

# Analysis of Health Survey for England (HSE) 2019

Candidate Numbers: 24226, 14925, 23766, 24170.

March 15, 2024

## **Abstract**

This report provides an analysis of data related to health, age, socio-economic factors and lifestyle habits in adults (from the age of 16) from the population in England, derived from the Health Survey for England 2019. Please note, no generative AI technology was utilized in the creation of this report or the subsequent data analysis presented herein. All findings and conclusions are derived from human-driven methodologies.

# Summary

In England, smoking, “vaping” and alcohol consumption are widespread habits among adults, particularly younger adults. In fact, it is estimated that around 13.9% of adults in England identified as cigarette smokers (ONS 2019a), and a similar study found that there were 10.9 alcohol-related deaths per 100,000 people living in England in 2019 (ONS 2019c). Thus, it is crucial to be aware of the predictors and consequences of these habits, which is the motivation for this analysis.

We conducted this analysis on a subset containing 8,204 adults living in England (ages 16+), who were interviewed in their homes about their demographics and their smoking and drinking habits as part of the Health Survey for England (HSE) 2019. 4,947 of these participants who consented to an at-home nurse visit had vitals such as blood pressure, weight and height measured. We investigated which features of a participant are associated with smoking habits, and which habits are associated with high blood pressure.

We found alcohol consumption to be the most prevalent lifestyle habit under investigation. Older males were found to drink the most, especially those aged 65-74. Also, smoking tended to be more prevalent in areas of greater deprivation, whilst alcohol consumption was found to be more prevalent in less deprived areas. Socio-economic factors played an important role in the likelihood of smoking, with females 25% less likely to be smokers than males. Ethnicity also played a great role, as being Black or Asian decreased the likelihood of smoking by more than 70% and 65%, respectively. Females were much less likely to experience hypertension in comparison to males. However, systolic BP levels were positively correlated with age and how often the respondent drinks alcohol, regardless of gender.

## Introduction

The mechanisms which drive lifestyle habits are part of a complex and ever-changing field in behavioural psychology. What is known, however, is that such habits are driven by cues and cravings (Anastasia Droungas 1995; Kambouropoulos 2009). The HSE 2019 captures demographic and socioeconomic characteristics that potentially make such cues more visible and cravings harder to resist. For example, living in a region where smoking is more prevalent or not having a husband/wife for an accountability partner to help you quit or resist these habits (ONS 2019b).

Throughout, we will use weighting variables to estimate prevalence of cigarette, e-cig, and alcohol usage in our sample, and compare it with the population of adults in England. We will then attempt to identify associations with a participant’s current smoking status, explore whether predictive modelling can be used to accurately identify smokers, and investigate which lifestyle habits are associated with increased values of systolic blood pressure. Such an understanding would be crucial in enabling healthcare providers to implement preventative care for people at risk of developing hypertension.

Table 1: Summary of the analysis variables used in our report. (D) indicates that a variable was derived.

Variable	Definition	Notation
SerialA	Respondent serial number	ID
wt_int	Weight of observation	w
Sex	Respondent sex (M/F)	s
(D) Age35g	Respondent age grouped in 5 year bands for 16+	$a^{(5)}$
(D) ag16g10	Respondent age grouped in 10 year bands for 16+	$a^{(10)}$
(D) topqual2	Highest level qualification	t
(D) marstatD	Marital status	ms
(D) qimd19	IMD score (a measure of deprivation)	q
(D) urban14b	Level of rurality	u
(D) origin2	Grouped ethnic category	o
(D) cigdyl_19	# of cigarettes smoked per day	$n^{cig}$
(D) cigsta3_19	Smoking status (Reg / Ex-Reg / Never-Reg)	cig
(D) NDPNow_19	Use of E-cigarettes and or NDPs	ndp
(D) DrinkYN_19	Alcohol consumed in the last 12 months	alc
dnofit_19	Frequency of alcohol consumed in the last 12 months (grouped)	f
(D) d7many3_19	# of days alcohol was consumed in the last week	$n^{alc}$
(D) GOR1	Government region office number	g
(D) BMIval	Valid BMI measurement (weight estimated if 130kg)	bmi
(D) omsysval	Valid mean systolic blood pressure	sbp
(D) smoking_status	Current smoking status (binary)	S
(D) age_estim	Estimated age	a

## Methodology

### Exploratory Analysis

The full HSE 2019 cross-sectional data set contains 10,299 observations across thousands of variables (NatCen Social Research and Health 2019), but we only studied patients over the age of 16 among key variables. Our subset included 8,204 participants and 19 variables, which are described in Table 1.

We found 36 pairs of observations with exactly equal variables (excluding ID variables and lab measurements), but we did not remove these from our analysis because the supporting documentation didn’t state a protocol for repeated visits. We assumed these were genuine observations from different participants and thus included them in our dataset. However, there were 3 pairs of observations that had duplicate lab variables also. Due to the high precision of measurements, we believed these to be true duplicates and removed one observation from each of the three pairs, resulting in a dataset of 8201 observations.

All variables were originally coded as numeric in our dataset, so we recoded these to factor variables accordingly. Some of the variables used for analysis were dichotomised by us for easier interpretation. We coded a binary variable indicating current smoking status, which served as a response variable. Also, we dichotomised urbanity into two levels being ‘Rural’ and ‘Urban’, and finally marital status into ‘Married’ and ‘Not Married’. We grouped respondent’s level of education into ‘Higher Education’, ‘Further Education’, ‘A-Level equiv.’, ‘GCSE equiv.’, ‘No qualification’ and ‘Foreign/other’, and we grouped alcohol frequency into “Frequent”, “Occasional” and

Table 2: Missing values in the training dataset

Variable	Missing Values	% Missing
omsysval	4036	61.5%
BMIVal	1519	23.2%
dnoft_19	1496	22.8%
cigdyal_19	57	0.869%
cigsta3_19	56	0.854%
NDPNow_19	53	0.808%
d7many3_19	52	0.793%
drinkYN_19	51	0.777%
topqual2	46	0.701%
origin2	29	0.442%
marstatD	1	0.0152%

“Rarely” based on the response to the question “How often did you consume alcohol in the last 12 months?”.

There is significant missingness in the lab values, as shown in Table 2. These can be explained by the fact additional consent was needed from the participant to allow a nurse visit. One variable with a substantial amount of missing data was the frequency of drinking in the last 12 months, which due to the retrospective and sensitive nature of the question could be explained by either recall bias or a participant’s refusal to answer. As a result, we decided to remove this variable from analysis in favour of a binary variable which simply states whether the respondent has drunk alcohol in the last 12 months. Education level, marital status and ethnicity are the only socioeconomic variables with any missing entries, with 46 observations (0.701% of the data) missing at least one of the three. These were not necessary for identifiability and as they are sensitive, we did not expect every participant to answer these questions. Therefore, we did not remove these observations.

We attempted to determine whether the two lab measurement variables contained potential outliers, and as can be seen in Figure 1, there were many outliers for BMI. We note that calculated BMI used self-estimated weight from the participant if this weight exceeded 130kg. Additionally, weight measurements were taken in inconsistent environments such as on different flooring in the participants’ homes which is known to impact measurement accuracy (Scientist 2002). We believed this could have explained the extremely high BMI values in the range of 54.82 to 73.49. We coded such values as missing in the analysis dataset, but kept the observation.

The readings for systolic blood pressure were collected by taking an average of three readings, each five minutes apart, performed by a trained nurse who had to declare each reading to be valid. For this reason, it seemed highly unlikely that these readings were mistakes, and so we included them in our analysis. The measurement device used has been validated for use in this environment (Yechiam Ostchega PhD 2009). We did, however, note that values of systolic blood pressure that were this consistently high ( $>140\text{mmHg}$ ) were indicators for hypertensive crisis, and so a part of our population may have serious underlying health conditions that could affect the generalisability of our findings (Association 2023).

Table 3: Estimates and 95% Confidence Intervals for % of Population

Habit	Estimate	C.I.
Drinking	80.4%	(79.5%, 81.2%)
Smoking	16.5%	(15.7%, 17.3%)
Smoking E-cigarettes	4.28%	(3.84%, 4.72%)

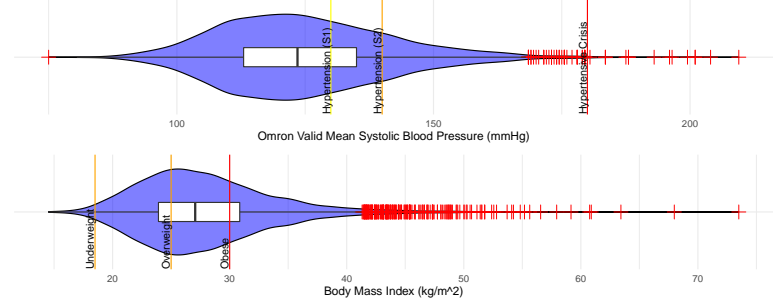


Figure 1: Distribution of BMI and Mean Systolic Blood Pressure

## Analysis

### What is the prevalence of drinking, smoking and E-cig usage?

To calculate the prevalence of these three habits we assumed each of the  $n$  observations,  $x_1, \dots, x_n$ , to be independent, identically distributed (iid) random variables (RVs) where  $x_i \sim \text{Bern}(p) \forall i = 1, \dots, n$  and  $p$  denotes the probability of a habit being present for the given observation.

We used the household-level weighting variable to calculate a weighted Maximum Likelihood Estimate (MLE) of  $p$ . That is, letting  $w_i$  denote the weight of the  $i^{\text{th}}$  observation, we altered the standard likelihood function of a Bernoulli distribution as below:

$$L(p|\mathbf{x}) = \prod_{i=1}^n (p^{x_i} (1-p)^{1-x_i})^{w_i}$$

From this, we calculated our weighted MLE as  $\hat{p} = \frac{\sum_{i=1}^n x_i w_i}{\sum_{i=1}^n w_i}$ . It can also be shown that the MLE has variance given by  $\text{var}(p) = \frac{p(1-p)}{\sum_{i=1}^n w_i}$ , which we estimated using  $\hat{p}$ . We used large sample properties of the MLE to get a normal approximation and estimated 95% confidence intervals for each habit, which are shown in Table 3.

We found that the e-cigarette usage among adults is relatively low, making it challenging to dissect any significant trends within the data. However, since drinking and smoking were so prevalent among the respondents, we had sufficient data to identify any potential trends even within smaller groups of our sample.

Next, we worked to uncover factors that may have associations with smoking prevalence. We started with the demographic factors of age and gender, and plotted the prevalence of these habits across the groups in Figure 2. The plot suggests a negative correlation between age and the prevalence of current cigarette smokers across both genders, with males having

higher drinking prevalence across nearly all age groups. Interestingly, we found in our data that the proportion of males who quit smoking in later-life was much greater than that of females (75yrs+; M: 50%, F: 31%).

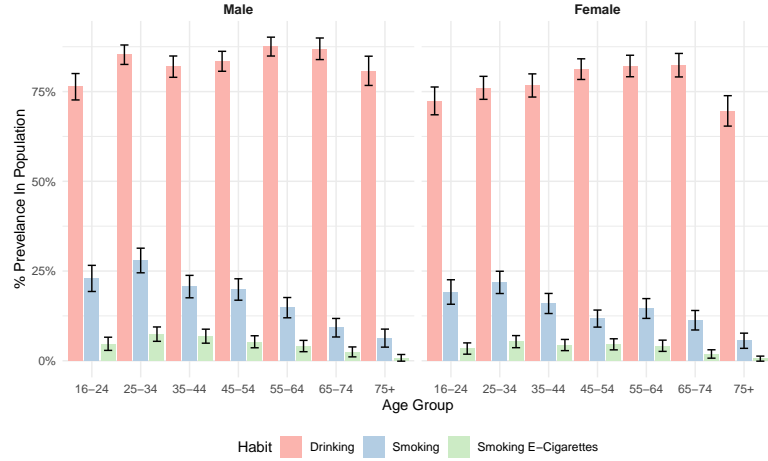


Figure 2: Estimation of the prevalence of habits by age group and gender.

During our analysis, we also observed a positive correlation between smoking and deprivation levels, with 25.76% of individuals in the most deprived quintile being smokers, a significant increase from 9.33% in the least deprived.

### How is smoking associated with socioeconomic factors and age?

We reserved 80% of our data to train the model, and used the remaining 20% to test the model, which was made possible due to the large size of our dataset. The training set contained 6561 observations and the test set contained 1640. Table 4 summarises the key variables between the test and train dataset to illustrate that both are representative of the whole dataset. This reduced the risk of overfitting and allowed us to test the predictive power of each model we proposed.

To develop a predictive model for identifying smokers, we used the binary variable for current smoking status as a response, with various socioeconomic and demographic factors as predictors. These were selected based on the associations suggested by our prior analysis.

Under the assumption that the age of participants were approximately uniformly distributed within their respective age bands, we defined the estimated age of the  $i^{th}$  observation as the midpoint of the participant's respective 5-year age group (taking the estimation of the 90+ category as 92.5), denoted  $a_i$ . We found that age may represent a somewhat quadratic effect on the probability of smoking, leading us to include a  $a_i^2$  term in our model. A comparison of these models are summarised in Table 5.

We selected the model based on AIC, which enabled us to balance model complexity and model fit (quantified by likelihood), which further reduced the risk of overfitting. This

Table 4: Comparison of the characteristics of the training set compared with the test set. P-values derived from two-sample t-test for means and proportion test for proportions.

Variable:Label	Proportion (%)		
	Test	Train	P Value
Sex:Male	43.5	45.1	0.24
Sex:Female	56.5	54.8	0.24
topqual2:No qualification	19.7	19.8	0.93
topqual2:GCSE equiv.	21.2	20.9	0.81
topqual2:A-Level equiv.	14.0	13.4	0.58
topqual2:Further Education	28.5	28.4	1.00
topqual2:Higher Education	15.9	16.1	0.89
topqual2:Foreign/other	0.7	1.3	0.08
marstatD:Married	55.4	51.9	0.01
marstatD:Not Married	44.6	48.1	0.01
urban14b:Urban	80.3	81.4	0.32
urban14b:Rural	19.7	18.6	0.32
origin2:White	85.9	85.8	1.00
origin2:Black	3.1	2.9	0.71
origin2:Asian	8.8	8.7	0.90
origin2:Multiple	1.4	1.7	0.50
origin2:Other	0.8	0.9	0.79
	<b>Mean</b>	<b>Mean</b>	
Age(Estimated)	51.2	51.0	0.80

Table 5: Comparison of selected model evaluations, wherein  $\eta_i = \text{logit}(\mu_i)$ , and Acc. refers to the model accuracy on test data.

Model	AIC	AUC	RMSE		Acc.
			Train	Test	
$\eta_i \sim a_i + a_i^2 + ms_i + q_i + u_i + o_i + t_i + s_i$	4987.3	0.759	0.342	0.324	0.861
$\eta_i \sim a_i + ms_i + q_i + u_i + o_i + t_i + s_i$	5036.4	0.740	0.343	0.327	0.861
$\eta_i \sim a_i^{(5)} + ms_i + q_i + u_i + o_i + t_i + s_i$	4994.3	0.759	0.341	0.325	0.861
$\eta_i \sim a_i^{(10)} + ms_i + q_i + u_i + o_i + t_i + s_i$	5005.9	0.753	0.342	0.324	0.860

model is:

$$S_i \stackrel{\text{ind}}{\sim} \text{Binomial}(\mu_i)$$

$$\text{logit}(\mu_i) = \beta_0 + \beta_1 a_i + \beta_2 a_i^2 + \beta_3 q_i + \alpha_j^{ms} + \alpha_k^u + \alpha_l^o + \alpha_n^t + \alpha_m^s$$

where:  $\iota = ijklmn$ ,

$j \in \{\text{Married, Not Married}\}$ ,

$k \in \{\text{Urban, Not Urban}\}$ ,

$l \in \{\text{White, } \dots, \text{Black}\}$ ,

$n \in \{\text{Male, Female}\}$

$m \in \{\text{Higher Education, } \dots, \text{No qualification}\}$

wherein:  $\alpha_{\text{Not Married}}^{ms} = \alpha_{\text{Not Urban}}^u = \alpha_{\text{White}}^o = \alpha_{\text{Male}}^t = \alpha_{\text{No qualification}}^s = 0$

To evaluate the predictive performance of our final model using our test data, Figure 3 demonstrates the predicted probability of each 'probability bin' against the mean actual outcomes. As we can see the calibration curve closely follows the line  $y = x$ , which is indicative of a well-calibrated model.

This model doesn't predict high values due to the low prevalence of smoking within the population. The highest likelihood our model predicts is 63.75%. All results should be interpreted with the population average in mind, which is ~16.5%. For example, if the model predicts the likelihood of an individual being a smoker at 33%, said individual is twice as more likely to be a smoker than the average. This effect can be seen in the following calibration chart.

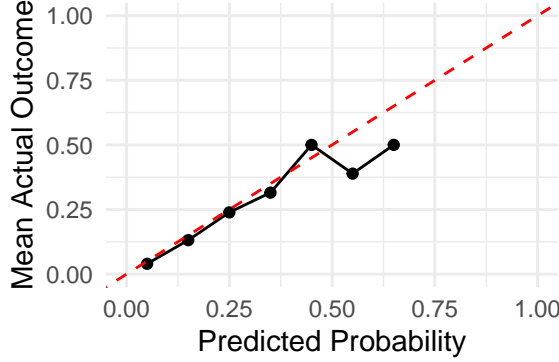


Figure 3: Calibration chart for Binomial model

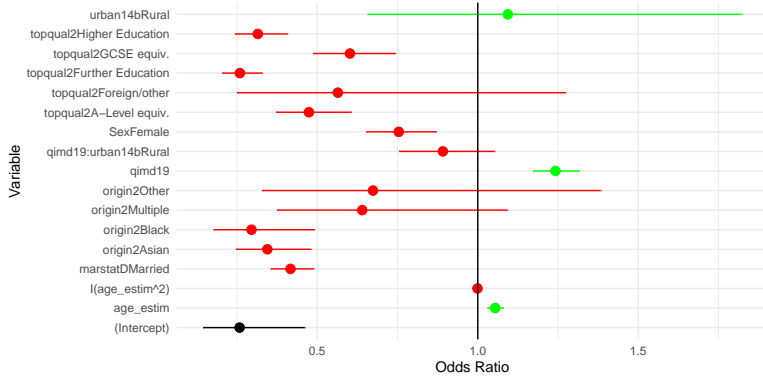


Figure 4: Forest plot of Odds Ratios for the binomial model

For interpretability, we decided to include a forest plot of the odds ratio for each variable in our model, as the log-odds can be difficult to interpret. The intercept 0.257 can be interpreted as the probability of a white, non-married, unqualified male living in urban area being a smoker, using the reference levels of each factor. On average, the green factors increase the probability and the red factors decrease the probability. Being Black, Asian, Married or a female significantly decreases this probability, so this is evidence of a socioeconomic effect.

### Which lifestyle habits are associated with systolic blood pressure?

We chose to model the systolic blood pressure (BP) using the Inverse-Normal distribution with a  $1/\mu_i^2$  link function. This distribution was a suitable choice because this variable has a lot of outliers, is right-skewed and strictly positive. We also tested a Log-Normal distribution, a Normal distribution with identity link and a Gamma distribution with an inverse link but found that the Inverse-Normal distribution was best able

to account for the higher values of BP in the Q-Q plots. The data we used to fit this model was the same training data set as the Binomial model. We filtered out *NA* values of BMI and BP as we only want to consider respondents that had lab measurements taken, leaving us with 2960 observations.

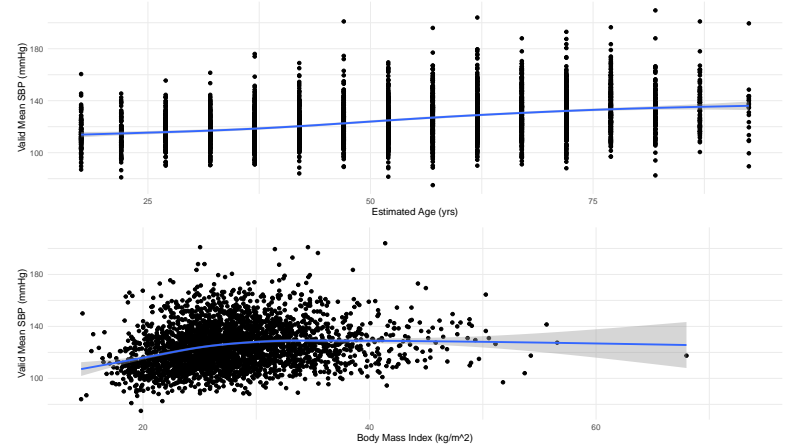


Figure 5: Relationship of BMI and Age with Mean Systolic Blood Pressure

We studied the relationship between both age and BMI with BP, shown in Figure 5. There was a quadratic relationship between BMI and BP so we decided to include a quadratic term in our model. To account for possible dependencies between variables, we tried fitting models with different interaction terms and found three that were significant. These were between BMI & Age (Medicine 2017), BMI & gender (K. et al. 2019) and Age & Sex (M. et al. 2019), which are supported by previous research.

To determine the predictors in the final model, we used a stepwise approach starting with the full model and eliminating variables based on AIC. The final model selected was

$$\begin{aligned} \text{sbp}_i &\stackrel{\text{ind}}{\sim} \text{IG}(\mu_i, \lambda) \\ (\mu_i)^{-2} &= \beta_0 + \beta_1 \text{cig}_i + \beta_2 a_i + \beta_3 \text{bmi}_i + \beta_4 \text{bmi}_i^2 + \beta_5 n_i^{\text{alc}} + \\ &\quad + \beta_6 \text{bmi}_i \alpha_n^s + \beta_7 a_i \alpha_n^s + \beta_8 \text{bmi}_i a_i + \alpha_j^{ms} + \alpha_k^u + \alpha_n^t \end{aligned}$$

where:  $i = ijkn$

$j \in \{\text{Married, Not Married}\},$

$k \in \{\text{Urban, Not Urban}\},$

$n \in \{\text{Male, Female}\}$

wherein:  $\alpha_{\text{Not Married}}^{ms} = \alpha_{\text{Not Urban}}^u = \alpha_{\text{Male}}^t = 0$

The Q-Q plot in Figure 6 shows a good fit to the distribution and, despite some deviation at both ends of the scale, the model captured the outliers more accurately than other distributions we tried. On our testing dataset we achieved a RMSE of 14.3mmHg, which was relatively small compared to the threshold for hypertensive crisis of 140mmHg+.

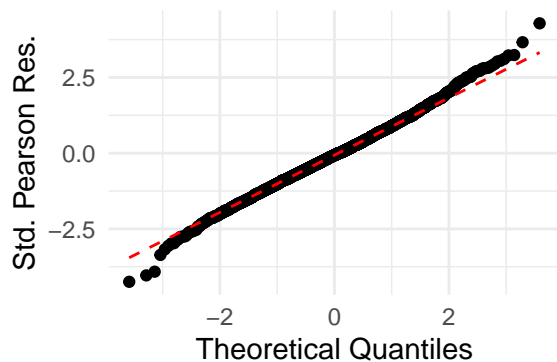


Figure 6: Q-Q plot of Inverse-Normal model residuals

We found that, on average, young females are predicted lower blood pressure than young men but as they age, their blood pressure increases at a faster rate. Controlling for other lifestyle choices included in the model, a female's blood pressure can be expected to exceed an equivalent male's at around 72 years old. Figure 7 shows this relationship. One of the biggest influences on BP was smoking status, and we found smoking was associated with higher blood pressure for both males and females, also shown in Figure 8. This effect is largest for males and a typical male can expect a 1.69% increase in BP by taking up smoking.

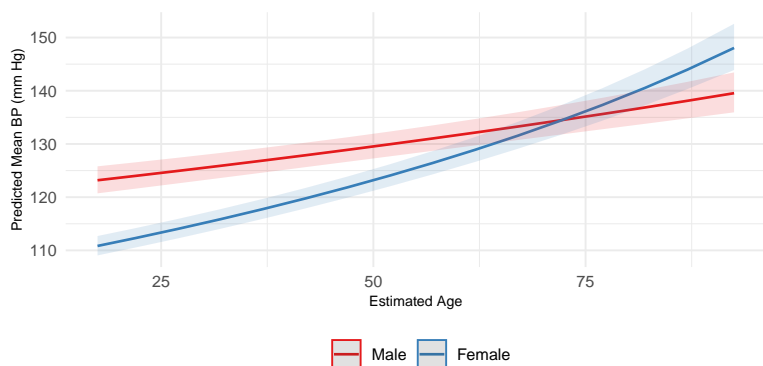


Figure 7: Marginal effects of Age on Systolic Blood Pressure

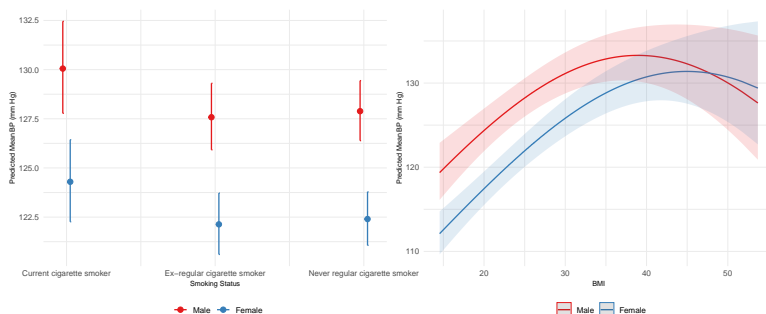


Figure 8: Marginal effects of Smoking Status and BMI on Systolic Blood Pressure

Drinking was also related to higher BP, and we observed a linear relationship in Figure 9 with the number of days that alcohol was consumed in the previous week. On average, one extra day of drinking led to an increase of 0.74 mmHg.

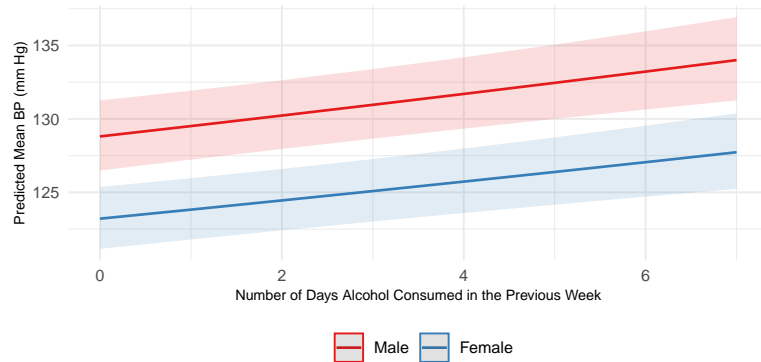


Figure 9: Marginal effects of alcohol consumption on Systolic Blood Pressure

## Results/Conclusion

Our analysis showed that alcohol consumption is highly common in our sample, with approximately 80.4% being consumers, while smoking and vaping rates were lower at 16.5% and 4.3%, respectively. Older males showed the highest tendencies for alcohol consumption, with 'frequent' drinking the most prevalent among males aged 65-74. Unlike 'frequent' alcohol consumption, the prevalence of smoking decreases with age. Individuals aged 16-24 appear to be the most susceptible to smoking cigarettes, indicating a significant issue within youth culture. Moreover, we found interesting, opposing associations between deprivation levels and smoking and drinking behaviours. Smoking prevalence seemed to increase as deprivation levels decreased, whereas alcohol consumption appeared to do the opposite.

To help us uncover the underlying socio-economic factors that may drive smoking prevalence, we looked at the most 'comparable' respondent in our dataset. This 'respondent' is a white, single male who has no qualifications and lives in an urban area. We found that the probability of our hypothetical respondent being a smoker is around 25%, which can be seen from Figure 4. The smoking prevalence estimate in our sample is significantly higher than that of the population of interest (16.5% vs. 13.9%, outside of the CI in Table 3). The only socio-economic variables to certainly increase this probability is the deprivation level our respondent falls under. The rest of the socio-economic variables we have access to, appear to decrease the likelihood of smoking. For example, if our respondent is any of the following: Asian, black, married, or went on to further education, the chances of them being a smoker is at-least halved.

If we had more data from current smokers, we may have been able to increase the accuracy of our final model for smoking status. Due to the cross-sectional format of the study, we are unable to comment on temporal changes in habits or establish a causal relationship. Furthermore, access to data from the whole of the UK could improve generalisability, but including participants with hypertensive crisis helped to generalise our findings to those who may have the most need for preventative care from healthcare providers.



## References

- al., Katsuya et. 2019. "The Investigation of Sex Differences in the Effect of Body Mass Index." <https://www.hindawi.com/journals/ijhy/2019/1360328/>.
- al., Martins et. 2019. "The effect of gender on age-related blood pressure changes and the prevalence of isolated systolic hypertension among older adults." <https://pubmed.ncbi.nlm.nih.gov/11605350/>.
- Anastasia Droungas, Anna Rose Childress, Ronald N. Ehrman. 1995. "Effect of smoking cues and cigarette availability on craving and smoking behavior." <https://www.sciencedirect.com/science/article/pii/S030646039500029C>.
- Association, American Heart. 2023. "Understanding Blood Pressure Readings." <https://www.heart.org/en/health-topics/high-blood-pressure/understanding-blood-pressure-readings>.
- GOV.UK. 2014. "Tax on shopping and services." <https://www.gov.uk/tax-on-shopping/beer-cigars>.
- Kambouropoulos, Nicolas. 2009. "' Cue reward salience' predicts craving in response to alcohol cues." <https://www.sciencedirect.com/science/article/abs/pii/S0191886908003127>.
- Medicine, National Library of. 2017. "Weight and Body Mass Index in Old Age: Do They Still Matter?" <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5704942/>.
- NatCen Social Research, Department of Epidemiology, University College London, and Public Health. 2019. "Health Survey for England." <http://doi.org/10.5255/UKDA-SN-8860-1>.
- ONS. 2019a. "Adult smoking habits in the UK: 2019." <https://shorturl.at/DMQXZ>.
- . 2019b. "Adult smoking habits in the UK: 2019." <https://shorturl.at/itPQU>.
- . 2019c. "Alcohol-specific deaths in the UK: registered in 2019." <https://shorturl.at/gqxY1>.
- . 2023. "Deprivation and the impact on smoking prevalence, England and Wales: 2017 to 2021." <https://shorturl.at/fuy48>.
- Scientist, New. 2002. "People weigh less on a hard surface." <https://www.newscientist.com/article/dn2462-people-weigh-less-on-a-hard-surface/>.
- Yechiam Ostchega PhD, Tatiana Nwankwo MS, RN. 2009. "Assessing the Validity of the Omron HEM-907XL Oscillometric Blood Pressure Measurement Device in a National Survey Environment." <https://onlinelibrary.wiley.com/doi/full/10.1111/j.1751-7176.2009.00199.x>.