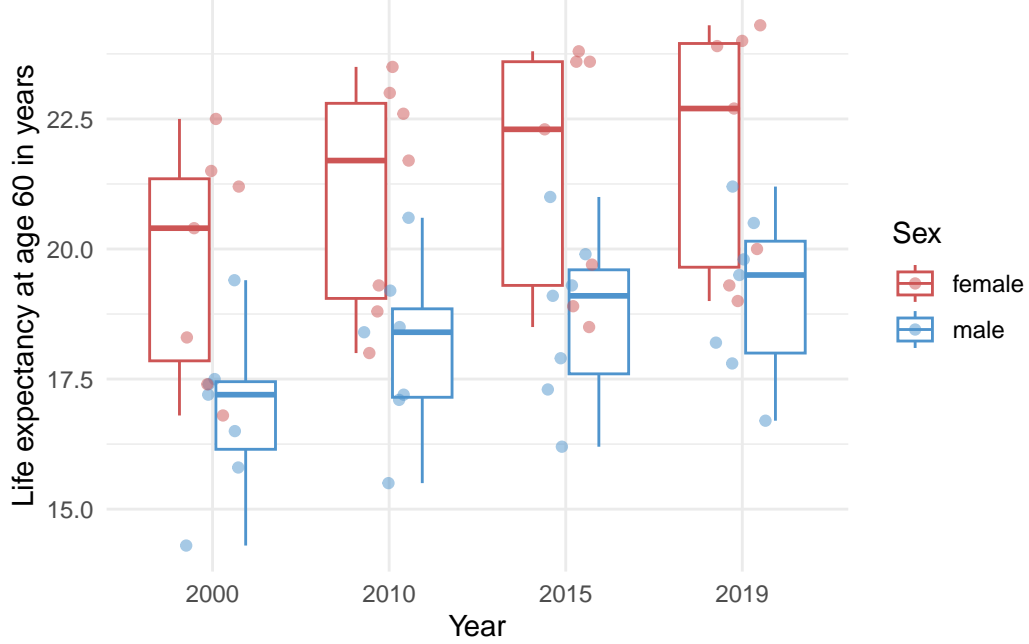


Figure 3: Life expectancy at age 60 across years for different sex

3 Model

Based on the data visualization, there are clear trends between prevalence of tobacco, alcohol consumption, sex, year, region, and life expectancy. So we will fit multiple linear regression models to discover how these factors contribute to life expectancy at age 60.

To justify our choice of including interaction terms based on observation, we will fit two models. Model 1 is the full model including interaction terms and model 2 is the reduced model without interaction terms. The full model is shown as

$$Y = \beta_0 + \beta_1 X_{tobacco} + \beta_2 D_{male} + \beta_3 X_{alcohol} + \beta_4 D_{year} + \beta_5 D_{region} + \beta_6 X_{tobacco} * D_{male} + \beta_7 X_{alcohol} * D_{male} + \epsilon$$

where

- Y is the dependent variable life expectancy at age 60
- β_0 represents the intercept of the model, which is the expected life expectancy when all other variables are zero
- β_1 represents the change in expected life expectancy for a one-unit change in the variable prevalence of tobacco $X_{tobacco}$ when other variables are held constant

- β_2 represents the average difference in expected life expectancy between male and female when other variables are held constant
- β_3 represents the change in expected life expectancy for a one-unit change in the variable alcohol consumption $X_{alcohol}$ when other variables are held constant
- β_4 represents a matrix of β for each year in the variable **year**
- β_5 represents a matrix of β for each region in the variable **region**
- β_6 represents the average difference in the change of expected life expectancy for a one-unit change in $X_{tobacco}$ for male and female
- β_7 represents the average difference in the change of expected life expectancy for a one-unit change in $X_{alcohol}$ for male and female
- $X_{tobacco}$ represents the variable **prevalence_of_tobacco**
- D_{male} represents the variable **sex** with baseline (0) to be 0
- $X_{alcohol}$ represents the variable **alcohol_consumption**
- D_{year} refers to the variable **year** with 3 dummy variables
- D_{region} refers to the variable **region** with 6 dummy variables
- ϵ is the random error

For the reduced model, we just removed the interaction terms, that is the predictors with coefficient β_6 and β_7 in the full model.

The linear model will generate the best estimates for parameters β_i and D_j that minimize the residual sum of squares (RSS). After getting the best fit of the model, we need to implement model validation to make sure the assumptions for the model hold, that is linearity, homoscedasticity of errors, independence of errors, and influential observations. The model validation is done in the Section 7. It seems that assumptions for both models holds.

4 Result

Table 3 shows a summary for both models with the listed values representing the estimates of parameters and values in brackets representing the standard error for the estimate of the parameter. The two models produced different but close estimates of parameters. It is noticeable that the direction of correlation between alcohol assumption and life expectancy is different in two models, positive correlation in the full model and negative in the reduced model. This is consistent with data visualization in Figure 2 that fitted dash lines within two sex groups show an increasing trend but green fitted line for the whole data shows a decreasing trend.

For the full model, according to the p-values for each parameters, the variable **sex** is not significant at significance level $\alpha = 0.05$ since it has a p-value of 0.15, indicating the variable **sex** does not have a significant effect on the expected life expectancy at age 60. However, since the interaction term between age and alcohol consumption is statistically significant with p-value smaller than the significance level $\alpha = 0.05$, we should still include the separate parameter **sex**.

For the reduced model, according to the p-values for each parameters, the variable `sex` is also insignificant with p-value equal 0.3, larger than that in the full model. This may be due to the interaction effect. Also, it is noticeable that the dummy variable for Eastern Mediterranean region is also statistically insignificant with p-value equal to 0.08 at significance level $\alpha = 0.05$, however, if we set the significant level at $\alpha = 0.1$, then the predictor becomes significant. It all depends on our choice of significance level.

To justify our choice of interaction terms and choose a better model, we will compare AIC, BIC, RMSE, and adjusted R^2 values for both models. We found that the rounded AIC for the full model is 45 with 64 for the reduced model, the rounded BIC for the full model is 77 with 92 for the reduced model, the RMSE for the full model is 0.27 with 0.33 for the reduced model, and the adjusted R^2 for the full model is 0.98 with 0.97 for the reduced model. Since we want AIC, BIC, and RMSE to be as small as possible while adjusted R^2 to be as large as possible, the full model with interaction terms could be a better fit for the data. It has higher predictive performance without loss of general interpretability. Generally, the prevalence of tobacco use has a negative linear relationship with the life expectancy at age 60, the alcohol per capita consumption has a positive linear relationship with the life expectancy at age 60 while sex itself does not have a significant effect on the life expectancy at age 60.

5 Discussion

5.1 Brief summary

In this paper, we have conducted a analysis of life expectancy at age 60 with several important predictors including prevalence of tobacco use, alcohol per capita consumption, sex at birth, year in 2000, 2010, 2015, 2019, and WHO regions. We have found potential linear relationships between each variable and life expectancy through data visualization. To quantitatively understand how these factors affect life expectancy at age 60, we have utilized multiple linear regression models. We have fitted two model, a full model with interaction effect between sex and prevalence of tobacco as well as the interaction effect between sex and alcohol consumption, and a reduced model with five independent variables but no interaction terms. The results show that the full model have a better fit with smaller root mean squared errors (RMSE), AIC, and BIC, and a larger adjusted R^2 .

5.2 Findings

According to Table 3 and Figure 1, when other variables are held constant, the increase of the prevalence of tobacco use will decrease the life expectancy at age 60. This is consistent with the fact that tobacco is a health risk factor for cardiovascular and respiratory diseases and thus people with a habit of use tobacco products not only cigarettes will have a lower life expectancy at age 60 (Doll et al. 2004). Combined with Table 1 and Table 2, we could find

Table 3: Summary of two linear regression models

	With interaction	Without interaction
(Intercept)	16.569 (0.351)	17.523 (0.189)
prevalence_of_tobacco	−0.059 (0.013)	−0.035 (0.010)
sexmale	−0.536 (0.368)	−0.436 (0.406)
alcohol_consumption	0.432 (0.164)	−0.225 (0.031)
as.factor(year)2010	0.956 (0.133)	1.104 (0.154)
as.factor(year)2015	1.404 (0.142)	1.591 (0.163)
as.factor(year)2019	1.761 (0.147)	1.879 (0.172)
regionAmericas	5.349 (0.288)	6.037 (0.235)
regionEastern Mediterranean	1.542 (0.349)	0.408 (0.233)
regionEurope	4.798 (0.455)	6.047 (0.328)
regionGlobal	3.761 (0.201)	3.884 (0.234)
regionSouth-East Asia	2.386 (0.278)	1.768 (0.290)
regionWestern Pacific	4.268 (0.209)	4.590 (0.235)
prevalence_of_tobacco \times sexmale	0.024 (0.012)	
sexmale \times alcohol_consumption	−0.509 (0.123)	
Num.Obs.	56	56
R ²	0.988	0.982
R ² Adj.	0.984	0.977
AIC	45.0	63.9
BIC	77.4	92.2
Log.Lik.	−6.510	−17.937
RMSE	0.27	0.33

that although the prevalence of tobacco in general decreases across years, it has an uneven region distribution and notably for Africa, it has the lowest prevalence of tobacco but relatively low life expectancy.

Other than tobacco products, the effect of alcohol on health reveals a complex relationship. The full model we choose suggested that the alcohol consumption is associated with increased life expectancy at age 60. However, in the reduced model, when we do not consider sex impact on alcohol consumption, it shows a negative association. This is consistent with the data visualization that shows a Simpson’s Paradox. Based on our justification of taking into account sex interaction with alcohol consumption, it indicates an overall positive association. Additionally, evidence-based study has found that modest drinking has been approved to be protective against some diseases, but in most cases, alcohol could be seen as a health risk factor (Liu et al. 2022). In our model, we did not specify alcohol consumption level, so it could be further investigated by creating a new dummy variable indicating different drinking levels.

According to Table 1 and Table 2, the values that are not consistent with our model result indicate that there might be other factors that influence the life expectancy at age 60 for a specific year or region. For different regions, the difference in economic development, medical treatments, and environmental conditions could be potential influencing factors (Rahman, Rana, and Khanam 2022). However, if we introduced these factors into our model, we have to consider the problem of multicollinearity that the independent variables correlates since these socio-economic factors and environmental factors could also contribute to the prevalence of tobacco and alcohol consumption.

5.3 Weakness and future work

One weakness of this paper is about data quality. For the life expectancy dataset, the limitation comes from lack of complete and reliable mortality data for some civil registrations. This may cause the estimation of life expectancy biased since additional model estimations were used instead of real data. For the prevalence of tobacco dataset, one limitation is that some countries or regions do not have reliable data and estimation using Bayesian models held many assumption that may not actually satisfy in real cases. Another limitation is the data source. Since the estimation was made by population-based surveys, self-report biases exist that people may tend to hide some smoking habits if they think it is “not good”. For alcohol consumption dataset, same problem for missing several data points for some regions or in a specific year. Additionally, the unrecorded consumption calculation and tourists consumption were calculated from model estimations and several assumption were made during the process that may not be actually valid.

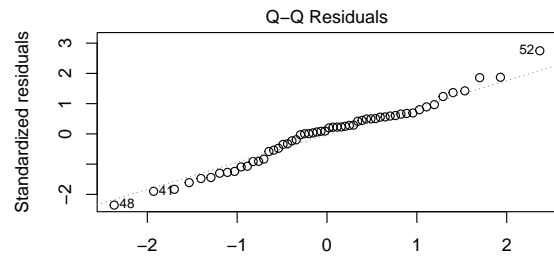
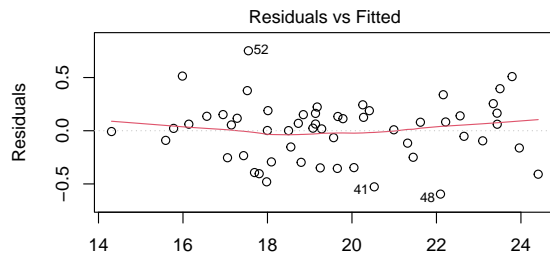
Another weakness is about the model. Since we decide to employ linear regression models based on data visualization, we assume linear relationships between variables, which may ignore some non-linear relationship.

For further studies, we could increase our sample size and split into training and testing data to better assess the model performance. We could also use cross validation to find a more robust estimate. Also, additional factors could be considered as well as more complex models that could capture potential non-linear relationships

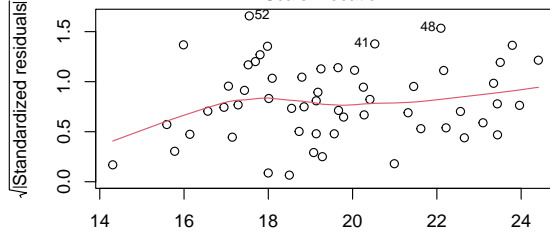
6 Reference

2020. *World Health Organization*. World Health Organization. <https://apps.who.int/gho/data/view.main.SDG2016LEXREGv?lang=en>.
- . 2022. *World Health Organization*. World Health Organization. <https://apps.who.int/gho/data/view.main.SDG3aWHOREGv?lang=en>.
- . 2023. *World Health Organization*. World Health Organization. <https://apps.who.int/gho/data/view.main.SDG35WHOREGv?lang=en>.
- Arel-Bundock, Vincent. 2022. “modelsummary: Data and Model Summaries in R.” *Journal of Statistical Software* 103 (1): 1–23. <https://doi.org/10.18637/jss.v103.i01>.
- Doll, Richard, Richard Peto, Jillian Boreham, and Isabelle Sutherland. 2004. “Mortality in Relation to Smoking: 50 Years’ Observations on Male British Doctors.” *BMJ* 328 (7455): 1519. <https://doi.org/10.1136/bmj.38142.554479.ae>.
- Firke, Sam. 2021. *Janitor: Simple Tools for Examining and Cleaning Dirty Data*. <https://CRAN.R-project.org/package=janitor>.
- Liu, Yen-Tze, June Han Lee, Min Kuang Tsai, James Cheng-Chung Wei, and Chi-Pang Wen. 2022. “The Effects of Modest Drinking on Life Expectancy and Mortality Risks: A Population-Based Cohort Study.” *Nature News*. Nature Publishing Group. <https://www.nature.com/articles/s41598-022-11427-x#citeas>.
- Müller, Kirill. 2020. *Here: A Simpler Way to Find Your Files*. <https://CRAN.R-project.org/package=here>.
- R Core Team. 2022. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Rahman, Mohammad Mafizur, Rezwanul Rana, and Rasheda Khanam. 2022. “Determinants of Life Expectancy in Most Polluted Countries: Exploring the Effect of Environmental Degradation.” *PloS One*. U.S. National Library of Medicine. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8782287/>.
- Roser, Max, Esteban Ortiz-Ospina, and Hannah Ritchie. 2013. “Life Expectancy.” *Our World in Data*. <https://ourworldindata.org/life-expectancy#:~:text=Life%20expectancy%20is%20the%20key,of%20death%20in%20a%20population>.
- Wei, Ying. 2021. “Colors in r.” University of Columbia. <http://www.stat.columbia.edu/~tzheng/files/Rcolor.pdf>.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain François, Garrett Golemund, et al. 2019. “Welcome to the tidyverse.” *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.
- World Health Organization. 2023. *World Health Organization*. [https://www.who.int/data/gho/data/indicators/indicator-details/GHO/life-expectancy-at-age-60-\(years\)](https://www.who.int/data/gho/data/indicators/indicator-details/GHO/life-expectancy-at-age-60-(years)).
- Xie, Yihui. 2014. “Knitr: A Comprehensive Tool for Reproducible Research in R.” In *Implementing Reproducible Computational Research*, edited by Victoria Stodden, Friedrich Leisch, and Roger D. Peng. Chapman; Hall/CRC. <http://www.crcpress.com/product/isbn/9781466561595>.

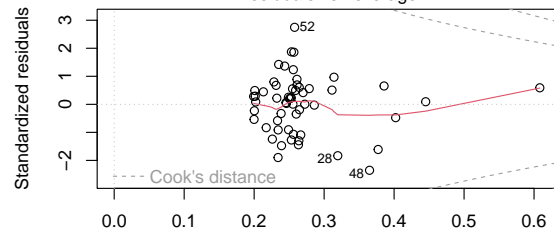
7 Appendix



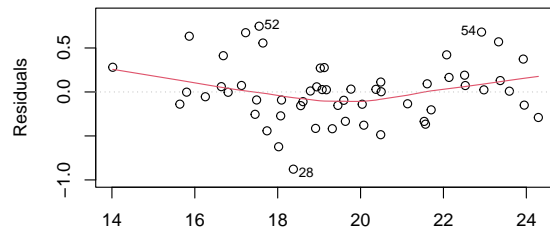
Fitted values
 $m(\text{life_expectancy} \sim \text{prevalence_of_tobacco} * \text{sex} + \text{alcohol_consumption})$
 Scale-Location



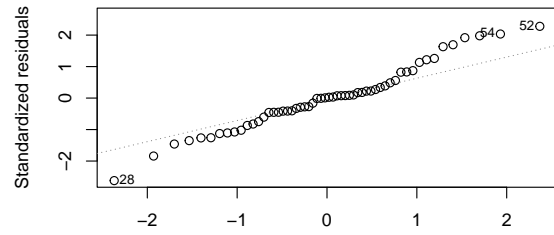
Theoretical Quantiles
 $m(\text{life_expectancy} \sim \text{prevalence_of_tobacco} * \text{sex} + \text{alcohol_consumption})$
 Residuals vs Leverage



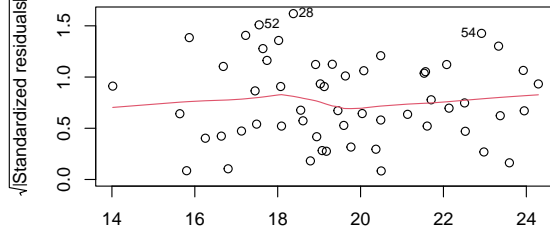
Fitted values
 $m(\text{life_expectancy} \sim \text{prevalence_of_tobacco} * \text{sex} + \text{alcohol_consumption})$
 Residuals vs Fitted



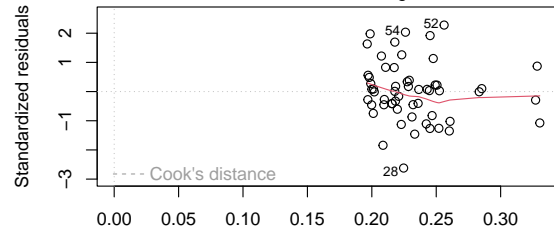
Leverage
 $m(\text{life_expectancy} \sim \text{prevalence_of_tobacco} * \text{sex} + \text{alcohol_consumption})$
 Q-Q Residuals



Fitted values
 $m(\text{life_expectancy} \sim \text{alcohol_consumption} + \text{prevalence_of_tobacco} + \text{region})$
 Scale-Location



Theoretical Quantiles
 $m(\text{life_expectancy} \sim \text{alcohol_consumption} + \text{prevalence_of_tobacco} + \text{region})$
 Residuals vs Leverage



Fitted values
 $m(\text{life_expectancy} \sim \text{alcohol_consumption} + \text{prevalence_of_tobacco} + \text{region})$