# Air Pollution and Respiratory Difficulties

P408 - Main Body Word Count: 1974

December 17, 2025

# Contents

## Declaration of AI use

In creating this work, I have used AI tools for the following tasks:

1. Literature search and ideation: not used

2. Spelling and grammar: Overleaf has in-built spell-check and AI (TexGPT). TexGPT highlights parts of the text in the editor, though I have not used it.

3. Improvements to writing: not used

4. Coding: I used Claude AI for research on the stargazer package.

5. Plotting, formatting etc: I used Claude AI and ChatGPT Edu for help with latex formatting. In particular, for altering font size in tables, as well as controlling the number of significant figures in the output.

## 1   Introduction

This report analyses a dataset from a survey on air pollution in a South Asian city. The report begins with exploratory data analysis, which summarises the variables in the dataset and displays some high level associations between respiratory difficulties (henceforth, RD) and age, gender, physical activity, and air pollution. The next section concerns model selection, beginning with a Bernoulli generalised linear regression model (henceforth, GLM) of RD on all the regressors mentioned above and interactions between exposure to air pollution and the remaining covariates. I then present the results of model selection, finding that all criteria (likelihood ratio tests, Wald tests, and Information Criteria) support dropping all the main effects of the regressors (with the exception of activity), while keeping all the interactions. Invoking the hierarchy principle, I expose the fragility of the model to non-substantive changes in the data, such as recoding the air pollution binary variable. Given the main aim of the model is interpretive, I conclude that the preferred model must include all baselines. I then present model diagnostics based on plots, s.a., residual plots, QQ plots, and leverage/ Cook's distance plots. While some potential outliers are identified, there is no apparent reason to reject them. The following section then provides an interpretation of the findings, including an evaluation of the odds ratios for each variable and the respective confidence intervals. The final section predicts the probability of developing RD for two new data points, one within-sample and one out-of-sample. The predictions are paired with the respective confidence interval and a brief discussion of out of sample reliability.

## 2   Exploratory Data Analysis

The dataset comprises 285 observations of one binary variable (Respiratory Difficulty) and four regressors (Activity, Gender, Air Pollution, and Age). The first three regressors are uncontroversially categorical (binary), while the last is treated as continuous. In fact, age takes only four values (25, 35, 45, 55) representing 10-year age bands, rather than the age of the participants. This could substantiate a categorical interpretation of age as a factor, which is shown in the appendix (Table 12). Since part of the model's objective is to predict the probability of RD for an out-of-sample individual (85 year old), age will be treated as continuous for the remainder of the analysis.

The following figures illustrate the distribution of the data. First, Figure 1 shows the marginal distribution of each variable and a histogram of age.

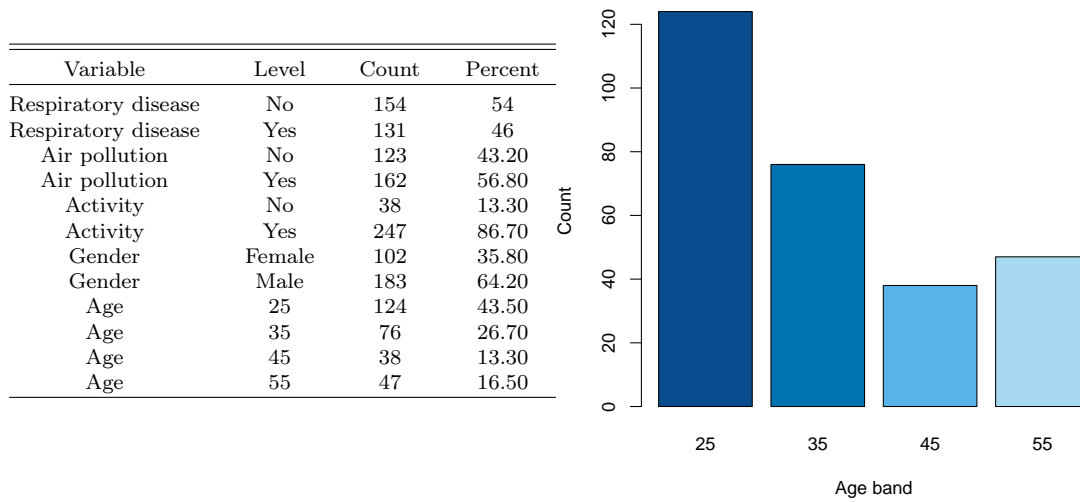| Variable | Level | Count | Percent |
|---|---|---|---|
| Respiratory disease | No | 154 | 54 |
| Respiratory disease | Yes | 131 | 46 |
| Air pollution | No | 123 | 43.20 |
| Air pollution | Yes | 162 | 56.80 |
| Activity | No | 38 | 13.30 |
| Activity | Yes | 247 | 86.70 |
| Gender | Female | 102 | 35.80 |
| Gender | Male | 183 | 64.20 |
| Age | 25 | 124 | 43.50 |
| Age | 35 | 76 | 26.70 |
| Age | 45 | 38 | 13.30 |
| Age | 55 | 47 | 16.50 |

Figure 1: Counts and proportion breakdown of each category, with visual representation for age

Table 1 then breaks down the dataset by the values of the regressors. Of 32 possible combinations of the regressors, only 24 are represented in the data, and the proportion of each subset experiencing RD varies widely. The proportion of subjects experiencing RD tends to be higher in those exposed to air pollution.

| Air Pollution | Activity | Gender | Age | Resp. diff. proportion | Count |
|---|---|---|---|---|---|
| No | No | Male | 25 | 0.75 | 4 |
| No | No | Male | 35 | 0.00 | 1 |
| No | No | Male | 45 | 1.00 | 1 |
| No | Yes | Female | 25 | 0.11 | 18 |
| No | Yes | Female | 35 | 0.17 | 12 |
| No | Yes | Female | 45 | 0.25 | 4 |
| No | Yes | Female | 55 | 0.33 | 3 |
| No | Yes | Male | 25 | 0.26 | 34 |
| No | Yes | Male | 35 | 0.16 | 19 |
| No | Yes | Male | 45 | 0.00 | 8 |
| No | Yes | Male | 55 | 0.05 | 19 |
| Yes | No | Female | 25 | 1.00 | 6 |
| Yes | No | Female | 35 | 1.00 | 2 |
| Yes | No | Male | 25 | 0.25 | 16 |
| Yes | No | Male | 35 | 0.20 | 5 |
| Yes | No | Male | 45 | 0.00 | 3 |
| Yes | Yes | Female | 25 | 0.83 | 18 |
| Yes | Yes | Female | 35 | 0.94 | 16 |
| Yes | Yes | Female | 45 | 1.00 | 5 |
| Yes | Yes | Female | 55 | 0.94 | 18 |
| Yes | Yes | Male | 25 | 0.46 | 28 |
| Yes | Yes | Male | 35 | 0.57 | 21 |
| Yes | Yes | Male | 45 | 0.65 | 17 |
| Yes | Yes | Male | 55 | 1.00 | 7 |

Table 1: Counts and proportion affected by respiratory difficulty for each combination of the regressors

The following barplots display the conditional probability of being affected by RD given the values of the regressors.

$$\mathbf{P}(\texttt{respd} = 1 \,|\, X = x) \quad \forall X \in \{\texttt{gender, activity, airpollution, age}\}$$
$$\mathbf{P}(\texttt{respd} = 0 \,|\, X = x) \quad \forall X \in \{\texttt{gender, activity, airpollution, age}\}$$

The visual shows that air pollution and gender are strongly associated with RD, with females

much more likely to suffer from RD, and similarly for subjects exposed to air pollution. The plot displays limited dependence on activity and a weak positive association between age and RD.
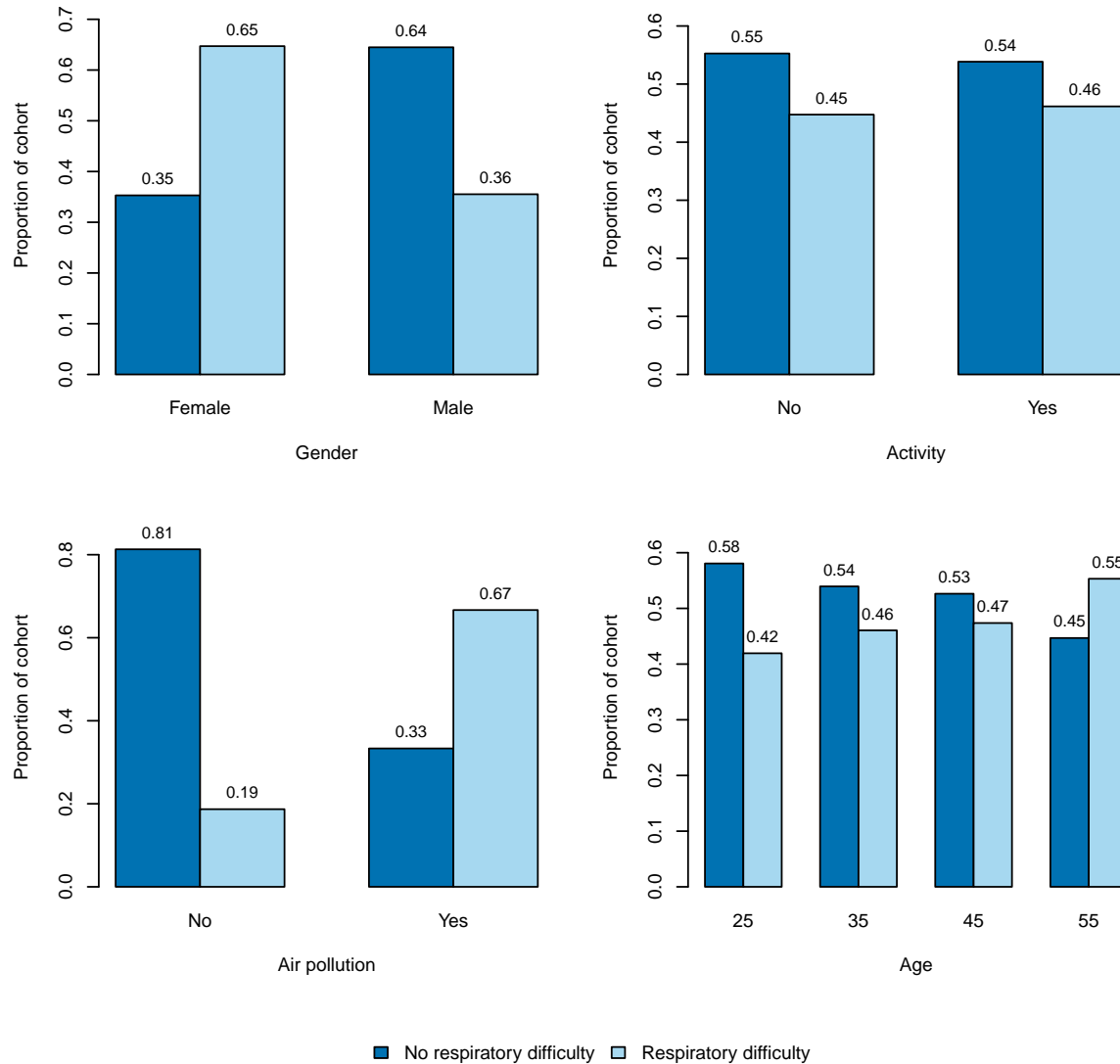


Figure 2: Incidence of respiratory difficulty by regressor levels

Since air pollution has the strongest intuitive link with RD, consider further conditioning the distributions on air pollution, so plotting:

$$\mathbf{P}(\texttt{respd} = 1 \,|\, X = x, \texttt{ air pollution} = i) \quad \forall X \in \{\texttt{gender, activity, age}\}, i \in \{0, 1\}$$
$$\mathbf{P}(\texttt{respd} = 0 \,|\, X = x, \texttt{ air pollution} = i) \quad \forall X \in \{\texttt{gender, activity, age}\}, i \in \{0, 1\}$$

This illuminates the structure of the dataset, showing that RD appears almost conditionally independent of gender when air pollution is absent, and displaying a strong association with activity both when present and when absent, as well as a strong positive association with age when present. This is in line with our intuition, since being physically active should improve health outcomes in general, but would also expose subjects more to pollutants, leading to diverging effects on RD depending on pollution. Similarly, we expect older people to be more vulnerable to the effects of air pollution.
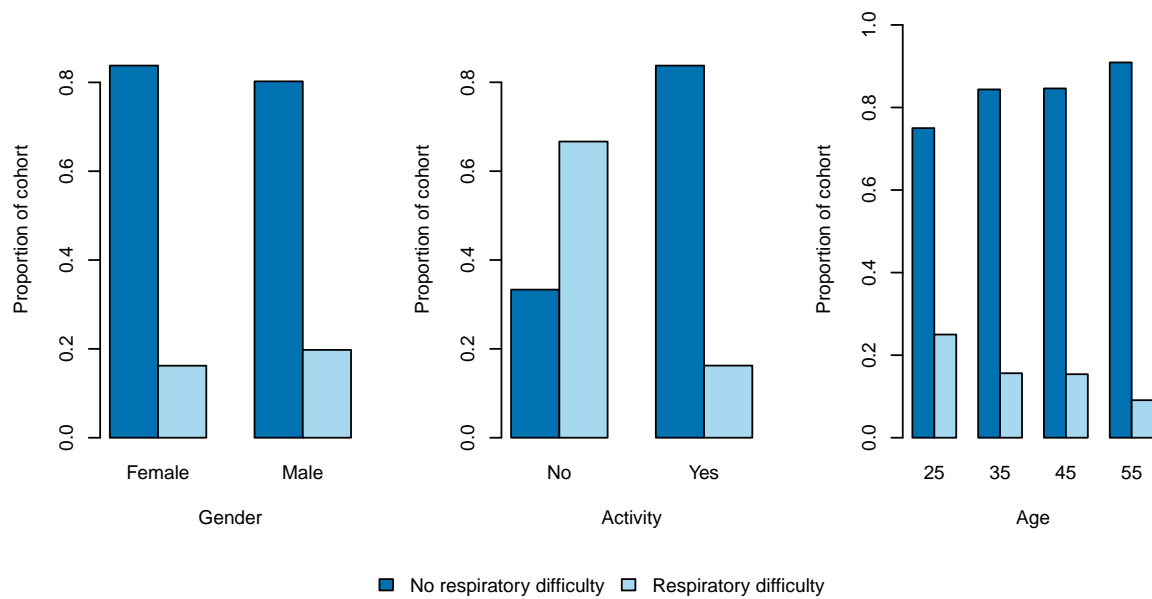
Figure 3: Incidence of respiratory difficulty in subjects not exposed to air pollution by regressor levels
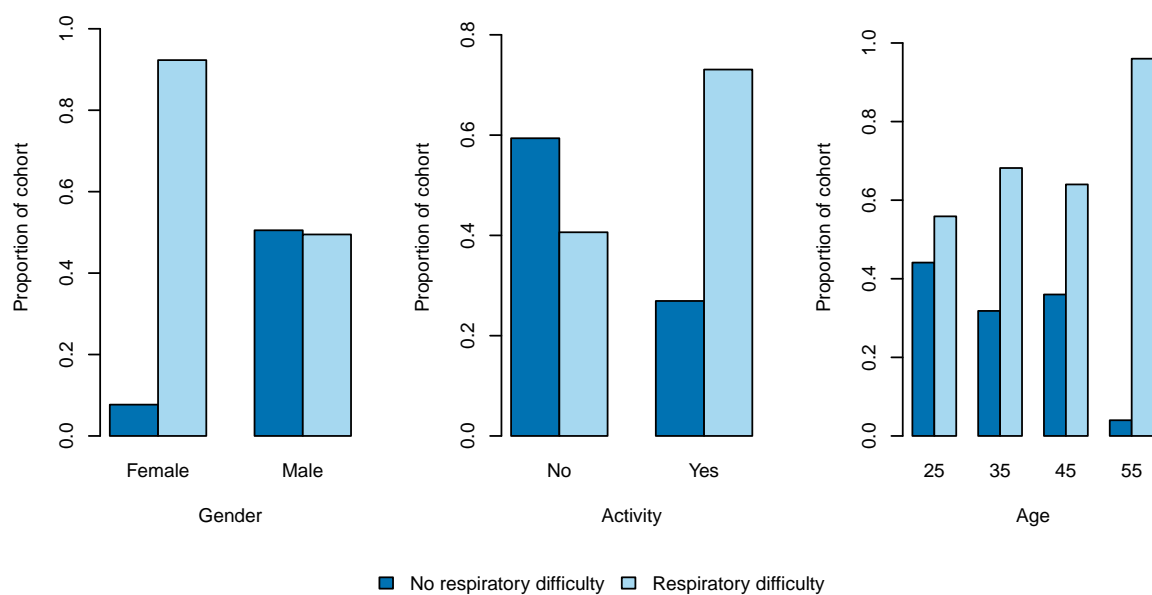


Figure 4: Incidence of respiratory difficulty in subjects exposed to air pollution by regressor levels

These plots are strongly informative, since they imply that interactions between air pollution and gender, activity, and age must be included to effectively capture the structure of the data.

# 3   Modelling

I begin with a Bernoulli GLM, regressing RD on all the regressors and including the interactions between air pollution and age, gender, and activity. This gives the linear predictor below, where the fitted value is $\mu = \sigma(\eta)$[1]

$$\eta = x^T\beta = \beta_0 + \beta_1\,\texttt{activity} + \beta_2\,\texttt{airpollution} + \beta_3\,\texttt{gender} + \beta_4\,\texttt{age}$$
$$+ \beta_5\,\texttt{activity} \times \texttt{airpollution} + \beta_6\,\texttt{gender} \times \texttt{airpollution}$$
$$+ \beta_7\,\texttt{age} \times \texttt{airpollution}$$

This is also the **chosen model**, as justified below. Table 2 then provides the coefficients and standard errors, as well as the p-values for the Wald statistic. All the interaction terms are highly significant under a Wald test, while all the main effects are insignificant.[2]

|  | *Dependent variable:* | |
| --- | --- | --- |
|  | Respiratory difficulty | p-val (Wald) |
| Air pollution (Yes) | −1.55 (1.50) | p = 0.31 |
| Age | −0.04 (0.02) | p = 0.16 |
| Male | 0.11 (0.55) | p = 0.85 |
| Active | −2.17* (0.93) | p = 0.02 |
| Air pollution × Age | 0.08* (0.03) | p = 0.02 |
| Air pollution × Male | −2.56*** (0.76) | p = 0.001 |
| Air pollution × Active | 3.17** (1.04) | p = 0.003 |
| Constant | 1.65 (1.24) | p = 0.19 |
| Observations | 285 | |
| Log Likelihood | −133.06 | |
| Akaike Inf. Crit. | 282.12 | |
| *Note:* | *p<0.05; **p<0.01; ***p<0.001 | |

Table 2: Logistic Regression of Respiratory difficulty

Since the Wald test and the likelihood ratio test (LRT)[3] are asymptotically equivalent, we expect to reject the hypotheses that (1) the model is useless and (2) interactions should be

---

[1]The Bernoulli GLM takes the following form (Dobson & Barnett, 2018 [1]): given the canonical link function $g$, we have:

$$y = \texttt{respd} \sim Bernoulli(\mu)$$
$$\eta = g(\mu(\texttt{respd})) = log\left(\frac{\mu(\texttt{respd})}{1 - \mu(\texttt{respd})}\right)$$
$$\Rightarrow \mu = g^{-1}(\eta) = \sigma(\eta) = \sigma(x^T\beta) = \frac{e^{x^T\beta}}{1 + e^{x^T\beta}}$$

where the linear predictor $\eta$ is defined as above.

[2]The Wald statistic is defined as follows:

$$w_i = \frac{\hat{\beta}_i}{SE(\hat{\beta}_i)} \overset{a}{\sim} N[0,1]$$

where $\hat{\beta}_i$ is the estimate obtained by iteratively reweighted least squares (IRLS) and $SE(\hat{\beta}_i)$ is the respective standard error. This is given by the respective entry of the inverse of the observed information, which for the Bernoulli model with canonical link above coincides with the Fisher information $J_{ii}^{-1} = I_{ii}^{-1}$.

[3]The LRT is derived as follows. Given the regressors (including the interactions) $\mathbf{x}$, we can partition the set and the coefficients into:

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_{(1)} \\ \mathbf{x}_{(2)} \end{pmatrix}; \quad \beta = \begin{pmatrix} \beta_{(1)} \\ \beta_{(2)} \end{pmatrix},$$

where $\beta_2$ is a $k$-dimensional vector. We then define the maximum likelihood estimator (MLE) of the restricted

dropped. Table 3 provides the results of the LRT, with the likelihood ratio statistic (LRS) and p-value given for each model restriction. The "drop" column indicates which regressor/interaction was removed in the test. The first test ("null") checks whether the model is useless, meaning it compares the model to the intercept-only model. The p-values in Table 3 provide the probability of obtaining such values of the LRT under the null.

| Full Model | Restricted Model | Drop | Likelihood Ratio Statistic | P values |
|---|---|---|---|---|
| glm1 | glm0 | All (null test) | 127.11 | < 0.001 |
| glm1 | glm2 | Age x Air Pollution | 6.79 | 0.01 |
| glm1 | glm3 | Activity x Air Pollution | 9.69 | 0.002 |
| glm1 | glm4 | Gender x Air Pollution | 12.95 | < 0.001 |
| glm1 | glm5 | Air Pollution, Age, Gender | 2.49 | 0.48 |
| glm1 | glm6 | Age | 2.25 | 0.13 |
| glm6 | glm7 | Gender | 0 | 1.00 |
| glm7 | glm8 | Activity | 7.09 | 0.01 |

Table 3: Results of nested Likelihood Ratio Tests

The LRTs confirm what was expected from the Wald tests, so we reject (1) that the model is useless ("null" test) and (2) any models that drop interactions at the 1% significance level. Similarly, the LRT agrees with the Wald test that none of the main effects except activity are significant. Moreover, the joint likelihood ratio test for the main effects strongly supports dropping them, with a p-value of 0.48. The last three tests also consider dropping age, gender, and activity sequentially, again showing that we cannot reject the null that the main effects are 0 with the sole exception of activity (as in the Wald test). This leads to the reduced model below, with estimates and Wald p-values in Table 4.

$$\eta = x^T\beta = \beta_0 + \beta_1\,\texttt{activity} + \beta_2\,\texttt{airpollution} \times \texttt{activity}$$
$$+ \beta_3\,\texttt{airpollution} \times \texttt{gender} + \beta_4\,\texttt{airpollution} \times \texttt{age}$$

| | Dependent variable: | |
|---|---|---|
| | Respiratory Difficulty | p-val (Wald) |
| Active | −2.04** (0.64) | p = 0.002 |
| Air pollution × Age | 0.04* (0.02) | p = 0.03 |
| Air pollution × Male | −2.53*** (0.51) | p < 0.0001 |
| Air pollution × Active | 3.00*** (0.74) | p = 0.0001 |
| Constant | 0.40 (0.59) | p = 0.51 |
| Observations | 285 | |
| Log Likelihood | −134.31 | |
| Akaike Inf. Crit. | 278.61 | |
| Note: | *p<0.05; **p<0.01; ***p<0.001 | |

Table 4: Logistic Regression of Respiratory Difficulty

This model (glm5 in Table 5) is not only the most parsimonious, but it also outperforms the initial model and all the other nested models that were tested under both information criteria, with little change in the value of the $R^2$ KL. This would provide strong evidence for selecting model 5.

---

model under the null $beta_2 = 0$ as $\widetilde{\beta}_{(1)}$, with maximised log-likelihood given by $\ell\left(\widetilde{\beta}_{(1)}; y\right)$. Similarly, the maximised log-likelihood of the unrestricted model is $\ell\left(\widehat{\beta}_{(1)}; y\right)$, giving the LRT statistic:

$$\Lambda(y) = 2\left(\ell(\widehat{\beta}; Y) - \ell\left(\tilde{\beta}_{(1)}; y\right)\right) \xrightarrow{d} \chi^2_k$$

The distribution converges to $\chi^2_k$ under the null, so we reject the null for large values of the statistic relative to the degrees of freedom $k$.

| Model | Description | Residual Deviance | AIC | BIC | $R^2$ Kullback Leibler |
|-------|-------------|-------------------|-----|-----|------------------------|
| glm1 | Full model | 266.12 | 282.12 | 311.34 | 0.32 |
| glm2 | Drop Age x Air Pollution | 272.92 | 286.92 | 312.48 | 0.31 |
| glm3 | Drop Activity x Air Pollution | 275.82 | 289.82 | 315.38 | 0.30 |
| glm4 | Drop Gender x Air Pollution | 279.08 | 293.08 | 318.64 | 0.29 |
| glm5 | Drop Air Pollution, Age, Gender | 268.61 | 278.61 | 296.88 | 0.32 |
| glm6 | Drop Age | 268.37 | 282.37 | 307.94 | 0.32 |
| glm7 | Drop Age, Gender | 268.37 | 280.37 | 302.29 | 0.32 |
| glm8 | Drop Activity, Age, Gender | 275.46 | 285.46 | 303.72 | 0.30 |

Table 5: Information Criteria, Deviances and $R^2$ KL

However, dropping the main effects represents a violation of the principle of marginality[4] which holds that we should not drop the main effect of a variable when including the interaction term. This is crucial here, given that:

1. The main purpose of the model is interpretation. Since the coefficient on the interaction $x_1 \times x_2$ indicates the change in the effect of $x_1$ on the log-odds when $x_2$ is present, interpretation is difficult when there is no main effect of $x_1$ (and vice versa)

2. The data is binary, so arbitrary encoding decisions can have serious effects on the significance of the coefficients.

The below exemplifies the second pitfall. The toy model here redefines[5] `airpollution =` "Yes" as 0 and "No" as 1, resulting in:

| | Dependent variable: | |
|---|---|---|
| | Respiratory difficulty | p-val (Wald) |
| Air pollution (No) | 1.55 (1.50) | p = 0.31 |
| Age | 0.05* (0.02) | p = 0.04 |
| Male | −2.45*** (0.52) | p = 0.0000 |
| Active | 1.00* (0.48) | p = 0.04 |
| Air pollution × Age | −0.08* (0.03) | p = 0.02 |
| Air pollution × Male | 2.56*** (0.76) | p = 0.001 |
| Air pollution × Active | −3.17** (1.04) | p = 0.003 |
| Constant | 0.11 (0.84) | p = 0.90 |
| Observations | 285 | |
| Log Likelihood | −133.06 | |
| Akaike Inf. Crit. | 282.12 | |
| *Note:* | *p<0.05; **p<0.01; ***p<0.001 | |

Table 6: Toy model – Logistic Regression of Respiratory difficulty

The LRT of the nested model without the main effects for age, gender, and air pollution now yields LRS 49.3, with p-value $< 0.001$, so we would **not** prefer the reduced model after this minor non-substantive re-encoding. This corroborates interpretive fragility, since reencoding `airpollution` now provides a model where `gender` alone is significant. Given the goal of interpretation, we must opt to keep the main effects despite statistical insignificance. The remaining parts of the analysis for the reduced model can be found in the appendix, and are pointed to in the relevant sections.

---

[4]Also known as the hierarchy principle

[5]The binary variables are encoded alphabetically, by the R default. This means that, in the context of regression, gender assigns 0 to F and 1 to M, while activity and air pollution assign 0 to No and 1 to Yes.

# 4 Model Diagnostics

This section assesses the quality of model fit using residual plots to determine if any large outliers may be skewing the analysis, and if patterns in the residuals display obvious signs of misfit.
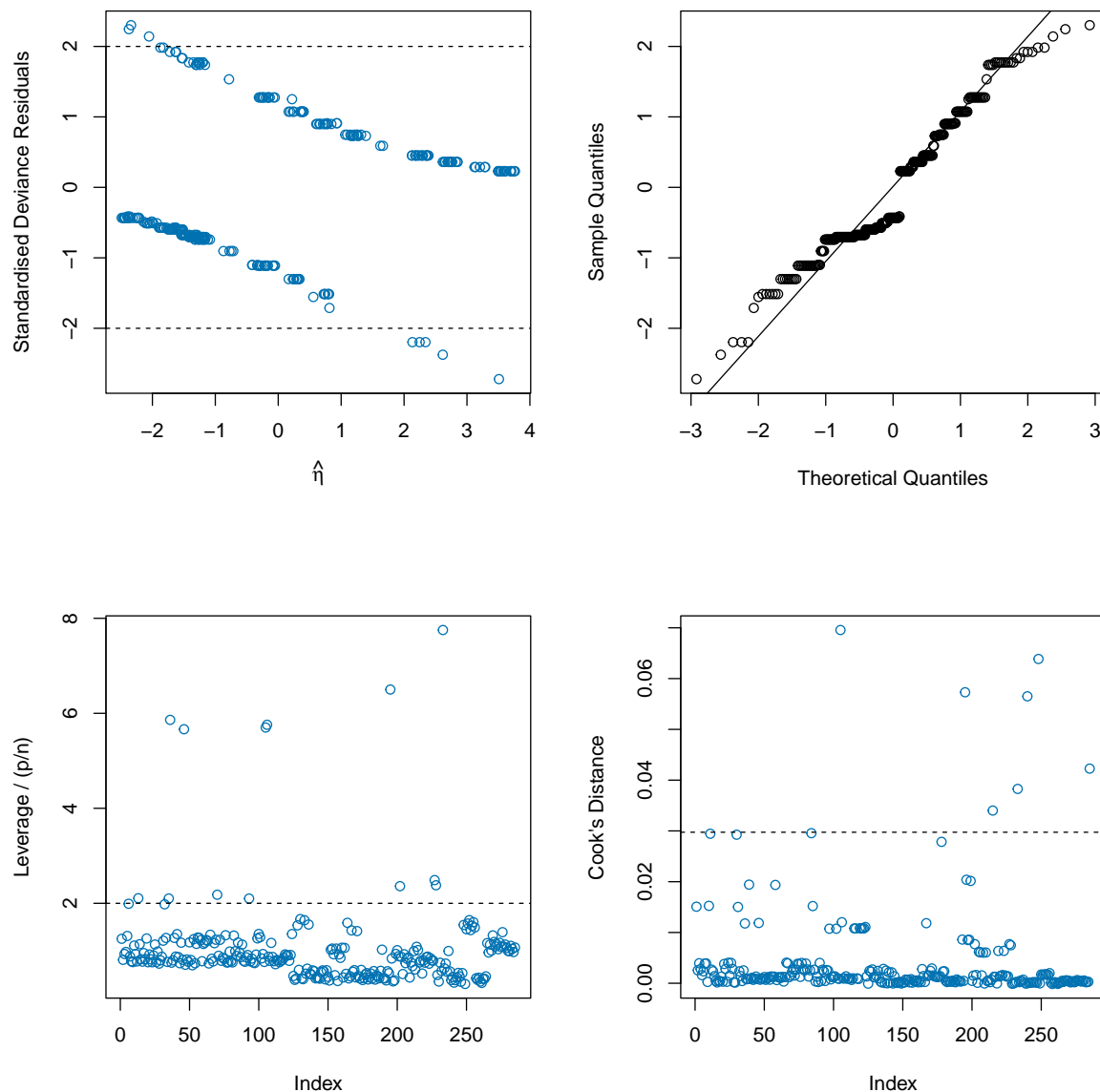


Figure 5: In order Left to Right, Top to Bottom: residuals vs fitted linear predictor, residual QQ-plot, leverage plot, Cook's distance plot

The first plot displays the values of the standardised deviance residuals against the fitted values of the linear predictor $\eta$. A plot indexing against the fitted response $\mu$ is in the appendix, but the linear predictor is preferable for the present analysis (Faraway, 2016 [2]). For a Bernoulli model, residual normality is not expected as for Poisson or Binomial with large counts. However, the standardisation means the residuals have approximately unit variance, allowing us to identify residuals outside of $(-2, 2)$ as potential outliers. The obvious pattern in the plot is due to the binary nature of dependent variable, where the higher cluster will be residuals of positive deviations (respd = 1) and vice versa for the lower cluster. Eight observations exceed the $\pm 2$

threshold: not particularly alarming for 285 observations.[6]

Again, since the model is Bernoulli we do not expect the residuals to be approximately normal. However, we expect the QQ plot to be approximately centred around 0, have symmetric tails, and display no strong curvature. This is satisfied by the plot, with the separation in the middle due to the binary nature of `respd`.

The third plot displays leverage divided by the ratio of parameters to sample size (here 8/285).[7] Large observations will have a value greater than $2p/n$, displayed by the values above the dashed line in the third plot.

Finally, the fourth plot displays Cook's Distance (CD) - a measure of the influence of observations on the estimates.[8] Heuristically, potential outliers exceed $8/(n-2p)$, where the denominator comes from the degrees of freedom of the numerator after two estimations of the coefficient vector, and the 8 in the numerator is a conservative empirical threshold. [9]

We then isolate these observations, with scaled leverage and Cook's Distance above the respective thresholds (2, 0.030).

| Air Pollution | Age | Respiratory Difficulty | Activity | Gender | Leverage | Cook's D |
|---|---|---|---|---|---|---|
| No | 25 | No | No | Male | 5.775 | 0.070 |
| No | 35 | No | No | Male | 6.445 | 0.057 |
| No | 45 | Yes | No | Male | 7.870 | 0.038 |

Table 7: Observations with High Leverage and High Cook's Distance

Notably, all three are inactive males not exposed to air pollution - a combination with very few observations (Table 1). This low count naturally produces higher leverage, and there is no apparent theoretical justification to exclude these observations from the dataset. Accordingly, we proceed conservatively and leave the sample intact.

---

[6]In fact, if normality did hold, given a cutoff of $\pm 2$ we would expect just under 5% to exceed this threshold, which would be approximately 14 values.

[7]The leverage of observation i indicates the entry $h_{ii}$ of the model's hat matrix $H$:

$$H = W^{1/2} X \left( X^T W X \right)^{-1} X^T W^{1/2}$$

where $W$ is the weight matrix from the IRLS algorithm. The trace of the matrix is $p$, by the cyclic property, so the average leverage value is $p/n$. Thus, the heuristic $2p/n$ indicates a value with leverage twice the mean. For the scaled leverage, the threshold is naturally 2.

[8]CD is defined as

$$C_i = \frac{\left( \widehat{\beta} - \widehat{\beta}_{-i} \right)^T \left( X^T W X \right) \left( \widehat{\beta} - \widehat{\beta}_{-i} \right)}{p\hat{\phi}} = \frac{\left( \widehat{\beta} - \widehat{\beta}_{-i} \right)^T \left( X^T W X \right) \left( \widehat{\beta} - \widehat{\beta}_{-i} \right)}{8}$$

where $\hat{\phi}$ would estimate the scaling parameter, which is 1 for the Bernoulli GLM. This represents the average change in the estimated coefficient vector when removing one observation, scaled by the covariance.

[9]The threshold is based partly on parallelism with linear models, where Cook's Distance follows an F(p, n - p) distribution which has 5% tail bound $\leq 8$ for any p, n > 3. Notice, however, that the distribution does not apply to GLMs.

## 5   Interpretation (Odds Ratios)

In order to comprehensively interpret the model, we must introduce several concepts. First, we define the odds for the model as

$$\text{odds} = \frac{\mu}{1-\mu} = \frac{e^{x^T\beta}}{1+e^{x^T\beta}} \times \left(\frac{1}{1+e^{x^T\beta}}\right)^{-1} = e^{x^T\beta}$$

and thus the log-odds as

$$\text{log-odds} = log\left(\frac{\mu}{1-\mu}\right) = x^T\beta$$

Then, the effect of a unit change in any of the regressors on the log odds will be given by the respective coefficient. If interactions are present, the effect is the sum of the main and interaction coefficients. For instance, with age we will have:

$$\log(\text{odds}(\texttt{age} = x+1)) - \log(\text{odds}(\texttt{age} = x)) = \hat{\beta}_4 + \hat{\beta}_7\texttt{airpollution}$$

which simplifies to $\hat{\beta}_4$ when air pollution is absent and $\hat{\beta}_4 + \hat{\beta}_7$ when present. This provides a blueprint to interpret the coefficients of the model. The coefficients on the main effects represent the additive effect on the log-odds of the dependent variable from a unit increase in the regressor (when the interaction is absent). Similarly, the coefficients on the interaction terms represent the extra additive effect on the log-odds of a unit increase in the regressor when air pollution is present. We can provide 95% confidence intervals for the coefficients, and thus the effect on the log-odds, as follows:

| Variable | Estimate | 95 % CI |
|---|---|---|
| Intercept | 1.65 | (-0.78, 4.09) |
| Air pollution (Yes) | -1.55 | (-4.48, 1.39) |
| Age | -0.04 | (-0.08, 0.01) |
| Male | 0.11 | (-0.96, 1.18) |
| Active | -2.17 | (-3.99, -0.35) |
| Air pollution × Age | 0.08 | (0.02, 0.15) |
| Air pollution × Male | -2.56 | (-4.04, -1.07) |
| Air pollution × Active | 3.17 | (1.13, 5.22) |

Table 8: Logistic Regression Coefficients with 95% Confidence Intervals

Notice that, for the baseline effects on age, gender, and air pollution, the effect is **not** statistically significant, so we cannot offer a strong directional interpretation. We can, however, provide a directional interpretation for activity, with the log-odds of RD between $-3.99$ and $-0.35$ lower for active individuals not exposed to air pollution. Conversely, individuals exposed to air pollution who are active will have log-odds of RD between 1.13 and 5.22 higher than active individuals who were not exposed to air pollution (interaction). Similar interpretations follow for the other coefficients.

While the additive effect on the log-odds is the most immediate interpretation, the effect on the odds is more intuitive. Following from the example of `age` above, we can find the effect on the Odds, i.e., the Odds Ratio (OR)

$$OR = \frac{\text{odds}(\texttt{age}+1)}{\text{odds}(\texttt{age})} = e^{(\log(\text{odds}(\texttt{age}=x+1)) - \log(\text{odds}(\texttt{age}=x)))}$$

$$= e^{(\hat{\beta}_4 + \hat{\beta}_7\texttt{airpollution})}$$

Since it affects the ratio, the exponentiated effect is **multiplicative**, rather than additive as in the effect on the log-odds.

We now compute confidence intervals and estimates for the ORs. The coefficients are correlated (the covariance matrix is not diagonal), so we must compute the standard errors for the change in the log-odds for every combination of the interactions using the delta method. The calculation is illustrated in the below for age (10-year increase) and further down for air pollution, since gender and activity are structurally similar to age, *mutatis mutandis*, while air pollution is more involved. After the derivation, two tables provide an exhaustive summary of all the multiplicative effects on the odds.[10]

Let the coefficient vector be $\beta = (\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5, \beta_6, \beta_7)^\top$, where $\beta_4$ is the age main effect and $\beta_7$ is the age×airpollution interaction (in line with the Modelling section above). Write $\widehat{V} = \widehat{\mathrm{Var}}(\widehat{\beta})$ for the estimated covariance matrix.

First, consider a 10-year increase in age when air pollution is absent.

$$\text{The effect on the log-odds is:} \quad \Delta_{\text{age}|\text{airpollution}=0} = 10 \times \beta_4,$$

so $\widehat{\mathrm{SE}}$ is the square root of the corresponding diagonal entry, scaled by the age increase:

$$\widehat{\mathrm{SE}}\big(\Delta_{\text{age}|\text{airpollution}=0}\big) = 10 \times \sqrt{\widehat{V}_{44}}.$$

Then, consider the effect of a 10-year increase in age when air pollution is present.

$$\text{The effect on the log-odds is:} \quad \Delta_{\text{age}|\text{airpollution}=1} = 10 \times (\beta_4 + \beta_7).$$

By the delta method by the symmetry of the covariance matrix,

$$\widehat{\mathrm{SE}}\big(\Delta_{\text{age}|\text{airpollution}=1}\big) = \widehat{\mathrm{Var}}\big(\Delta_{\text{age}|\text{airpollution}=1}\big) = 10 \times \sqrt{\widehat{V}_{44} + \widehat{V}_{77} + 2\widehat{V}_{47}},$$

where $\widehat{V}_{jk}$ denotes the $(j,k)$-entry of $\widehat{V}$.

For any estimated linear effect $\widehat{\Delta}$ with estimated standard error $\widehat{\mathrm{SE}}(\widehat{\Delta})$, a (Wald) (1-$\alpha$) CI for the log-odds change is

$$\widehat{\Delta} \pm z_{1-\alpha/2}\,\widehat{\mathrm{SE}}(\widehat{\Delta}).$$

and thus we can derive the CI for the OR as

$$\mathrm{CI}_{\mathrm{OR}} = \exp\Big(\widehat{\Delta} \pm z_{1-\alpha/2}\,\widehat{\mathrm{SE}}(\widehat{\Delta})\Big).$$

The same derivation applies to gender and activity, without the 10x factor. The decision to use 10-year increments for age comes from the fact that age is meant to represent 10-year bands, and also for ease of interpretation since one year has a relatively small effect.

| Variable | Air Pollution | Odds Ratio | 95 % CI |
|---|---|---|---|
| Age (10-year increase) | Present | 1.58 | (1.04, 2.40) |
| | Absent | 0.70 | (0.43, 1.14) |
| Gender (Male vs Female) | Present | 0.09 | (0.03, 0.24) |
| | Absent | 1.11 | (0.38, 3.25) |
| Activity (Active vs Inactive) | Present | 2.72 | (1.07, 6.93) |
| | Absent | 0.11 | (0.02, 0.70) |

Table 9: Odds Ratios with 95% CI: age, gender, activity

Table 9 summarises the results: the odds of experiencing RD are between 1.04x and 2.40x higher, with 95% confidence, when increasing the age of a subject exposed to air pollution by 10 years. Conversely, if the subject is not exposed to air pollution, age is ambiguous (insignificant main effect), with the odds multiplicative factor between 0.43x and 1.14x for a 10-year increase.

---

[10]refer to the R-code in the Appendix for the implementation

Similarly, the odds of RD for men are much lower than women when exposed to air pollution, at just 3-24%, whereas we find no significant effect when pollution is absent (0.38-3.25, very wide interval). More strikingly, when air pollution is present, the odds of RD are much higher among active subjects, between 1.07 and 6.93, whereas they are much lower when air pollution is absent, at just 2-70% of the inactive population.[11]

Now, we can exemplify the effect of air pollution when all interactions are present, with gender = 1 (male), activity = 1 (active) and age = 35.

$$\Delta_{\text{airpollution}|g=1,a=1,\text{age}=35} = \beta_2 + \beta_5 \cdot (\text{activity} = 1) + \beta_6 \cdot (\text{gender} = 1) + \beta_7 \cdot (\text{age} = 35)$$
$$= \beta_2 + \beta_5 + \beta_6 + 35\,\beta_7.$$

Let $c = (0, 0, 1, 0, 0, 1, 1, 35)^\top$. Then,

$$\widehat{\text{Var}}\big(\Delta_{\text{air}|g=1,a=1,\text{age}=35}\big) = c^\top \widehat{V}\, c, \quad \widehat{\text{SE}}\big(\Delta_{\text{air}|g=1,a=1,\text{age}=35}\big) = \sqrt{c^\top \widehat{V}\, c}.$$

Expanding in covariance notation:

$$\widehat{\text{Var}}(\Delta) = \widehat{V}_{22} + \widehat{V}_{55} + \widehat{V}_{66} + 35^2 \widehat{V}_{77}$$
$$+ 2\big(\widehat{V}_{25} + \widehat{V}_{26} + 35\widehat{V}_{27} + \widehat{V}_{56} + 35\widehat{V}_{57} + 35\widehat{V}_{67}\big).$$

The standard error and confidence intervals for the odds ratio follow as above. Table 10 breaks down all the effects of air pollution on the odds by the values of the other regressors. Some particularly interesting effects are exposure to air pollution for 25-year-old inactive men, where the odds of developing RD are **lower** (0.02x to 0.87x) than when air pollution is absent, and active females of all ages, where air pollution substantially increases the odds of RD, ranging from 9.74x-154.25x for 25-year-old women, to 61.10x-3205.61x for 55-year-old women.

| Gender | Activity | Age | Odds Ratio | 95 % OR Confidence Interval |
|--------|----------|-----|------------|------------------------------|
| M | Active | 25 | 3.01 | (1.17, 7.76) |
| M | Active | 35 | 6.77 | (3.16, 14.52) |
| M | Active | 45 | 15.24 | (5.38, 43.20) |
| M | Active | 55 | 34.33 | (7.28, 161.87) |
| F | Active | 25 | 38.75 | (9.74, 154.25) |
| F | Active | 35 | 87.27 | (23.14, 329.03) |
| F | Active | 45 | 196.52 | (41.31, 934.80) |
| F | Active | 55 | 442.55 | (61.10, 3205.61) |
| M | Inactive | 25 | 0.13 | (0.02, 0.87) |
| M | Inactive | 35 | 0.28 | (0.04, 1.94) |
| M | Inactive | 45 | 0.64 | (0.08, 5.32) |
| M | Inactive | 55 | 1.44 | (0.12, 16.98) |
| F | Inactive | 25 | 1.62 | (0.16, 16.54) |
| F | Inactive | 35 | 3.65 | (0.35, 38.36) |
| F | Inactive | 45 | 8.22 | (0.64, 105.01) |
| F | Inactive | 55 | 18.50 | (1.04, 328.00) |

Table 10: Odds Ratios with 95% CI: air pollution

Of course, we should refrain from causal interpretations of the results, since it seems highly unlikely that exposure to air pollution should cause a decline in RD in any demographic.

---

[11]This is in line with the interpretation provided in the EDA section.

# 6   Prediction

Finally, we estimate the probability of developing RD for a 45-year-old woman who was active during an extreme air pollution event, and an 85-year-old inactive man who was exposed to air pollution. First, notice that the first datapoint is an in-sample prediction, since from Table 1 we see only five individuals that fall into this description, all of whom developed RD. Conversely, the 85-year-old man is an out-of-sample prediction, since we do not have any observations above the age of 55. Therefore, we expect the predictions in the second case to be less accurate (wider confidence interval) and less reliable, given extreme extrapolation.

| Activity | Age | Gender | Air Pollution | Predicted Probability | 95 % CI |
|----------|-----|--------|---------------|-----------------------|---------|
| Yes | 45 | Female | Yes | 0.960 | (0.889, 0.986) |
| No | 85 | Male | Yes | 0.826 | (0.293, 0.982) |

Table 11: Predicted Probabilities of Respiratory Disease - Full Model

As expected, the predicted probability for the first observation is high, in line with the prevalence of RD in that demographic, with a relatively tight confidence interval. Again as expected, the second interval is very wide, offering limited guidance for the probability of RD. The predicted value is high at 0.826, in line with the marginal effect of 0.04 (-0.04 + 0.08, see Table 8) from an additional year of age during an extreme air pollution event. However, this out-of-sample prediction must be treated cautiously, since it relies on substantial extrapolation beyond the observed age range.

# 7  Appendix

The table below summarises the result of the full GLM with all main effects and interactions with air pollution, but treating age as a categorical variable. This choice is more reflective of the data, since age represents 10-year bands rather than the actual age of the participants. This makes the continuous interpretation more fragile, since the coefficient on age indicates the effect of an additional year of age on the log-odds of RD, which requires some level of extrapolation. Moreover, an unwanted effect of the continuous treatment is that the fitted values for, say, 24 and 26 will be different, when they could presumably be the same, since may be in the same 10-year band. That is also the reason the paper displays the effect of a 10-year increase in age, rather than a 1-year increase. While more reflective of the data, the model below is also more restrictive, since it does not allow extrapolation for ages beyond the upper limit of 55. Also, notice that the model below includes 4 more coefficients, which leads to higher AIC (288.6 vs 282.1).

|  | Dependent variable: |
| --- | --- |
|  | Respiratory difficulty |
| Air pollution (Yes) | 0.44 (1.21) |
| Age 35 | −0.50 (0.60) |
| Age 45 | −0.69 (0.88) |
| Age 55 | −1.02 (0.82) |
| Male | 0.10 (0.55) |
| Active | −2.17* (0.94) |
| Air pollution × Age 35 | 0.87 (0.76) |
| Air pollution × Age 45 | 1.23 (1.03) |
| Air pollution × Age 55 | 3.22* (1.37) |
| Air pollution × Male | −2.46** (0.76) |
| Air pollution × Active | 3.15** (1.05) |
| Constant | 0.81 (1.06) |
| Observations | 285 |
| Log Likelihood | −132.32 |
| Akaike Inf. Crit. | 288.63 |
| Note: | *p<0.05; **p<0.01; ***p<0.001 |

Table 12: Logistic Regression of Respiratory difficulty, age as factor

The plot below displays the standardised residuals against the predicted responses for the chosen model (glm1 or the "full" model). We do not see any differences in terms of outlying observations from the plot, while the two clusters appear non-monotonic due to the sigmoid transformation.
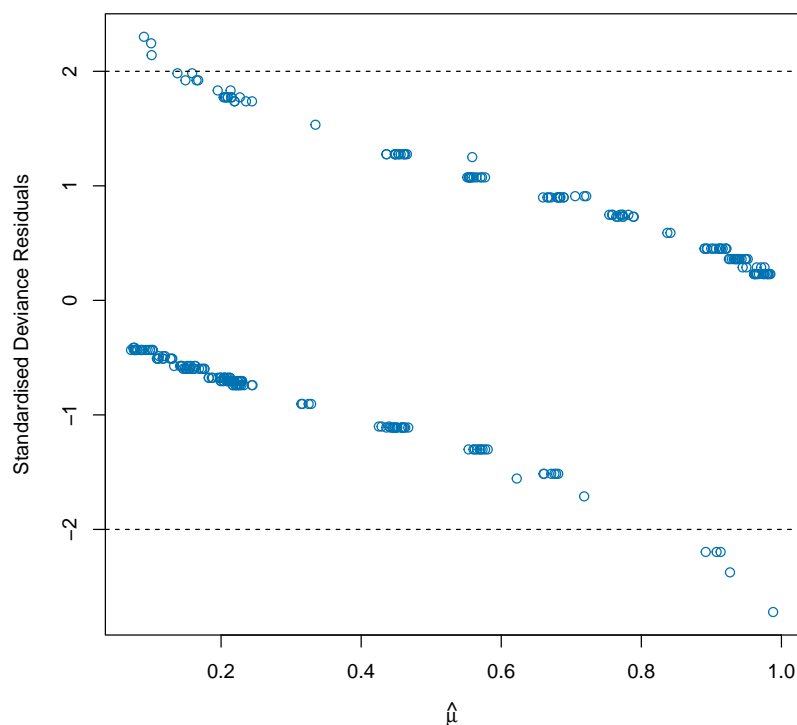
Figure 6: Plot of standardised deviance residuals vs predicted responses

The following plots show the model diagnostics plots for the reduced model (glm5). Interestingly, we see fewer outlying observations in the standardised residual plot. Notice that this may be partly due to fewer fitted values for a more constrained model. The plot uses jitter to display the concentration of values, but of course the total number of fitted values will at most be 18, since we only have 18 possible combinations of the regressors in then model and thus 36 possible values of the residuals. Therefore, fewer outliers will be expected. The comments in the model diagnostics section otherwise extend to the remaining plots. The QQ plot is somewhat more anomalous but still merely reflects the limited number of residual values.
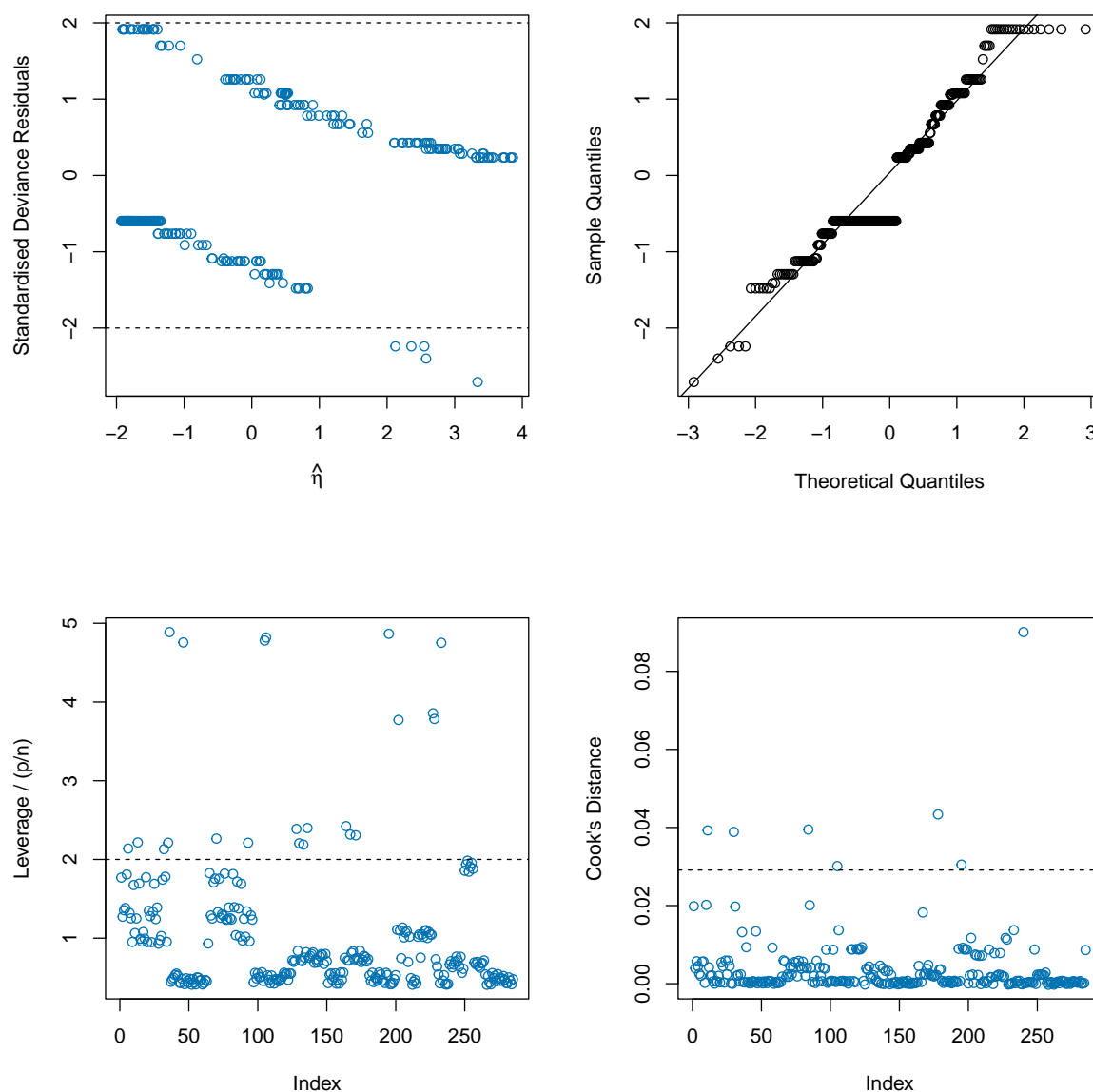
Figure 7: In order Left to Right, Top to Bottom: residuals vs fitted linear regressor, residual QQ-plot, leverage plot, Cook's distance plot

Finally, the table below illustrates the odds ratio and respective confidence intervals for the reduced model. These are broadly in line with the full model, but notice that this model has much lower variance for high values of age. This is because the standard error of age is larger in the full model (0.03 vs 0.02), which has a minor effect when age is small but increases exponentially with age.

| Gender | Activity | Age | Odds Ratio | 95% OR Confidence Interval |
|--------|----------|-----|------------|----------------------------|
| M | Active | 25 | 4.42 | (1.39, 14.10) |
| M | Active | 35 | 6.63 | (2.22, 19.84) |
| M | Active | 45 | 9.96 | (3.17, 31.25) |
| M | Active | 55 | 14.95 | (4.11, 54.44) |
| F | Active | 25 | 55.59 | (19.38, 159.46) |
| F | Active | 35 | 83.46 | (28.42, 245.11) |
| F | Active | 45 | 125.30 | (37.35, 420.34) |
| F | Active | 55 | 188.12 | (45.36, 780.16) |
| M | Inactive | 25 | 0.22 | (0.07, 0.68) |
| M | Inactive | 35 | 0.33 | (0.08, 1.29) |
| M | Inactive | 45 | 0.49 | (0.10, 2.54) |
| M | Inactive | 55 | 0.74 | (0.11, 5.16) |
| F | Inactive | 25 | 2.76 | (1.13, 6.73) |
| F | Inactive | 35 | 4.15 | (1.19, 14.44) |
| F | Inactive | 45 | 6.23 | (1.25, 30.96) |
| F | Inactive | 55 | 9.35 | (1.32, 66.38) |

Table 13: Odds Ratios with 95% CI: air pollution, under the reduced model (glm5)

Below is the complete R code used for this analysis.

```r
data <- read.csv("airpollution.csv", stringsAsFactors = TRUE)
packages <- c("ggplot2", "xtable", "stargazer", "rsq", "emmeans")
install.packages(packages)
library(xtable)
library(ggplot2)
library(stargazer)
library(rsq)
library(emmeans)

#### DEFACTORED BECAUSE OF WHAT SAID IN Qe ABOUT OOS AGE PREDICT
# Convert age to factors, since they represent age ranges
# data$age <- as.factor(data$age)
attach(data) #attach lets us call the columns directly, making code more
    concise
summary(data)


################################################################################
### A - Exploratory data analysis
################################################################################

# summary table with counts and proportions for each variable
make_summary <- function(var, name) {
  tbl <- table(var)
  data.frame(
    Variable = name,
    Level = names(tbl),
    Count = as.numeric(tbl),
    Percent = round(100 * as.numeric(prop.table(tbl)), 1)
  )
}

summary_table <- rbind(
  make_summary(respd, "Respiratory Difficulty"),
  make_summary(airpollution, "Air pollution"),
  make_summary(activity, "Activity"),
  make_summary(gender, "Gender"),
  make_summary(age, "Age")
)


```

```r
41 stargazer::stargazer(summary_table,
42                        type = "latex",
43                        summary = FALSE,
44                        rownames = FALSE,
45                        digits = 2,
46                        font.size = "scriptsize",
47                        out = "tables/summary_table.tex")
48
49 # proportion with respd == "Yes"
50 prop_tbl <- aggregate(I(respd == "Yes") ~ airpollution + activity + gender +
       age,
51                        data = data, FUN = mean)
52 names(prop_tbl)[5] <- "PropYes"
53
54 # count for each combination
55 count_tbl <- aggregate(respd ~ airpollution + activity + gender + age,
56                        data = data, FUN = length)
57 names(count_tbl)[5] <- "Count"
58
59 prop_tbl <- merge(prop_tbl, count_tbl,
60                   by = c("airpollution", "activity", "gender", "age"))
61
62 prop_tbl <- prop_tbl[order(prop_tbl$airpollution, prop_tbl$activity, prop_tbl$
       gender, prop_tbl$age), ]
63 prop_tbl$PropYes <- sprintf("%.2f", prop_tbl$PropYes)
64
65 names(prop_tbl) <- c("Air Pollution", "Activity", "Gender", "Age",
66                      "Proportion", "Count")
67
68 stargazer::stargazer(prop_tbl,
69                        type = "latex",
70                        summary = FALSE,
71                        rownames = FALSE,
72                        digits = 2,
73                        font.size = "scriptsize",
74                        out = "tables/combo_prop_table.tex")
75
76
77 # PLOTS
78 plot_names <- list("Air pollution", "Age", "Respiratory Difficulty", "Activity"
       , "Gender")
79 names(plot_names) <- names(data)
80
81 blue_palette <- c("#084C8D", "#0072B2", "#56B4E9", "#A6D8F0")
82
83 # Bar plot for age
84 pdf("figs/age_barplot_cex1.pdf", width = 6, height = 6)
85 par(mfrow = c(1,1), cex = 1.1)
86 barplot(table(age),
87         col = blue_palette,
88         xlab = "Age band",
89         ylab = "Count")
90 dev.off()
91
92 # Plot proportions with respiratory Difficulty for binary categories
93 prop_plot <- function(x, y, text = FALSE) {
94   x_label <- plot_names[[deparse(substitute(x))]]
95   y_label <- plot_names[[deparse(substitute(y))]]
96   table <- table(x, y)
97   print(table)
98   prop <- prop.table(table, 2)
99   print(prop)
100  bp <- barplot(prop,
```

```r
101            col = blue_palette[c(2, 4)],
102            xlab = y_label,
103            ylab = "Proportion of cohort",
104            beside = TRUE,
105            ylim = c(0, max(prop) + 0.1))
106   if(text == TRUE) {text(x = bp,
107         y = prop,
108         labels = signif(prop, 2),
109         pos = 3,
110         cex = 0.9)}
111 }
112
113 pdf("figs/prop_plots_initial.pdf", width = 7, height = 7)
114 layout(matrix(c(1, 2, 3, 4, 5, 5), nrow = 3, byrow = TRUE),
115        heights = c(1, 1, 0.3))  # Bottom row is shorter
116
117 par(mar = c(4, 4, 2, 2))
118 prop_plot(respd, gender, text = TRUE)
119 prop_plot(respd, activity, text = TRUE)
120 prop_plot(respd, airpollution, text = TRUE)
121 prop_plot(respd, age, text = TRUE)
122
123 par(mar = c(0, 0, 0, 0))
124 plot.new()
125 legend("center",
126        legend = c("No respiratory difficulty", "Respiratory difficulty"),
127        fill = blue_palette[c(2, 4)],
128        bty = "n",
129        horiz = TRUE)
130
131 dev.off()
132
133 # Now breaking the data down by airpollution = yes or no, showing that gender
134 # is irrelevant when there is no air pollution, but very relevant when there is
135 # whereas activity seems to have opposite effects with/ without air pollution
136
137 data_pollutionyes <- subset(data, subset = airpollution == "Yes")
138 data_pollutionno <- subset(data, subset =airpollution == "No")
139
140 plot_names <- list(
141   respd                      = "Respiratory difficulty",
142   gender                     = "Gender",
143   activity                   = "Activity",
144   age                        = "Age",
145   "data_pollutionno$respd"   = "Respiratory difficulty",
146   "data_pollutionno$gender"  = "Gender",
147   "data_pollutionno$activity" = "Activity",
148   "data_pollutionno$age"     = "Age",
149   "data_pollutionyes$respd"  = "Respiratory difficulty",
150   "data_pollutionyes$gender" = "Gender",
151   "data_pollutionyes$activity" = "Activity",
152   "data_pollutionyes$age"    = "Age"
153 )
154
155 pdf("figs/barplots_noairpollution.pdf", width = 7, height = 4)
156 par(mar = c(5.1, 4.1, 4.1, 2.1))
157 layout(matrix(c(1, 2, 3, 4, 4, 4), nrow = 2, byrow = TRUE),
158        heights = c(1, 0.1))  # Bottom row is shorter
159
160 prop_plot(data_pollutionno$respd, data_pollutionno$gender)
161 prop_plot(data_pollutionno$respd, data_pollutionno$activity)
162 prop_plot(data_pollutionno$respd, data_pollutionno$age)
163
```

```r
164 par(mar = c(0, 0, 0, 0))
165 plot.new()
166
167 legend("center",
168        legend = c("No respiratory difficulty", "Respiratory difficulty"),
169        fill = blue_palette[c(2, 4)],
170        bty = "n",
171        horiz = TRUE)
172 dev.off()
173
174 pdf("figs/barplots_yesairpollution.pdf", width = 7, height = 4)
175 par(mar = c(5.1, 4.1, 4.1, 2.1))
176 layout(matrix(c(1, 2, 3, 4, 4, 4), nrow = 2, byrow = TRUE),
177        heights = c(1, 0.1))  # Bottom row is shorter
178
179 prop_plot(data_pollutionyes$respd, data_pollutionyes$gender)
180 prop_plot(data_pollutionyes$respd, data_pollutionyes$activity)
181 prop_plot(data_pollutionyes$respd, data_pollutionyes$age)
182
183 par(mar = c(0, 0, 0, 0))
184 plot.new()
185
186 legend("center",
187        legend = c("No respiratory difficulty", "Respiratory difficulty"),
188        fill = blue_palette[c(2, 4)],
189        bty = "n",
190        horiz = TRUE)
191
192 dev.off()
193
194 ################################################################################
195 ### B - MODELLING
196 ################################################################################
197
198 # Baseline model: all interactions between airpollution and other explanatory
       variables
199
200 resp.glm1 <- glm(respd ~ airpollution * (age + gender + activity),
201                  data = data,
202                  family = binomial)
203
204 summary(resp.glm1)
205
206 anova(resp.glm1) #initial anova suggests all interactions are significant
207 # thus by the hierarchy principle we should keep all the lower order terms too
208 # by the wald we should get rid of all baselines except activity. Gender is
       especially insignificant
209 # by the LR gender is extremely significant and activity insignificant.
210 # this bc anova sequential and goes gender first, then activity and
       interactions. see glm1.1 below
211
212 LRT <- function(baseline, restricted) {
213   LRT <- restricted$deviance - baseline$deviance
214   p_val <- 1 - pchisq(LRT, restricted$df.residual - baseline$df.residual)
215   out <- list(
216     LRT = LRT,
217     p_value = p_val
218   )
219   return(out)
220 }
221
222 # likelihood ratio for model useless:
223
```

```r
224 resp.glm0 <- glm(respd ~ 1, data = data, family = binomial)
225
226 L1 <- LRT(resp.glm1, resp.glm0) # approximately 0 --> the model is not useless
227
228 # Consider dropping any of the interactions
229 # All effects are significant, do not drop any. Agrees with Wald
230
231 resp.glm2 <- glm(respd ~ airpollution * (gender + activity) + age,
232                  data = data,
233                  family = binomial)
234
235 L2 <- LRT(resp.glm1, resp.glm2)
236 resp.glm3 <- glm(respd ~ airpollution * (gender + age) + activity,
237                  data = data,
238                  family = binomial)
239 L3 <- LRT(resp.glm1, resp.glm3)
240
241 resp.glm4 <- glm(respd ~ airpollution * (activity + age) + gender,
242                  data = data,
243                  family = binomial)
244 L4 <- LRT(resp.glm1, resp.glm4)
245
246 # Dropping the baseline effects except activity. This is suggested by Wald
247 data$agexpollution <- data$age*(data$airpollution == "Yes")
248 data$genderxpollution <- as.integer(data$gender == "Male" & data$airpollution
        == "Yes")
249 data$activityxpollution <- as.integer(data$activity == "Yes" & data$
        airpollution == "Yes")
250
251 resp.glm5 <- glm(respd ~ activity + agexpollution + genderxpollution +
        activityxpollution,
252                  data = data,
253                  family = binomial)
254 L5 <- LRT(resp.glm1, resp.glm5) # p-value is 0.478, suggesting we should drop
        all baselines!
255
256 # drop age baseline
257 resp.glm6 <- glm(respd ~ airpollution*(gender + activity) + agexpollution,
258                  data = data,
259                  family = binomial)
260
261 L6 <- LRT(resp.glm1, resp.glm6)
262
263 # drop gender baseline
264 resp.glm7 <- glm(respd ~ airpollution*(activity) + agexpollution +
        genderxpollution,
265                  data = data,
266                  family = binomial)
267 summary(resp.glm7)
268
269 L7 <- LRT(resp.glm6, resp.glm7)
270
271 # drop activity main effect (against Wald)
272 resp.glm8 <- glm(respd ~ airpollution + activityxpollution + agexpollution +
        genderxpollution,
273                  data = data,
274                  family = binomial)
275 summary(resp.glm8)
276 L8 <- LRT(resp.glm7, resp.glm8) # reject
277
278 # finally, drop airpollution main effect
279 # since all the coefficients are statistically significant, we should choose
        final model resp.glm5
```

```r
280 # for final verification , see a plot of the AICs and a table of the R^2 KL for
        all the models
281
282 # why we should keep the baseline interaction - the hierarchy principle
283
284 data$airpollutionprime <- ifelse (( data$airpollution == "Yes") , 0, 1)
285 resp.glmairpolltoy <- glm(respd ~ airpollutionprime * (age + gender + activity)
        ,
286                              data = data ,
287                              family = binomial )
288
289 resp.glmairpolltoy_nointeractions <- glm(respd ~ airpollutionprime +
        airpollutionprime:gender +
290                                          airpollutionprime:activity +
        airpollutionprime:age ,
291                                          data = data ,
292                                          family = binomial )
293 summary ( resp.glmairpolltoy_nointeractions )
294 LRT( resp.glmairpolltoy , resp.glmairpolltoy_nointeractions )
295
296 chosen_model <- resp.glm1
297 reduced_model <- resp.glm5
298
299 ########### TABLES ###############
300 stargazer ( resp.glm1 ,
301           type = "latex",
302           title = "Logistic Regression of Respiratory difficulty",
303           dep.var.labels = "Respiratory difficulty",
304           covariate.labels = c(
305             "Air pollution (Yes)",
306             "Age",
307             "Male",
308             "Active",
309             "Air pollution    Age",
310             "Air pollution    Male",
311             "Air pollution    Active"
312           ),
313           digits = 2,
314           star.cutoffs = c(0.05 , 0.01 , 0.001) ,
315           no.space = TRUE ,
316           float = TRUE ,
317           font.size = "scriptsize",
318           report = "vc*sp",
319           single.row = TRUE ,
320           out = "tables/resp.glm1.tex")
321
322 stargazer ( resp.glm5 ,
323           type = "latex",
324           title = "Logistic Regression of Respiratory Difficulty",
325           dep.var.labels = "Respiratory Difficulty",
326           covariate.labels = c(
327             "Active",
328             "Air pollution    Age",
329             "Air pollution    Male",
330             "Air pollution    Active"
331           ),
332           digits = 2,
333           star.cutoffs = c(0.05 , 0.01 , 0.001) ,
334           no.space = TRUE ,
335           float = TRUE ,
336           single.row = TRUE ,
337           font.size = "scriptsize",
338           report = "vc*sp",
```

```
339            out = "tables/resp.glm5.tex")
340
341 likelihood_ratios <- round(c(L1$LRT, L2$LRT, L3$LRT, L4$LRT, L5$LRT, L6$LRT, L7
        $LRT, L8$LRT), 2)
342 p_values <- round(c(L1$p_value, L2$p_value, L3$p_value, L4$p_value, L5$p_value,
         L6$p_value, L7$p_value, L8$p_value), 3)
343 full_model <- c("glm1", "glm1", "glm1", "glm1", "glm1", "glm1", "glm6", "glm7")
344 restricted_model <- c("glm0","glm2","glm3","glm4","glm5","glm6","glm7","glm8")
345 drop_vector <- c(
346   "All (null test)",
347   "Age x Air Pollution",
348   "Activity x Air Pollution",
349   "Gender x Air Pollution",
350   "Air Pollution, Age, Gender",
351   "Age",
352   "Gender",
353   "Activity"
354 )
355
356 likelihood_ratio_table <- data.frame(
357   "Full Model" = full_model,
358   "Restricted Model" = restricted_model,
359   "Drop:" = drop_vector,
360   "Likelihood Ratio Statistic" = likelihood_ratios,
361   "P-values" = p_values)
362
363 stargazer(likelihood_ratio_table,
364           type = "latex",
365           summary = FALSE,
366           rownames = FALSE,
367           digits = 2,
368           font.size = "scriptsize",
369           out = "tables/likelihood_ratio_table.tex")
370
371 models <- list(resp.glm1, resp.glm2, resp.glm3, resp.glm4, resp.glm5, resp.glm6
        , resp.glm7, resp.glm8)
372 model_names <- c("glm1", "glm2", "glm3", "glm4", "glm5", "glm6", "glm7", "glm8"
        )
373 aics <- sapply(models, AIC)
374 bics <- sapply(models, BIC)
375 deviances <- sapply(models, deviance)
376 rsq_vals <- round(sapply(models, rsq.kl), 3)
377
378 model_desc <- c(
379   "Full model",
380   "Drop Age x Air Pollution",
381   "Drop Activity x Air Pollution",
382   "Drop Gender x Air Pollution",
383   "Drop Air Pollution, Age, Gender",
384   "Drop Age",
385   "Drop Age, Gender",
386   "Drop Activity, Age, Gender"
387 )
388
389 aic_bic_etc_table <- data.frame(
390   "Model" = model_names,
391   "Description" = model_desc,
392   "Residual Deviance" = deviances,
393   "AIC" = aics,
394   "BIC" = bics,
395   "R^2 Kullback-Leibler" = rsq_vals
396 )
397
```

```
398  stargazer(aic_bic_etc_table,
399            type = "latex",
400            font.size = "scriptsize",
401            summary = FALSE,
402            rownames =FALSE,
403            digits = 2,
404            out = "tables/aic_bic_etc_table.tex")
405
406  stargazer(resp.glmairpolltoy,
407            type = "latex",
408            title = "Toy model - Logistic Regression of Respiratory difficulty",
409            dep.var.labels = "Respiratory difficulty",
410            covariate.labels = c(
411              "Air pollution (No)",
412              "Age",
413              "Male",
414              "Active",
415              "Air pollution    Age",
416              "Air pollution    Male",
417              "Air pollution    Active"
418            ),
419            digits = 2,
420            star.cutoffs = c(0.05, 0.01, 0.001),
421            no.space = TRUE,
422            float = TRUE,
423            font.size = "scriptsize",
424            report = "vc*sp",
425            single.row = TRUE,
426            out = "tables/toymodel.tex")
427
428
429
430  ###############################################################################
431  ### C - Model diagnostics
432  ###############################################################################
433
434  # Residual plots
435
436  set.seed(30)
437
438  pdf("figs/residual_plots.pdf", width = 9.5, height = 10)
439  par(mfrow = c(2, 2), margin(5, 5, 3, 3), cex = 1)
440
441  plot(jitter(predict(chosen_model, type = 'link'), 10),
442                     jitter(rstandard(chosen_model), 0), col = blue_palette
       [2],
443                     xlab = expression(hat(eta)),
444                     ylab = "Standardised Deviance Residuals")
445  abline(a = -2, b = 0, lty = 2)
446  abline(a = 2, b = 0, lty = 2)
447
448  # QQ plot
449
450  qqnorm(rstandard(chosen_model), main = "")
451  qqline(rstandard(chosen_model))
452
453  # Leverage
454
455  p <- chosen_model$df.null - chosen_model$df.residual + 1
456  n <- nrow(data)
457
458  plot(jitter(influence(chosen_model)$hat/(p/n), 200), ylab = "Leverage / (p/n)",
459       col = blue_palette[2])
```

```r
460 abline(2, 0, lty = 2)
461
462 # Cook's distance
463
464 plot(jitter(cooks.distance(chosen_model), 100), ylab = "Cook's Distance",
465       col = blue_palette[2])
466 abline(8/(n - 2 * p), 0, lty = 2)
467 dev.off()
468
469 leverage <- influence(chosen_model)$hat
470 cooks_d <- cooks.distance(chosen_model)
471
472 # Thresholds
473 leverage_threshold <- 2 * (p / n)
474 cooks_threshold <- 8 / (n - 2 * p)
475
476 # Problematic observations
477 high_leverage <- which(leverage > leverage_threshold)
478 high_cooks <- which(cooks_d > cooks_threshold)
479
480 # Summary table of problematic observations
481 problematic <- intersect(high_leverage, high_cooks)
482 diag_table <- data.frame(
483   data[problematic, 1:5],
484   Leverage = round(leverage[problematic], 4) / (p/n),
485   Cooks_D = round(cooks_d[problematic], 4)
486 )
487 print(diag_table)
488
489 colnames(diag_table) <- c("Air Pollution","Age", "Respiratory Difficulty", "
        Activity","Gender", "Leverage", "Cook's D")
490
491 stargazer(diag_table,
492           summary = FALSE,
493           rownames = FALSE,
494           type = "latex",
495           title = "Observations with High Leverage and High Cook's Distance",
496           font.size = "scriptsize",
497           float = TRUE,
498           out = "tables/diag_table.tex")
499
500 #############################################################################
501 ### D - INTERPRETATION
502 #############################################################################
503
504 beta <- coef(chosen_model)
505 se <- sqrt(diag(vcov(chosen_model)))
506 z <- qnorm(0.975)
507
508 # Calculate CIs
509 ci_lower <- beta - z * se
510 ci_upper <- beta + z * se
511
512 # Build table
513 ci_table <- data.frame(
514   Variable = c(
515     "Intercept",
516     "Air pollution (Yes)",
517     "Age",
518     "Male",
519     "Active",
520     "Air pollution     Age",
521     "Air pollution     Male",
```

```r
522        "Air pollution    Active"
523      ),
524      Estimate = sprintf("%.2f", beta),
525      CI = sprintf("(%.2f, %.2f)", ci_lower, ci_upper)
526  )
527
528  colnames(ci_table) <- c("Variable", "Estimate", "95\\% CI")
529
530  stargazer(ci_table,
531            summary = FALSE,
532            rownames = FALSE,
533            type = "latex",
534            title = "Logistic Regression Coefficients with 95\\% Confidence
      Intervals",
535            font.size = "scriptsize",
536            float = TRUE,
537            out = "tables/ci_table.tex")
538
539
540
541  odds_ratio_and_ci <- function(variable) {
542    beta <- coef(chosen_model)
543    V <- vcov(chosen_model)
544    interaction <- paste("airpollutionYes", variable, sep = ":")
545
546    # Air pollution present
547    logor_airpollution1 <- beta[variable] + beta[interaction]
548    variance_1 <- V[variable, variable] + V[interaction, interaction] + 2 * V[
      variable, interaction]
549    se_1 <- sqrt(variance_1)
550    logor_airpollution1_ci <- logor_airpollution1 + 1.96 * se_1 * c(-1, 1)
551
552    # Air pollution absent
553    logor_airpollution0 <- beta[variable]
554    variance_0 <- V[variable, variable]
555    se_0 <- sqrt(variance_0)
556    logor_airpollution0_ci <- logor_airpollution0 + 1.96 * se_0 * c(-1, 1)
557
558    # For age: multiply log-odds by 10
559    if (variable == "age") {
560      logor_airpollution1 <- logor_airpollution1 * 10
561      logor_airpollution1_ci <- logor_airpollution1_ci * 10
562      logor_airpollution0 <- logor_airpollution0 * 10
563      logor_airpollution0_ci <- logor_airpollution0_ci * 10
564    }
565
566    data.frame(
567      Air_Pollution = c("Present", "Absent"),
568      Odds_Ratio = sprintf("%.2f", exp(c(logor_airpollution1, logor_airpollution0
      ))),
569      CI = c(
570        sprintf("(%.2f, %.2f)", exp(logor_airpollution1_ci[1]), exp(logor_
      airpollution1_ci[2])),
571        sprintf("(%.2f, %.2f)", exp(logor_airpollution0_ci[1]), exp(logor_
      airpollution0_ci[2]))
572      )
573    )
574  }
575
576  variable_names <- names(beta)[3:5]
577  variable_nicenames <- c("Age", "Gender", "Activity")
578
579  # calculate using the delta method
```

```r
580  odds_ratio_and_ci_airpollution <- function(age, gender, activity) {
581    beta <- coef(chosen_model)
582    V <- vcov(chosen_model)
583    ap <- "airpollutionYes"
584    apag <- "airpollutionYes:age"
585    apge <- "airpollutionYes:genderMale"
586    apac <- "airpollutionYes:activityYes"
587    var <- V[ap, ap] +
588      sum(c(V[apag, apag], V[apge, apge], V[apac, apac]) * c(age, gender,
         activity)**2) +
589      2 * sum(c(V[ap, apag], V[ap, apge], V[ap, apac]) * c(age, gender, activity)
         ) +
590      2 * age * sum(c(V[apag, apge], V[apag, apac]) * c(gender, activity)) +
591      2 * gender * activity * V[apge, apac]
592    se <- sqrt(var)
593    log_odds_ratio <- sum(beta[c(ap, apag, apge, apac)] * c(1, age, gender,
         activity))
594    log_odds_ci <- log_odds_ratio + 1.96 * c(-1, 1) * se
595    odds_ratio <- exp(log_odds_ratio)
596    odds_ci <- exp(log_odds_ci)
597    out <- c(
598      "Odds Ratio" = round(odds_ratio, 2),
599      "Odds Ratio CI" = round(odds_ci, 2)
600    )
601    return(out)
602  }
603
604  ######### TABLES #############
605
606  # OTHER ODDS
607
608  combined_table <- rbind(
609    cbind(Variable = "Age (10-year increase)", odds_ratio_and_ci("age")),
610    cbind(Variable = "Gender (Male vs Female)", odds_ratio_and_ci("genderMale")),
611    cbind(Variable = "Activity (Active vs Inactive)", odds_ratio_and_ci("
         activityYes"))
612  )
613  combined_table$Variable[c(2,4,6)] <- ""  # Remove duplicate labels
614
615  colnames(combined_table) <- c("Variable", "Air Pollution", "Odds Ratio", "95 %
         CI")
616
617  stargazer(combined_table,
618            summary = FALSE,
619            rownames = FALSE,
620            type = "latex",
621            font.size = "scriptsize",
622            digits = 2,
623            out = "tables/combined_odds_ratios_table.tex")
624
625  # AIRPOLLUTION ODDS
626
627  combos <- expand.grid(
628    age = c(25, 35, 45, 55),
629    gender = c(1, 0),
630    activity = c(1, 0)
631  )
632
633  results <- t(mapply(
634    gender = combos$gender,
635    activity = combos$activity,
636    age = combos$age,
637    odds_ratio_and_ci_airpollution
```

```r
638 ))
639
640 odds_table <- cbind(combos, results)
641 odds_table$gender <- ifelse(odds_table$gender == 1, "\texttt{M}", "\texttt{F}")
642 odds_table$activity <- ifelse(odds_table$activity == 1, "Active", "Inactive")
643 odds_table$CI <- sprintf("(%.2f, %.2f)", odds_table[,5], odds_table[,6])
644
645 odds_table <- odds_table[, c("gender", "activity", "age", "Odds Ratio", "CI")]
646 colnames(odds_table) <- c("Gender", "Activity", "Age", "Odds Ratio", "95\\% OR
       Confidence Interval")
647
648 stargazer(odds_table,
649           type = "latex",
650           font.size = "scriptsize",
651           summary = FALSE,
652           rownames =FALSE,
653           digits = 2,
654           out = "tables/odds_table.tex")
655
656
657 ###############################################################################
658 ### E - Prediction
659 ###############################################################################
660 # new prediction data frames
661 predict_1 <- data.frame(
662   activity = "Yes",
663   age = 45,
664   gender = "Female",
665   airpollution = "Yes"
666 )
667
668 predict_2 <- data.frame(
669   activity = "No",
670   age = 85,
671   gender = "Male",
672   airpollution = "Yes"
673 )
674
675 # linear regression first (for conf. int's)
676 pred1_link <- predict(chosen_model, newdata = predict_1, type = "link", se.fit
       = TRUE)
677 pred2_link <- predict(chosen_model, newdata = predict_2, type = "link", se.fit
       = TRUE)
678
679 get_prob_ci <- function(pred_link) {
680   fit <- pred_link$fit
681   se <- pred_link$se.fit
682   ci_link <- fit + 1.96 * se * c(-1, 1)
683   prob <- plogis(fit)
684   ci_prob <- plogis(ci_link)
685   return(c(prob, ci_prob))
686 }
687
688 prob1 <- get_prob_ci(pred1_link)
689 prob2 <- get_prob_ci(pred2_link)
690
691 pred_table <- data.frame(
692   Activity = c("Yes", "No"),
693   Age = c(45, 85),
694   Gender = c("Female", "Male"),
695   Air_Pollution = c("Yes", "Yes"),
696   Probability = sprintf("%.3f", c(prob1[1], prob2[1])),
697   CI = c(
```

```
698      sprintf("(%.3f, %.3f)", prob1[2], prob1[3]),
699      sprintf("(%.3f, %.3f)", prob2[2], prob2[3])
700   )
701 )
702
703 colnames(pred_table) <- c("Activity", "Age", "Gender", "Air Pollution",
704                           "Predicted Probability", "95 % CI")
705
706 stargazer(pred_table,
707           summary = FALSE,
708           rownames = FALSE,
709           type = "latex",
710           title = "Predicted Probabilities of Respiratory Difficulty - Full
     Model",
711           font.size = "scriptsize",
712           out = "tables/predicted_probs.tex")
713
714
715 # new prediction data frames - reduced model
716 predict_1_rm <- data.frame(
717   activity = "Yes",
718   agexpollution = 45,
719   genderxpollution = 0,
720   activityxpollution = 1
721 )
722
723 predict_2_rm <- data.frame(
724   activity = "No",
725   agexpollution = 85,
726   genderxpollution = 1,
727   activityxpollution = 0
728 )
729
730 # linear regression first (for confidence intervals)
731 pred1_link_rm <- predict(reduced_model, newdata = predict_1_rm, type = "link",
     se.fit = TRUE)
732 pred2_link_rm <- predict(reduced_model, newdata = predict_2_rm, type = "link",
     se.fit = TRUE)
733
734 get_prob_ci(pred1_link_rm)  # Female, 45, active, pollution
735 get_prob_ci(pred2_link_rm)  # Male, 85, inactive, pollution
736
737 # notice that while we can get confidence intervals for the *probabilities*, we
     cannot get
738 # prediction intervals for the probabilities, since the realisation will always
     jsut be 0 or 1 for resp
739 # so prediction intervals are meaningless here
740
741 ##############################################################################
742 ### APPENDIX CODE
743 ##############################################################################
744
745 ### THE FACTOR MODEL ###
746
747 data$agef <- as.factor(data$age)
748
749 respd.glmf <- glm(respd ~ airpollution *(agef + gender + activity),
750                   data = data,
751                   family = binomial)
752 respd.glmf.1 <- glm(respd ~ airpollution *(agef + activity + gender),
753                     data = data,
754                     family = binomial)
755 respd.glmf.2 <- glm(respd ~ airpollution *(activity + agef + gender),
```

```
756                          data = data ,
757                          family = binomial )
758  summary ( respd.glmf )
759  anova ( respd.glmf )
760  anova ( respd.glmf .1)
761  anova ( respd.glmf .2)
762
763  stargazer ( respd.glmf ,
764           type = "latex",
765           title = "Logistic Regression of Respiratory difficulty",
766           dep.var.labels = "Respiratory difficulty",
767           covariate.labels = c(
768             "Air pollution (Yes)",
769             "Age 35",
770             "Age 45",
771             "Age 55",
772             "Male",
773             "Active",
774             "Air pollution    Age 35",
775             "Air pollution    Age 45",
776             "Air pollution    Age 55",
777             "Air pollution    Male",
778             "Air pollution    Active"
779           ),
780           digits = 2,
781           star.cutoffs = c(0.05, 0.01, 0.001),
782           no.space = TRUE ,
783           float = TRUE ,
784           font.size = "scriptsize",
785           single.row = TRUE ,
786           out = "tables/resp.glmf.tex")
787
788  # in all cases in the factor model we see that the baseline age effect is not
        significant
789  # this provides further evidence for dropping
790  # activity is significant when added first, gender is always highly significant
        in the factor model
791
792
793  ### extra plot - residuals vs fitted values
794
795  pdf ("figs/residual_vs_mu.pdf")
796  plot ( jitter ( predict ( chosen_model , type = 'response'), 10) ,
797       jitter ( rstandard ( chosen_model ), 0) , col = blue_palette [2] ,
798       xlab = expression ( hat (mu)) ,
799       ylab = "Standardised Deviance Residuals")
800  abline (a = -2, b = 0, lty = 2)
801  abline (a = 2, b = 0, lty = 2)
802  dev.off ()
803
804
805  ######### FOR THE REDUCED MODEL ##############
806
807  ##### DIAGNOSTICS #######
808
809
810  # residual plots
811
812  set.seed (30)
813
814  pdf ("figs/residual_plots_reducedmodel.pdf", width = 9.5, height = 10)
815  par ( mfrow = c(2, 2), margin (5, 5, 3, 3), cex = 1)
816
```

```r
817 plot(jitter(predict(reduced_model, type = 'link'), 10),
818      jitter(rstandard(reduced_model), 0), col = blue_palette[2],
819      xlab = expression(hat(eta)),
820      ylab = "Standardised Deviance Residuals")
821 abline(a = -2, b = 0, lty = 2)
822 abline(a = 2, b = 0, lty = 2)
823
824 # qq plot
825
826 qqnorm(rstandard(reduced_model), main = "")
827 qqline(rstandard(reduced_model))
828
829 # leverage
830
831 p <- reduced_model$df.null - reduced_model$df.residual + 1
832 n <- nrow(data)
833
834 plot(jitter(influence(reduced_model)$hat/(p/n), 200), ylab = "Leverage / (p/n)"
        ,
835      col = blue_palette[2])
836 abline(2, 0, lty = 2)
837
838 # cook's distance
839
840 plot(jitter(cooks.distance(reduced_model), 100), ylab = "Cook's Distance",
841      col = blue_palette[2])
842 abline(8/(n - 2 * p), 0, lty = 2)
843 dev.off()
844
845 # ACTIVITY
846 # odds ratio depends on air pollution
847
848 beta <- coef(reduced_model)
849 V <- vcov(reduced_model)
850
851 se_sum <- function(b1_name, b2_name) {
852   b_sum <- beta[b1_name] + beta[b2_name]
853   var_sum <- V[b1_name, b1_name] + V[b2_name, b2_name] + 2 * V[b1_name, b2_name
        ]
854   return(se_sum <- sqrt(var_sum))
855 }
856
857 OR_activity_1 <- exp (beta[activity])
858 or_act_ap0 <- exp(reduced_model$coefficients[2])
859
860 CI_or_act_ap1 <- exp(summary(reduced_model)$coefficients[2,1] + summary(reduced
      _model)$coefficients[5,1]
861                      + 1.96 * se_sum("activityYes", "activityxpollution") *
862                        c(-1, 1))
863
864 CI_or_act_ap0 <- exp(summary(reduced_model)$coefficients[2,1]
865                      + 1.96 * summary(reduced_model)$coefficients[2,2] *
866                        c(-1, 1))
867
868 # Genderxairpollution
869 # odds ratio again depends on air pollution, but there is no term for airpo = 0
870
871 or_gender_ap1 <- exp(reduced_model$coefficients[4])
872 ci_or_gender_ap1 <- exp(summary(reduced_model)$coefficients[4,1]
873                        + 1.96 * summary(reduced_model)$coefficients[4,2] *
874                          c(-1, 1))
875
876 # age --> measured for 10 years
```

```r
or_age_ap1 <- exp(reduced_model$coefficients[3] * 10)
ci_or_age_ap1 <- exp(summary(reduced_model)$coefficients[3,1] * 10
                     + 1.96 * 10 *summary(reduced_model)$coefficients[3,2] *
                       c(-1, 1))

# airpollution
# start with age = 25
act1_gender1 <- function(age) {
  or_ap_a1g1 <- exp(reduced_model$coefficients[3]*age +
                      reduced_model$coefficients[4] +
                      reduced_model$coefficients[5])

  se_sum_three <- sqrt(se_sum("genderxpollution", "activityxpollution") +
                         age ** 2 * V["agexpollution", "agexpollution"] +
                         2 * age * (V["agexpollution", "genderxpollution"] +
                                      V["agexpollution", "activityxpollution"])
  )

  CI_or_ap_a1g1 <- exp(reduced_model$coefficients[3]* age +
                         reduced_model$coefficients[4] +
                         reduced_model$coefficients[5] +
                         1.96 *
                         se_sum_three * c(-1, 1))
  output <- c(odds_ratio = or_ap_a1g1,
              conf_int = CI_or_ap_a1g1)
  return(output)
}

act1g1confints <- lapply(c(25,35,45,55), act1_gender1)

act1_gender0 <- function(age) {
  or_ap_a1g1 <- exp(reduced_model$coefficients[3]*age +
                      reduced_model$coefficients[5])

  se_sum_three <- sqrt(V["activityxpollution", "activityxpollution"]+
                         age ** 2 * V["agexpollution", "agexpollution"] +
                         2 * age * V["agexpollution", "activityxpollution"]
  )

  CI_or_ap_a1g1 <- exp(reduced_model$coefficients[3]* age +
                         reduced_model$coefficients[5] +
                         1.96 *
                         se_sum_three * c(-1, 1))
  output <- c(odds_ratio = or_ap_a1g1,
              conf_int = CI_or_ap_a1g1)
  return(output)
}

act1g0confints <- lapply(c(25,35,45,55), act1_gender0)

act0_gender1 <- function(age) {
  or_ap_a1g1 <- exp(reduced_model$coefficients[3]*age +
                      reduced_model$coefficients[4])

  se_sum_three <- sqrt(V["genderxpollution", "genderxpollution"]+
                         age ** 2 * V["agexpollution", "agexpollution"] +
                         2 * age * V["agexpollution", "genderxpollution"]
  )

  CI_or_ap_a1g1 <- exp(reduced_model$coefficients[3]* age +
                         reduced_model$coefficients[4] +
                         1.96 *
                         se_sum_three * c(-1, 1))
```

```
940    output  <- c( odds_ratio = or_ap_a1g1 ,
941                  conf_int = CI_or_ap_a1g1)
942    return ( output )
943  }
944
945  act0g1confints <- lapply (c (25 ,35 ,45 ,55) , act0_gender1 )
946
947  act0_gender0 <- function ( age ) {
948    or_ap_a1g1 <- exp ( reduced_model$coefficients [3]* age )
949
950    se_sum_three <- sqrt ( age ** 2 * V["agexpollution", "agexpollution"])
951
952    CI_or_ap_a1g1 <- exp ( reduced_model$coefficients [3]* age +
953                          1.96 *
954                          se_sum_three * c ( -1 , 1))
955    output  <- c( odds_ratio = or_ap_a1g1 ,
956                  conf_int = CI_or_ap_a1g1)
957    return ( output )
958  }
959
960  act0g0confints <- lapply (c (25 ,35 ,45 ,55) , act0_gender0 )
```

# References

[1] Annette J. Dobson and Adrian G. Barnett. *An Introduction to Generalized Linear Models.* Chapman and Hall/CRC Texts in Statistical Science. CRC Press, Boca Raton, FL, 4 edition, 2018.

[2] Julian J. Faraway. *Extending the Linear Model with R.* Chapman and Hall/CRC, Boca Raton, FL, 2 edition, 2016.