

Code for the manuscript titled Causal association between systolic
blood pressure and quality-adjusted life years: a mendelian
randomisation study.

Tamrat Befekadu Abebe*

April 10, 2025

*Centre for Medicine Use and Safety, Faculty of Pharmacy and Pharmaceutical Sciences, Monash University,
tamrat.abebe@monash.edu

Contents

1	Acknowledgement	4
2	Preface	5
3	Introduction	6
4	Steps	6
4.1	Step 1: Working on Phenotype data	7
4.1.1	Main dataset import	7
4.1.2	Renaming variables	10
4.1.3	Exclusion criteria	14
4.1.4	Hospital admission data	14
4.1.4.1	Working on ICD-10 diagnosis codes	18
4.1.4.2	Working on ICD-9 diagnosis codes	34
4.1.5	Setting up the data for HRQoL prediction	45
4.1.5.1	Predicting Utility	46
4.1.6	Working on EQ-5D data collected by UK Biobank	59
4.2	Step 2: Working on genotype data	63
4.2.1	Preparing our data	63
4.2.2	Association of Genetic Variants with blood pressure traits	64
4.2.3	LD Clumping	67
4.2.4	Selecting the genetic variants from the BGEN file	70
4.2.5	Association of Genetic Variants with Quality-Adjusted Life Years	71
4.2.6	Data harmonisation	73
4.2.7	Calculating polygenic risk scores	74
4.3	Step 3: Combining Phenotype and Genotype data	75
4.4	Step 4: Main analysis	76
4.5	Step 5: Sensitivity analyses	78
4.5.1	Untreated populaton	78
4.5.2	Two-sample MR	80
4.5.3	Sub-group analysis	83
4.5.4	Non-linear MR	85
4.5.4.1	Rerun the main MR within the fifty quantiles of PRS-free SBP	85
4.5.4.2	Sub-group analysis mainly stratified by PRS-free SBP	86
4.5.4.3	Select observations for all participants in each SBP	88
4.5.4.4	Estimate the mean SBP for each quantile	88
4.5.4.5	Merge the previous datasets for VWLS analyses and plots	89
4.5.5	EQ-5D index from UK Biobank survey	92
4.6	Step 6: Secondary analysis	94
4.6.1	Multivariable Mendelian Randomisation	94
4.6.1.1	Working on the DBP data	94
4.6.1.2	Clumping SNPs in LD	94
4.6.1.3	Selecting the genetic variants from the BGEN file	95
4.6.1.4	Association of Genetic Variants with Quality-Adjusted Life Years	97

4.6.1.5	Data Harmonisation	98
4.6.1.6	Combining the genetic data	98
4.6.1.7	Working on SNPs primarily associated with SBP and seconarily associated with DBP	100
4.6.1.8	Working on SNPs primarily associated with DBP and seconarily associated with SBP	102
4.6.1.9	Combining Phenotype and Genotype data	103
4.6.1.10	Analysis	104
4.7	Step 7: Tables and Figures	106
4.7.1	Tables	106
4.7.1.1	Table 1: Background characteristics	106
4.7.1.2	Table 2: Main analysis	110
4.7.1.3	Table 3: Sensitivity analyses - No history of antihypertensive med- ication cohort	111
4.7.1.4	Table 4: Sensitivity analyses - Two-sample MR	112
4.7.1.5	Table 5: Sensitivity analyses - EQ-5D-index from UK Biobank survey	113
4.7.1.6	Table 6: Secondary analysis	114
4.7.2	Figures	115
4.7.2.1	Figure 1	115
4.7.2.2	Figure 2	116
4.7.2.3	Figure 3	117
4.7.2.4	Figure 4	120
4.7.2.5	Figure 5	122

1 Acknowledgement

Firstly, I would like to express my deepest gratitude to **Professor Zanfina Ademi** for her exceptional supervision, unwavering support, and for giving me the opportunity to work on this project. Her guidance has been instrumental throughout the entire journey.

I am also sincerely thankful to **Dr. Jedidiah I. Morton** for his consistent encouragement, insightful advice, and for continually inspiring me to expand my skillset.

My heartfelt thanks go to **Dr. Padraig Dixon** for his invaluable input, particularly in guiding the analysis of the data, which significantly enhanced the quality of this work.

Lastly, I would like to thank **Dr. Jenni Ilomaki** for her helpful feedback and constructive comments on the overall project.

I gratefully acknowledge **Mr. Adam Livori** for being a source of inspiration to work in TeXdoc and LaTeX, and for the time he generously invested in guiding me through the essential materials for setting up this documentation.

In addition, I acknowledge the work of **Sean Harrison and colleagues**[1], from whom the current Stata code was adapted. GitHub link: <https://github.com/seanharrison-bristol/Robust-causal-inference-for-long-term-policy-decisions>.

This study is supported by the National Health and Medical Research Council Ideas Grants Application ID: 2012582. The funder had no input into the design of the study or decision to submit for publication.

2 Preface

This document presented the code and workflow for the manuscript titled *Causal association between systolic blood pressure and quality-adjusted life years: a mendelian randomisation study*. It detailed the data preparation (cleaning) that was performed using a dataset provided by the UK Biobank. This study was approved by the Research Ethics Committee (REC reference for UK Biobank is 11/NW/0382) of National Committee North West-Haydock, National Health Service, UK.

To generate this document, the Stata package `texdoc` was used, which is available at: <http://repec.sowi.unibe.ch/stata/texdoc/>.

Our code is also available at: <https://github.com/tamrat-works/mr-sbp-qalys>.

3 Introduction

Hypertension has long been recognised as a major risk factor for heart disease due to both oxidative and mechanical stress exerted on the arterial wall[2]. A recent study reported that for every 10 mmHg increase in systolic blood pressure (SBP), there was a 53% higher risk of atherosclerotic cardiovascular disease[3]. Mendelian randomisation (MR) studies have also demonstrated a lifetime causal association between SBP and cardiovascular disease[4, 5].

However, a gap remained in the literature regarding the lifetime causal association between SBP and health-related quality of life (HRQoL). One possible approach to demonstrate this causal association could have been a randomised controlled trial (RCT). However, the feasibility of conducting such trials, along with the limited generalisability of their findings, may hinder their applicability. An alternative approach was to conduct an observational study, which is often more cost-effective and may yield findings generalisable to the broader population. Nonetheless, observational studies carry inherent limitations, such as confounding and reverse causation[6].

Mendelian randomisation offered a potential solution to these limitations by using genetic variants as instrumental variables for modifiable traits (i.e., risk factors) associated with outcomes. These outcomes could include clinical conditions (e.g., coronary artery disease) or HRQoL.

This document was developed to showcase the application of MR techniques to investigate the lifetime association between SBP and HRQoL using data from the UK Biobank. HRQoL data were sourced from[7]. The following steps were undertaken to conduct the study.

4 Steps

- Step 1: Working on Phenotype data
- Step 2: Working on Genotype data
- Step 3: Combining Phenotype and Genotype data
- Step 4: Main analysis
- Step 5: Sensitivity analyses
- Step 6: Secondary analysis
- Step 7: Tables and Figures

4.1 Step 1: Working on Phenotype data

4.1.1 Main dataset import

The UK Biobank main dataset contained phenotype data for participants enrolled in the study. Briefly, more than 500,000 individuals were recruited across 22 centres in the UK between 2006 and 2010[8]. Hospital admission data and primary care data (general practice) for UK Biobank participants were linked to Hospital Episode Statistics (HES) in England and Scottish Morbidity Records (SMR) in Scotland up to 31 October 2022, and to the Patient Episode Database for Wales (PEDW) up to 26 May 2022.

Importing the entire dataset into STATA required significant time and computational resources. Therefore, it was more efficient to first select the key variables relevant to the study using the data dictionary.

For those primarily using the UK Biobank Research Analysis Platform (UKB RAP), data extraction was carried out through DNAnexus's JupyterLab environment, specifically using Spark JupyterLab. The following lines of Python code were used to extract the necessary variables.

```
1
2 # Building cohorts using Spark JupyterLab
3
4 # Folders
5 exome_folder = 'Population level exome OQFE variants, PLINK format -
    interim 450k release'
6 exome_field_id = '23149'
7 output_dir = '/Data/'
8
9 # Import important variables
10 import os
11
12 # Set environment variable before importing pyspark
13 os.environ['PYARROW_IGNORE_TIMEZONE'] = '1'
14
15 # Import necessary libraries
16 import pyspark
17 from pyspark.sql import SparkSession
18 from pyspark.sql.functions import col
19 import dxdpy
20 import dxddata
21 import pandas as pd
22 import re
23
24 # Initialize Spark
25 # Spark initialization (Done only once; do not rerun this cell unless you
    select Kernel -> Restart kernel).
26 sc = pyspark.SparkContext()
27 spark = pyspark.sql.SparkSession(sc)
28
29 # Automatically discover dispensed dataset ID and load the dataset
30 dispensed_dataset = dxdpy.find_one_data_object(
31     typename="Dataset",
```

```

32     name="app*.dataset",
33     folder="/",
34     name_mode="glob"
35 )
36 dispensed_dataset_id = dispensed_dataset["id"]
37 dataset = dxdata.load_dataset(id=dispensed_dataset_id)
38
39 dataset.entities
40
41 participant = dataset['participant']
42
43 main_cohort = dxdata.load_cohort("/cohort/pheno")
44
45 field_names = [
46     'eid', 'p31', 'p22001', 'p21022', 'p21003_i0', 'p738_i0', 'p22019', '
47     p22021', 'p53_i0', 'p40000_i0', 'p22018',
48     'p22011_a0', 'p22011_a1', 'p22011_a2', 'p22011_a3', 'p22011_a4', '
49     p22012_a0', 'p22012_a1', 'p22012_a2',
50     'p22012_a3', 'p22012_a4', 'p22013_a0', 'p22013_a1', 'p22013_a2', '
51     p22013_a3', 'p22013_a4', 'p22020',
52     'p21000_i0', 'p48_i0', 'p49_i0', 'p50_i0', 'p54_i0', 'p4079_i0_a0', '
53     p4080_i0_a0', 'p4080_i0_a1',
54     'p93_i0_a0', 'p93_i0_a1', 'p4079_i0_a1', 'p94_i0_a0', 'p94_i0_a1', '
55     p20117_i0', 'p20160_i0', 'p21001_i0',
56     'p21002_i0', 'p22000', 'p22007', 'p22008', 'p22003', 'p22027', 'p22004
57     ', 'p40007_i0', 'p26201_a0',
58     'p26201_a1', 'p26201_a2', 'p26201_a3', 'p22009_a1', 'p22009_a2', '
59     p22009_a3', 'p22009_a4', 'p22009_a5',
60     'p22009_a6', 'p22009_a7', 'p22009_a8', 'p22009_a9', 'p22009_a10', '
61     p22009_a11', 'p22009_a12', 'p22009_a13',
62     'p22009_a14', 'p22009_a15', 'p22009_a16', 'p22009_a17', 'p22009_a18',
63     'p22009_a19', 'p22009_a20',
64     'p22009_a21', 'p22009_a22', 'p22009_a23', 'p22009_a24', 'p22009_a25',
65     'p22009_a26', 'p22009_a27',
66     'p22009_a28', 'p22009_a29', 'p22009_a30', 'p22009_a31', 'p22009_a32',
67     'p22009_a33', 'p22009_a34',
68     'p22009_a35', 'p22009_a36', 'p22009_a37', 'p22009_a38', 'p22009_a39',
69     'p22009_a40', 'p20002_i0_a0',
70     'p20002_i0_a1', 'p20002_i0_a2', 'p20002_i0_a3', 'p20002_i0_a4', '
71     p20002_i0_a5', 'p20002_i0_a6',
72     'p20002_i0_a7', 'p20002_i0_a8', 'p20002_i0_a9', 'p20002_i0_a10', '
73     p20002_i0_a11', 'p20002_i0_a12',
74     'p20002_i0_a13', 'p20002_i0_a14', 'p20002_i0_a15', 'p20002_i0_a16', '
75     p20002_i0_a17', 'p20002_i0_a18',
76     'p20002_i0_a19', 'p20002_i0_a20', 'p20002_i0_a21', 'p20002_i0_a22', '
77     p20002_i0_a23', 'p20002_i0_a24',
78     'p20002_i0_a25', 'p20002_i0_a26', 'p20002_i0_a27', 'p20002_i0_a28', '
79     p20002_i0_a29', 'p20002_i0_a30',
80     'p20002_i0_a31', 'p20002_i0_a32', 'p20002_i0_a33', 'p120098', 'p120099
81     ', 'p120100', 'p120101', 'p120102',
82     'p120103', 'p120128', 'p29150', 'p29151', 'p29152', 'p29153', 'p29154'

```



```

        , 'p29155', 'p29206', 'fid',
65     'p6153_i0', 'p6153_i0_a1', 'p6153_i0_a2', 'p6153_i0_a3', 'p6177_i0', '
        p6177_i0_a1', 'p6177_i0_a2',
66     'p6138_i0', 'p34_i0_a0', 'p189_i0_a0', 'p30690_i0', 'p30691_i0', '
        p30760_i0', 'p30761_i0', 'p30780_i0',
67     'p30781_i0', 'p30870_i0', 'p30871_i0'
68 ]
69
70 df_main_cohort = participant.retrieve_fields(names=field_names, engine=
    dxdata.connect())
71
72 print("Initial columns:", df_main_cohort.columns)
73
74 # Rename columns for better readability
75 df_main_cohort = df_main_cohort.withColumnRenamed("eid", "IID")
76
77 # Add FID column -- required input format for regenie
78 print(type(df_main_cohort))
79
80 df_main_cohort = df_main_cohort.withColumn('FID', col('IID'))
81
82 df_main_cohort.show()
83
84 df_main_cohort_pandas = df_main_cohort.toPandas()
85
86 df_main_cohort_pandas.shape
87 df_main_cohort_pandas.p31.value_counts()
88
89 print("Initial columns:", df_main_cohort_pandas.columns)
90
91 # Get WES
92 path_to_family_file = f'/mnt/project/Bulk/Exome sequences/{exome_folder}/
    ukb{exome_field_id}_c1_b0_v1.fam'
93 plink_fam_df = pd.read_csv(path_to_family_file, delimiter='\s', dtype='
    object',
94                             names=['FID', 'IID', 'Father ID', 'Mother ID', '
    sex', 'Pheno'], engine='python')
95
96 # Intersect the phenotype file and the 450K WES .fam file
97 # to generate phenotype DataFrame for the 450K participants
98 main_wes_450k_df = df_main_cohort_pandas.join(plink_fam_df.set_index('IID'
    ), on='IID', rsuffix='_fam', how='inner')
99
100 # Drop unuseful columns from .fam file
101 main_wes_450k_df.drop(
102     columns=['FID_fam', 'Father ID', 'Mother ID', 'sex_fam', 'Pheno'], axis
    =1, inplace=True, errors='ignore'
103 )
104
105 print(type(main_wes_450k_df))
106

```

```

107 pheno_IDs = main_wes_450k_df[["IID", "FID"]]
108
109 print(pheno_IDs)
110
111 pheno_IDs_main = df_main_cohort_pandas[["IID", "FID"]]
112 print(pheno_IDs_main)
113
114 # Write phenotype files to a TSV file
115 main_wes_450k_df.to_csv('main_wes_450k.phe', sep='\t', na_rep='NA', index=
    False, quoting=3)
116 main_wes_450k_df.to_csv('main_wes_450k.csv', sep='\t', na_rep='NA', index=
    False, quoting=3, escapechar='\\')
117 df_main_cohort_pandas.to_csv('main_cohort.csv', sep='\t', na_rep='NA',
    index=False, quoting=3, escapechar='\\')
118 pheno_IDs.to_csv('pheno_id_450k.phe', sep='\t', na_rep='NA', index=False,
    quoting=3)
119 pheno_IDs.to_csv('pheno_id_450k.csv', sep='\t', na_rep='NA', index=False,
    quoting=3, escapechar='\\')
120
121 # Write phenotype files to a TSV file (for the main (full) cohort)
122 df_main_cohort_pandas.to_csv('main_cohort.csv', sep='\t', na_rep='NA',
    index=False, quoting=3, escapechar='\\')
123 pheno_IDs_main.to_csv('pheno_id_main.csv', sep='\t', na_rep='NA', index=
    False, quoting=3, escapechar='\\')
124
125 %bash -s "$output_dir"
126 dx upload main_wes_450k.phe -p --path $1 --brief
127
128 %bash -s "$output_dir"
129 dx upload pheno_id_450k.phe -p --path $1 --brief
130
131 %bash -s "$output_dir"
132 dx upload main_wes_450k.csv -p --path $1 --brief
133
134 %bash -s "$output_dir"
135 dx upload pheno_id_450k.csv -p --path $1 --brief
136
137 %bash -s "$output_dir"
138 dx upload main_cohort.csv -p --path $1 --brief
139
140 %bash -s "$output_dir"
141 dx upload pheno_id_main.csv -p --path $1 --brief

```

4.1.2 Renaming variables

The important phenotype variables were selected and stored in the **main_cohort.csv** dataset. The next step involved downloading the CSV file to the local machine. The downloaded data were saved to the **stata_sbp_input** file path. Once this was completed, the next step was to prepare the data for analysis. For participants prescribed antihypertensive medications, 15 mmHg and 10 mmHg were added to their baseline systolic blood pressure (SBP) and diastolic blood pressure (DBP) measurements, respectively[9].

```

import delimited "$stata_sbp_input\main_cohort.csv", clear
*File paths
global data_source "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evangelou
> \new_sbp_snps\stata\data_source"
global snps "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evangelou\new_sb
> p_snps\stata\snps"
global sbp_snps "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evangelou\ne
> w_sbp_snps\stata\sbp_snps"
global dbp_snps "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evangelou\ne
> w_sbp_snps\stata\dbp_snps"
global sbp_dbp_snps "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evangelou
> \new_sbp_snps\stata\sbp_dbp_snps"
global dx_data_sbp "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evangelou
> \new_sbp_snps\stata\dx_data\sbp\output"
global dx_data_dbp "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evangelou
> \new_sbp_snps\stata\dx_data\dbp\output"
global stata_sbp_input "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evang
> elou\new_sbp_snps\stata\stata_sbp_input"
global stata_sbp_output "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evan
> gelou\new_sbp_snps\stata\stata_sbp_output"
global stata_sbp_result "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evan
> gelou\new_sbp_snps\stata\stata_sbp_result"
global stata_dbp_input "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evang
> elou\new_sbp_snps\stata\stata_dbp_input"
global stata_dbp_output "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evan
> gelou\new_sbp_snps\stata\stata_dbp_output"
global stata_dbp_result "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evan
> gelou\new_sbp_snps\stata\stata_dbp_result"
global plot_png "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evangelou\ne
> w_sbp_snps\stata\stata_sbp_plot\png"
global plot_pdf "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evangelou\ne
> w_sbp_snps\stata\stata_sbp_plot\pdf"
global hesin_data "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evangelou\
> new_sbp_snps\stata\hesin_data"
global r_data "C:\Users\tabe0010\OneDrive - Monash University\MR_backup_file\Articles\Evangelou\new_
> sbp_snps\stata\r_codes\r_data"
*rename the variables

rename (iid p31 p22001 p21022 p21003_i0 p738_i0 p22019 p22021 p53_i0 p40000_i0 p22018 p22011_a0 p22
> 011_a1 p22011_a2 p22011_a3 p22011_a4 p22012_a0 p22012_a1 p22012_a2 p22012_a3 p22012_a4 p22013_a0 p
> 22013_a1 p22013_a2 p22013_a3 p22013_a4 p22020 p21000_i0 p48_i0 p49_i0 p50_i0 p54_i0 p4079_i0_a0 p4
> 080_i0_a0 p4080_i0_a1 p93_i0_a0 p93_i0_a1 p4079_i0_a1 p94_i0_a0 p94_i0_a1 p20117_i0 p20160_i0 p210
> 01_i0 p21002_i0 p22000 p22007 p22008 p22003 p22027 p22004 p40007_i0 p26201_a0 p26201_a1 p26201_a2
> p26201_a3 p22009_a1 p22009_a2 p22009_a3 p22009_a4 p22009_a5 p22009_a6 p22009_a7 p22009_a8 p22009_a
> 9 p22009_a10 p22009_a11 p22009_a12 p22009_a13 p22009_a14 p22009_a15 p22009_a16 p22009_a17 p22009_a
> 18 p22009_a19 p22009_a20 p22009_a21 p22009_a22 p22009_a23 p22009_a24 p22009_a25 p22009_a26 p22009_
> a27 p22009_a28 p22009_a29 p22009_a30 p22009_a31 p22009_a32 p22009_a33 p22009_a34 p22009_a35 p22009
> _a36 p22009_a37 p22009_a38 p22009_a39 p22009_a40 p20002_i0_a0 p20002_i0_a1 p20002_i0_a2 p20002_i0_
> a3 p20002_i0_a4 p20002_i0_a5 p20002_i0_a6 p20002_i0_a7 p20002_i0_a8 p20002_i0_a9 p20002_i0_a10 p20
> 002_i0_a11 p20002_i0_a12 p20002_i0_a13 p20002_i0_a14 p20002_i0_a15 p20002_i0_a16 p20002_i0_a17 p20
> 002_i0_a18 p20002_i0_a19 p20002_i0_a20 p20002_i0_a21 p20002_i0_a22 p20002_i0_a23 p20002_i0_a24 p20
> 002_i0_a25 p20002_i0_a26 p20002_i0_a27 p20002_i0_a28 p20002_i0_a29 p20002_i0_a30 p20002_i0_a31 p20
> 002_i0_a32 p20002_i0_a33 p120098 p120099 p120100 p120101 p120102 p120103 p120128 p29150 p29151 p29
> 152 p29153 p29154 p29155 p29206 fid p6153_i0 p6153_i0_a1 p6153_i0_a2 p6153_i0_a3 p6177_i0 p6177_i0
> _a1 p6177_i0_a2 p6138_i0 p34_i0_a0 p189_i0_a0 p30690_i0 p30691_i0 p30760_i0 p30761_i0 p30780_i0 p3
> 0781_i0 p30870_i0 p30871_i0) (iid n_31_0_0 n_22001_0_0 n_21022_0_0 n_21003_0_0 n_738_0_0 n_22019_0
> _0 n_22021_0_0 n_53_0_0 n_40000_0_0 n_22018_0_0 n_22011_0_0 n_22011_0_1 n_22011_0_2 n_22011_0_3 n_
> 22011_0_4 n_22012_0_0 n_22012_0_1 n_22012_0_2 n_22012_0_3 n_22012_0_4 n_22013_0_0 n_22013_0_1 n_22
> 013_0_2 n_22013_0_3 n_22013_0_4 n_22020_0_0 n_21000_0_0 n_48_0_0 n_49_0_0 n_50_0_0 n_54_0_0 n_4079
> _0_0 n_4080_0_0 n_4080_0_1 n_93_0_0 n_93_0_1 n_4079_0_1 n_94_0_0 n_94_0_1 n_20117_0_0 n_20160_0_0
> n_21001_0_0 n_21002_0_0 n_22000_0_0 n_22007_0_0 n_22008_0_0 n_22003_0_0 n_22027_0_0 n_22004_0_0 n_
> 40007_0_0 n_26201_0_0 n_26201_0_1 n_26201_0_2 n_26201_0_3 n_22009_0_1 n_22009_0_2 n_22009_0_3 n_22
> 009_0_4 n_22009_0_5 n_22009_0_6 n_22009_0_7 n_22009_0_8 n_22009_0_9 n_22009_0_10 n_22009_0_11 n_22
> 009_0_12 n_22009_0_13 n_22009_0_14 n_22009_0_15 n_22009_0_16 n_22009_0_17 n_22009_0_18 n_22009_0_1
> 9 n_22009_0_20 n_22009_0_21 n_22009_0_22 n_22009_0_23 n_22009_0_24 n_22009_0_25 n_22009_0_26 n_220
> 09_0_27 n_22009_0_28 n_22009_0_29 n_22009_0_30 n_22009_0_31 n_22009_0_32 n_22009_0_33 n_22009_0_34
> n_22009_0_35 n_22009_0_36 n_22009_0_37 n_22009_0_38 n_22009_0_39 n_22009_0_40 n_20002_0_0 n_20002

```

```

> _0_1 n_20002_0_2 n_20002_0_3 n_20002_0_4 n_20002_0_5 n_20002_0_6 n_20002_0_7 n_20002_0_8 n_20002_0
> _9 n_20002_0_10 n_20002_0_11 n_20002_0_12 n_20002_0_13 n_20002_0_14 n_20002_0_15 n_20002_0_16 n_20
> 002_0_17 n_20002_0_18 n_20002_0_19 n_20002_0_20 n_20002_0_21 n_20002_0_22 n_20002_0_23 n_20002_0_2
> 4 n_20002_0_25 n_20002_0_26 n_20002_0_27 n_20002_0_28 n_20002_0_29 n_20002_0_30 n_20002_0_31 n_200
> 02_0_32 n_20002_0_33 n_120098 n_120099 n_120100 n_120101 n_120102 n_120103 n_120128 n_29150 n_2915
> 1 n_29152 n_29153 n_29154 n_29155 n_29206 fid n_6153_0_0 n_6153_0_1 n_6153_0_2 n_6153_0_3 n_6177_0
> _0 n_6177_0_1 n_6177_0_2 n_6138_0_0 n_34_0_0 n_189_0_0 n_30690_0_0 n_30691_0_0 n_30760_0_0 n_30761
> _0_0 n_30780_0_0 n_30781_0_0 n_30870_0_0 n_30871_0_0)

*Let's destrung the variables
foreach j of varlist n_* {
  tostring `j', replace force
  replace `j' = "88" if `j' == "NA"
  destrung `j', replace
  capture replace `j'=. if `j' == 88
}

*work on date variables and rename them.
gen ts_53_0_0 = date(n_53_0_0, "YMD")
gen ts_40000_0_0 = date(n_40000_0_0, "YMD")
gen ts_120128 = cclock(n_120128, "YMDhms")
gen ts_29206 = date(n_29206, "YMD")
format ts_53_0_0 ts_40000_0_0 ts_29206 %td
format ts_120128 %tc
replace ts_120128 = dofc(ts_120128)
format ts_120128 %td

order ts_53_0_0, a(n_53_0_0)
order ts_40000_0_0, a(n_40000_0_0)
order ts_120128, a(n_120128)
order ts_29206, a(n_29206)

drop n_53_0_0 n_40000_0_0 n_120128 n_29206

*Work on systolic blood pressure measurements to take the mean of measurements (n_4080_* are automa
> tice measurements while n_93_* are manual measurements)
gen phe_sbp=.
replace phe_sbp = (n_4080_0_0 + n_4080_0_1)/2 if n_4080_0_0 !=. & n_4080_0_1 !=.
replace phe_sbp = ( n_93_0_0 + n_93_0_1)/2 if n_93_0_0 !=. & n_93_0_1 !=. & n_4080_0_0 ==. & n_4080
> _0_1 ==.
replace phe_sbp = (n_93_0_0 + n_4080_0_1)/2 if n_93_0_0 !=. & n_93_0_1 ==. & n_4080_0_0 ==. & n_408
> 0_0_1 !=.
replace phe_sbp = (n_93_0_1 + n_4080_0_0)/2 if n_93_0_0 ==. & n_93_0_1 !=. & n_4080_0_0 !=. & n_408
> 0_0_1 ==.

*the following measurements did not have any observations
*replace phe_sbp = (n_93_0_1 + n_4080_0_1)/2 if n_93_0_0 ==. & n_93_0_1 !=. & n_4080_0_0 ==. & n_40
> 80_0_1 !=.
*replace phe_sbp = (n_93_0_0 + n_4080_0_0)/2 if n_93_0_0 !=. & n_93_0_1 ==. & n_4080_0_0 !=. & n_40
> 80_0_1 ==.

order phe_sbp, a(n_93_0_1)

*Work on diastolic blood pressure measurements to take the mean of measurements (n_4079_* are auto
> matic measurements while n_94_* are manual measurements)
gen phe_dbp=.
replace phe_dbp = (n_4079_0_0 + n_4079_0_1)/2 if n_4079_0_0 !=. & n_4079_0_1 !=.
replace phe_dbp = ( n_94_0_0 + n_94_0_1)/2 if n_94_0_0 !=. & n_94_0_1 !=. & n_4079_0_0 ==. & n_4079
> _0_1 ==.
replace phe_dbp = (n_94_0_0 + n_4079_0_1)/2 if n_94_0_0 !=. & n_94_0_1 ==. & n_4079_0_0 ==. & n_407
> 9_0_1 !=.
replace phe_dbp = (n_94_0_1 + n_4079_0_0)/2 if n_94_0_0 ==. & n_94_0_1 !=. & n_4079_0_0 !=. & n_407
> 9_0_1 ==.

*the following measurements did not have any observations
*replace phe_dbp = (n_94_0_1 + n_4079_0_1)/2 if n_94_0_0 ==. & n_94_0_1 !=. & n_4079_0_0 ==. & n_40
> 79_0_1 !=.
*replace phe_dbp = (n_94_0_0 + n_4079_0_0)/2 if n_94_0_0 !=. & n_94_0_1 ==. & n_4079_0_0 !=. & n_40

```

```

> 79_0_1 ==.

order phe_4bp, a(n_94_0_1)

*Work on SBP (add 15 mmHg) for those prescribed antihypertensive medications
gen phe_dbp_adj = phe_dbp+15 if n_6153_0_0 == 2 | n_6153_0_1 == 2 | n_6153_0_2 == 2 | n_6153_0_3 ==
> 2 | n_6177_0_0 == 2 | n_6177_0_1 == 2 | n_6177_0_2 == 2

replace phe_dbp_adj = phe_dbp if phe_dbp_adj ==.
order(phe_dbp_adj), a(phe_dbp)

*Work on DBP (add 10 mmHg) for those prescribed antihypertensive medications
gen phe_dbp_adj = phe_dbp+10 if n_6153_0_0 == 2 | n_6153_0_1 == 2 | n_6153_0_2 == 2 | n_6153_0_3 ==
> 2 | n_6177_0_0 == 2 | n_6177_0_1 == 2 | n_6177_0_2 == 2

replace phe_dbp_adj = phe_dbp if phe_dbp_adj ==.
order(phe_dbp_adj), a(phe_dbp)
save "$stata_sbp_input\main_data.dta", replace

```

4.1.3 Exclusion criteria

The revised phenotype variables were stored in the **Phenotype.data** dataset. The subsequent step involved excluding participants based on a set of criteria. These exclusion criteria included:

- Sex mismatch between genetic sex and reported sex
- Sex chromosome aneuploidy
- Outliers for heterozygosity or missing rate
- Ethnicity: non-White British participants
- Participants without SBP values
- Related participants (based on kinship)

```
use "$stata_sbp_input\main_data.dta", clear
gen id_phe = iid
order fid id_phe, a(iid)
*Sex mismatch between genetic sex and reported
replace n_22001_0_0 = n_31_0_0 if n_22001_0_0 ==.
drop if n_31_0_0 != n_22001_0_0 // sex mismatch between genetic sex and reported sex 372 obs delete
> d
drop if n_22019_0_0 == 1 // Sex chromosome aneuploidy 470 obs deleted
drop if n_22027_0_0 == 1 // Outliers for heterozygosity or missing rate 963 obs deleted
drop if n_21000_0_0 != 1001 // Non white British 59584 obs deleted
drop if phe_sbp ==. // Participants without SBP values 1,229 obs deleted
keep if n_22021_0_0 == 0 // 148,621 observations deleted
save "$stata_sbp_output\part_1.dta",replace
keep id_phe ts_53_0_0
save "$stata_sbp_input\date_attending.dta", replace
keep id_phe
save "$stata_sbp_input\id_list.dta", replace
```

The dataset labeled **part_1.dta** contained the variables after applying the exclusion criteria. A separate dataset containing only the date of attending the assessment centre was stored in **date_attending.dta**.

4.1.4 Hospital admission data

The **hesin_diag** dataset was downloaded from the UKB RAP using the following Python code.

```
1
2 #Building cohorts using Spark JupyterLab
3 #Import important variables
4
5 import os
6
7 # Set environment variable before importing pyspark
8
9 os.environ['PYARROW_IGNORE_TIMEZONE'] = '1'
10
11 # Import necessary libraries
12 import pyspark
13 from pyspark import SparkConf, SparkContext
14 from pyspark.sql import SparkSession
```

```

15 from pyspark.sql.functions import col
16 import dxdy
17 import dxdata
18 import pandas as pd
19 import re
20
21 # Initialize SparkConf with the necessary configurations
22 conf = SparkConf() \
23     .setAppName("HESIN Data Analysis") \
24     .set("spark.kryoserializer.buffer.max", "1g")
25
26 # Initialize Spark
27 # Spark initialization (Done only once; do not rerun this cell unless you
    select Kernel -> Restart kernel).
28
29 sc = pyspark.SparkContext()
30 spark = pyspark.sql.SparkSession(sc)
31
32
33 print(f"Kryo serializer buffer max size set to: {conf.get('spark.
    kryoserializer.buffer.max')}")
34
35
36 # Automatically discover dispensed dataset ID and load the dataset
37 dispensed_dataset = dxdy.find_one_data_object(
38     typename="Dataset",
39     name="app*.dataset",
40     folder="/",
41     name_mode="glob")
42 dispensed_dataset_id = dispensed_dataset["id"]
43 dataset = dxdata.load_dataset(id=dispensed_dataset_id)
44
45 dataset.entities
46
47 participant = dataset['hesin_diag']
48
49 print(type(participant))
50
51 help(participant)
52
53 field_names = ['eid', 'ins_index', 'arr_index', 'level', 'diag_icd9', '
    diag_icd9_nb', 'diag_icd10', 'diag_icd10_nb']
54
55 df_hesin_diag = participant.retrieve_fields(names=field_names, engine=
    dxdata.connect())
56
57 print("Initial columns:", df_hesin_diag.columns)
58
59 print(type(df_hesin_diag))
60
61 df_hesin_diag.show(5)

```

```

62
63 df_hesin_diag.count()
64
65 df_hesin_diag = df_hesin_diag.repartition(10)
66
67 df_hesin_diag_main=df_hesin_diag.toPandas()
68
69 print(type(df_hesin_diag_main))
70
71 print(df_hesin_diag_main)
72
73 df_hesin_diag_main.to_csv('hesin_diag_main.csv', index=False)
74
75 %%bash
76 dx upload hesin_diag_main.csv --dest project-Gkz56gjJx5g1zB269F0ybP63:/
    Data/

```

The **hesin_main** dataset, which included the date of diagnosis for participants admitted to hospitals, was also downloaded. The following Python code was used to create the necessary file from the UKB RAP.

```

1  #Building cohorts using Spark JupyterLab
2  #Import important variables
3
4  import os
5
6  # Set environment variable before importing pyspark
7
8  os.environ['PYARROW_IGNORE_TIMEZONE'] = '1'
9
10 # Import necessary libraries
11 import pyspark
12 from pyspark import SparkConf, SparkContext
13 from pyspark.sql import SparkSession
14 from pyspark.sql.functions import col
15 import dxdpy
16 import dxdata
17 import pandas as pd
18 import re
19
20 # Initialize SparkConf with the necessary configurations
21 conf = SparkConf() \
22     .setAppName("HESIN Data Analysis") \
23     .set("spark.kryoserializer.buffer.max", "1g")
24
25 # Initialize Spark
26 # Spark initialization (Done only once; do not rerun this cell unless you
    select Kernel -> Restart kernel).
27
28 sc = pyspark.SparkContext()
29 spark = pyspark.sql.SparkSession(sc)
30

```



```

31
32 print(f"Kryo serializer buffer max size set to: {conf.get('spark.
    kryoserIALIZER.buffer.max')}")
33
34
35 # Automatically discover dispensed dataset ID and load the dataset
36 dispensed_dataset = dxpy.find_one_data_object(
37     typename="Dataset",
38     name="app*.dataset",
39     folder="/",
40     name_mode="glob")
41 dispensed_dataset_id = dispensed_dataset["id"]
42 dataset = dxdata.load_dataset(id=dispensed_dataset_id)
43
44
45 dataset.entities
46
47 participant = dataset['hesin']
48
49 print(type(participant))
50
51 help(participant)
52
53 field_names = ['eid', 'ins_index', 'dsource', 'epistart', 'epiend', '
    epidur', 'admidate', 'disdate']
54
55
56 df_hesin = participant.retrieve_fields(names=field_names, engine=dxdata.
    connect())
57
58 print("Initial columns:", df_hesin.columns)
59
60
61 print(type(df_hesin))
62
63
64 df_hesin.show(5)
65
66 df_hesin.count()
67
68 df_hesin = df_hesin.repartition(10)
69
70 df_hesin_main=df_hesin.toPandas()
71
72 print(type(df_hesin_main))
73
74 print(df_hesin_main)
75
76 df_hesin_main.to_csv('hesin_main.csv', index=False)
77
78 %%bash

```

```
79 dx upload hesin_main.csv --dest project-Gkz56gjJx5g1zB269F0ybP63:/Data/
```

Save the **hesin_diag_main.csv** and **hesin_main.csv** files to the **hesin_data** file path on the local machine.

The next step involved working with the ICD codes. The **hesin_diag_main.csv** file contained both ICD-9 and ICD-10 diagnosis codes. The analysis first focused on the ICD-10 codes before proceeding to the ICD-9 codes. In addition to diagnosis codes, the dataset also included variables labeled **instance** (**ins_index**), **array** (**arr_index**), and **level**.

Instance indicates how many occasions participants have measurements performed. There are three categories:

- **Singular:** only one instance can be present, for example sex or year-of-birth
- **Defined:** more than one instance may be present, and each instance represents a fixed identifiable set of results across all participants
- **Variable:** more than one instance may be present, however there is no correspondence between (say) the 3rd instance for one participant and the 3rd instance for another

Array describes whether there are multiple data items for a given participant instance. There are two categories:

- **Single:** only one data item is present for each participant, for instance the answer to "What is your favourite colour of the rainbow?"
- **Multiple:** more than one data item may be present for each participant, for instance the answer to "Which colours of the rainbow do you like?"

level describes whether the diagnosis is primary (1) or secondary (2)

4.1.4.1 Working on ICD-10 diagnosis codes

The next line of stata codes will do the following tasks:

- filter ICD-10 codes
- merge with the **hesin_main.dta** dataset (we need to convert the csv file to dta for convenience)
- identify admission or episode start date
- identify discharge or episode end date (if required)
- merge with **data_attending.dta** file which contains participants' date of attending the UK Biobank assessment centre
- create 240 dummy variables for the conditions we would like to work on the next step.

```
*Let's save the hesin_main.csv in stata file format (hesin_main.dta)

import delimited "$hesin_data\hesin_main.csv", clear
save "$hesin_data\hesin_main.dta", replace
```

```

* Let's work on the ICD 10 codes
import delimited "$hesin_data\hesin_diag_main.csv", clear

drop diag_icd9_nb diag_icd10_nb diag_icd9
egen idlong = concat(eid ins_index arr_index level), format(%20.0g) p(",")
drop eid ins_index level
reshape wide diag_icd10, i(idlong) j(arr_index)
split idlong, p(,)
rename idlong1 eid
rename idlong2 ins_index
rename idlong3 arr_index
rename idlong4 level
destring eid, replace
destring ins_index, replace
destring arr_index, replace
destring level, replace
drop idlong
order eid ins_index arr_index level diag*
merge m:1 eid ins_index using "$hesin_data\hesin_main.dta"
drop if _merge == 2
drop _merge*

drop if diag_icd100 == "" & diag_icd101 == "" & diag_icd102 == "" & diag_icd103 == "" & diag_icd104
> == "" & diag_icd105 == "" & diag_icd106 == "" & diag_icd107 == "" & diag_icd108 == "" & diag_icd1
> 09 == "" & diag_icd1010 == "" & diag_icd1011 == "" & diag_icd1012 == "" & diag_icd1013 == "" & dia
> g_icd1014 == "" & diag_icd1015 == "" & diag_icd1016 == "" & diag_icd1017 == "" & diag_icd1018 == "
> " & diag_icd1019 == "" & diag_icd1020 == ""

gen epistart_1 = date(epistart, "YMD")
gen epiend_1 = date(epiend, "YMD")
gen admdate_1 = date(admdate, "YMD")
gen disdate_1 = date(disdate, "YMD")

format epistart_1 epiend_1 admdate_1 disdate_1 %td
drop epistart epiend admdate disdate
rename epistart_1 epistart
rename epiend_1 epiend
rename admdate_1 admdate
rename disdate_1 disdate
replace epistart = admdate if epistart ==.
replace epiend = disdate if epiend ==.
drop admdate disdate
rename eid id_phe
merge m:1 id_phe using "$stata_sbp_input\date_attending.dta", keep(3) nogen
order ts_53_0_0, a(epistart)
rename ts_53_0_0 date_baseline
rename epistart date_episode

*Generate health condition phenotypes
forvalues i = 1/240 {
    gen v`i' = .
    format v`i' %td
}

save "$hesin_data\hesin_icd10_dates.dta", replace

```

The **hesin_icd10_dates.dta** dataset was large, and processing it in full required several hours to complete the analysis. To improve efficiency, the dataset was divided into smaller subsets with reduced sample sizes. Within each subset, the date of diagnosis was identified for 240 predefined comorbidities.

```

local j = 250000
local k = 1
forval i = 1/39 {
    di `i' " " `k' " " `j'

```

```

        use "$hesin_data\hesin_icd10_dates.dta", clear
        keep if _n >= `k' & _n <= `j'
        save "$hesin_data\hesin_icd10_dates_`i'.dta", replace

        local k = `j'+1
        local j = `j'+250000

    }

*Loop through HES diagnosis fields to find relevant ICD10 codes

forval j = 1/39 {
    use "$hesin_data\hesin_icd10_dates_`j'.dta",clear
    forval i = 0/20 {

        dis "Looking through `i' of 20 diagnosis fields for ICD10 codes of data `j' "

        qui replace v1 = date_episode if strpos(diag_icd10`i',"B20") > 0
        qui replace v2 = date_episode if strpos(diag_icd10`i',"D86") > 0
        qui replace v3 = date_episode if strpos(diag_icd10`i',"C14") > 0
        qui replace v4 = date_episode if strpos(diag_icd10`i',"C16") > 0
        qui replace v5 = date_episode if strpos(diag_icd10`i',"C18") > 0
        qui replace v6 = date_episode if strpos(diag_icd10`i',"C22") > 0
        qui replace v7 = date_episode if strpos(diag_icd10`i',"C33") > 0
        qui replace v7 = date_episode if strpos(diag_icd10`i',"C34") > 0
        qui replace v8 = date_episode if strpos(diag_icd10`i',"C43") > 0
        qui replace v9 = date_episode if strpos(diag_icd10`i',"C44") > 0
        qui replace v10 = date_episode if strpos(diag_icd10`i',"C50") > 0
        qui replace v11 = date_episode if strpos(diag_icd10`i',"C55") > 0
        qui replace v12 = date_episode if strpos(diag_icd10`i',"C53") > 0
        qui replace v13 = date_episode if strpos(diag_icd10`i',"C61") > 0
        qui replace v14 = date_episode if strpos(diag_icd10`i',"C67") > 0
        qui replace v15 = date_episode if strpos(diag_icd10`i',"C64") > 0
        qui replace v15 = date_episode if strpos(diag_icd10`i',"C65") > 0
        qui replace v15 = date_episode if strpos(diag_icd10`i',"C66") > 0
        qui replace v15 = date_episode if strpos(diag_icd10`i',"C68") > 0
        qui replace v16 = date_episode if strpos(diag_icd10`i',"C71") > 0
        qui replace v17 = date_episode if strpos(diag_icd10`i',"C76") > 0
        qui replace v18 = date_episode if strpos(diag_icd10`i',"C80") > 0
        qui replace v19 = date_episode if strpos(diag_icd10`i',"C81") > 0
        qui replace v20 = date_episode if strpos(diag_icd10`i',"C82") > 0
        qui replace v20 = date_episode if strpos(diag_icd10`i',"C83") > 0
        qui replace v20 = date_episode if strpos(diag_icd10`i',"C84") > 0
        qui replace v20 = date_episode if strpos(diag_icd10`i',"C85") > 0
        qui replace v20 = date_episode if strpos(diag_icd10`i',"C96") > 0
        qui replace v21 = date_episode if strpos(diag_icd10`i',"C95") > 0
        qui replace v22 = date_episode if strpos(diag_icd10`i',"D12") > 0
        qui replace v22 = date_episode if strpos(diag_icd10`i',"D13") > 0
        qui replace v22 = date_episode if strpos(diag_icd10`i',"D20") > 0
        qui replace v23 = date_episode if strpos(diag_icd10`i',"D21") > 0
        qui replace v24 = date_episode if strpos(diag_icd10`i',"D22") > 0
        qui replace v24 = date_episode if strpos(diag_icd10`i',"D23") > 0
        qui replace v25 = date_episode if strpos(diag_icd10`i',"D24") > 0
        qui replace v26 = date_episode if strpos(diag_icd10`i',"D25") > 0
        qui replace v27 = date_episode if strpos(diag_icd10`i',"D36") > 0
        qui replace v28 = date_episode if strpos(diag_icd10`i',"D45") > 0
        qui replace v28 = date_episode if strpos(diag_icd10`i',"D46") > 0
        qui replace v28 = date_episode if strpos(diag_icd10`i',"D47") > 0
        qui replace v28 = date_episode if strpos(diag_icd10`i',"D48") > 0
        qui replace v29 = date_episode if strpos(diag_icd10`i',"D49") > 0
        qui replace v30 = date_episode if strpos(diag_icd10`i',"E01") > 0
        qui replace v30 = date_episode if strpos(diag_icd10`i',"E04") > 0
        qui replace v31 = date_episode if strpos(diag_icd10`i',"E05") > 0
        qui replace v32 = date_episode if strpos(diag_icd10`i',"E02") > 0
        qui replace v32 = date_episode if strpos(diag_icd10`i',"E03") > 0
        qui replace v33 = date_episode if strpos(diag_icd10`i',"E06") > 0
        qui replace v34 = date_episode if strpos(diag_icd10`i',"E07") > 0
    }
}

```

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]


```

qui replace v235 = date_episode if strpos(diag_icd10`i`, "Z85") > 0
qui replace v236 = date_episode if strpos(diag_icd10`i`, "Z720") > 0 // changed from "Z72.0"
qui replace v236 = date_episode if strpos(diag_icd10`i`, "Z91") > 0
qui replace v237 = date_episode if strpos(diag_icd10`i`, "Z865") > 0 // changed from "Z86.5"
qui replace v238 = date_episode if strpos(diag_icd10`i`, "Z94") > 0
qui replace v239 = date_episode if strpos(diag_icd10`i`, "Z95") > 0
qui replace v239 = date_episode if strpos(diag_icd10`i`, "Z96") > 0
qui replace v239 = date_episode if strpos(diag_icd10`i`, "Z97") > 0
qui replace v240 = date_episode if strpos(diag_icd10`i`, "Z49") > 0

*Label variables
label variable v1 "Human immunodeficiency virus (hiv) disease"
label variable v2 "Sarcoidosis"
label variable v3 "Malignant neoplasm of other and ill-defined sites within the lip oral cavity and
> pharynx"
label variable v4 "Malignant neoplasm of stomach"
label variable v5 "Malignant neoplasm of colon"
label variable v6 "Malignant neoplasm of liver"
label variable v7 "Malignant neoplasm of trachea bronchus and lung"
label variable v8 "Malignant melanoma of skin"
label variable v9 "Other malignant neoplasm of skin"
label variable v10 "Malignant neoplasm of female breast"
label variable v11 "Malignant neoplasm of uterus-part unspecified"
label variable v12 "Malignant neoplasm of cervix uteri"
label variable v13 "Malignant neoplasm of prostate"
label variable v14 "Malignant neoplasm of bladder"
label variable v15 "Malignant neoplasm of kidney and other and unspecified urinary organs"
label variable v16 "Malignant neoplasm of brain"
label variable v17 "Malignant neoplasm of other and ill-defined sites"
label variable v18 "Malignant neoplasm without specification of site"
label variable v19 "Hodgkin's disease"
label variable v20 "Other malignant neoplasms of lymphoid and histiocytic tissue"
label variable v21 "Leukemia of unspecified cell type"
label variable v22 "Benign neoplasm of other parts of digestive system"
label variable v23 "Other benign neoplasm of connective and other soft tissue"
label variable v24 "Benign neoplasm of skin"
label variable v25 "Benign neoplasm of breast"
label variable v26 "Uterine leiomyoma"
label variable v27 "Benign neoplasm of other and unspecified sites"
label variable v28 "Neoplasm of uncertain behavior of other and unspecified sites and tissues"
label variable v29 "Neoplasms Of Unspecified Nature"
label variable v30 "Simple and unspecified goiter"
label variable v31 "Thyrototoxicosis with or without goiter"
label variable v32 "Acquired hypothyroidism"
label variable v33 "Thyroiditis"
label variable v34 "Other disorders of thyroid"
label variable v35 "Diabetes mellitus"
label variable v36 "Other disorders of pancreatic internal secretion"
label variable v37 "Disorders of parathyroid gland"
label variable v38 "Disorders of the pituitary gland and its hypothalamic control"
label variable v39 "Disorders of adrenal glands"
label variable v40 "Ovarian dysfunction"
label variable v41 "Other endocrine disorders"
label variable v42 "Disorders of carbohydrate transport and metabolism"
label variable v43 "Disorders of lipid metabolism"
label variable v44 "Gout"
label variable v45 "Disorders of mineral metabolism"
label variable v46 "Other and unspecified disorders of metabolism"
label variable v47 "Overweight, obesity and other hyperalimentation"
label variable v48 "Disorders involving the immune mechanism"
label variable v49 "Iron deficiency anemias"
label variable v50 "Hereditary hemolytic anemias"
label variable v51 "Other and unspecified anemias"
label variable v52 "Coagulation defects"
label variable v53 "Purpura and other hemorrhagic conditions"
label variable v54 "Diseases of white blood cells"
label variable v55 "Other diseases of blood and blood-forming organs"

```

label variable v56 "Persistent mental disorders due to conditions classified elsewhere"
 label variable v57 "Schizophrenic disorders"
 label variable v58 "Episodic mood disorders"
 label variable v59 "Delusional disorders"
 label variable v60 "Other nonorganic psychoses"
 label variable v61 "Anxiety, dissociative and somatoform disorders"
 label variable v62 "Personality disorders"
 label variable v63 "Sexual and gender identity disorders"
 label variable v64 "Alcohol dependence syndrome"
 label variable v65 "Drug dependence"
 label variable v66 "Nondependent abuse of drugs"
 label variable v67 "Physiological malfunction arising from mental factors"
 label variable v68 "Special symptoms or syndromes not elsewhere classified"
 label variable v69 "Adjustment reaction"
 label variable v70 "Specific nonpsychotic mental disorders due to brain damage"
 label variable v71 "Depressive disorder not elsewhere classified"
 label variable v72 "Disturbance of conduct not elsewhere classified"
 label variable v73 "Hyperkinetic syndrome of childhood"
 label variable v74 "Specific delays in development"
 label variable v75 "Unspecified mental retardation"
 label variable v76 "Other cerebral degenerations"
 label variable v77 "Parkinson's disease"
 label variable v78 "Other extrapyramidal disease and abnormal movement disorders"
 label variable v79 "Other diseases of spinal cord"
 label variable v80 "Disorders of the autonomic nervous system"
 label variable v81 "Multiple sclerosis"
 label variable v82 "Infantile cerebral palsy"
 label variable v83 "Other paralytic syndromes"
 label variable v84 "Epilepsy"
 label variable v85 "Migraine"
 label variable v86 "Cataplexy and narcolepsy"
 label variable v87 "Other conditions of brain"
 label variable v88 "Other and unspecified disorders of the nervous system"
 label variable v89 "Facial nerve disorders"
 label variable v90 "Nerve root and plexus disorders"
 label variable v91 "Mononeuritis of upper limb and mononeuritis multiplex"
 label variable v92 "Mononeuritis of lower limb and unspecified site"
 label variable v93 "Hereditary and idiopathic peripheral neuropathy"
 label variable v94 "Muscular dystrophies and other myopathies"
 label variable v95 "Disorders of the globe"
 label variable v96 "Retinal detachments and defects"
 label variable v97 "Other retinal disorders"
 label variable v98 "Glaucoma"
 label variable v99 "Cataract"
 label variable v100 "Disorders of refraction and accommodation"
 label variable v101 "Visual disturbances"
 label variable v102 "Blindness and low vision"
 label variable v103 "Keratitis"
 label variable v104 "Disorders of conjunctiva"
 label variable v105 "Other disorders of eyelids"
 label variable v106 "Disorders of lacrimal system"
 label variable v107 "Disorders of optic nerve and visual pathways"
 label variable v108 "Other disorders of eye"
 label variable v109 "Other disorders of tympanic membrane"
 label variable v110 "Hearing loss"
 label variable v111 "Essential hypertension"
 label variable v112 "Acute myocardial infarction"
 label variable v113 "Other acute and subacute forms of ischemic heart disease"
 label variable v114 "Old myocardial infarction"
 label variable v115 "Angina pectoris"
 label variable v116 "Other forms of chronic ischemic heart disease"
 label variable v117 "Chronic pulmonary heart disease"
 label variable v118 "Other diseases of endocardium"
 label variable v119 "Cardiomyopathy"
 label variable v120 "Conduction disorders"
 label variable v121 "Cardiac dysrhythmias"
 label variable v122 "Heart failure"
 label variable v123 "Ill-defined descriptions and complications of heart disease"

label variable v124 "Occlusion and stenosis of precerebral arteries"
 label variable v125 "Transient cerebral ischemia"
 label variable v126 "Acute but ill-defined cerebrovascular disease"
 label variable v127 "Other and ill-defined cerebrovascular disease"
 label variable v128 "Late effects of cerebrovascular disease"
 label variable v129 "Atherosclerosis"
 label variable v130 "Aortic aneurysm and dissection"
 label variable v131 "Other aneurysm"
 label variable v132 "Other peripheral vascular disease"
 label variable v133 "Arterial embolism and thrombosis"
 label variable v134 "Polyarteritis nodosa and allied conditions"
 label variable v135 "Other disorders of arteries and arterioles"
 label variable v136 "Phlebitis and thrombophlebitis"
 label variable v137 "Other venous embolism and thrombosis"
 label variable v138 "Varicose veins of lower extremities"
 label variable v139 "Hemorrhoids"
 label variable v140 "Noninfectious disorders of lymphatic channels"
 label variable v141 "Other disorders of circulatory system"
 label variable v142 "Chronic sinusitis"
 label variable v143 "Chronic disease of tonsils and adenoids"
 label variable v144 "Allergic rhinitis"
 label variable v145 "Chronic bronchitis"
 label variable v146 "Emphysema"
 label variable v147 "Asthma"
 label variable v148 "Chronic airway obstruction not elsewhere classified"
 label variable v149 "Coal workers' pneumoconiosis"
 label variable v150 "Other diseases of lung"
 label variable v151 "Other diseases of respiratory system"
 label variable v152 "Diseases of hard tissues of teeth"
 label variable v153 "Dentofacial anomalies including malocclusion"
 label variable v154 "Other diseases and conditions of the teeth and supporting structures"
 label variable v155 "Diseases of esophagus"
 label variable v156 "Gastric ulcer"
 label variable v157 "Gastritis and duodenitis"
 label variable v158 "Disorders of function of stomach"
 label variable v159 "Other disorders of stomach and duodenum"
 label variable v160 "Inguinal hernia"
 label variable v161 "Other hernia of abdominal cavity without mention of obstruction or gangrene"
 label variable v162 "Regional enteritis"
 label variable v163 "Ulcerative enterocolitis"
 label variable v164 "Other and unspecified noninfectious gastroenteritis and colitis"
 label variable v165 "Diverticula of intestine"
 label variable v166 "Functional digestive disorders not elsewhere classified"
 label variable v167 "Other disorders of intestine"
 label variable v168 "Chronic liver disease and cirrhosis"
 label variable v169 "Other disorders of liver"
 label variable v170 "Cholelithiasis"
 label variable v171 "Other disorders of gallbladder"
 label variable v172 "Other disorders of biliary tract"
 label variable v173 "Diseases of pancreas"
 label variable v174 "Intestinal malabsorption"
 label variable v175 "Nephritis and nephropathy not specified as acute or chronic"
 label variable v176 "Renal failure unspecified"
 label variable v177 "Calculus of kidney and ureter"
 label variable v178 "Other disorders of kidney and ureter"
 label variable v179 "Other disorders of bladder"
 label variable v180 "Other disorders of urethra and urinary tract"
 label variable v181 "Hyperplasia of prostate"
 label variable v182 "Other disorders of prostate"
 label variable v183 "Disorders of penis"
 label variable v184 "Benign mammary dysplasias"
 label variable v185 "Other disorders of breast"
 label variable v186 "Endometriosis"
 label variable v187 "Genital prolapse"
 label variable v188 "Noninflammatory disorders of ovary fallopian tube and broad ligament"
 label variable v189 "Disorders of uterus not elsewhere classified"
 label variable v190 "Noninflammatory disorders of cervix"
 label variable v191 "Noninflammatory disorders of vagina"


```

label variable v192 "Disorders of menstruation and other abnormal bleeding from female genital tract
> "
label variable v193 "Menopausal and postmenopausal disorders"
label variable v194 "Female infertility"
label variable v195 "Psoriasis and similar disorders"
label variable v196 "Pruritus and related conditions"
label variable v197 "Corns and callosities"
label variable v198 "Other hypertrophic and atrophic conditions of skin"
label variable v199 "Diseases of nail"
label variable v200 "Diseases of hair and hair follicles"
label variable v201 "Chronic ulcer of skin"
label variable v202 "Diffuse diseases of connective tissue"
label variable v203 "Rheumatoid arthritis and other inflammatory polyarthropathies"
label variable v204 "Osteoarthritis and allied disorders"
label variable v205 "Other and unspecified arthropathies"
label variable v206 "Internal derangement of knee"
label variable v207 "Other derangement of joint"
label variable v208 "Other and unspecified disorders of joint"
label variable v209 "Ankylosing spondylitis and other inflammatory spondylopathies"
label variable v210 "Spondylosis and allied disorders"
label variable v211 "Intervertebral disc disorders"
label variable v212 "Other disorders of cervical region"
label variable v213 "Other and unspecified disorders of back"
label variable v214 "Polymyalgia rheumatica"
label variable v215 "Peripheral enthesopathies and allied syndromes"
label variable v216 "Osteomyelitis periostitis and other infections involving bone"
label variable v217 "Osteitis deformans and osteopathies associated with other disorders classified
> elsewhere"
label variable v218 "Other disorders of bone and cartilage"
label variable v219 "Acquired deformities of toe"
label variable v220 "Curvature of spine"
label variable v221 "Spina bifida"
label variable v222 "Other congenital anomalies of heart"
label variable v223 "Other congenital anomalies of circulatory system"
label variable v224 "Congenital anomalies of urinary system"
label variable v225 "Certain congenital musculoskeletal deformities"
label variable v226 "Other congenital musculoskeletal anomalies"
label variable v227 "Chromosomal anomalies"
label variable v228 "Other and unspecified congenital anomalies"
label variable v229 "Senility without psychosis"
label variable v230 "Intracranial injury of other and unspecified nature"
label variable v231 "Late effects of other and unspecified injuries"
label variable v232 "Crushing injury of lower limb"
label variable v233 "Spinal cord injury without evidence of spinal bone injury"
label variable v234 "Asymptomatic human immunodeficiency virus (hiv) infection status"
label variable v235 "Personal history of malignant neoplasm"
label variable v236 "Other personal history presenting hazards to health"
label variable v237 "Mental and behavioral problems"
label variable v238 "Organ or tissue replaced by transplant"
label variable v239 "Organ or tissue replaced by other means"
label variable v240 "Encounter for dialysis and dialysis catheter care"

*drop diag* date_episode
}

save "$hesin_data\hesin_icd10_complete_`j`.dta", replace

}

*Append the icd10_complete data
use "$hesin_data\hesin_icd10_complete_1.dta", clear
forval j=2/39{
    append using "$hesin_data\hesin_icd10_complete_`j`.dta"
}

save "$hesin_data\hesin_icd10_complete.dta",replace

```

The necessary **date_episodes** were identified for each of the 240 comorbidities and saved as **hesin_icd10_complete.dta**.

4.1.4.2 Working on ICD-9 diagnosis codes

Next, we followed the same approach to work on the ICD-9 codes as we did for the ICD-10 codes.

```
import delimited "$hesin_data\hesin_diag_main.csv", clear
drop diag_icd9_nb diag_icd10_nb diag_icd10
egen idlong = concat(eid ins_index arr_index level), format(%20.0g) p(",")
drop eid ins_index level
reshape wide diag_icd9, i(idlong) j(arr_index)
split idlong, p(,)
rename idlong1 eid
rename idlong2 ins_index
rename idlong3 arr_index
rename idlong4 level
destring eid, replace
destring ins_index, replace
destring arr_index, replace
destring level, replace
drop idlong
order eid ins_index arr_index level diag*
merge m:1 eid ins_index using "$hesin_data\hesin_main.dta"
drop if _merge == 2
drop _merge*
drop if diag_icd90 == "" & diag_icd91 == "" & diag_icd92 == "" & diag_icd93 == "" & diag_icd94 == ""
> & diag_icd95 == "" & diag_icd96 == "" & diag_icd97 == "" & diag_icd98 == "" & diag_icd99 == "" &
> diag_icd910 == "" & diag_icd911 == "" & diag_icd912 == "" & diag_icd913 == "" & diag_icd914 == ""
> & diag_icd915 == "" & diag_icd916 == "" & diag_icd917 == "" & diag_icd918 == "" & diag_icd919 == ""
> " & diag_icd920 == ""

gen epistart_1 = date(epistart, "YMD")
gen epiend_1 = date(epiend, "YMD")
gen admidate_1 = date(admidate, "YMD")
gen disdate_1 = date(disdate, "YMD")
format epistart_1 epiend_1 admidate_1 disdate_1 %td
drop epistart epiend admidate disdate
rename epistart_1 epistart
rename epiend_1 epiend
rename admidate_1 admidate
rename disdate_1 disdate
replace epistart = admidate if epistart ==.
replace epiend = disdate if epiend ==.
drop admidate disdate
rename eid id_phe
merge m:1 id_phe using "$stata_sbp_input\date_attending.dta", keep(3) nogen
order ts_53_0_0, a(epistart)
rename ts_53_0_0 date_baseline
rename epistart date_episode
*Generate health condition phenotypes
forvalues i = 1/240 {
    gen v`i' = .
    format v`i' %td
}
save "$hesin_data\hesin_icd9_dates.dta", replace
use "$hesin_data\hesin_icd9_dates.dta", clear
forvalues i = 0/20 {
    dis "Looking through `i' of 20 diagnosis fields for ICD9 codes"
    qui replace v1 = date_episode if substr(diag_icd9`i',1,3) == "042"
    qui replace v2 = date_episode if substr(diag_icd9`i',1,3) == "135"
    qui replace v3 = date_episode if substr(diag_icd9`i',1,3) == "149"
```

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]


```

        append using "$hesin_data\hesin_icd9_complete.dta"

        *Replace all v`i` values with the minimum, so when duplicates are dropped, any instance of t
        > he ICD code is kept
        forvalues i = 1/240 {
            dis "Sorting through `i` of 240 variables"
            qui {
                bysort id_phe: egen x = min(v`i`)
            }
            qui replace v`i` = x
            qui drop x
        }

        *Keep only a single row per participant
        duplicates drop id_phe, force

        save "$hesin_data\hesin_icd_complete.dta",replace

```

4.1.5 Setting up the data for HRQoL prediction

We compiled the phenotype data, including both ICD-9 and ICD-10 diagnosis dates. The next step was to prepare the data for HRQoL prediction based on the types of comorbidities that a participant might have developed.

The next stata codes will execute:

- Merge the **part_1.dta** dataset with **hesin_icd_compelet.dta** dataset
- Generate variables labeled "z" to make sure if a study participant have developed a condition before the utility day.
- Generate variables labeled "ncc" to count the number of comorbidities a particiapant may have
- Rename some variables
- Reclassify the "qualification" variable for the purpose of analysis

```

use "$stata_sbp_output\part_1.dta", clear
merge 1:1 id_phe using "$hesin_data\hesin_icd_complete.dta", nogen keep(1 3)

*Generate dummy variables for each condition
*These will be 1 if the participant had the condition before the utility day
forvalues i = 1/240 {
    gen z`i` = 0
}

*Generate dummy variables for number of comorbidities
forvalues i = 2/10 {
    gen ncc`i` = 0
}

*Rename some variables
rename n_31_0_0 sex
rename n_54_0_0 centre
rename n_189_0_0 eco_tdi
rename n_738_0_0 eco_household_income
rename n_6138_0_0 eco_qualifications
rename n_20160_0_0 ever_smoked
rename n_21001_0_0 phe_bmi
rename n_21003_0_0 age
forvalues i = 1/40 {
    rename n_22009_0_`i` pc`i`
}

replace date_baseline = ts_53_0_0 if date_baseline == .

```

```

drop ts_53_0_0
rename ts_40000_0_0 date_death
*Make qualifications into dummy variables
forvalues i = 1/3 {
    gen qual_`i' = 0
}
*Assume everyone who preferred not to answer and missing responses have no qualifications
*High school = A levels, O levels, CSEs & GCSEs
replace qual_1 = 1 if eco_qualifications == 2 | eco_qualifications == 3 | eco_qualifications == 4
*Other degree = NVQs and other
replace qual_2 = 1 if eco_qualifications == 5 | eco_qualifications == 6
*Bachelor degree = College or university degree
replace qual_3 = 1 if eco_qualifications == 1
*Not coding ma_phd
*Remove all participants who died before entry
qui drop if date_death < date_baseline
save "$stata_sbp_output\part_2a.dta", replace

```

Dividing the part_2a.dta Dataset into Small Blocks

The **part_2a.dta** dataset was large, making utility prediction slow. To expedite the analysis, the following steps were taken:

First, the data were divided into 12 blocks, each containing 25,000 observations. Then, the utility was predicted for each block.

4.1.5.1 Predicting Utility

The initial plan was to predict the daily utilities for the participants. However, running the *for loop* for daily predictions would have taken a significant amount of time. Therefore, a monthly utility prediction was chosen by converting the dates into month values. The following steps were taken for the prediction:

- Step 1: Changed variables with date values to monthly values.
- Step 2: Determined the length of the follow-up window for total QALYs (denoted as *fu*).
- Step 3: Created a *for loop* to predict utility using coefficients sourced from Sullivan et al.[7] for the covariates adjusted for prediction and comorbidities.
- Step 4: Predicted the utilities for the participants, considering:
 - The 240 conditions - the main prediction model
 - The four major conditions: cancer, cardiovascular disease, cerebrovascular disease, and diabetes

N.B. Since the Stata code was long and somewhat complicated, comments were added throughout to improve the interpretability of the code.

N.B. Since the Stata command used a *local macro*, the entire command was run at once.

```

use "$stata_sbp_output\part_2a.dta", clear
foreach j of varlist date* v* epiend {
    replace `j' = mofd(`j')
    format `j' %tm
}

```



```

save "$stata_sbp_output\part_2a_months.dta", replace
    local j = 25000
    local k = 1
forval i = 1/12 {
    di `i' " " `k' " " `j'
    use part_2a_months, clear

    keep if _n >= `k' & _n<= `j'
    save "$stata_sbp_output\part_2a_months_`i'.dta", replace

    local k = `j'+1
    local j = `j'+25000
}

*Let's look at the maximum follow-up dates by data source (HES, SMR, and PEDW)
use "$stata_sbp_output\part_2a_months.dta", clear
bysort dsources: su date_episode
*The maximum follow up dates for hospitalization were:
    * HES: upto 31 Oct 2022 (753 months)
    * SMR: upto 31 Oct 2022 (753 months)
    * PEDW: upto 26 May 2022 (748 months)
    * For participant without hospital reported data follow-up assumed until 31 Oct 2022 (753 mo
> nths)
*****
*Let's work on follow-up reported by HES, SMR or no hospital data reported
*****
forval x = 1/12 {

    use "$stata_sbp_output\part_2a_months_`x'.dta", clear

    keep if dsources != "PEDW"

*Find out how long the follow-up window should be for total QALYs (`fu')
*This is the longest amount of follow-up, i.e. the first entry to 31/10/2022
su date_baseline
qui su date_baseline
local min = r(min)
local fu = 753 -`min'
gen months = 0

forvalues i = 0/`fu' {
    *Utility variables are the estimated utility for each day after baseline (1 is the end of ba
> seline). For convenience we converted each date to monthly format
    *Start the utility as baseline, adjusted for age
    gen utility_`i' = 0.9630907 + ((age+0.5+`i'/12)*-0.0002747) + (sex*0.0010046) /// constant +
> age (half a year added to baseline age, plus days after baseline/365.25) + sex
    + 0.0396568 /// constant added for income (mid income), as this isn't mappable
    + (qual_1*0.0028418) + (qual_2*0.0056836) + (qual_3*0.0060444) // qualifications

    *Update the dummy variables to 1 if the participant had the condition on or before the utili
> ty day
    forvalues v = 1/240 {
        qui replace z`v' = 1 if v`v' <= (date_baseline + `i')
    }

    *Generate d4_utility just using the 4 main conditions (stroke, heart disease, cancer, type 2
> diabetes)
    *Cancer = v3-v21, diabetes = v35, cardiovascular disease = v112-v123, stroke = v124-v128
    {
        gen d4_utility_`i' = utility_`i'
        qui replace d4_utility_`i' = d4_utility_`i' ///
        + z3*-0.0278367 /// Cancer
        + z4*-0.0705565 ///
        + z5*-0.0673908 ///
        + z6*-0.0929452 ///
    }
}

```

```

+ z7*-0.1192427 ///
+ z8*-0.0020176 ///
+ z9*-0.0134124 ///
+ z10*-0.0194279 ///
+ z11*-0.1132353 ///
+ z12*-0.0513159 ///
+ z13*-0.049392 ///
+ z14*-0.056553 ///
+ z15*-0.0479982 ///
+ z16*-0.0414255 ///
+ z17*-0.0858906 ///
+ z18*-0.0335121 ///
+ z19*-0.0584308 ///
+ z20*-0.0099273 ///
+ z21*-0.0500121 ///
+ z35*-0.0714349 /// Diabetes
+ z113*-0.0866826 /// Cardiovascular disease
+ z114*-0.0367975 ///
+ z115*-0.0854255 ///
+ z116*-0.0626527 ///
+ z117*-0.0387891 ///
+ z118*-0.0288829 ///
+ z119*-0.1556403 ///
+ z120*-0.0866662 ///
+ z121*-0.0383929 ///
+ z122*-0.1166656 ///
+ z123*-0.0867575 ///
+ z124*-0.0348978 /// Cerebrovascular disease
+ z125*-0.0330382 ///
+ z126*-0.1170501 ///
+ z127*-0.0310476 ///
+ z128*-0.0731964

*Account for number of comorbidities
egen ncc = rowtotal(z3-z21 z35 z112-z123 z124-z128)
forvalues j = 2/10 {
    qui replace ncc`j' = 0 // reset the ncc counters to 0 each time, so if switching fro
> m 2 to 3 comorbidities, you don't double count the nccs
    qui replace ncc`j' = 1 if ncc == `j'
}
*ncc10 if 10 or more, so account for that
qui replace ncc10 = 1 if ncc > 10 & ncc < .
drop ncc

*Update the utility variable for number of comorbidities
qui replace d4_utility_`i' = d4_utility_`i'+(ncc2*-0.0528484) ///
+(ncc3*-0.0415352) ///
+(ncc4*-0.0202969) ///
+(ncc5*0.0083033) ///
+(ncc6*0.0408673) ///
+(ncc7*0.0668729) ///
+(ncc8*0.1158895) ///
+(ncc9*0.1344392) ///
+(ncc10*0.183614)

*Account for people who have died: set utility to 0 if death is on or before the utility day
qui replace d4_utility_`i' = 0 if date_death <= (date_baseline + `i')

*Account for utilities after the end date
qui replace d4_utility_`i' = . if date_baseline + `i' > 753

*End of d4_utility code
}

*Account for number of comorbidities
egen ncc = rowtotal(z*)
forvalues j = 2/10 {
    qui replace ncc`j' = 0 // reset the ncc counters to 0 each time, so if switching fro

```



```

> m 2 to 3 comorbidities, you don't double count the nccs
    qui replace ncc`j` = 1 if ncc == `j`
  }
  *ncc10 if 10 or more, so account for that
  qui replace ncc10 = 1 if ncc > 10 & ncc < .
  drop ncc

  *Update the utility variable for number of comorbidities
  qui replace utility_`i` = utility_`i`+(ncc2*-0.0528484) ///
  +(ncc3*-0.0415352) ///
  +(ncc4*-0.0202969) ///
  +(ncc5*0.0083033) ///
  +(ncc6*0.0408673) ///
  +(ncc7*0.0668729) ///
  +(ncc8*0.1158895) ///
  +(ncc9*0.1344392) ///
  +(ncc10*0.183614)

  *Main equation code
  {
  qui replace utility_`i` = utility_`i` ///
  + z1*-0.0844983 ///
  + z2*-0.0957076 ///
  + z3*-0.0278367 ///
  + z4*-0.0705565 ///
  + z5*-0.0673908 ///
  + z6*-0.0929452 ///
  + z7*-0.1192427 ///
  + z8*-0.0020176 ///
  + z9*-0.0134124 ///
  + z10*-0.0194279 ///
  + z11*-0.1132353 ///
  + z12*-0.0513159 ///
  + z13*-0.049392 ///
  + z14*-0.056553 ///
  + z15*-0.0479982 ///
  + z16*-0.0414255 ///
  + z17*-0.0858906 ///
  + z18*-0.0335121 ///
  + z19*-0.0584308 ///
  + z20*-0.0099273 ///
  + z21*-0.0500121 ///
  + z22*0.0035188 ///
  + z23*-0.0427221 ///
  + z24*-0.0005495 ///
  + z25*-0.0012877 ///
  + z26*0.0025671 ///
  + z27*-0.0457617 ///
  + z28*0.0021036 ///
  + z29*-0.0356286 ///
  + z30*-0.0320323 ///
  + z31*-0.0345585 ///
  + z32*-0.0064666 ///
  + z33*0.0048511 ///
  + z34*-0.0471758 ///
  + z35*-0.0714349 ///
  + z36*-0.0041212 ///
  + z37*-0.1202905 ///
  + z38*-0.043536 ///
  + z39*-0.1555965 ///
  + z40*0.0251649 ///
  + z41*-0.0302843 ///
  + z42*0.0374719 ///
  + z43*-0.0065353 ///
  + z44*-0.0400731 ///
  + z45*-0.0094726 ///
  + z46*-0.0299147 ///
  + z47*-0.0708597 ///

```

```

+ z48*-0.0807239 ///  

+ z49*-0.035499 ///  

+ z50*-0.1858287 ///  

+ z51*-0.0218789 ///  

+ z52*-0.1290684 ///  

+ z53*-0.007524 ///  

+ z54*0.001099 ///  

+ z55*-0.0606417 ///  

+ z56*-0.0679451 ///  

+ z57*-0.1129732 ///  

+ z58*-0.1269103 ///  

+ z59*-0.3272071 ///  

+ z60*-0.0860614 ///  

+ z61*-0.093641 ///  

+ z62*-0.087199 ///  

+ z63*0.0075985 ///  

+ z64*-0.0314791 ///  

+ z65*-0.0514317 ///  

+ z66*-0.0284652 ///  

+ z67*-0.0537585 ///  

+ z68*-0.0367998 ///  

+ z69*-0.0634499 ///  

+ z70*-0.1445047 ///  

+ z71*-0.1123433 ///  

+ z72*-0.1196535 ///  

+ z73*-0.0020176 ///  

+ z74*-0.025434 ///  

+ z75*-0.2098308 ///  

+ z76*-0.2165659 ///  

+ z77*-0.227641 ///  

+ z78*-0.0968456 ///  

+ z79*-0.1496606 ///  

+ z80*-0.3915088 ///  

+ z81*-0.2114183 ///  

+ z82*-0.2459921 ///  

+ z83*-0.5262339 ///  

+ z84*-0.0398664 ///  

+ z85*-0.0438483 ///  

+ z86*-0.0128252 ///  

+ z87*-0.0855975 ///  

+ z88*-0.1495273 ///  

+ z89*-0.0761332 ///  

+ z90*-0.1432425 ///  

+ z91*-0.0757759 ///  

+ z92*-0.0653636 ///  

+ z93*-0.1491389 ///  

+ z94*-0.3675736 ///  

+ z95*-0.0296473 ///  

+ z96*-0.0282526 ///  

+ z97*-0.0357925 ///  

+ z98*-0.0384753 ///  

+ z99*-0.0334403 ///  

+ z100*0.0014598 ///  

+ z101*-0.0408324 ///  

+ z102*-0.0642188 ///  

+ z103*0.0206044 ///  

+ z104*-0.0002747 ///  

+ z105*-0.0116543 ///  

+ z106*-0.0253964 ///  

+ z107*-0.1127119 ///  

+ z108*-0.0143062 ///  

+ z109*-0.0003608 ///  

+ z110*-0.0217701 ///  

+ z111*-0.0460375 ///  

+ z112*-0.0625727 ///  

+ z113*-0.0866826 ///  

+ z114*-0.0367975 ///  

+ z115*-0.0854255 ///  


```

```

+ z116*-0.0626527 ///  

+ z117*0.0387891 ///  

+ z118*-0.0288829 ///  

+ z119*-0.1556403 ///  

+ z120*-0.0866662 ///  

+ z121*-0.0383929 ///  

+ z122*-0.1166656 ///  

+ z123*-0.0867575 ///  

+ z124*-0.0348978 ///  

+ z125*-0.0330382 ///  

+ z126*-0.1170501 ///  

+ z127*-0.0310476 ///  

+ z128*-0.0731964 ///  

+ z129*-0.0364444 ///  

+ z130*-0.0351324 ///  

+ z131*-0.0983128 ///  

+ z132*0.0012497 ///  

+ z133*-0.0390077 ///  

+ z134*-0.0101898 ///  

+ z135*-0.0408994 ///  

+ z136*-0.0849659 ///  

+ z137*-0.0646133 ///  

+ z138*-0.0105162 ///  

+ z139*-0.0049371 ///  

+ z140*-0.0625365 ///  

+ z141*-0.0633589 ///  

+ z142*-0.0021036 ///  

+ z143*0.0125577 ///  

+ z144*-0.0012877 ///  

+ z145*-0.0443551 ///  

+ z146*-0.1090566 ///  

+ z147*-0.0463398 ///  

+ z148*-0.133609 ///  

+ z149*-0.1262862 ///  

+ z150*-0.0776235 ///  

+ z151*-0.0372242 ///  

+ z152*-0.0021036 ///  

+ z153*-0.0473768 ///  

+ z154*-0.001099 ///  

+ z155*-0.0438595 ///  

+ z156*-0.0568783 ///  

+ z157*-0.0439011 ///  

+ z158*-0.043736 ///  

+ z159*-0.068438 ///  

+ z160*-0.018698 ///  

+ z161*-0.0641628 ///  

+ z162*-0.1218071 ///  

+ z163*-0.0608405 ///  

+ z164*-0.0752718 ///  

+ z165*-0.0574825 ///  

+ z166*-0.0726647 ///  

+ z167*-0.0516125 ///  

+ z168*-0.0831554 ///  

+ z169*-0.0955822 ///  

+ z170*-0.0583732 ///  

+ z171*-0.0277284 ///  

+ z172*-0.1269329 ///  

+ z173*-0.1753948 ///  

+ z174*0.0367826 ///  

+ z175*-0.0013305 ///  

+ z176*-0.1103664 ///  

+ z177*-0.0232626 ///  

+ z178*-0.1005789 ///  

+ z179*-0.0900796 ///  

+ z180*-0.0053797 ///  

+ z181*-0.0022102 ///  

+ z182*-0.0408141 ///  

+ z183*-0.015596 ///  


```

```

+ z184*-0.0019232 ///
+ z185*-0.003297 ///
+ z186*-0.0659663 ///
+ z187*-0.0097629 ///
+ z188*-0.0024727 ///
+ z189*-0.0206471 ///
+ z190*-0.0205704 ///
+ z191*-0.044975 ///
+ z192*-0.0013737 ///
+ z193*-0.0088572 ///
+ z194*0.0010046 ///
+ z195*-0.0037521 ///
+ z196*-0.0479614 ///
+ z197*-0.0482219 ///
+ z198*0.003297 ///
+ z199*-0.0335269 ///
+ z200*0.0011933 ///
+ z201*-0.0705302 ///
+ z202*-0.0832538 ///
+ z203*-0.1659431 ///
+ z204*-0.1144509 ///
+ z205*-0.1179321 ///
+ z206*-0.0692656 ///
+ z207*-0.0330332 ///
+ z208*-0.0796054 ///
+ z209*-0.0280394 ///
+ z210*-0.1168881 ///
+ z211*-0.1442472 ///
+ z212*-0.057222 ///
+ z213*-0.0865975 ///
+ z214*-0.0932371 ///
+ z215*-0.071455 ///
+ z216*-0.1250332 ///
+ z217*-0.1696578 ///
+ z218*-0.0362949 ///
+ z219*-0.0594539 ///
+ z220*-0.0809203 ///
+ z221*-0.271901 ///
+ z222*-0.0181141 ///
+ z223*-0.2049198 ///
+ z224*-0.020492 ///
+ z225*-0.0186058 ///
+ z226*-0.0366582 ///
+ z227*0.0105347 ///
+ z228*-0.0257964 ///
+ z229*-0.2136477 ///
+ z230*-0.122754 ///
+ z231*-0.0045418 ///
+ z232*-0.0019232 ///
+ z233*-0.1567364 ///
+ z234*-0.0667603 ///
+ z235*-0.0347814 ///
+ z236*-0.0575569 ///
+ z237*-0.0946193 ///
+ z238*-0.1420118 ///
+ z239*-0.0651164 ///
+ z240*-0.0414367 ///
}

*Account for people who have died: set utility to 0 if death is on or before the utility day
qui replace utility_`i` = 0 if date_death <= (date_baseline + `i`)

*Account for utilities after the end date
qui replace utility_`i` = . if date_baseline + `i` > 753

*Add one to the count of days (this is equivalent to adding 0.03285 months)
qui replace months = months + 1 if date_baseline + `i` <= 753

```

```

        dis "Completed day `i` of `fu`"
    }

    *Utilities are created for all participants for all days between their registration date and the fol
    > low-up window
    *Total QALYs are the sum of all utilities, divided by months of follow-up (accounting for the number
    > of utility observations)
    egen qaly_hes = rowtotal(utility*)
    replace qaly_hes = qaly_hes/months
    label variable qaly_hes "Average Utility to October 2022 (i.e. QALYs per year) [HES]"

    *And for d4_utility
    egen d4_qaly_hes = rowtotal(d4_utility*)
    replace d4_qaly_hes = d4_qaly_hes/months
    label variable d4_qaly_hes "Average D4 Utility to October 2022 (i.e. QALYs per year) [HES]"

    drop z* ncc* utility* v* d4_utility*

    compress

    save "$stata_sbp_output\part_3a_753_`x`.dta", replace
}

```

The QALYs were predicted over an average follow-up period and the data were saved into **part_3a_753_x.dta** for England and Scotland, where **x** represents the data block (ranging from 1 to 12). Next, the analysis was focused on the Wales data.

```

*****
*Let's work on follow-up reported by PEDW
*****
forval x = 1/12 {

    use "$stata_sbp_output\part_2a_months_`x`.dta", clear

    keep if dsources == "PEDW"

    *Find out how long the follow-up window should be for total QALYs (`fu`)
    *This is the longest amount of follow-up, i.e. the first entry to 26/05/2022.

    su date_baseline
    qui su date_baseline
    local min = r(min)
    local fu = 748 - `min`
    gen months = 0

    forvalues i = 0/`fu` {
        *Utility variables are the estimated utility for each day after baseline (1 is the end of ba
        > seline). For convenience we converted each date to monthly format
        *Start the utility as baseline, adjusted for age
        gen utility_`i` = 0.9630907 + ((age+0.5+`i`/12)*-0.0002747) + (sex*0.0010046) /// constant +
        > age (half a year added to baseline age, plus days after baseline/365.25) + sex
        + 0.0396568 /// constant added for income (mid income), as this isn't mappable
        + (qual_1*0.0028418) + (qual_2*0.0056836) + (qual_3*0.0060444) // qualifications

        *Update the dummy variables to 1 if the participant had the condition on or before the utili
        > ty day
        forvalues v = 1/240 {
            qui replace z`v` = 1 if v`v` <= (date_baseline + `i`)
        }

        *Generate d4_utility just using the 4 main conditions (stroke, heart disease, cancer, type 2
        > diabetes)
        *Cancer = v3-v21, diabetes = v35, cardiovascular disease = v112-v123, stroke = v124-v128
        {
            gen d4_utility_`i` = utility_`i`
            qui replace d4_utility_`i` = d4_utility_`i` ///
            + z3*-0.0278367 /// Cancer
            + z4*-0.0705565 ///
            + z5*-0.0673908 ///
            + z6*-0.0929452 ///
            + z7*-0.1192427 ///

```

```

+ z8*-0.0020176 ///
+ z9*-0.0134124 ///
+ z10*-0.0194279 ///
+ z11*-0.1132353 ///
+ z12*-0.0513159 ///
+ z13*-0.049392 ///
+ z14*-0.056553 ///
+ z15*-0.0479982 ///
+ z16*-0.0414255 ///
+ z17*-0.0858906 ///
+ z18*-0.0335121 ///
+ z19*-0.0584308 ///
+ z20*-0.0099273 ///
+ z21*-0.0500121 ///
+ z35*-0.0714349 /// Diabetes
+ z113*-0.0866826 /// Cardiovascular disease
+ z114*-0.0367975 ///
+ z115*-0.0854255 ///
+ z116*-0.0626527 ///
+ z117*-0.0387891 ///
+ z118*-0.0288829 ///
+ z119*-0.1556403 ///
+ z120*-0.0866662 ///
+ z121*-0.0383929 ///
+ z122*-0.116656 ///
+ z123*-0.0867575 ///
+ z124*-0.0348978 /// Cerebrovascular disease
+ z125*-0.0330382 ///
+ z126*-0.1170501 ///
+ z127*-0.0310476 ///
+ z128*-0.0731964

*Account for number of comorbidities
egen ncc = rowtotal(z3-z21 z35 z112-z123 z124-z128)
forvalues j = 2/10 {
    qui replace ncc`j' = 0 // reset the ncc counters to 0 each time, so if switching fro
> m 2 to 3 comorbidities, you don't double count the nccs
    qui replace ncc`j' = 1 if ncc == `j'
}
*ncc10 if 10 or more, so account for that
qui replace ncc10 = 1 if ncc > 10 & ncc < .
drop ncc

*Update the utility variable for number of comorbidities
qui replace d4_utility_`i' = d4_utility_`i' + (ncc2*-0.0528484) ///
+ (ncc3*-0.0415352) ///
+ (ncc4*-0.0202969) ///
+ (ncc5*0.0083033) ///
+ (ncc6*0.0408673) ///
+ (ncc7*0.0668729) ///
+ (ncc8*0.1158895) ///
+ (ncc9*0.1344392) ///
+ (ncc10*0.183614)

*Account for people who have died: set utility to 0 if death is on or before the utility day
qui replace d4_utility_`i' = 0 if date_death <= (date_baseline + `i')

*Account for utilities after the end date
qui replace d4_utility_`i' = . if date_baseline + `i' > 748

*End of d4_utility code
}

*Account for number of comorbidities
egen ncc = rowtotal(z*)
forvalues j = 2/10 {
    qui replace ncc`j' = 0 // reset the ncc counters to 0 each time, so if switching fro
> m 2 to 3 comorbidities, you don't double count the nccs

```

```

        qui replace ncc`j` = 1 if ncc == `j`
    }
    *ncc10 if 10 or more, so account for that
    qui replace ncc10 = 1 if ncc > 10 & ncc < .
    drop ncc

    *Update the utility variable for number of comorbidities
    qui replace utility_`i` = utility_`i`+(ncc2*-0.0528484) ///
    +(ncc3*-0.0415352) ///
    +(ncc4*-0.0202969) ///
    +(ncc5*0.0083033) ///
    +(ncc6*0.0408673) ///
    +(ncc7*0.0668729) ///
    +(ncc8*0.1158895) ///
    +(ncc9*0.1344392) ///
    +(ncc10*0.183614)

    *Main equation code
    {
    qui replace utility_`i` = utility_`i` ///
    + z1*-0.0844983 ///
    + z2*-0.0957076 ///
    + z3*-0.0278367 ///
    + z4*-0.0705565 ///
    + z5*-0.0673908 ///
    + z6*-0.0929452 ///
    + z7*-0.1192427 ///
    + z8*-0.0020176 ///
    + z9*-0.0134124 ///
    + z10*-0.0194279 ///
    + z11*-0.1132353 ///
    + z12*-0.0513159 ///
    + z13*-0.049392 ///
    + z14*-0.056553 ///
    + z15*-0.0479982 ///
    + z16*-0.0414255 ///
    + z17*-0.0858906 ///
    + z18*-0.0335121 ///
    + z19*-0.0584308 ///
    + z20*-0.0099273 ///
    + z21*-0.0500121 ///
    + z22*0.0035188 ///
    + z23*-0.0427221 ///
    + z24*-0.0005495 ///
    + z25*-0.0012877 ///
    + z26*0.0025671 ///
    + z27*-0.0457617 ///
    + z28*0.0021036 ///
    + z29*-0.0356286 ///
    + z30*-0.0320323 ///
    + z31*-0.0345585 ///
    + z32*-0.0064666 ///
    + z33*0.0048511 ///
    + z34*-0.0471758 ///
    + z35*-0.0714349 ///
    + z36*-0.0041212 ///
    + z37*-0.1202905 ///
    + z38*-0.043536 ///
    + z39*-0.1555965 ///
    + z40*0.0251649 ///
    + z41*-0.0302843 ///
    + z42*0.0374719 ///
    + z43*-0.0065353 ///
    + z44*-0.0400731 ///
    + z45*-0.0094726 ///
    + z46*-0.0299147 ///
    + z47*-0.0708597 ///
    + z48*-0.0807239 ///

```

```

+ z49*-0.035499 ///
+ z50*-0.1858287 ///
+ z51*-0.0218789 ///
+ z52*-0.1290684 ///
+ z53*-0.007524 ///
+ z54*0.001099 ///
+ z55*-0.0606417 ///
+ z56*-0.0679451 ///
+ z57*-0.1129732 ///
+ z58*-0.1269103 ///
+ z59*-0.3272071 ///
+ z60*-0.0860614 ///
+ z61*-0.093641 ///
+ z62*-0.087199 ///
+ z63*0.0075985 ///
+ z64*-0.0314791 ///
+ z65*-0.0514317 ///
+ z66*-0.0284652 ///
+ z67*-0.0537585 ///
+ z68*-0.0367998 ///
+ z69*-0.0634499 ///
+ z70*-0.1445047 ///
+ z71*-0.1123433 ///
+ z72*-0.1196535 ///
+ z73*-0.0020176 ///
+ z74*-0.025434 ///
+ z75*-0.2098308 ///
+ z76*-0.2165659 ///
+ z77*-0.227641 ///
+ z78*-0.0968456 ///
+ z79*-0.1496606 ///
+ z80*-0.3915088 ///
+ z81*-0.2114183 ///
+ z82*-0.2459921 ///
+ z83*-0.5262339 ///
+ z84*-0.0398664 ///
+ z85*-0.0438483 ///
+ z86*-0.0128252 ///
+ z87*-0.0855975 ///
+ z88*-0.1495273 ///
+ z89*-0.0761332 ///
+ z90*-0.1432425 ///
+ z91*-0.0757759 ///
+ z92*-0.0653636 ///
+ z93*-0.1491389 ///
+ z94*-0.3675736 ///
+ z95*-0.0296473 ///
+ z96*-0.0282526 ///
+ z97*-0.0357925 ///
+ z98*-0.0384753 ///
+ z99*-0.0334403 ///
+ z100*0.0014598 ///
+ z101*-0.0408324 ///
+ z102*-0.0642188 ///
+ z103*0.0206044 ///
+ z104*-0.0002747 ///
+ z105*-0.0116543 ///
+ z106*-0.0253964 ///
+ z107*-0.1127119 ///
+ z108*-0.0143062 ///
+ z109*-0.0003608 ///
+ z110*-0.0217701 ///
+ z111*-0.0460375 ///
+ z112*-0.0625727 ///
+ z113*-0.0866826 ///
+ z114*-0.0367975 ///
+ z115*-0.0854255 ///
+ z116*-0.0626527 ///

```



```

+ z117*0.0387891 ///
+ z118*-0.0288829 ///
+ z119*-0.1556403 ///
+ z120*-0.0866662 ///
+ z121*-0.0383929 ///
+ z122*-0.1166656 ///
+ z123*-0.0867575 ///
+ z124*-0.0348978 ///
+ z125*-0.0330382 ///
+ z126*-0.1170501 ///
+ z127*-0.0310476 ///
+ z128*-0.0731964 ///
+ z129*-0.0364444 ///
+ z130*-0.0351324 ///
+ z131*-0.0983128 ///
+ z132*0.0012497 ///
+ z133*-0.0390077 ///
+ z134*-0.0101898 ///
+ z135*-0.0408994 ///
+ z136*-0.0849659 ///
+ z137*-0.0646133 ///
+ z138*-0.0105162 ///
+ z139*-0.0049371 ///
+ z140*-0.0625365 ///
+ z141*-0.0633589 ///
+ z142*-0.0021036 ///
+ z143*0.0125577 ///
+ z144*-0.0012877 ///
+ z145*-0.0443551 ///
+ z146*-0.1090566 ///
+ z147*-0.0463398 ///
+ z148*-0.133609 ///
+ z149*-0.1262862 ///
+ z150*-0.0776235 ///
+ z151*-0.0372242 ///
+ z152*-0.0021036 ///
+ z153*-0.0473768 ///
+ z154*-0.001099 ///
+ z155*-0.0438595 ///
+ z156*-0.0568783 ///
+ z157*-0.0439011 ///
+ z158*-0.043736 ///
+ z159*-0.068438 ///
+ z160*-0.018698 ///
+ z161*-0.0641628 ///
+ z162*-0.1218071 ///
+ z163*-0.0608405 ///
+ z164*-0.0752718 ///
+ z165*-0.0574825 ///
+ z166*-0.0726647 ///
+ z167*-0.0516125 ///
+ z168*-0.0831554 ///
+ z169*-0.0955822 ///
+ z170*-0.0583732 ///
+ z171*-0.0277284 ///
+ z172*-0.1269329 ///
+ z173*-0.1753948 ///
+ z174*0.0367826 ///
+ z175*-0.0013305 ///
+ z176*-0.1103664 ///
+ z177*-0.0232626 ///
+ z178*-0.1005789 ///
+ z179*-0.0900796 ///
+ z180*-0.0053797 ///
+ z181*-0.0022102 ///
+ z182*-0.0408141 ///
+ z183*-0.015596 ///
+ z184*-0.0019232 ///

```

```

+ z185*-0.003297 ///
+ z186*-0.0659663 ///
+ z187*-0.0097629 ///
+ z188*-0.0024727 ///
+ z189*-0.0206471 ///
+ z190*-0.0205704 ///
+ z191*-0.044975 ///
+ z192*-0.0013737 ///
+ z193*-0.0088572 ///
+ z194*0.0010046 ///
+ z195*-0.0037521 ///
+ z196*-0.0479614 ///
+ z197*-0.0482219 ///
+ z198*0.003297 ///
+ z199*-0.0335269 ///
+ z200*0.0011933 ///
+ z201*-0.0705302 ///
+ z202*-0.0832538 ///
+ z203*-0.1659431 ///
+ z204*-0.1144509 ///
+ z205*-0.1179321 ///
+ z206*-0.0692656 ///
+ z207*-0.0330332 ///
+ z208*-0.0796054 ///
+ z209*-0.0280394 ///
+ z210*-0.1168881 ///
+ z211*-0.1442472 ///
+ z212*-0.057222 ///
+ z213*-0.0865975 ///
+ z214*-0.0932371 ///
+ z215*-0.071455 ///
+ z216*-0.1250332 ///
+ z217*-0.1696578 ///
+ z218*-0.0362949 ///
+ z219*-0.0594539 ///
+ z220*-0.0809203 ///
+ z221*-0.271901 ///
+ z222*-0.0181141 ///
+ z223*-0.2049198 ///
+ z224*-0.020492 ///
+ z225*-0.0186058 ///
+ z226*-0.0366582 ///
+ z227*0.0105347 ///
+ z228*-0.0257964 ///
+ z229*-0.2136477 ///
+ z230*-0.122754 ///
+ z231*-0.0045418 ///
+ z232*-0.0019232 ///
+ z233*-0.1567364 ///
+ z234*-0.0667603 ///
+ z235*-0.0347814 ///
+ z236*-0.0575569 ///
+ z237*-0.0946193 ///
+ z238*-0.1420118 ///
+ z239*-0.0651164 ///
+ z240*-0.0414367 ///
}

*Account for people who have died: set utility to 0 if death is on or before the utility day
qui replace utility_`i` = 0 if date_death <= (date_baseline + `i`)

*Account for utilities after the end date
qui replace utility_`i` = . if date_baseline + `i` > 748

*Add one to the count of days (this is equivalent to adding 0.03285 months)
qui replace months = months + 1 if date_baseline + `i` <= 748

dis "Completed day `i` of `fu`"

```

```

}
*Utilities are created for all participants for all days between their registration date and the fol
> low-up window
*Total QALYs are the sum of all utilities, divided by months of follow-up (accounting for the number
> of utility observations)
egen qaly_hes = rowtotal(utility*)
replace qaly_hes = qaly_hes/months
label variable qaly_hes "Average Utility to May 2022 (i.e. QALYs per year) [HES]"

*And for d4_utility
egen d4_qaly_hes = rowtotal(d4_utility*)
replace d4_qaly_hes = d4_qaly_hes/months
label variable d4_qaly_hes "Average D4 Utility to May 2022 (i.e. QALYs per year) [HES]"
drop z* ncc* utility* v* d4_utility*
compress

save "$stata_sbp_output\part_3a_748_`x`.dta", replace
}

```

The QALYs were also predicted over an average follow-up period and the data were saved into **part_3a_748_x.dta** for Wales, where **x** represents the data block (ranging from 1 to 12). Next, we merged these two datasets.

```

use "$stata_sbp_output\part_3a_753_1.dta", clear
forval i == 2/12 {
    append using "$stata_sbp_output\part_3a_753_`i`.dta"
}
save "$stata_sbp_output\part_3a_753.dta", replace

use "$stata_sbp_output\part_3a_748_1.dta", clear
forval i == 2/12 {
    append using "$stata_sbp_output\part_3a_748_`i`.dta"
}
save "$stata_sbp_output\part_3a_748.dta", replace
append using "$stata_sbp_output\part_3a_753.dta"

tabstat qaly_hes d4_qaly_hes, statistics(mean sd median p25 p75 min max)
replace qaly_hes = 1 if qaly_hes >1
replace d4_qaly_hes = 1 if d4_qaly_hes >1
save "$stata_sbp_output\part_3a.dta", replace

```

4.1.6 Working on EQ-5D data collected by UK Biobank

Web based **EQ-5D-5L** questionnaires were administered to the UK Biobank participants as part of the chronic pain (administered in 2019–20) and mental well-being (administered in 2022–23) surveys. We calculated the **EQ-5D index** using the **UK tariffs** (i.e., value set) for each survey[10]. Once the EQ-5D-indexes were calculated, we took the average EQ-5D-index for participants who had EQ-5D data for both surveys (124,830). The remaining participants had EQ-5D data for either chronic pain survey (42,281) or mental well-being survey (44,707); hence, the average EQ-5D index was not calculated. The next line of codes will prepare our data for EQ-5D index calculation, then apply the UK tariffs.

```

use "$stata_sbp_output\part_3a.dta", clear
*Work on EQ_5D data collected by UKB
rename (n_120098 n_120099 n_120100 n_120101 n_120102) (mobility selfcare activity pain anxiety)
rename (n_29150 n_29151 n_29152 n_29153 n_29154) (mobility_1 selfcare_1 activity_1 pain_1 anxiety_1
> )

foreach j in mobility selfcare activity pain anxiety {
    replace `j` = 1 if `j` == -521
    replace `j` = 2 if `j` == -522
}

```

```

        replace `j` = 3 if `j` == -523
        replace `j` = 4 if `j` == -524
        replace `j` = 5 if `j` == -525
    }
    foreach j in mobility_1 selfcare_1 activity_1 pain_1 anxiety_1 {
        replace `j` = `j`+1 if `j` !=.
    }

    *Work on the health states
    egen health_states = concat(mobility_1 selfcare_1 activity_1 pain_1 anxiety_1)
    tostring health_states, replace force
    replace health_states = "88" if health_states == "....."
    destring health_states, replace
    replace health_states = . if health_states == 88
    order health_states, a(anxiety_1)

    egen health_states_1 = concat(mobility_1 selfcare_1 activity_1 pain_1 anxiety_1)
    tostring health_states_1, replace force
    replace health_states_1 = "88" if health_states_1 == "....."
    destring health_states_1, replace
    replace health_states_1 = . if health_states_1 == 88
    order health_states_1, a(anxiety_1)

    /*
    Computing EQ-5D-5L index values with STATA using the English (ENG) Devlin value set
    Version 1.1 (Updated 01/12/2020)

    The variables for the 5 dimensions of the EQ-5D-5L descriptive system should be named `mobility`,
    `selfcare`, `activity`, `pain`, and `anxiety`. If they are given different names the syntax code
    below will not work properly. The 5 variables should contain the values for the different dimensions
    >
    in the EQ-5D health profile (i.e. 1, 2, 3, 4 or 5). The variable `EQindex` contains the values of th
    > e
    EQ-5D-5L index values on the basis of the ENG set of weights.
    You can copy and paste the syntax below directly into a STATA syntax window.
    */
    *****
    *STATA syntax code for the computation of index*
    *values with the English value set*
    *****

    gen disut_mo= .
    replace disut_mo= 0 if missing(disut_mo) ///
    & mobility == 1
    replace disut_mo= 0.058 if missing(disut_mo) ///
    & mobility == 2
    replace disut_mo= 0.076 if missing(disut_mo) ///
    & mobility == 3
    replace disut_mo= 0.207 if missing(disut_mo) ///
    & mobility == 4
    replace disut_mo= 0.274 if missing(disut_mo) ///
    & mobility == 5

    gen disut_sc= .
    replace disut_sc= 0 if missing(disut_sc) ///
    & selfcare == 1
    replace disut_sc= 0.050 if missing(disut_sc) ///
    & selfcare == 2
    replace disut_sc= 0.080 if missing(disut_sc) ///
    & selfcare == 3
    replace disut_sc= 0.164 if missing(disut_sc) ///
    & selfcare == 4
    replace disut_sc= 0.203 if missing(disut_sc) ///
    & selfcare == 5

    gen disut_ua= .
    replace disut_ua= 0 if missing(disut_ua) ///

```

```

& activity == 1
replace disut_ua= 0.050 if missing(disut_ua) ///
& activity == 2
replace disut_ua= 0.063 if missing(disut_ua) ///
& activity == 3
replace disut_ua= 0.162 if missing(disut_ua) ///
& activity == 4
replace disut_ua= 0.184 if missing(disut_ua) ///
& activity == 5

gen disut_pd= .
replace disut_pd= 0 if missing(disut_pd) ///
& pain == 1
replace disut_pd= 0.063 if missing(disut_pd) ///
& pain == 2
replace disut_pd= 0.084 if missing(disut_pd) ///
& pain == 3
replace disut_pd= 0.276 if missing(disut_pd) ///
& pain == 4
replace disut_pd= 0.335 if missing(disut_pd) ///
& pain == 5

gen disut_ad= .
replace disut_ad= 0 if missing(disut_ad) ///
& anxiety == 1
replace disut_ad= 0.078 if missing(disut_ad) ///
& anxiety == 2
replace disut_ad= 0.104 if missing(disut_ad) ///
& anxiety == 3
replace disut_ad= 0.285 if missing(disut_ad) ///
& anxiety == 4
replace disut_ad= 0.289 if missing(disut_ad) ///
& anxiety == 5

gen disut_total=disut_mo+disut_sc+disut_ua+disut_pd+disut_ad
gen EQindex=.
replace EQindex=1-disut_total
replace EQindex=round(EQindex,.001)
order disut_mo disut_sc disut_ua disut_pd disut_ad disut_total EQindex, b(mobility_1)

gen disut_mo_1= .
replace disut_mo_1= 0 if missing(disut_mo_1) ///
& mobility_1 == 1
replace disut_mo_1= 0.058 if missing(disut_mo_1) ///
& mobility_1 == 2
replace disut_mo_1= 0.076 if missing(disut_mo_1) ///
& mobility_1 == 3
replace disut_mo_1= 0.207 if missing(disut_mo_1) ///
& mobility_1 == 4
replace disut_mo_1= 0.274 if missing(disut_mo_1) ///
& mobility_1 == 5

gen disut_sc_1= .
replace disut_sc_1= 0 if missing(disut_sc_1) ///
& selfcare_1 == 1
replace disut_sc_1= 0.050 if missing(disut_sc_1) ///
& selfcare_1 == 2
replace disut_sc_1= 0.080 if missing(disut_sc_1) ///
& selfcare_1 == 3
replace disut_sc_1= 0.164 if missing(disut_sc_1) ///
& selfcare_1 == 4
replace disut_sc_1= 0.203 if missing(disut_sc_1) ///
& selfcare_1 == 5

gen disut_ua_1= .
replace disut_ua_1= 0 if missing(disut_ua_1) ///
& activity_1 == 1
replace disut_ua_1= 0.050 if missing(disut_ua_1) ///
& activity_1 == 2

```

```

replace disut_ua_1= 0.063 if missing(disut_ua_1) ///
& activity_1 == 3
replace disut_ua_1= 0.162 if missing(disut_ua_1) ///
& activity_1 == 4
replace disut_ua_1= 0.184 if missing(disut_ua_1) ///
& activity_1 == 5

gen disut_pd_1= .
replace disut_pd_1= 0 if missing(disut_pd_1) ///
& pain_1 == 1
replace disut_pd_1= 0.063 if missing(disut_pd_1) ///
& pain_1 == 2
replace disut_pd_1= 0.084 if missing(disut_pd_1) ///
& pain_1 == 3
replace disut_pd_1= 0.276 if missing(disut_pd_1) ///
& pain_1 == 4
replace disut_pd_1= 0.335 if missing(disut_pd_1) ///
& pain_1 == 5

gen disut_ad_1= .
replace disut_ad_1= 0 if missing(disut_ad_1) ///
& anxiety_1 == 1
replace disut_ad_1= 0.078 if missing(disut_ad_1) ///
& anxiety_1 == 2
replace disut_ad_1= 0.104 if missing(disut_ad_1) ///
& anxiety_1 == 3
replace disut_ad_1= 0.285 if missing(disut_ad_1) ///
& anxiety_1 == 4
replace disut_ad_1= 0.289 if missing(disut_ad_1) ///
& anxiety_1 == 5

gen disut_total_1=disut_mo_1+disut_sc_1+disut_ua_1+disut_pd_1+disut_ad_1
gen EQindex_1=.
replace EQindex_1=1-disut_total_1
replace EQindex_1 =round(EQindex_1,.001)
order disut_mo_1 disut_sc_1 disut_ua_1 disut_pd_1 disut_ad_1 disut_total_1 EQindex_1, a(health_state
> s_1)

rename n_120103 EQ_VAS
rename n_29155 EQ_VAS_1
rename ts_120128 qol_date
rename ts_29206 qol_date_1
tabstat EQindex EQindex_1 EQ_VAS EQ_VAS_1, statistics(mean sd median p25 p75 min max)
count if (EQindex !=. & EQindex_1 !=.)
count if (EQindex !=. & EQindex_1 ==.)
count if (EQindex ==. & EQindex_1 !=.)

save "$stata_sbp_output\part_3b.dta", replace

```

We have now completed the first step. We will work on the **genotype data**.

4.2 Step 2: Working on genotype data

4.2.1 Preparing our data

To run a mendelian randomisation (MR) analysis, a type of instrumental variable analysis, to estimate the causal association of SBP (i.e., exposure) with QALYs (i.e., outcome) using genetic variants as instruments, we need genetic variants that are strongly associated with the exposure (i.e., trait) of interest. In this case, we need genetic variants (also known as single nucleotide polymorphisms) that are associated with SBP. There are a number of assumptions genetic variants should meet to run an MR analysis[6, 11, 12].

- **relevance:** the variant is associated with the exposure
- **exchangeability:** the variant is not associated with the outcome via a confounding pathway
- **exclusion restriction:** the variant does not affect the outcome directly, only possibly indirectly via exposure

For this we look into a study conducted by Evangelou E. et al[13]. This genome-wide association study (GWAS) included data from the UK Biobank, International Consortium for Blood Pressure (ICBP), the US Million Veteran Program (MVP) and Estonian Genome Centre, University of Tartu (EGCUT). The UK Biobank and ICBP data were used for discovery meta-analysis, while the MVP and EGCUT data were used for replication meta-analysis. A combined meta-analysis was also performed using all data sources[13]. The GWAS study has identified 535 novel loci (1 variant per locus) that have reached the significant threshold to one of the blood pressure traits[13]. The criteria for significance threshold for genetic variants with a specific trait was based on one-stage or two-stage analysis design set by the GWAS study[13]. SNPs were filtered to meet criteria of genotype missingness below 0.015 and minor allele frequency above 0.01[13]. They were also tested for Hardy-Weinberg equilibrium and linkage disequilibrium within the GWAS[13]. Linkage disequilibrium (LD) was calculated for all variants within a 500kb window on either side of the reference SNP[13]. Variants in linkage disequilibrium with the reference SNP, reaching an r^2 threshold of 0.1 or higher, were identified[13].

The GWAS study also included 92 sentinel SNPs previously known but replicated for the first time and 357 SNPs already known and validated to have association with one of the blood pressure traits.

Among the 984 SNPs, 282 were primarily related to SBP[13]. After LD clumping (r^2 0.001 and 10,000 kb window) and removing ambiguous SNPs, 181 genetic instruments were candidates for building polygenic risk scores (PRS). If a sentinel SNP was not available in the UK Biobank, a proxy SNP was substituted ($r^2 \geq 0.8$)[13]. PRS was constructed by summing the effects of the 181 SNPs on SBP, each weighted by its effect size derived from non-UK Biobank cohorts (the ICBP meta-analysis and replication meta-analysis). The ICBP and replication meta-analyses were selected to avoid cohort overlap with the UK Biobank population.

We also need diastolic blood pressure for the multivariable MR (MVMR)analysis, an extension of standard MR analysis. The GWAS study reported sentinel SNPs association with primary and secondary blood pressure traits[13]. For the 187 sentinel SNPs primarily associated with SBP (after LD clumping), we also identified association with DBP as a secondary trait. Similarly, we also identified 208 sentinel SNPs (after LD clumping) primarily associated with DBP that were

also linked to SBP. The combined 395 SNPs were candidates for the MVMR analysis. After the exclusion of ambiguous and missing effect size SNPs, 384 and 382 SNPs were used to construct PRS for SBP and DBP, respectively. Effect estimates for the SNPs were sourced from either ICBP or replication meta-analysis.

4.2.2 Association of Genetic Variants with blood pressure traits

The following stata code will help us to format the data suitable for the next analysis.

```
*****
*Let's work on known SNPs
*****
import excel "$data_source\known_bp_snps.xlsx", firstrow clear
save "$snps\known_bp_snps.dta",replace
import excel "$data_source\known_bp_snps_association.xlsx", firstrow clear
save "$snps\known_bp_snps_association.dta",replace
merge 1:1 rsID using "$snps\known_bp_snps.dta"
replace Trait = "SBP" if rsID == "rs2076328"
replace Trait = "PP" if rsID == "rs2157597"
replace Trait = "PP" if rsID == "rs28427409"
replace Trait = "SBP" if rsID == "rs6783086"
replace Trait = "DBP" if rsID == "rs73030266"
replace Trait = "DBP" if rsID == "rs73091767"
replace Trait = "PP" if rsID == "rs7480089"
replace Trait = "DBP" if rsID == "rs7777128"
replace Trait = "DBP" if rsID == "rs7810028"
replace Trait = "PP" if rsID == "rs9479200"

count if SNP_Gwsig == "FALSE" // not reached GWAS significant level for the primary trait (16)
count if SNP_Gwsig == "n/a" // no locus coverage, rare variant (MAF<1%), low frequency variant (MAF
> 1%-5%), or not in Haplotype Reference Consortium (28)
drop if _merge == 2
drop _merge*
* the snps included in the analysis are low frequency or common varinats, and LD pruned.
preserve
keep if Trait == "SBP"
gen status = "known"
gen effect_source = "icbp"
replace P_min = Pmin_ICBP_UKBmeta if P_min == ""
drop *PP Pmin_ICBP_UKBmeta

foreach var of varlist EAF_SBP Beta_SBP se_SBP P_SBP EAF_DBP Beta_DBP se_DBP P_DBP MAF P_min {
    destring `var`, replace
}

save "$sbp_snps\known_snps_sbp.dta",replace
export delimited "$sbp_snps\known_snps_sbp.csv", replace
restore
keep if Trait == "DBP"
gen status = "known"
gen effect_source = "icbp"
replace P_min = Pmin_ICBP_UKBmeta if P_min == ""
drop *PP Pmin_ICBP_UKBmeta

foreach var of varlist EAF_SBP Beta_SBP se_SBP P_SBP EAF_DBP Beta_DBP se_DBP P_DBP MAF P_min {
    destring `var`, replace
}

save "$dbp_snps\known_snps_dbp.dta",replace
export delimited "$dbp_snps\known_snps_dbp.csv", replace
*****
*Let's work on novel (1stage210 and 2stage325) previously known but replicated
*for the first time SNPs (non) associated with SBP
```

```

*****
import excel "$data_source\1stage_snps_secondary_effect.xlsx", firstrow clear
rename (SBPICBPBETA SBPICBPSE SBPICBPPVAL DBPICBPBETA DBPICBPSE DBPICBPPVAL) (SBP_ICBP_BETA SBP_ICBP
> _SE SBP_ICBP_PVAL DBP_ICBP_BETA DBP_ICBP_SE DBP_ICBP_PVAL)
foreach x in SBP_ICBP_BETA SBP_ICBP_SE SBP_ICBP_PVAL DBP_ICBP_BETA DBP_ICBP_SE DBP_ICBP_PVAL {
    destring `x`, replace
}
save "$snps\one_stage_snps_secondary_effect.dta",replace
import excel "$data_source\2stage_snps_secondary_effect.xlsx", firstrow clear
rename (SBPrepBETA SBPrepSE SBPrepPVAL DBPrepBETA DBPrepSE DBPrepPVAL) (SBP_rep_BETA SBP_rep_SE SBP_
> rep_PVAL DBP_rep_BETA DBP_rep_SE DBP_rep_PVAL)
foreach x in SBP_rep_BETA SBP_rep_SE SBP_rep_PVAL DBP_rep_BETA DBP_rep_SE DBP_rep_PVAL {
    destring `x`, replace
}
save "$snps\two_stage_snps_secondary_effect.dta",replace

import excel "$sbp_snps\novel_sbp_snps.xlsx", firstrow clear

gen rsID = ""
gen chrpos = ""
gen A1 = ""
gen A2 = ""
gen EAF_SBP = .
gen Beta_SBP = .
gen se_SBP = .
gen P_SBP = .
gen EAF_DBP = .
gen Beta_DBP = .
gen se_DBP = .
gen P_DBP = .
gen MAF = .
gen P_min = .
gen SNP_Gwsig = ""
gen status = ""
gen effect_source = ""

replace rsID = rsID_proxy_SNP
replace chrpos = ChrPos_RepSNP
replace A1 = A1_rep if Type == "2stage325" | Type == "non"
replace A1 = A1_icbp if Type == "1stage210"
replace A2 = A2_rep if Type == "2stage325" | Type == "non"
replace A2 = A2_icbp if Type == "1stage210"
replace EAF_SBP = EAF_rep if Type == "2stage325" | Type == "non"
replace EAF_SBP = EAF_icbp if Type == "1stage210"
replace Beta_SBP = BETA_rep if Type == "2stage325" | Type == "non"
replace Beta_SBP = BETA_icbp if Type == "1stage210"
replace se_SBP = SE_rep if Type == "2stage325" | Type == "non"
replace se_SBP = SE_icbp if Type == "1stage210"
replace P_SBP = P_rep if Type == "2stage325" | Type == "non"
replace P_SBP = P_icbp if Type == "1stage210"
replace MAF = MAF_uk
replace P_min = P_comb if Type == "2stage325" | Type == "non"
replace P_min = P_disc if Type == "1stage210"
replace SNP_Gwsig = "TRUE"
replace status = "novel" if Type == "1stage210" | Type == "2stage325"
replace status = "replicated" if Type == "non"
replace effect_source = "replicaton_study" if Type == "2stage325" | Type == "non"
replace effect_source = "icbp" if Type == "1stage210"
keep rsID chrpos A1 A2 EAF_SBP Beta_SBP se_SBP P_SBP EAF_DBP Beta_DBP se_DBP P_DBP MAF Trait P_min
> SNP_Gwsig status effect_source Type
order rsID chrpos A1 A2 EAF_SBP Beta_SBP se_SBP P_SBP EAF_DBP Beta_DBP se_DBP P_DBP MAF Trait P_min
> SNP_Gwsig status effect_source Type

save "$sbp_snps\novel_rep_sbp_snps.dta", replace
export delimited "$sbp_snps\novel_rep_sbp_snps.csv", replace

```

```

use "$sbp_snps\novel_rep_sbp_snps.dta", clear
append using "$sbp_snps\known_snps_sbp.dta"
merge 1:1 rsID using "$snps\one_stage_snps_secondary_effect.dta"
drop if _merge == 2
drop _merge*

merge 1:1 rsID using "$snps\two_stage_snps_secondary_effect.dta"
drop if _merge == 2
drop _merge*

replace Beta_DBP = DBP_ICBP_BETA if Beta_DBP==. & Type == "1stage210"
replace se_DBP = DBP_ICBP_SE if se_DBP ==. & Type == "1stage210"
replace P_DBP = DBP_ICBP_PVAL if P_DBP ==. & Type == "1stage210"

replace Beta_DBP = DBP_rep_BETA if Beta_DBP==. & Type == "2stage325"
replace se_DBP = DBP_rep_SE if se_DBP ==. & Type == "2stage325"
replace P_DBP = DBP_rep_PVAL if P_DBP ==. & Type == "2stage325"

replace Beta_DBP = DBP_rep_BETA if Beta_DBP==. & Type == "non"
replace se_DBP = DBP_rep_SE if se_DBP ==. & Type == "non"
replace P_DBP = DBP_rep_PVAL if P_DBP ==. & Type == "non"
replace EAF_DBP = EAF_SBP if EAF_DBP ==.

drop a1_icbp a2_icbp SBP_ICBP_BETA SBP_ICBP_SE SBP_ICBP_PVAL DBP_ICBP_BETA DBP_ICBP_SE DBP_ICBP_PV
> AL a1_rep a2_rep SBP_rep_BETA SBP_rep_SE SBP_rep_PVAL DBP_rep_BETA DBP_rep_SE DBP_rep_PVAL

save "$sbp_snps\all_sbp_snps.dta", replace
export delimited "$sbp_snps\all_sbp_snps.csv", replace

*****
*Let's work on novel (1stage210 and 2stage325) previously known but replicated
*for the first time SNPs (non) associated with DBP
*****
import excel "$dbp_snps\novel_dbp_snps.xlsx", firstrow clear

gen rsID = ""
gen chrpos = ""
gen A1 = ""
gen A2 = ""
gen EAF_SBP = .
gen Beta_SBP = .
gen se_SBP = .
gen P_SBP = .
gen EAF_DBP = .
gen Beta_DBP = .
gen se_DBP = .
gen P_DBP = .
gen MAF = .
gen P_min = .
gen SNP_Gwsig = ""
gen status = ""
gen effect_source = ""

replace rsID = rsID_proxy_SNP
replace chrpos = ChrPos_RepSNP
replace A1 = A1_rep if Type == "2stage325" | Type == "non"
replace A1 = A1_icbp if Type == "1stage210"
replace A2 = A2_rep if Type == "2stage325" | Type == "non"
replace A2 = A2_icbp if Type == "1stage210"
replace EAF_DBP = EAF_rep if Type == "2stage325" | Type == "non"
replace EAF_DBP = EAF_icbp if Type == "1stage210"
replace Beta_DBP = BETA_rep if Type == "2stage325" | Type == "non"
replace Beta_DBP = BETA_icbp if Type == "1stage210"
replace se_DBP = SE_rep if Type == "2stage325" | Type == "non"
replace se_DBP = SE_icbp if Type == "1stage210"
replace P_DBP = P_rep if Type == "2stage325" | Type == "non"
replace P_DBP = P_icbp if Type == "1stage210"
replace MAF = MAF_uk

```

```

replace P_min = P_comb if Type == "2stage325" | Type == "non"
replace P_min = P_disc if Type == "1stage210"
replace SNP_Gwsig = "TRUE"
replace status = "novel" if Type == "1stage210" | Type == "2stage325"
replace status = "replicated" if Type == "non"
replace effect_source = "replicaton_study" if Type == "2stage325" | Type == "non"
replace effect_source = "icbp" if Type == "1stage210"
keep rsID chrpos A1 A2 EAF_SBP Beta_SBP se_SBP P_SBP EAF_DBP Beta_DBP se_DBP P_DBP MAF Trait P_min
> SNP_Gwsig status effect_source Type
order rsID chrpos A1 A2 EAF_SBP Beta_SBP se_SBP P_SBP EAF_DBP Beta_DBP se_DBP P_DBP MAF Trait P_min
> SNP_Gwsig status effect_source Type

save "$dbp_snps\novel_rep_dbp_snps.dta", replace
export delimited "$dbp_snps\novel_rep_dbp_snps.csv", replace

use "$dbp_snps\novel_rep_dbp_snps.dta", clear
append using "$dbp_snps\known_snps_dbp.dta"
merge 1:1 rsID using "$snps\one_stage_snps_secondary_effect.dta"
drop if _merge == 2
drop _merge*

merge 1:1 rsID using "$snps\two_stage_snps_secondary_effect.dta"
drop if _merge == 2
drop _merge*

replace Beta_SBP = SBP_ICBP_BETA if Beta_SBP==. & Type == "1stage210"
replace se_SBP = SBP_ICBP_SE if se_SBP==. & Type == "1stage210"
replace P_SBP = SBP_ICBP_PVAL if P_SBP==. & Type == "1stage210"

replace Beta_SBP = SBP_rep_BETA if Beta_SBP==. & Type == "2stage325"
replace se_SBP = SBP_rep_SE if se_SBP==. & Type == "2stage325"
replace P_SBP = SBP_rep_PVAL if P_SBP==. & Type == "2stage325"

replace Beta_SBP = SBP_rep_BETA if Beta_SBP==. & Type == "non"
replace se_SBP = SBP_rep_SE if se_SBP==. & Type == "non"
replace P_SBP = SBP_rep_PVAL if P_SBP==. & Type == "non"
replace EAF_SBP = EAF_DBP if EAF_SBP==.

drop a1_icbp a2_icbp SBP_ICBP_BETA SBP_ICBP_SE SBP_ICBP_PVAL DBP_ICBP_BETA DBP_ICBP_SE DBP_ICBP_PV
> AL a1_rep a2_rep SBP_rep_BETA SBP_rep_SE SBP_rep_PVAL DBP_rep_BETA DBP_rep_SE DBP_rep_PVAL

save "$dbp_snps\all_dbp_snps.dta", replace
export delimited "$dbp_snps\all_dbp_snps.csv", replace

*****
*Let's combine all SBP and DBP SNPs
*****
use "$sbp_snps\all_sbp_snps.dta", clear
append using "$dbp_snps\all_dbp_snps.dta"

save "$sbp_dbp_snps\all_sbp_dbp_snps.dta"
export delimited "$sbp_dbp_snps\all_sbp_dbp_snps.csv", replace

```

4.2.3 LD Clumping

The previous Stata codes provided CSV files for the next step: LD clumping using the **TwoSampleMR** and **ieugwas** packages from R. First, an Application Programming Interface (API) was set up to access the **IEU GWAS** database. Then, the SNPs were clumped.

```

1 #This is to set up the API
2 # Get the location of your .Renviron file
3 renviro_path <- Sys.getenv("R_ENVIRON_USER")

```

```

4
5 # If the .Renviron file doesn't exist, create it
6 if (renviron_path == "") {
7   renviron_path <- file.path(Sys.getenv("HOME"), ".Renviron")
8   file.create(renviron_path)
9 }
10
11 # Append your token to the .Renviron file
12 cat("OPENGWAS_JWT=<Add Your Token Here>", file = renviron_path, append =
    TRUE)
13
14 # Print the path to verify
15 print(renviron_path)
16
17
18 library(ieugwasr)
19
20 # Check if the token is loaded
21 jwt <- ieugwasr::get_opengwas_jwt()
22 if (nzchar(jwt)) {
23   cat("Token is recognized:\n", jwt, "\n")
24 } else {
25   cat("Token is not recognized. Check your .Renviron file.\n")
26 }
27
28
29 # Retrieve user information
30 user_info <- ieugwasr::user()
31
32 # Check if user information is retrieved
33 if (!is.null(user_info)) {
34   print("Token is working. User information:")
35   print(user_info)
36 } else {
37   print("Token is not working. Check your token and internet connection.")
38 }
39
40 #####
41
42 #Load the necessary packages
43
44 library(MRPracticals)
45 library(TwoSampleMR)
46 library(ieugwasr)
47 library(MRInstruments)
48 library(dplyr)
49 library(readxl)
50
51 vignette("MRBase")
52
53 #####

```

```

54
55 #Run the analysis from here onwards
56 setwd("C:/Users/tabe0010/OneDrive - Monash University/MR_backup_file/
    Articles/Evangelou/new_sbp_snps/stata/sbp_snps")
57 getwd()
58
59
60
61 sbp_data<-read.csv("all_sbp_snps.csv")
62 sbp_data[c("Chromosome", "Position")]<-do.call(rbind, strsplit(sbp_data$
    chrpos, ":"))
63 sbp_data$Chromosome<-as.numeric(sbp_data$Chromosome)
64 sbp_data$Position<-as.numeric(sbp_data$Position)
65
66
67 sbp_data_2<-sbp_data%>%select(Chromosome, Position, rsID, Beta_SBP, se_SBP
    , A1, A2, EAF_SBP, P_min, Trait, Beta_DBP, se_DBP, EAF_DBP)%>%mutate(id
    .exposure = "icbp_rep")
68 colnames(sbp_data_2)
69
70 colnames(sbp_data_2)<-c("chr.exposure", "pos.exposure", "SNP", "beta.
    exposure", "se.exposure", "effect_allele.exposure", "other_allele.
    exposure", "eaf.exposure", "pval.exposure", "exposure", "Beta_DBP", "se
    _DBP", "EAF_DBP", "id.exposure")
71 head(sbp_data_2)
72
73 sbp_data_2<-sbp_data_2[order(sbp_data_2$chr.exposure),]
74
75 clumped_sbp_data_2 <- clump_data(sbp_data_2,
76                                clump_kb = 10000, # Clumping window
                                (10,000 kb)
77                                clump_r2 = 0.001, # LD threshold (r2 <
                                0.001)
78                                pop = "EUR") # European LD reference
79
80 sbp_snplist<-clumped_sbp_data_2%>%select(SNP)
81
82 #sbp_effect_list<-clumped_sbp_data_2%>%select(SNP,effect_allele.exposure,
    beta.exposure)
83
84
85
86 write.csv(clumped_sbp_data_2, "sbp_exposure.csv", row.names = FALSE)
87 write.table(sbp_snplist, "C:/Users/tabe0010/OneDrive - Monash University/
    MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dx_data/sbp/data/
    sbp_snplist.txt", row.names = FALSE, col.names = FALSE, quote = FALSE, sep
    = " ")

```

We have clumped the SNPs, and the necessary genetic variants primarily associated with SBP were selected, totaling 187 SNPs. The **snps_snplist.txt** file was then used to select the required genetic variants from the UK Biobank imputed BGEN file.

4.2.4 Selecting the genetic variants from the BGEN file

PLINK2 was used via the **Swiss Army Knife** to select the necessary SNPs from the UK Biobank. The UKB RAP served as the platform to access the SNPs from the UK Biobank. Alternatively, a *bash* command was used on the local machine to run the process. First, log in with your UKB RAP credentials using the Command Prompt (on a Windows machine). Then, open Git Bash, select your project, and run the analysis. A job request was sent, and users were notified when the process was complete.

```
#####
#Run the following PLINK2 codes on Git Bash
#####

#Login through Command Prompt on your Windows machine
dx login

#Run the code below on Git Bash

dx select --level VIEW
#select your project
# make sure your sbp_snplist.txt file is uploaded to the UKB RAP.
#select the "instance type" you want: this makes sure you have enough computation
  power (CPU and GPU).
#In the command below, I put chromosome 1 to 22 to loop through all autosomal
  chromosomes just to show the code. But in actuality, I put two chromosomes at a
  time. This makes sure I have enough computational space and if there is any
  error, I could adjust the code.

# Loop over chromosomes 1 to 22 and process each one with the SNP list
run_merge=""
for chr in {1..22}; do
  run_merge+="cp /mnt/project/Bulk/Imputation/UKB\ imputation\ from\ genotype/
    ukb22828_c${chr}_b0_v3.bgen .; "
  run_merge+="cp /mnt/project/Bulk/Imputation/UKB\ imputation\ from\ genotype/
    ukb22828_c${chr}_b0_v3.sample .; "
  run_merge+="plink2 --bgen ukb22828_c${chr}_b0_v3.bgen ref-first --sample
    ukb22828_c${chr}_b0_v3.sample --extract sbp_snplist.txt --make-pgen --
    autosome-xy --out ukb22828_c${chr}_v3; "
done

dx run swiss-army-knife -iin="project-GpbQqBjJb7jb1vQjf8ZxVpVY:/SBP_data/sbp_txt/
  sbp_snplist.txt" -icmd="{run_merge}" --tag="Step1" --instance-type "
  mem1_ssd1_v2_x36" --destination="project-GpbQqBjJb7jb1vQjf8ZxVpVY:/SBP_data/
  sbp_geno_data/" --brief --yes

#####
#Run the following PLINK2 codes via Swiss Army Knife on UKB RAP platform
#####
#Make sure you have uploaded the sbp_merge_list.txt to UKB RAP file path.
#The merge list should have a sigle column list containing the following text, "
  ukb22828_cx_v3" (without the quotations). The column wil have 22 rows for each
  autosomal chromosomes. Replace "x" with 1-22.

#merging the pgen files
```



```

#Execute the command on Swiss Army Knife interface
#inputs are the plink files for the chromosomes and the txt file for the merging
  chromosomes
plink2 --pmerge-list sbp_merge_list.txt pfile --make-pgen --out
  ukb22828_c1_22_v3_sbp_merged

#Calculate the allele dosage
#creating a .raw file for participants with the number of effect allele (0, 1, or
  2) for each snp
# Input the for code below is the merged plink files (pfiles)

plink2 --pfile ukb22828_c1_22_v3_sbp_merged --export A --out ukb22828_sbp_alleles

#Calculate allele frequency

plink2 --pfile ukb22828_c1_22_v3_sbp_merged --freq --out ukb22828_sbp_allele_freq

```

The final PLINK output files, **ukb22828_sbp_alleles.raw** and **ukb22828_sbp_allele_freq.afreq**, contained the allele dosage and allele frequency for each of the 187 SNPs. These files were downloaded to the local machine and saved to the `$dx_data_sbp` file path.

4.2.5 Association of Genetic Variants with Quality-Adjusted Life Years

We then worked on the association between the genetic variants and QALYs. The QALYs were regressed on the allele dosages, adjusting for age, sex, and the first 10 genetic principal components to account for population stratification.

```

*****
*SNP-QALYs association
*****

import delimited "$dx_data_sbp\ukb22828_sbp_alleles.raw",clear
keep iid rs*
rename iid id_phe
merge 1:1 id_phe using "$stata_sbp_input\id_list.dta", keep(3) nogen
save "$stata_sbp_input\snp_alleles_sbp.dta", replace

use "$stata_sbp_output\part_3b.dta", clear
merge 1:1 id_phe using "$stata_sbp_input\snp_alleles_sbp.dta", keep(1 3) nogen
*gen imputation = .
gen snp = ""
gen effect_allele = ""
gen eaf = .
gen outcome = ""
gen beta = .
gen se = .
gen variance = .
gen p = .
gen n = .
local i = 1
local outcomes = "qaly_hes"

foreach outcome in `outcomes' {

    qui regress `outcome' rs* age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10

    foreach snp of varlist rs* {
        local snpx = substr("`snp'",1,length("`snp'")-2)
        *qui replace imputation = `imputation' in `i'
    }
}

```

```

        qui replace snp = "`snp'" in `i'
        qui replace outcome = "`outcome'" in `i'
        qui replace beta = _b[`snp'] if snp == "`snp'" & outcome == "`outcome'"
        qui replace se = _se[`snp'] if snp == "`snp'" & outcome == "`outcome'"
        qui sum `snp'
        qui replace eaf = r(mean)/2 if snp == "`snp'"
        local effect_allele = upper(substr("`snp'",length("`snp'"),1))
        qui replace effect_allele = "`effect_allele'" if snp == "`snp'"
        local i = `i'+1
    }

    *Ns
    qui sum `outcome'
    qui replace n = r(N) if outcome == "`outcome'"
}

keep snp-n
keep if snp != ""
qui replace variance = se^2
qui replace p = 2*normal(-abs(beta/se))
save "$stata_sbp_result\Results_snp_qalys.dta", replace
use "$stata_sbp_result\Results_snp_qalys.dta", clear
keep if outcome == "qaly_hes"
save "$stata_sbp_result\Results_snp_qaly_hes.dta"

*****
*merge with allele frequency data
*****
import delimited "$dx_data\ukb22828_sbp_allele_freq.afreq",clear
rename id snp
rename alt other_allele
merge 1:1 snp using "$stata_sbp_result\Results_snp_qaly_hes.dta"
drop chrom ref _merge* alt_freqs obs_ct
export delimited using "$sbp_snps\snp_qaly_hes.csv",replace

```

4.2.6 Data harmonisation

We now had the SNP-exposure and SNP-outcome association data. The next task was to harmonize these two datasets. For this, we continued working in the previous R environment. Data harmonization ensured that the effect alleles between the two datasets were properly aligned and removed any ambiguous SNPs. Ambiguous SNPs were those with A/T or C/G allele pairs, as they could create strand alignment issues due to their complementarity, making it difficult to determine the correct effect direction.

```
1 # we will continue working in R environment
2 # clumped_dbp_data_2 contains the SNP-SBP association
3
4 # Upload the SNP-outcome csv file: snp_qaly_hes.csv
5 qaly_hes_data<-read.csv("C:/Users/tabe0010/OneDrive - Monash University/MR
  _backup_file/Articles/Evangelou/new_sbp_snps/stata/sbp_snps/snp_qaly_
  hes.csv")
6
7 # rename the column names
8 colnames(qaly_hes_data)<-c("SNP", "other_allele.outcome", "effect_allele.
  outcome", "eaf.outcome", "outcome", "beta.outcome", "se.outcome", "
  variance.outcome", "pval.outcome", "samplesize.outcome")
9
10 #add outcome ID
11 qaly_hes_data$id.outcome = "ukb" # "ukb" added here just a reminder the
  SNP-outcome association is from UKB cohort.
12
13 #Harmonise the data
14 harmonise_data <- harmonise_data(
15   exposure_dat = clumped_sbp_data_2,
16   outcome_dat = qaly_hes_data
17 )
18
19 #save the file
20 write.csv(harmonise_data, "C:/Users/tabe0010/OneDrive - Monash University/
  MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/sbp_snps/snp_sbp_
  qaly_hes_harmonised.csv", row.names = FALSE)
21
22 #select the SNP-exposure associatoin column for the next analysis
23 sbp_effect_list<-harmonise_data%>%select(SNP,effect_allele.exposure,beta.
  exposure)
24
25 #Filter ambiguous to exclude for the next analysis
26 sbp_snplist_exclude<-harmonise_data%>%filter(mr_keep == "FALSE")%>%select(
  SNP)
27
28 #save both files in a .txt format
29 write.table(sbp_effect_list, "C:/Users/tabe0010/OneDrive - Monash
  University/MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dx_data
  /sbp/data/sbp_effect_list.txt",row.names = FALSE, col.names =FALSE,
  quote = FALSE,sep = " ")
30 write.table(sbp_snplist_exclude, "C:/Users/tabe0010/OneDrive - Monash
  University/MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dx_data
```

```
/sbp/data/sbp_snplist_exclude.txt", row.names = FALSE, col.names =FALSE  
,quote = FALSE,sep = " ")
```

We have harmonised our data and created .txt files for the next analysis.

4.2.7 Calculating polygenic risk scores

Using the **sbp_effect_list.txt** file and excluding ambiguous SNPs based on the **sbp_snplist_exclude.txt** file, the **polygenic risk score (PRS)** was calculated. PRS is a weighted sum of risk alleles across many genetic variants (SNPs). In our study, it was calculated across 181 SNPs after excluding the ambiguous ones.

The two .txt files were uploaded to the UKB RAP platform, and PLINK2 was used via the Swiss Army Knife to calculate the PRS.

```
# Upload sbp_effect_list.txt and sbp_snplist_exclude.txt files to UKB RAP  
# We will also use the merged plink files we created before  
# Use the two .txt files and the merged plink files to run the following code to  
  calculate the PRS via Swiss Army Knife  
  
plink2 --pfile ukb22828_c1_22_v3_sbp_merged --score sbp_effect_list.txt cols=+  
  scoresums --exclude sbp_snplist_exclude.txt --out ukb22828_sbp_prs
```

We have now calculated the PRS for SNPs effect on SBP. Download the **ukb22828_sbp_prs.sscore** file to your local machine and save them to the **\$dx_data_sbp** file path.

4.3 Step 3: Combining Phenotype and Genotype data

Now we will combine our phenotype data with PRS. The **part_3b.dta** contains the phenotype data for our cohort while the **ukb22828_sbp_prs.sscore** contains the PRS for each participant in the UK Biobank. The next stata line of codes will merge the two datasets. In addition prepare our data for the main and sensitivity analyses.

```
import delimited "$dx_data_sbp\ukb22828_sbp_prs.sscore", clear
gen id_phe = iid // IID: Individual ID
save "$dx_data_sbp\ukb22828_sbp_prs.sscore.dta", replace

use "$dx_data_sbp\ukb22828_sbp_prs.sscore.dta", clear
merge 1:1 id_phe using "$stata_sbp_output\part_3b.dta"
keep if _merge == 3
drop _merge*
rename score1_sum prs_sbp
drop fid iid allele_ct named_allele_dosage_sum score1_avg

* Create a genotype array type variable for the next analysis
gen geno_array = .
drop if n_22000_0_0 == .
replace geno_array = 0 if n_22000_0_0 < 0
replace geno_array = 1 if n_22000_0_0 > 0

*Estimate prs-free SBP for the next analysis
su prs_sbp
local mean = r(mean)
reg phe_sbp_adj prs_sbp
gen gf_sbp = phe_sbp_adj - _b[prs_sbp]*prs_sbp + _b[prs_sbp]*`mean', a(phe_sbp_adj)

*50 xtiles of prs-free SBP
qui xtile cat_gf_sbp = gf_sbp, nq(50)

save "$stata_sbp_output\part_4a.dta", replac
```

We have prepared our dataset for the main and sensitivity analyses.

4.4 Step 4: Main analysis

We estimated the causal association of SBP with QALY using **MR technique**[14, 11, 6]. Specifically, **two-stage least square (2SLS)** run by regressing the exposure variable (SBP) on the PRS for SBP at the first stage followed by regressing the outcome variable (QALYs) on the predicted SBP from the first stage. Age, sex, UK Biobank assessment centre, genotyping array, and the first 10 genetic principal components for population stratification were used as covariates in the model. F-statistics was used to assess for weak instrument bias. Outputs of the model interpreted as change in QALYs caused by a 1 mmHg increase in SBP over an average year of follow-up. For convenience, the final output was presented as percentage change in QALY per 10 mmHg increase in SBP.

We also performed **multivariable linear regression model** fitting the QALY outcome on SBP exposure data adjusting for age, sex, assessment centre, genotyping array and the first 10 genetic principal components. Then compared the estimate from 2SLS with the multivariable linear regression model and test for presence of **endogeneity (Hausman test)**[15, 16, 17]. A low p value in the Hausman test indicated difference in the estimates between the 2SLS and multivariable linear regression model.

```
use "$stata_sbp_output\part_4a.dta", clear

*****
*Main analysis
*****
*Create table
gen outcome = ""
gen type = ""
*gen imputation = .
gen n = .
gen beta = .
gen variance = .
gen se = .
gen double p = .
gen double p_endog = .
gen f_stat = .

local x = 1

foreach var in qaly_hes {
    dis "Outcome = `var'"

    *MR analysis
    ivreg2 `var' (phe_sbp_adj = prs_sbp) age i.sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 p
    > c10 i.centre i.geno_array, robust endog(phe_sbp_adj)

    matrix a = e(b)
    matrix b = e(V)
    local beta = a[1,1]
    local variance = b[1,1]

    local n = e(N)
    local f_stat = e(widstat)
    local p_endog = e(estatp)

    replace outcome = "`var'" in `x'
    replace type = "Main Analysis MR" in `x'
    foreach z in beta variance n p_endog f_stat {
        replace `z' = ``z'' in `x'
    }

    local x = `x' + 1

    *Linear regression
    reg `var' phe_sbp_adj age i.sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre i.
```

```

> geno_array

        matrix a = e(b)
        matrix b = e(V)
        local reg_beta = a[1,1]
        local reg_variance = b[1,1]

        local reg_n = e(N)

        *Linear regression estimates
        replace outcome = "`var'" in `x'
        replace type = "Multivariable Adjusted" in `x'
        *qui replace imputation = `j' in `x'
        foreach z in beta variance n {
            replace `z' = `reg_`z'' in `x'
        }

        local x = `x' + 1

    }

keep outcome-f_stat
keep if outcome != ""
replace outcome = "QALYs per year (with 240 comorbidities)" if outcome == "qaly_hes"
qui replace se = sqrt(var)
qui replace p = 2*normal(-abs(beta/se))
sort outcome type
*save "Result_sbp_exclusive_table.dta", replace
save "$stata_sbp_result\Result_sbp_table.dta", replace

*****
*SBP variations explained by PRS
*****

clear
set obs 1
gen r2 = .

        preserve
        use "$stata_sbp_output\part_4a.dta", clear

        corr phe_sbp_adj prs_sbp
        local r2 = r(rho)^2
        restore
        replace r2 = `r2' in 1

su r2, d
save "$stata_sbp_result\R2_value_sbp.dta", replace

```

4.5 Step 5: Sensitivity analyses

To test the robustness of the main analyses outcome, a number of sensitivity analyses were performed.

- **Untreated population:** Rerun the main analysis for the cohort without antihypertensive medications
- **Two-sample MR:** Run a number of summary level MR analyses
- **Sub-group analysis:** Stratified by age, sex, PRS-free SBP
- **Non-linear MR**
- **EQ-5D index from UK Biobank survey**

4.5.1 Untreated population

We repeated the main analysis excluding participants prescribed with antihypertensive medications.

```
use "$stata_sbp_output\part_4a.dta", clear
gen sbp_treat = .
replace sbp_treat = 1 if n_6153_0_0 == 2 | n_6153_0_1 == 2 | n_6153_0_2 == 2 | n_6153_0_3 == 2 | n_6177_0_0 == 2 | n_6177_0_1 == 2 | n_6177_0_2 == 2
replace sbp_treat = 2 if sbp_treat == .
label define treatlbl 1 "medication" 2 "no medication"
label values sbp_treat treatlbl
keep if sbp_treat == 2
*Create table
gen outcome = ""
gen type = ""
gen n = .
gen beta = .
gen variance = .
gen se = .
gen double p = .
gen double p_endog = .
gen f_stat = .
local x = 1
*local outcomes = "cost qaly qaly_cost_20k"

foreach var in qaly_hes {
    dis "Outcome = `var'"

    *MR analysis
    ivreg2 `var' (phe_sbp_adj = prs_sbp) age i.sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 p
    > c10 i.centre i.geno_array, robust endog(phe_sbp_adj)

    matrix a = e(b)
    matrix b = e(V)
    local beta = a[1,1]
    local variance = b[1,1]

    local n = e(N)
    local f_stat = e(widstat)
    local p_endog = e(estatp)

    replace outcome = "`var'" in `x'
    replace type = "Main Analysis MR" in `x'
}
```



```

        foreach z in beta variance n p_endog f_stat {
            replace `z' = ``z'' in `x'
        }

        local x = `x' + 1

        *Linear regression
        reg `var' phe_sbp_adj age i.sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre i.

> geno_array

        matrix a = e(b)
        matrix b = e(V)
        local reg_beta = a[1,1]
        local reg_variance = b[1,1]

        local reg_n = e(N)

        *Linear regression estimates
        replace outcome = ``var' in `x'
        replace type = "Multivariable Adjusted" in `x'
        foreach z in beta variance n {
            replace `z' = `reg_`z'' in `x'
        }

        local x = `x' + 1

    }

    keep outcome-f_stat
    keep if outcome != ""
    replace outcome = "QALYs per year (with 240 comorbidities)" if outcome == "qaly_hes"
    qui replace se = sqrt(var)
    qui replace p = 2*normal(-abs(beta/se))
    sort outcome type
    save "$stata_sbp_result\Result_sbp_exclusive_table_no_medication.dta", replace

```

4.5.2 Two-sample MR

To test the robustness of the main analysis estimate, 2SLS, we also performed a two-sample MR (i.e., summary level MR)[11]. The first sample data was sourced from the ICBP and replication meta-analyses which included the effect of the 181 SNPs on SBP[13]. For the second sample, we estimated the effect of the 181 SNPs on QALYs using the UK Biobank population by regressing the QALYs on the 181 SNPs allele dosage (the number of effect allele each participant has) adjusting for age, sex and the first 10 genetic principal components[8, 18]. Inverse variance weighting (IVW), MR Egger, weighted median and weighted mode analyses were performed using the summary level data[11].

- IVW
 - IVW method combines the individual ratio estimates from multiple genetic instruments (i.e., SNPs) into a single overall estimate, using a weighted average approach, where the weights are the inverse of the variance of the SNP-outcome association[11].
 - Presence of heterogeneity among the instruments was tested using Cochran’s Q statistics, and a statistically significant p-value indicates evidence of heterogeneity[19].
 - To quantify the degree of heterogeneity, I^2 was also reported, with 0% indicating no observed heterogeneity, $0\% < I^2 \leq 25\%$ representing low heterogeneity, $25\% < I^2 \leq 50\%$ indicating moderate heterogeneity, $50\% < I^2 \leq 75\%$ denoting substantial heterogeneity, and $I^2 > 75\%$ reflecting considerable heterogeneity[19].
- MR Egger
 - The MR Egger method used to test for and account for directional pleiotropy, a potential source of bias in causal estimates[11]. It’s an extension of the IVW approach, with additional flexibility to model horizontal (i.e., directional) pleiotropy, which occurs when genetic variants (SNPs) influence the outcome through pathways other than the exposure of interest[11]. The main difference between MR Egger and IVW is that MR Egger includes an intercept term in the regression model. This allows for testing and adjusting for directional pleiotropy[11].
- Weighted median
 - The weighted median method is useful when some of the genetic variants (SNPs) used as instrumental variables may be invalid due to pleiotropy or other issues[11]. This method gives consistent causal estimate if at least 50
- Weighted mode
 - The weighted mode method estimates the causal effect by finding the most common (modal) value of SNP-specific causal estimates, with more weight given to SNPs that are more precise[11]. This method is robust to invalid instruments and can produce valid estimates even when most instruments are invalid[11].

```
*****
*Prepare the data
*****
*import the harmonised data
```

```

import delimited "$sbp_snps\snp_sbp_qaly_harmonised.csv", clear
    rename betaexposure beta_exposure
    rename seexposure se_exposure
    rename betaoutcome beta
    rename seoutcome se
    rename samplesizeoutcome n

*Make all the exposure betas positive
    replace beta = -beta if beta_exposure < 0
    replace beta_exposure = -beta_exposure if beta_exposure < 0
    count if beta_exposure < 0
    sort outcome snp
    keep snp beta_exposure se_exposure beta se outcome pvaloutcome pval exposure n

save "$stata_sbp_output\MR_data_sbp.dta", replace

*****
*Run the Two sample MR analyses
*****
use "$stata_sbp_output\MR_data_sbp.dta", clear

gen out = ""
gen ivw = .
gen ivw_se = .
gen ivw_p = .
gen egger_slope = .
gen egger_slope_se = .
gen egger_slope_p = .
gen egger_cons = .
gen egger_cons_se = .
gen egger_cons_p = .
gen double heterogeneity_p = .
gen median = .
gen median_se = .
gen median_p = .
gen mode = .
gen mode_se = .
gen mode_p = .

qui levelsof outcome, local(outcome)
local i = 1

foreach out in `outcome' {
    *MR robust takes the outcome first
    replace out = "`out'" in `i'
    mregger beta beta_exposure [aw=1/(se^2)] if outcome == "`out'", ivw heterogi
    if !_rc {
        replace heterogeneity_p = r(pval) in `i'
    }
    else {
        mregger beta beta_exposure [aw=1/(se^2)] if outcome == "`out'", ivw
    }
    replace ivw = _b[beta_exposure] in `i'
    replace ivw_se = _se[beta_exposure] in `i'

    mregger beta beta_exposure [aw=1/(se^2)] if outcome == "`out'"
    replace egger_slope = _b[slope] in `i'
    replace egger_slope_se = _se[slope] in `i'
    replace egger_cons = _b[_cons] in `i'
    replace egger_cons_se = _se[_cons] in `i'
    mrmedian beta se beta_exposure se_exposure if outcome == "`out'"
    replace median = _b[beta] in `i'
    replace median_se = _se[beta] in `i'
    mrmodal beta se beta_exposure se_exposure if outcome == "`out'"
    replace mode = _b[beta] in `i'
    replace mode_se = _se[beta] in `i'

    local i = `i' + 1
}

foreach var of varlist ivw egger_slope egger_cons median mode {

```

```

        replace `var`_p = 2*normal(-abs(`var`/`var`_se))
    }
    keep out-mode_p
    keep if out != ""
    rename out outcome
    *Make things look better
    replace outcome = "QALYs per year (HES only)" if outcome == "qaly_hes"
    sort outcome
    save "$stata_sbp_result\Results_table_sensitivity_sbp.dta", replace

```

4.5.3 Sub-group analysis

Sub-group analysis was performed by rerunning the main MR stratified by age categories (< 50 years, 50-54 years, 55-59 years, 60-64 years, and 65+ years), sex (Male and Female), and PRS-free SBP categories (<120 mmHg, 120-139 mmHg and 140+ mmHg). PRS-free SBP was estimated by first regressing observed SBP (i.e., mean SBP) on the PRS for SBP and then predicting each participant's SBP as if they had the average PRS for SBP[1].

```
*****
*Sex-, SBP category- and Age-Specific Analyses
*****
{
  use "$stata_sbp_output\part_4a.dta", clear

  *Mark participants depending on sex, age and SBP level on genetic-free SBP
  gen all = 1

  *replace age = age+38
  gen age_cat1 = 1 if age < 50
  gen age_cat2 = 1 if age < 55 & age >= 50
  gen age_cat3 = 1 if age < 60 & age >= 55
  gen age_cat4 = 1 if age < 65 & age >= 60
  gen age_cat5 = 1 if age >= 65

  gen sbp_cat1 = 1 if gf_sbp < 120
  gen sbp_cat2 = 1 if gf_sbp >= 120 & gf_sbp < 140
  gen sbp_cat3 = 1 if gf_sbp >= 140

  *Create table
  gen outcome = ""
  gen type = ""

  *gen imputation = .
  gen sensitivity = ""

  foreach sex in all male female {
    gen n_`sex' = .
    gen beta_`sex' = .
    gen variance_`sex' = .
    gen se_`sex' = .
    gen double p_`sex' = .
    gen double p_endog_`sex' = .
    gen f_stat_`sex' = .
  }

  foreach sex in all male female {
    local x = 1
    foreach sens in all age_cat1 age_cat2 age_cat3 age_cat4 age_cat5 sbp_cat1 sbp_cat2 sbp_cat3 {
      foreach var in qaly_hes {
        if "`sex'" == "all" {
          qui ivreg2 `var' (phe_sbp_adj = prs_sbp) age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.
          > centre geno_array if `sens' == 1, robust endog(phe_sbp_adj)
        }
        else if "`sex'" == "female" {
          qui ivreg2 `var' (phe_sbp_adj = prs_sbp) age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.
          > centre geno_array if `sens' == 1 & sex == 0, robust endog(phe_sbp_adj)
        }
        else {
          qui ivreg2 `var' (phe_sbp_adj = prs_sbp) age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.
          > centre geno_array if `sens' == 1 & sex == 1, robust endog(phe_sbp_adj)
        }
      }

      matrix a = e(b)
      matrix b = e(V)
      local beta = a[1,1]
      local variance = b[1,1]

      local n = e(N)
      local f_stat = e(widstat)
      local p_endog = e(estatp)

      *MR estimates
    }
  }
}
```

```

qui replace outcome = "`var'" in `x'
qui replace type = "Main Analysis MR" in `x'
qui replace sensitivity = "`sens'" in `x'
foreach z in beta variance n p_endog f_stat {
  qui replace `z'`sex' = ``z'' in `x'
}

local x = `x' + 1

*Linear regression
if "`sex'" == "all" {
  qui reg `var' phe_sbp_adj age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre geno_arr
  > ay if `sens' == 1
}
else if "`sex'" == "female" {
  qui reg `var' phe_sbp_adj age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre geno_arr
  > ay if `sens' == 1 & sex == 0
}
else {
  qui reg `var' phe_sbp_adj age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre geno_arr
  > ay if `sens' == 1 & sex == 1
}

matrix a = e(b)
matrix b = e(V)
local beta = a[1,1]
local variance = b[1,1]
local n = e(N)

*Linear regression estimates
qui replace outcome = "`var'" in `x'
qui replace type = "Multivariable Adjusted" in `x'
qui replace sensitivity = "`sens'" in `x'
foreach z in beta variance n {
  qui replace `z'`sex' = ``z'' in `x'
}

local x = `x' + 1
}
}

keep outcome-f_stat-female
drop if outcome == ""
sort outcome type sensitivity
foreach sex in all male female {
  qui replace se_`sex' = sqrt(variance_`sex')
  qui replace p_`sex' = 2*normal(-abs(beta_`sex'/se_`sex'))
}

*Altman-Bland/Fisher tests
*Male-female
gen double p_sex = 2*normal(-abs((beta_female-beta_male)/sqrt(se_female^2+se_male^2)))

save "$stata_sbp_result\Results_sbp_subgroup.dta", replace
}

```

4.5.4 Non-linear MR

Non-linear MR was performed by running the main MR within fifty quantiles of PRS-free SBP, estimating quantile specific **local average causal effects**. These local average estimates were used in the **variance weighted least squares (VWLS)** models to determine whether there was a stable or incremental change in the effect of SBP on QALYs as SBP increased. Both linear and cubic models (with respect to the mean PRS-free SBP in each quantile) used to describe the shape of the effect of the increase in SBP over the range of PRS-free SBP values. We followed the next steps to execute the non-linear MR analysis.

4.5.4.1 Rerun the main MR within the fifty quantiles of PRS-free SBP

```
use "$stata_sbp_output\part_4a.dta", clear
gen outcome = ""
gen type = ""
gen gf_sbp_cat = .
foreach sex in all male female {
    local x = 1
    gen n_`sex' = .
    gen beta_`sex' = .
    gen variance_`sex' = .
    gen se_`sex' = .
    gen double p_`sex' = .
    foreach var in qaly_hes {
        forvalues k = 1/50 {
            display "working on `k' of 50"
            if "`sex'" == "all" {
                qui ivreg2 `var' (phe_sbp_adj = prs_sbp) age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centres
                > geno_array if cat_gf_sbp == `k', robust
            }
            else if "`sex'" == "female" {
                qui ivreg2 `var' (phe_sbp_adj = prs_sbp) age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centres
                > geno_array if sex == 0 & cat_gf_sbp == `k', robust
            }
            else {
                qui ivreg2 `var' (phe_sbp_adj = prs_sbp) age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centres
                > geno_array if sex == 1 & cat_gf_sbp == `k', robust
            }
            matrix a = e(b)
            matrix b = e(V)
            local beta = a[1,1]
            local variance = b[1,1]
            local se = sqrt(b[1,1])
            local p = 2*normal(-abs(`beta'/`se'))
            local n = e(N)
            *MR estimates
            qui replace outcome = "`var'" in `x'
            qui replace type = "Main Analysis MR" in `x'
            qui replace gf_sbp_cat = `k' in `x'
            foreach z in beta se variance p n {
                qui replace `z'_`sex' = ``z'' in `x'
            }
            local x = `x' + 1
            *Linear regression
            if "`sex'" == "all" {
                qui reg `var' phe_sbp_adj age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centres geno_array if ca
                > t_gf_sbp == `k'
            }
        }
    }
}
```

```

else if "`sex'" == "female" {
qui reg `var' phe_sbp_adj age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre geno_array if se
> x == 0 & cat_gf_sbp == `k'
}

else {
qui reg `var' phe_sbp_adj age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre geno_array if se
> x == 1 & cat_gf_sbp == `k'
}

matrix a = e(b)
matrix b = e(V)
local beta = a[1,1]
local variance = b[1,1]
local se = sqrt(b[1,1])
local p = 2*normal(-abs(`beta'/'se'))
local n = e(N)

*Linear regression estimates
qui replace outcome = "`var'" in `x'
qui replace type = "Multivariable Adjusted" in `x'
qui replace gf_sbp_cat = `k' in `x'
*qui replace imputation = `j' in `x'

foreach z in beta se variance p n {
qui replace `z'_'sex' = ``z'' in `x'
}

local x = `x' + 1
}
}
}

keep outcome-p_female
drop if outcome == ""
*sort outcome imputation type gf_sbp_cat
sort outcome type gf_sbp_cat

*Altman-Bland/Fisher tests
*Male-female
gen double p_sex = 2*normal(-abs((beta_female-beta_male)/sqrt(se_female^2+se_male^2)))
save "$stata_sbp_result\Results_sensitivity_sbp_nl_50q.dta", replace

```

4.5.4.2 Sub-group analysis mainly stratified by PRS-free SBP

For each PRS-free SBP category, we further stratified by age-group and sex.

```

{
use "$stata_sbp_output\part_4a.dta", clear

su gf_sbp
gen sbp_cat =.

replace sbp_cat = 1 if gf_sbp <120
replace sbp_cat = 2 if gf_sbp >= 120 & gf_sbp <140
replace sbp_cat = 3 if gf_sbp >= 140

label define sbplbl 1 "<120 mmHg" 2 "120-139 mmHg" 3 "140+ mmHg"
label values sbp_cat sbplbl

*keep if sbp_cat == 1
*Mark participants depending on sex, age and overweight on genetic-free SBP
gen all = 1

gen age_cat1 = 1 if age < 50
gen age_cat2 = 1 if age < 55 & age >= 50
gen age_cat3 = 1 if age < 60 & age >= 55
gen age_cat4 = 1 if age < 65 & age >= 60
gen age_cat5 = 1 if age >=65

gen sbp_cat1 = 1 if gf_sbp <120
gen sbp_cat2 = 1 if gf_sbp >=120 & gf_sbp <=139
gen sbp_cat3 = 1 if gf_sbp >=140

*Create table

```



```

gen outcome = ""
gen type = ""
gen sbp_category = ""

*gen imputation = .
gen sensitivity = ""

foreach sex in all male female {
    gen n_`sex' = .
    gen beta_`sex' = .
    gen variance_`sex' = .
    gen se_`sex' = .
    gen double p_`sex' = .
    gen double p_endog_`sex' = .
    gen f_stat_`sex' = .
}

foreach sex in all male female {
    local x = 1
    levelsof sbp_cat, local(sbp)

    foreach g in `sbp' {
        foreach sens in all age_cat1 age_cat2 age_cat3 age_cat4 age_cat5 {
            foreach var in qaly_hes {
                display "Group = `g', Sex = `sex', Sensitivity `sens', and Outcome = `var'"

                if "`sex'" == "all" {
                    qui ivreg2 `var' (phe_sbp_adj = prs_sbp) age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.
> centre geno_array if `sens' == 1 &
                    sbp_cat == `g', robust endog(phe_sbp_adj)
                }
                else if "`sex'" == "female" {
                    qui ivreg2 `var' (phe_sbp_adj = prs_sbp) age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.cent
> e geno_array if `sens' == 1 &
                    sex == 0 & sbp_cat == `g', robust endog(phe_sbp_adj)
                }
                else {
                    qui ivreg2 `var' (phe_sbp_adj = prs_sbp) age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre
> geno_array if `sens' == 1 &
                    sex == 1 & sbp_cat == `g', robust endog(phe_sbp_adj)
                }
            }

            matrix a = e(b)
            matrix b = e(V)
            local beta = a[1,1]
            local variance = b[1,1]

            local n = e(N)
            local f_stat = e(widstat)
            local p_endog = e(estatp)

            *MR estimates
            qui replace outcome = "`var'" in `x'
            qui replace type = "Main Analysis MR" in `x'
            qui replace sbp_category = "`g'" in `x'
            qui replace sensitivity = "`sens'" in `x'

            foreach z in beta variance n p_endog f_stat {
                qui replace `z'_`sex' = ``z'' in `x'
            }

            local x = `x' + 1

            *Linear regression
            if "`sex'" == "all" {
                qui reg `var' phe_sbp_adj age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre geno_array if `
> sens' == 1 &
                sbp_cat == `g'
            }
            else if "`sex'" == "female" {
                qui reg `var' phe_sbp_adj age pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre geno_array if `sens'
> == 1 &
                sex == 0 & sbp_cat == `g'
            }
            else {

```

```

qui reg `var` phe_sbp_adj age pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre geno_array if `sens`
> == 1 &
sex == 1 & sbp_cat == `g`
}

matrix a = e(b)
matrix b = e(V)
local beta = a[1,1]
local variance = b[1,1]
local n = e(N)

*Linear regression estimates
qui replace outcome = "`var`" in `x`
qui replace type = "Multivariable Adjusted" in `x`
qui replace sbp_category = "`g`" in `x`
qui replace sensitivity = "`sens`" in `x`

foreach z in beta variance n {
  qui replace `z`_`sex` = ``z`` in `x`
}

local x = `x` + 1

}
}
}
}

keep outcome-f_stat_female
drop if outcome == ""
sort outcome type sbp_category sensitivity
foreach sex in all male female {
  qui replace se_`sex` = sqrt(variance_`sex`)
  qui replace p_`sex` = 2*normal(-abs(beta_`sex`/se_`sex`))
}

save "$stata_sbp_result\Results_sbp_subgroup_1.dta", replace
}

```

4.5.4.3 Select observations for all participants in each SBP

```

use "$stata_sbp_result\Results_sbp_subgroup_1.dta", clear
replace sbp_category = "Normal" if sbp_category == "1"
replace sbp_category = "Pre_hypertension" if sbp_category == "2"
replace sbp_category = "Hypertension" if sbp_category == "3"
keep if sensitivity == "all"
drop if type == "Multivariable Adjusted"
keep outcome beta* sbp_category

replace sbp_category = lower(sbp_category)
foreach var of varlist beta* {
  rename `var` `var`_
}

reshape wide beta*, i(outcome) j(sbp_category) string

save "$stata_sbp_result\Sensitivity_analysis_graph_data.dta", replace

```

4.5.4.4 Estimate the mean SBP for each quantile

```

use "$stata_sbp_output\Part_4a.dta", clear
su gf_sbp
gen touse = 1
gen gf_sbp_cat = _n in 1/50
foreach sex in all male female {
  gen gf_sbp_mean_`sex` = .
  replace touse = 1
  if "`sex`" == "male" {
    replace touse = . if sex == 0
  }
}

```

```

    }
    if "`sex'" == "female" {
        replace touse = . if sex == 1
    }
    forvalues i = 1/50 {
        qui su gf_sbp if cat_gf_sbp == `i' & touse == 1
        qui replace gf_sbp_mean_`sex' = r(mean) in `i'
    }
}
keep if gf_sbp_cat != .
keep gf_sbp_cat gf_sbp_mean*

    histogram gf_sbp_mean_all, bin(20) normal percent color(red) xlabel(100(5)200) title("SBP Di
> stribution")
    count if gf_sbp_mean_all >180

save "$stata_sbp_result\Sensitivity_analysis_PRS_free_SBP_data.dta", replace

```

4.5.4.5 Merge the previous datasets for VWLS analyses and plots

We run the analyses here

```

use "$stata_sbp_result\Results_sensitivity_sbp_nl_50q.dta", clear
merge m:1 gf_sbp_cat using "$stata_sbp_result\Sensitivity_analysis_PRS_free_SBP_data.dta", nogen
sort outcome type gf_sbp_cat
merge m:1 outcome using "$stata_sbp_result\Sensitivity_analysis_graph_data.dta", nogen

replace outcome = "QALYs per year" if outcome == "qaly_hes"
*Meta-regress using linear regression
encode outcome, gen(outcome2)
encode type, gen(type2)
foreach sex in all male female {
    gen l95_`sex' = beta_`sex' - se_`sex'*1.96
    gen u95_`sex' = beta_`sex' + se_`sex'*1.96

    gen gf_sbp_mean_`sex'_2 = gf_sbp_mean_`sex'`^2
    gen gf_sbp_mean_`sex'_3 = gf_sbp_mean_`sex'`^3

    gen b0_`sex' = .
    gen b1_`sex' = .
    gen b2_`sex' = .
    gen b3_`sex' = .

    gen se0_`sex' = .
    gen se1_`sex' = .
    gen se2_`sex' = .
    gen se3_`sex' = .
}
local obs = c(N)+1
local obs2 = c(N) + 10000
set obs `obs2'
qui gen n = _n
gen test = 1 in `obs'/'`obs2'
foreach sex in all male female {
    qui replace gf_sbp_mean_`sex' = 100 + (n-`obs')/100 in `obs'/'`obs2'
    qui replace gf_sbp_mean_`sex'_2 = (100 + (n-`obs')/100)^2 in `obs'/'`obs2'
    qui replace gf_sbp_mean_`sex'_3 = (100 + (n-`obs')/100)^3 in `obs'/'`obs2'
}

forvalues outcome = 1/1 {
    qui su n if outcome2 == `outcome'
    local outcome_label = outcome[r(min)]

    foreach sex in all male female {

        local xtitle = "SBP (mmHg)"

```

```

        if "`sex'" == "male" {
            local xtitle = "SBP (mmHg) (male)"
        }
        if "`sex'" == "female" {
            local xtitle = "SBP (mmHg) (female)"
        }
    }

    vwls beta_`sex' gf_sbp_mean_`sex' gf_sbp_mean_`sex'_2 gf_sbp_mean_`sex'_3 if outcome
> 2 == `outcome' & type2 == 1, sd(se_`sex')
    local b1 = _b[gf_sbp_mean_`sex']
    local b2 = _b[gf_sbp_mean_`sex'_2]
    local b3 = _b[gf_sbp_mean_`sex'_3]
    local cons = _b[_cons]

    foreach x in normal pre_hypertension hypertension {
        qui su n if outcome2 == `outcome'
        local `x' = beta_`sex'_'x'[r(min)]
    }

    scatter beta_`sex' gf_sbp_mean_`sex' if outcome2 == `outcome' & type == "Main Analys
> is MR" & gf_sbp_mean_`sex' <=160 || rcap l95_`sex' u95_`sex' gf_sbp_mean_`sex' if outcome2 == `ou
> tcome' & type == "Main Analysis MR" & gf_sbp_mean_`sex' <=160 || ///
        function y = `cons' + `b1'*x + `b2'*x^2 + `b3'*x^3, range(100 160) || function y = `
> normal', range (100 119.9) lcolor(navy) lpattern(dash) || ///
        function y = `pre_hypertension', range(120 139.9) lcolor(navy) lpattern(dash) || fun
> ction y = `hypertension', range(140 160) lcolor(navy) lpattern(dash) ///
        xtitle("`xtitle'") ytitle("Effect of a unit increase in SBP on `outcome_label'", siz
> e(small)) legend(off) xscale(range(100 160))

    graph export "$plot_png\`outcome_label' [`sex'].png", as(png) width(1200) replace

        *Just trend line and 95% CI
        *Maybe create a small dataset with fixed values to remove

        predict x if test == 1
        predict x_se if test == 1, stdp
        gen x_l95 = x - 1.96*x_se
        gen x_u95 = x + 1.96*x_se

        twoway rarea x_l95 x_u95 gf_sbp_mean_`sex' if test == 1 & gf_sbp_mean_`sex' <=160, l
> color(green%50) color(green%50) || line x gf_sbp_mean_`sex' if test == 1 & gf_sbp_mean_`sex' <=160
> , lcolor(dkgreen) ///
        legend(off) xtitle("`xtitle'") ytitle("Effect of a unit increase in SBP on `outcome_
> label'", size(small)) ///
        plotregion(fcolor(white)) graphregion(fcolor(white)) xline(120,lcolor(maroon%80) lpa
> ttern(-)) xline(140,lcolor(maroon%80) lpattern(-)) yline(0,lcolor(gs4))

        graph export "$plot_png\`outcome_label' [`sex'] v2.png", as(png) width(1200) replace

        drop x-x_u95
    }
}

*Figures
su n if outcome2 == 1
local outcome_label = outcome[r(min)]
local xtitle = "SBP (mmHg)"

vwls beta_all gf_sbp_mean_all gf_sbp_mean_all_2 gf_sbp_mean_all_3 if outcome2 == 1 & type2 == 1, sd(
> se_all)
local b1 = _b[gf_sbp_mean_all]
local b2 = _b[gf_sbp_mean_all_2]
local b3 = _b[gf_sbp_mean_all_3]
local cons = _b[_cons]

predict x if test == 1
predict x_se if test == 1, stdp
gen x_l95 = x - 1.96*x_se
gen x_u95 = x + 1.96*x_se

twoway rarea x_l95 x_u95 gf_sbp_mean_all if test == 1 & gf_sbp_mean_all<=160, lcolor(green%50) color

```

```

> (green%50) || line x gf_sbp_mean_all if test == 1 & gf_sbp_mean_all<=160, lcolor(dkgreen) ///
legend(off) xtitle("`xtitle'") ytitle("Effect of a unit increase in SBP on `outcome_label'" "{&uarr}"
> SBP leads to {&darr} QALYs {&uarr} SBP leads to {&uarr} QALYs", size(small)) ///
plotregion(fcolor(white)) graphregion(fcolor(white)) xline(120,lcolor(maroon%80) lpattern(-)) xline(
> 140,lcolor(maroon%80) lpattern(-)) yline(0,lcolor(gs4)) ///
xscale(range (100 160))

graph export "$plot_png\nl_MR.png", as(png) width(1200) replace
drop x-x_u95

*Meta-regression
forvalues outcome = 1/2 {
    forvalues type = 1/2 {
        foreach sex in all male female {
            qui vwls beta_`sex' gf_sbp_mean_`sex' gf_sbp_mean_`sex'_2 gf_sbp_mean_`sex'_
> 3 if outcome2 == `outcome' & type2 == `type', sd(se_`sex')

            qui replace b1_`sex' = _b[gf_sbp_mean_`sex'] if outcome2 == `outcome' & type
> 2 == `type'
            qui replace b2_`sex' = _b[gf_sbp_mean_`sex'_2] if outcome2 == `outcome' & ty
> pe2 == `type'
            qui replace b3_`sex' = _b[gf_sbp_mean_`sex'_3] if outcome2 == `outcome' & ty
> pe2 == `type'

            qui replace se1_`sex' = _se[gf_sbp_mean_`sex'] if outcome2 == `outcome' & ty
> pe2 == `type'
            qui replace se2_`sex' = _se[gf_sbp_mean_`sex'_2] if outcome2 == `outcome' &
            qui replace se3_`sex' = _se[gf_sbp_mean_`sex'_3] if outcome2 == `outcome' &
> type2 == `type'
            qui replace se3_`sex' = _se[gf_sbp_mean_`sex'_3] if outcome2 == `outcome' &
> type2 == `type'

            qui vwls beta_`sex' gf_sbp_mean_`sex' if outcome2 == `outcome' & type2 == `t
> ype', sd(se_`sex')

            qui replace b0_`sex' = _b[gf_sbp_mean_`sex'] if outcome2 == `outcome' & type
> 2 == `type'
            qui replace se0_`sex' = _se[gf_sbp_mean_`sex'] if outcome2 == `outcome' & ty
> pe2 == `type'
        }
    }
}

keep outcome type b0* b1* b2* b3* se0* se1* se2* se3*
duplicates drop

foreach sex in all male female {
    forvalues i = 0/3 {
        qui gen p`i'_`sex' = 2*normal(-abs(b`i'_`sex'/se`i'_`sex'))
    }
}

order outcome type *0* *_all *_male *_female
save "$stata_sbp_result\Results_table_sensitivity_Stata_vwls_raw.dta", replace
keep outcome type *0*
foreach x in b se p {
    foreach sex in all male female {
        rename `x'0_`sex' `x'1_`sex'
    }
}

gen model = "Linear"

save "$stata_sbp_result\Results_table_sensitivity_Stata_vwls_append.dta", replace
use "$stata_sbp_result\Results_table_sensitivity_Stata_vwls_append.dta", clear
append using "$stata_sbp_result\Results_table_sensitivity_Stata_vwls_raw.dta"
drop *0*
replace model = "Cubic" if model == ""
order b1* se1* p1* b2* se2* p2* b3* se3* p3*
order outcome type model *_all *_male *_female
save "$stata_sbp_result\Results_table_sensitivity_Stata_vwls.dta", replace

```

4.5.5 EQ-5D index from UK Biobank survey

We performed an additional sensitivity analysis using the EQ-5D-5L survey data from the UK Biobank[20, 21] to assess whether the direction of the effect observed in the main analysis would be similar when using direct EQ-5D-5L data from UK Biobank. Web-based EQ-5D-5L questionnaires were administered to the UK Biobank participants as part of the chronic pain (administered in 2019–20) and mental well-being (administered in 2022–23) surveys[20, 21]. There were 167,111 participants responded to the chronic pain survey[20] and 169,537 participants responded to the mental well-being survey[21]. We calculated the EQ-5D index using the UK tariffs (i.e., value set) for each survey[10]. Once the EQ-5D-indexes were calculated, we took the average of both EQ-5D-indexes for participants who had EQ-5D data for both surveys (124,830). The remaining participants had EQ-5D data for either chronic pain survey (42,281) or mental well-being survey (44,707); hence, the average EQ-5D index was not calculated. In total, there were 211,818 participants with EQ-5D index data from the combined survey. Of these participants, 128,635 had met the initial inclusion criteria and were included in this sensitivity analysis.

We rerun 2SLS analysis by regressing the SBP trait on the genetic risk scores for SBP on the first stage following by fitting the EQ-5D index on the predicted SBP values from the first stage. We adjusted for age, sex, UK Biobank assessment centre, genotyping array, and the first 10 genetic principal components for population stratification. Outputs from the 2SLS model were presented as the effect of 10 mmHg increase in SBP on percentage change in EQ-5D index.

```
*UKB Utilities
use "$stata_sbp_output\part_4a.dta", clear

gen ukb_utility = .
replace ukb_utility = (EQindex_1 + EQindex)/2 if EQindex !=. & EQindex_1 !=.
replace ukb_utility = EQindex if EQindex !=. & EQindex_1 ==.
replace ukb_utility = EQindex_1 if EQindex ==. & EQindex_1 !=.

*Create table
gen outcome = ""
gen type = ""

*gen imputation = .
gen n = .
gen beta = .
gen variance = .
gen se = .
gen double p = .
gen double p_endog = .
gen f_stat = .

*Number of imputations
*local m = 100

local x = 1
*local outcomes = "cost qaly qaly_cost_20k"

foreach var in ukb_utility {
    dis "Outcome = `var'"

    *MR analysis
    ivreg2 `var' (phe_sbp_adj = prs_sbp) age i.sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 p
    > c10 i.centre i.geno_array, robust endog(phe_sbp_adj)

    matrix a = e(b)
    matrix b = e(V)
    local beta = a[1,1]
    local variance = b[1,1]

    local n = e(N)
```

```

        local f_stat = e(widstat)
        local p_endog = e(estatp)

        replace outcome = "`var'" in `x'
        replace type = "Main Analysis MR" in `x'
        *qui replace imputation = `j' in `x'
        foreach z in beta variance n p_endog f_stat {
            replace `z' = ``z'' in `x'
        }

        local x = `x' + 1

        *Linear regression
        reg `var' phe_sbp_adj age i.sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre i.

> geno_array

        matrix a = e(b)
        matrix b = e(V)
        local reg_beta = a[1,1]
        local reg_variance = b[1,1]

        local reg_n = e(N)

        *Linear regression estimates
        replace outcome = "`var'" in `x'
        replace type = "Multivariable Adjusted" in `x'
        foreach z in beta variance n {
            replace `z' = `reg_`z'' in `x'
        }

        local x = `x' + 1

    }

    keep outcome-f_stat
    keep if outcome != ""
    replace outcome = "EQ-5D-index" if outcome == "ukb_utility"
    qui replace se = sqrt(var)
    qui replace p = 2*normal(-abs(beta/se))
    sort outcome type
    save "$stata_sbp_result\Result_sbp_table_ukb_utility.dta", replace

```

4.6 Step 6: Secondary analysis

4.6.1 Multivariable Mendelian Randomisation

We performed a multivariable MR(MVMR), an extension of a simple MR that accounts for multiple exposures that are potentially related and may influence the outcome of interest[22, 23, 24]. For this analysis, SBP and DBP were considered as exposures. The mean DBP calculated the same way as the mean SBP described above (see exposure and covariate section). For participants who reported taking antihypertensive medications, a 10 mmHg add to the mean DBP measurement[13, 5, 9, 25].

The GWAS study reported sentinel SNPs association with primary and secondary traits[13]. For the 187 sentinel SNPs primarily associated with SBP (after LD clumping), we also identified association with DBP as a secondary trait. Similarly, we also identified 208 sentinel SNPs (after LD clumping) primarily associated with DBP that were also linked to SBP. The combined 395 SNPs were candidates for the MVMR analysis. After the exclusion of ambiguous and missing effect size SNPs, 384 and 382 SNPs were used to construct PRS for SBP and DBP, respectively. Effect estimates for the SNPs were sourced from either ICBP or replication meta-analysis.

4.6.1.1 Working on the DBP data

We had already worked on known and validated SNPs, known but validated (i.e., replicated) for the first time, and novel SNPs associated with the DBP trait. We continued to work on these SNPs.

These SNPs were clumped using the **TwoSampleMR** and **ieugwasr** packages in the **R environment**.

4.6.1.2 Clumping SNPs in LD

The previous stata codes provided csv files for the next step, LD clumping using **TwoSampleMR** and **ieugwas** packages. First, we need to set up an Application Programming Interface (API) to access the **IEU GWAS** database. Then we clump the SNPs in LD.

```
1 # We will continue working on in the R environment.
2 # We clump the SNPs in LD.
3
4 dbp_data<-read.csv("C:/Users/tabe0010/OneDrive - Monash University/MR_
    backup_file/Articles/Evangelou/new_sbpmnps/stata/dbp_snpms/all_dbp_snpms
    .csv")
5 dbp_data[c("Chromosome", "Position")]<-do.call(rbind, strsplit(dbp_data$
    chrpos, ":"))
6 dbp_data$Chromosome<-as.numeric(dbp_data$Chromosome)
7 dbp_data$Position<-as.numeric(dbp_data$Position)
8
9 dbp_data_2<-dbp_data%>%select(Chromosome, Position, rsID, Beta_DBP, se_DBP
    , A1, A2, EAF_DBP, P_min, Trait, Beta_SBP, se_SBP, EAF_SBP)%>%mutate(id
    .exposure = "icbp_rep")
10 colnames(dbp_data_2)
11
12 colnames(dbp_data_2)<-c("chr.exposure", "pos.exposure", "SNP", "beta.
    exposure", "se.exposure", "effect_allele.exposure", "other_allele.
```



```

    exposure", "eaf.exposure", "pval.exposure", "exposure", "Beta_SBP", "se
    _SBP", "EAF_SBP", "id.exposure")
13 head(dbp_data_2)
14
15 dbp_data_2<-dbp_data_2[order(dbp_data_2$chr.exposure),]
16
17 clumped_dbp_data_2 <- clump_data(dbp_data_2,
18                                clump_kb = 10000,  # Clumping window
                                (10,000 kb)
19                                clump_r2 = 0.001,  # LD threshold (r2 <
                                0.001)
20                                pop = "EUR")  # European LD reference
21
22 dbp_snplist<-clumped_dbp_data_2%>%select(SNP)
23 #sbp_effect_list<-clumped_sbp_data_2%>%select(SNP,effect_allele.exposure,
    beta.exposure)
24
25
26
27 write.csv(clumped_dbp_data_2, "C:/Users/tabe0010/OneDrive - Monash
    University/MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dbp_
    snps/dbp_exposure.csv", row.names = FALSE)
28 write.table(dbp_snplist, "C:/Users/tabe0010/OneDrive - Monash University/
    MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dx_data/dbp/data/
    dbp_snplist.txt",row.names = FALSE, col.names =FALSE,quote = FALSE,sep
    = " ")

```

4.6.1.3 Selecting the genetic variants from the BGEN file

PLINK2 was used via the **Swiss Army Knife** to select the necessary SNPs from the UK Biobank. The UKB RAP served as the platform to access the SNPs from the UK Biobank. Alternatively, a *bash* command was used on the local machine to run the process. First, users logged in with their UKB RAP credentials using the Command Prompt (on a Windows machine). Then, Git Bash was opened, the project was selected, and the analysis was run. A job request was sent, and users were notified when it was completed.

```

#####
#Run the following PLINK2 codes on Git Bash
#####

#Login through Command Prompt on your Windows machine
dx login

#Run the code below on Git Bash

dx select --level VIEW
#select your project
# make sure your sbp_snplist.txt file is uploaded to the UKB RAP.
#select the "instance type" you want: this makes sure you have enough computation
    power (CPU and GPU).
#In the command below, I put chromosome 1 to 22 to loop through all autosomal
    chromosomes just to show the code. But in actuality, I put two chromosomes at a

```

```

    time. This makes sure I have enough computational space and if there is any
    error, I could adjust the code.

# Loop over chromosomes 1 to 22 and process each one with the SNP list
run_merge=""
for chr in {1..22}; do
    run_merge+="cp /mnt/project/Bulk/Imputation/UKB\ imputation\ from\ genotype/
    ukb22828_c${chr}_b0_v3.bgen .; "
    run_merge+="cp /mnt/project/Bulk/Imputation/UKB\ imputation\ from\ genotype/
    ukb22828_c${chr}_b0_v3.sample .; "
    run_merge+="plink2 --bgen ukb22828_c${chr}_b0_v3.bgen ref-first --sample
    ukb22828_c${chr}_b0_v3.sample --extract dbp_snplist.txt --make-pgen --
    autosome-xy --out ukb22828_c${chr}_v3; "
done

dx run swiss-army-knife -iin="project-GpbQqBjJb7jb1vQjf8ZxVpVY:/DBP_data/dbp_txt/
dbp_snplist.txt" -icmd="{run_merge}" --tag="Step1" --instance-type "
mem1-ssd1-v2-x36" --destination="project-GpbQqBjJb7jb1vQjf8ZxVpVY:/DBP_data/
dbp_geno_data/" --brief --yes

#####
#Run the following PLINK2 codes via Swiss Army Knife on UKB RAP platform
#####
#Make sure you have uploaded the sbp_merge_list.txt to UKB RAP file path.
#The merge list should have a single column list containing the following text, "
    ukb22828_cx_v3" (without the quotations). The column will have 22 rows for each
    autosomal chromosomes. Replace "x" with 1-22.

#merging the pgen files

#Execute the command on Swiss Army Knife interface
#inputs are the plink files for the chromosomes and the txt file for the merging
    chromosomes
plink2 --pmerge-list dbp_merge_list.txt pfile --make-pgen --out
    ukb22828_c1_22_v3_dbp_merged

#Calculate the allele dosage
#creating a .raw file for participants with the number of effect allele (0, 1, or
    2) for each snp
# Input the for code below is the merged plink files (pfiles)

plink2 --pfile ukb22828_c1_22_v3_dbp_merged --export A --out ukb22828_dbp_alleles

#Calculate allele frequency

plink2 --pfile ukb22828_c1_22_v3_dbp_merged --freq --out ukb22828_dbp_allele_freq

```

The last PLINK output files, **ukb22828_dbp_alleles.raw** and **ukb22828_dbp_allele_freq.afreq**, contain the allele dosage and allele frequency for each of the 208 SNPs. Download these files to your local machine and save them to the **\$dx_data_dbp** file path.

4.6.1.4 Association of Genetic Variants with Quality-Adjusted Life Years

We then worked on the association between the genetic variants and QALYs. The QALYs were regressed on the allele dosages, adjusting for age, sex, and the first 10 genetic principal components to account for population stratification.

```
*****
*SNP-QALYs association
*****
import delimited "$dx_data_dbp\ukb22828_dbp_alleles.raw",clear

keep iid rs*
rename iid id_phe
merge 1:1 id_phe using "$stata_sbp_input\id_list.dta", keep(3) nogen
save "$stata_dbp_input\snp_alleles_dbp.dta", replace

use "$stata_sbp_output\part_3b.dta", clear
merge 1:1 id_phe using "$stata_dbp_input\snp_alleles_dbp.dta", keep(1 3) nogen
*gen imputation = .
gen snp = ""
gen effect_allele = ""
gen eaf = .
gen outcome = ""
gen beta = .
gen se = .
gen variance = .
gen p = .
gen n = .
local i = 1
local outcomes = "qaly_hes"

    foreach outcome in `outcomes' {

        qui regress `outcome' rs* age sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10

        foreach snp of varlist rs* {
            local snpx = substr("`snp'",1,length("`snp'")-2)
            *qui replace imputation = `imputation' in `i'
            qui replace snp = "`snpx'" in `i'
            qui replace outcome = "`outcome'" in `i'
            qui replace beta = _b[`${snp}'] if snp == "`snpx'" & outcome == "`outcome'"
            qui replace se = _se[`${snp}'] if snp == "`snpx'" & outcome == "`outcome'"
            qui sum `snp'
            qui replace eaf = r(mean)/2 if snp == "`snpx'"
            local effect_allele = upper(substr("`snp'",length("`snp'"),1))
            qui replace effect_allele = "`effect_allele'" if snp == "`snpx'"
            local i = `i'+1
        }

        *Ns
        qui sum `outcome'
        qui replace n = r(N) if outcome == "`outcome'"
    }

keep snp-n
keep if snp != ""
qui replace variance = se^2
qui replace p = 2*normal(-abs(beta/se))
save "$stata_dbp_result\Results_snp_qalys_dbp.dta", replace
```

```

use "$stata_dbp_result\Results_snp_qalys_dbp.dta", clear
keep if outcome == "qaly_hes"
save "$stata_dbp_result\Results_snp_qaly_hes_dbp.dta"

*****
*merge with allele frequency data
*****
import delimited "$dx_data_dbp\ukb22828_dbp_allele_freq.afreq",clear
rename id snp
rename alt other_allele
merge 1:1 snp using "$stata_dbp_result\Results_snp_qaly_hes_dbp.dta"
drop chrom ref _merge* alt_freqs obs_ct
export delimited using "$dbp_snps\snp_qaly_hes_dbp.csv",replace

```

4.6.1.5 Data Harmonisation

We now had the SNP-exposure and SNP-outcome association data. The next task was to harmonize these two datasets. For this, we continued working in the previous R environment.

```

1 qaly_hes_data_dbp<-read.csv("C:/Users/tabe0010/OneDrive - Monash
  University/MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dbp_
  snps/snp_qaly_hes_dbp.csv")
2 colnames(qaly_hes_data_dbp)<-c("SNP", "other_allele.outcome", "effect_
  allele.outcome", "eaf.outcome", "outcome", "beta.outcome", "se.outcome"
  , "variance.outcome", "pval.outcome", "samplesize.outcome")
3 qaly_hes_data_dbp$id.outcome = "ukb"
4
5 harmonise_data_dbp <- harmonise_data(
6   exposure_dat = clumped_dbp_data_2,
7   outcome_dat = qaly_hes_data_dbp
8 )
9
10
11 dbp_effect_list<-harmonise_data_dbp%>%select(SNP,effect_allele.exposure,
  beta.exposure)
12 dbp_snplist_exclude<-harmonise_data_dbp%>%filter(mr_keep == "FALSE")%>%
  select(SNP)
13
14 write.table(dbp_effect_list, "C:/Users/tabe0010/OneDrive - Monash
  University/MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dx_data
  /dbp/data/dbp_effect_list.txt",row.names = FALSE, col.names =FALSE,
  quote = FALSE,sep = " ")
15 write.table(dbp_snplist_exclude, "C:/Users/tabe0010/OneDrive - Monash
  University/MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dx_data
  /dbp/data/dbp_snplist_exclude.txt", row.names = FALSE, col.names =FALSE
  ,quote = FALSE,sep = " ")

```

4.6.1.6 Combining the genetic data

We had now identified the genetic variants primarily associated with SBP and secondarily with DBP, and vice versa. These datasets were then combined.

```

1
2 #We will continue working on the previous environment
3

```

```

4  #Let's combine the SNPs that are LD clumped and associated with SBP and DBP
   to create the SNP list (this is the total sets of SNPs)
5
6  sbp_dbp_snplist<-rbind(sbp_snplist, dbp_snplist)
7
8  #Let's work on combining SNPs primarily associated with SBP and seconarily
   associated with DBP
9
10 temp_sbp<-harmonise_data%>%select(chr.exposure, pos.exposure, SNP, effect_
   allele.exposure, other_allele.exposure, beta.exposure, se.exposure, eaf.
   exposure)
11 temp_sbp_2<-harmonise_data_dbp%>%select(chr.exposure, pos.exposure, SNP,
   effect_allele.exposure, other_allele.exposure, Beta_SBP, se_SBP, EAF_SBP
   )%>%rename(beta.exposure = Beta_SBP, se.exposure = se_SBP, eaf.exposure
   = EAF_SBP)
12 temp_sbp_append<-rbind(temp_sbp, temp_sbp_2)
13 temp_sbp_append<-temp_sbp_append[order(temp_sbp_append$chr.exposure, temp_
   sbp_append$pos.exposure),]
14 temp_sbp_append$beta.exposure<-as.numeric(temp_sbp_append$beta.exposure)
15 temp_sbp_append$se.exposure<-as.numeric(temp_sbp_append$se.exposure)
16 n_distinct(temp_sbp_append$SNP)
17 temp_sbp_na<-temp_sbp_append%>%filter(is.na(beta.exposure))%>%select(SNP)
18 mvmr_sbp_snplist_exclude<-rbind(sbp_snplist_exclude, dbp_snplist_exclude,
   temp_sbp_na)
19 mvmr_sbp_effect_list<-temp_sbp_append%>%select(SNP, effect_allele.exposure
   , beta.exposure)
20
21 #Let's work on combining SNPs primarily associated with DBP and seconarily
   associated with SBP
22
23 temp_dbp<-harmonise_data_dbp%>%select(chr.exposure, pos.exposure, SNP,
   effect_allele.exposure, other_allele.exposure, beta.exposure, se.exposure,
   eaf.exposure)
24 temp_dbp_2<-harmonise_data%>%select(chr.exposure, pos.exposure, SNP,
   effect_allele.exposure, other_allele.exposure, Beta_DBP, se_DBP, EAF_DBP
   )%>%rename(beta.exposure = Beta_DBP, se.exposure = se_DBP, eaf.exposure
   = EAF_DBP)
25 temp_dbp_append<-rbind(temp_dbp, temp_dbp_2)
26 temp_dbp_append<-temp_dbp_append[order(temp_dbp_append$chr.exposure, temp_
   dbp_append$pos.exposure),]
27 temp_dbp_append$beta.exposure<-as.numeric(temp_dbp_append$beta.exposure)
28 temp_dbp_append$se.exposure<-as.numeric(temp_dbp_append$se.exposure)
29 temp_dbp_na<-temp_dbp_append%>%filter(is.na(beta.exposure))%>%select(SNP)
30 mvmr_dbp_snplist_exclude<-rbind(sbp_snplist_exclude, dbp_snplist_exclude,
   temp_dbp_na)
31 mvmr_dbp_effect_list<-temp_dbp_append%>%select(SNP, effect_allele.exposure
   , beta.exposure)
32
33 #Save the files
34

```

```

35 write.table(sbp_dbp_snplist, "C:/Users/tabe0010/OneDrive - Monash
    University/MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dx_data
    /sbp_dbp/data/sbp_dbp_snplist.txt",row.names = FALSE, col.names =FALSE,
    quote = FALSE,sep = " ")
36 write.table(mvmr_sbp_snplist_exclude, "C:/Users/tabe0010/OneDrive - Monash
    University/MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dx_
    data/sbp_dbp/data/mvmr_sbp_snplist_exclude.txt",row.names = FALSE, col.
    names =FALSE,quote = FALSE,sep = " ")
37 write.table(mvmr_dbp_snplist_exclude, "C:/Users/tabe0010/OneDrive - Monash
    University/MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dx_
    data/sbp_dbp/data/mvmr_dbp_snplist_exclude.txt",row.names = FALSE, col.
    names =FALSE,quote = FALSE,sep = " ")
38 write.table(mvmr_sbp_effect_list, "C:/Users/tabe0010/OneDrive - Monash
    University/MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dx_data
    /sbp_dbp/data/mvmr_sbp_effect_list.txt",row.names = FALSE, col.names =
    FALSE,quote = FALSE,sep = " ")
39 write.table(mvmr_dbp_effect_list, "C:/Users/tabe0010/OneDrive - Monash
    University/MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/dx_data
    /sbp_dbp/data/mvmr_dbp_effect_list.txt",row.names = FALSE, col.names =
    FALSE,quote = FALSE,sep = " ")
40 write.csv(temp_sbp_append, "C:/Users/tabe0010/OneDrive - Monash University
    /MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/sbp_dbp_snps/mvmr
    _sbp.csv", row.names = FALSE)
41 write.csv(temp_dbp_append, "C:/Users/tabe0010/OneDrive - Monash University
    /MR_backup_file/Articles/Evangelou/new_sbp_snps/stata/sbp_dbp_snps/mvmr
    _dbp.csv", row.names = FALSE)

```

4.6.1.7 Working on SNPs primarily associated with SBP and seconarily associated with DBP

We calculated the PRS, allele frequency, and allele dosage for these SNPs using the `sbp_dbp_snplist.txt`, `mvmr_sbp_snplist_exclude.txt`, and `mvmr_sbp_effect_list.txt` files.

We had 395 total SNPs; of these, 382 SNPs were associated with SBP as a primary or secondary trait, and 384 SNPs were associated with DBP as a primary or secondary trait. The remaining SNPs were either ambiguous SNPs or had missing effect sizes.

```

#####
#Run the following PLINK2 codes on Git Bash
#####

#Login through Command Prompt on your Windows machine
dx login

#Run the code below on Git Bash

dx select --level VIEW
#select your project
# make sure your sbp_dbp_snplist.txt file is uploaded to the UKB RAP.
#select the "instance type" you want: this makes sure you have enough computation
    power (CPU and GPU).
#In the command below, I put chromosome 1 to 22 to loop through all autosomal
    chromosomes just to show the code. But in actuality, I put two chromosomes at a

```

```

    time. This makes sure I have enough computational space and if there is any
    error, I could adjust the code.

# Loop over chromosomes 1 to 22 and process each one with the SNP list

run_merge=""
for chr in {1..22}; do
    run_merge+="cp /mnt/project/Bulk/Imputation/UKB\ imputation\ from\ genotype/
    ukb22828_c${chr}_b0_v3.bgen .; "
    run_merge+="cp /mnt/project/Bulk/Imputation/UKB\ imputation\ from\ genotype/
    ukb22828_c${chr}_b0_v3.sample .; "
    run_merge+="plink2 --bgen ukb22828_c${chr}_b0_v3.bgen ref-first --sample
    ukb22828_c${chr}_b0_v3.sample --extract sbp_dbp_snplist.txt --exclude
    mvmr_sbp_snplist_exclude.txt --make-pgen --autosome-xy --out ukb22828_c${
    chr}_v3; "
done

dx run swiss-army-knife -iin="project-GpbQqBjJb7jb1vQjf8ZxVpVY:/SBP_DBP_data/
mvmr_sbp_txt/sbp_dbp_snplist.txt" -iin="project-GpbQqBjJb7jb1vQjf8ZxVpVY:/
SBP_DBP_data/mvmr_sbp_txt/mvmr_sbp_snplist_exclude.txt" -icmd="{run_merge}" --
tag="mvmr_chr1_22" --instance-type "mem1-ssd1_v2_x36" --destination="project-
GpbQqBjJb7jb1vQjf8ZxVpVY:/SBP_DBP_data/mvmr_sbp_genotype_data/mvmr_sbp_chromosomes/
" --brief --yes

#####
#Run the following PLINK2 codes via Swiss Army Knife on UKB RAP platform
#####
#Make sure you have uploaded the sbp_dbp_merge_list.txt to UKB RAP file path.
#The merge list should have a single column list containing the following text, "
    ukb22828_cx_v3" (without the quotations). The column will have 22 rows for each
    autosomal chromosomes. Replace "x" with 1-22.

#merging the pgen files

#Execute the command on Swiss Army Knife interface
#inputs are the plink files for the chromosomes and the txt file for the merging
    chromosomes
plink2 --pmerge-list sbp_dbp_merge_list.txt pfile --make-pgen --out
    ukb22828_c1_22_v3_mvmr_sbp_merged

#PRS mvmr_SBP
plink2 --pfile ukb22828_c1_22_v3_mvmr_sbp_merged --score mvmr_sbp_effect_list.txt
    cols=+scoresums --out ukb22828_mvmr_sbp_prs

#Calculate the allele dosage
#creating a .raw file for participants with the number of effect allele (0, 1, or
    2) for each snp
# Input the for code below is the merged plink files (pfiles)

plink2 --pfile ukb22828_c1_22_v3_mvmr_sbp_merged --export A --out
    ukb22828_mvmr_sbp_alleles

#Calculate allele frequency

plink2 --pfile ukb22828_c1_22_v3_mvmr_sbp_merged --freq --out

```

4.6.1.8 Working on SNPs primarily associated with DBP and seconarily associated with SBP

We calculated the PRS, allele frequency, and allele dosage for these SNPs using the `sbp_dbp_snplist.txt`, `mvnr_dbp_snplist_exclude.txt`, and `mvnr_dbp_effect_list.txt` files.

```
#####
#Run the following PLINK2 codes on Git Bash
#####

#Login through Command Prompt on your Windows machine
dx login

#Run the code below on Git Bash

dx select --level VIEW
#select your project
# make sure your sbp_dbp_snplist.txt file is uploaded to the UKB RAP.
#select the "instance type" you want: this makes sure you have enough computation
  power (CPU and GPU).
#In the command below, I put chromosome 1 to 22 to loop through all autosomal
  chromosomes just to show the code. But in actuality, I put two chromosomes at a
  time. This makes sure I have enough computational space and if there is any
  error, I could adjust the code.

# Loop over chromosomes 1 to 22 and process each one with the SNP list

run_merge=""
for chr in {1..22}; do
  run_merge+="cp /mnt/project/Bulk/Imputation/UKB\ imputation\ from\ genotype/
    ukb22828_c${chr}_b0_v3.bgen .; "
  run_merge+="cp /mnt/project/Bulk/Imputation/UKB\ imputation\ from\ genotype/
    ukb22828_c${chr}_b0_v3.sample .; "
  run_merge+="plink2 --bgen ukb22828_c${chr}_b0_v3.bgen ref-first --sample
    ukb22828_c${chr}_b0_v3.sample --extract sbp_dbp_snplist.txt --exclude
    mvnr_dbp_snplist_exclude.txt --make-pgen --autosome-xy --out ukb22828_c${
    chr}_v3; "
done

dx run swiss-army-knife -iin="project-GpbQqBjJb7jb1vQjf8ZxVpVY:/SBP_DBP_data/
  mvnr_dbp_txt/sbp_dbp_snplist.txt" -iin="project-GpbQqBjJb7jb1vQjf8ZxVpVY:/
  SBP_DBP_data/mvnr_dbp_txt/mvnr_dbp_snplist_exclude.txt" -icmd="{run_merge}" --
  tag="mvnr_chr1_22" --instance-type "mem1-ssd1_v2_x36" --destination="project-
  GpbQqBjJb7jb1vQjf8ZxVpVY:/SBP_DBP_data/mvnr_dbp_genotype_data/mvnr_dbp_chromosomes/
  " --brief --yes

#####
#Run the following PLINK2 codes via Swiss Army Knife on UKB RAP platform
#####
#Make sure you have uploaded the sbp_dbp_merge_list.txt to UKB RAP file path.
#The merge list should have a single column list containing the following text, "
  ukb22828_cx_v3" (without the quotations). The column will have 22 rows for each
```



```

    autosomal chromosomes. Replace "x" with 1-22.

#merging the pgen files

#Execute the command on Swiss Army Knife interface
#inputs are the plink files for the chromosomes and the txt file for the merging
  chromosomes
plink2 --pmerge-list sbp_dbp_merge_list.txt pfile --make-pgen --out
      ukb22828_c1_22_v3_mvmmr_dbp_merged

#PRS mvmmr_DBP
plink2 --pfile ukb22828_c1_22_v3_mvmmr_dbp_merged --score mvmmr_dbp_effect_list.txt
      cols=+scoresums --out ukb22828_mvmmr_dbp_prs

#Calculate the allele dosage
#creating a .raw file for participants with the number of effect allele (0, 1, or
  2) for each snp
# Input the for code below is the merged plink files (pfiles)

plink2 --pfile ukb22828_c1_22_v3_mvmmr_dbp_merged --export A --out
      ukb22828_mvmmr_dbp_alleles

#Calculate allele frequency

plink2 --pfile ukb22828_c1_22_v3_mvmmr_dbp_merged --freq --out
      ukb22828_mvmmr_dbp_allele_freq

```

We have now calculated the PRS for SNPs effect on SBP. Download the **ukb22828_sbp_prs.sscore** file to your local machine and save them to the `$dx_data_sbp` file path.

4.6.1.9 Combining Phenotype and Genotype data

We then combined our phenotype data with the PRS. The **part_3b.dta** file contained the phenotype data for our cohort, while the **ukb22828_mvmmr_sbp_prs.sscore** and **ukb22828_mvmmr_dbp_prs.sscore** files contained the PRS for each participant in the UK Biobank. The next Stata code merged the two datasets and prepared the data for the MVMR analysis. We saved the data to the **part_4a.dta** Stata file that we had previously created.

```

import delimited "$dx_data_sbp\ukb22828_mvmmr_sbp_prs.sscore", clear
gen id_phe = iid // IID: Individual ID
save "$dx_data_sbp\ukb22828_mvmmr_sbp_prs.sscore.dta", replace

import delimited "$dx_data_sbp_dbp\ukb22828_mvmmr_dbp_prs.sscore", clear
gen id_phe = iid // IID: Individual ID
save "$dx_data_sbp_dbp\ukb22828_mvmmr_dbp_prs.sscore.dta", replace

use "$stata_sbp_output\part_3b.dta", clear

merge 1:1 id_phe using "$dx_data_sbp_dbp\ukb22828_mvmmr_sbp_prs.sscore.dta"
keep if _merge == 3
drop _merge*
rename score1_sum prs_mvmmr_sbp
drop fid iid allele_ct named_allele_dosage_sum score1_avg

merge 1:1 id_phe using "$dx_data_sbp_dbp\ukb22828_mvmmr_dbp_prs.sscore.dta"
keep if _merge == 3
drop _merge*

```

```

rename score1_sum prs_mvmmr_dbp
drop fid iid allele_ct named_allele_dosage_sum score1_avg

save "$stata_sbp_output\part_4a.dta", replace

```

4.6.1.10 Analysis

To estimate the direct effect of SBP on QALYs conditional on DBP, we employed the difference method approach [24]. The total effect (β_1^*) of SBP on QALYs was estimated in the main analysis. The direct effect (β_1) of SBP on QALYs was estimated by considering both SBP and DBP traits. At the first stage, we regressed the exposure trait (SBP) on the polygenic risk score (PRS) for SNPs that were associated with SBP (i.e., 382 SNPs). We repeated the same approach for the DBP exposure by regressing it on the PRS for SNPs that were associated with DBP (i.e., 384 SNPs). In the second stage, the outcome (QALYs) was regressed on the predicted values for SBP and DBP from the first-stage models. For the direct effect models, age, sex, UK Biobank assessment centre, genotyping array, and the first 10 genetic principal components for population stratification were used as covariates. The indirect effect was calculated as the difference between the total effect estimate and the direct effect estimate, i.e., $\beta_1^* - \beta_1$.

```

*MVMMR
use "$stata_sbp_output\part_4a.dta", clear
*Create table
gen outcome = ""
gen type = ""
gen n = .
gen beta = .
gen variance = .
gen se = .
gen double p = .
*gen double p_endog = .
gen f_stat = .
local x = 1

    foreach var in qaly_hes {
        dis "Outcome = `var'"

        *MR analysis (main)
        ivreg2 `var' (phe_sbp_adj = prs_sbp) age i.sex pc1 pc2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 p
> c10 i.centre i.geno_array, robust

        matrix a = e(b)
        matrix b = e(V)
        local beta_total = a[1,1]
        local variance_total = b[1,1]

        local n_total = e(N)
        local f_stat_total = e(widstat)
        *local p_endog_total = e(estatp)

        replace outcome = "`var'" in `x'
        replace type = "Total effect" in `x'
        *qui replace imputation = `j' in `x'
        foreach z in beta variance n f_stat {
            replace `z' = ``z'_total' in `x'
        }

        local x = `x' + 1

        *MV MR
        ivreg2 `var' (phe_sbp_adj phe_dbp_adj = prs_mvmmr_sbp prs_mvmmr_dbp) age i.sex pc1 pc
> 2 pc3 pc4 pc5 pc6 pc7 pc8 pc9 pc10 i.centre i.geno_array, robust

        matrix a = e(b)

```

```

matrix b = e(V)
local beta_direct = a[1,1]
local variance_direct = b[1,1]

local n_direct = e(N)
local f_stat_direct = e(widstat)
*local p_endog = e(estatp)

replace outcome = "`var'" in `x'
replace type = "Direct effect" in `x'
*qui replace imputation = `j' in `x'
foreach z in beta variance n f_stat {
    replace `z' = ``z'_direct' in `x'
}

local x = `x' + 1

}

keep outcome-f_stat
keep if outcome != ""
replace outcome = "QALYs per year (with 240 comorbidities)" if outcome == "qaly_hes"
qui replace se = sqrt(var)
qui replace p = 2*normal(-abs(beta/se))
sort outcome type
save "$stata_sbp_result\Result_mvmmr_sbp_table.dta", replace

```

4.7 Step 7: Tables and Figures

The codes that were used to generate the main results' tables and figures are provided here.

4.7.1 Tables

4.7.1.1 Table 1: Background characteristics

```
{
use "$stata_sbp_output\part_4a.dta", clear
gen Variable = ""
gen All = ""
gen N_All = ""
gen Men = ""
gen N_Men = ""
gen Women = ""
gen N_Women = ""
order Variable-N_Women, first

gen years = months/12
gen qual_0 = 0
replace qual_0 = 1 if qual_1 == 0 & qual_2 == 0 & qual_3 == 0
gen death = 0

*bysort dsources: su date_death
replace death = 1 if date_death <= 753 & dsources != "PEDW"
replace death = 1 if date_death <= 748 & dsources == "PEDW"

qui replace Variable = "N" in 1
qui replace Variable = "Age at recruitment, years [Median (IQR)]" in 2
qui replace Variable = "Body Mass Index, kg/m2 [Median (IQR)]" in 3
qui replace Variable = "Systolic blood pressure, mmHg [Median (IQR)]" in 4
qui replace Variable = "Years of follow-up [Median (IQR)]" in 5
qui replace Variable = "Death [N (%)]*" in 6
qui replace Variable = "Qualification: None [N (%)]" in 7
qui replace Variable = "Qualification: A levels, O level, GCSE or CSE [N (%)]" in 8
qui replace Variable = "Qualification: NVQ or other [N (%)]" in 9
qui replace Variable = "Qualification: College or university degree [N (%)]" in 10
qui replace Variable = "Average QALYs per year, HES only [Median (IQR)]*" in 11

local i = 1

*Number of participants
{
qui sum sex
local x = r(N)
local x: dis %9.0fc `x'
local x = strtrim("`x'")
local x1 = r(N)*r(mean)
local x1: dis %9.0fc `x1'
local x1 = strtrim("`x1'")
local x2 = r(N)*(1-r(mean))
local x2: dis %9.0fc `x2'
local x2 = strtrim("`x2'")
qui replace All = "`x'" in `i'
qui replace Men = "`x1'" in `i'
qui replace Women = "`x2'" in `i'
local i = `i' + 1
}

*Age, SBP, systolic blood pressure, follow-up
foreach var of varlist age phe_bmi phe_sbp years {
sum `var',d
local median = r(p50)
local median: dis %9.2f `median'
local median = strtrim("`median'")
local N = r(N)
local N: dis %9.0fc `N'
local N = strtrim("`N'")
```

```

local p25 = r(p25)
local p25: dis %9.2f `p25`
local p25 = strtrim("`p25`")
local p75 = r(p75)
local p75: dis %9.2f `p75`
local p75 = strtrim("`p75`")
local x = "`median` (`p25` to `p75`)"
qui replace All = "`x`" in `i`
qui replace N_All = "`N`" in `i`

sum `var` if sex == 1,d
local median = r(p50)
local median: dis %9.2f `median`
local median = strtrim("`median`")
local N = r(N)
local N: dis %9.0fc `N`
local N = strtrim("`N`")
local p25 = r(p25)
local p25: dis %9.2f `p25`
local p25 = strtrim("`p25`")
local p75 = r(p75)
local p75: dis %9.2f `p75`
local p75 = strtrim("`p75`")
local x = "`median` (`p25` to `p75`)"
qui replace Men = "`x`" in `i`
qui replace N_Men = "`N`" in `i`

sum `var` if sex == 0,d
local median = r(p50)
local median: dis %9.2f `median`
local median = strtrim("`median`")
local N = r(N)
local N: dis %9.0fc `N`
local N = strtrim("`N`")
local p25 = r(p25)
local p25: dis %9.2f `p25`
local p25 = strtrim("`p25`")
local p75 = r(p75)
local p75: dis %9.2f `p75`
local p75 = strtrim("`p75`")
local x = "`median` (`p25` to `p75`)"
qui replace Women = "`x`" in `i`
qui replace N_Women = "`N`" in `i`

local i = `i` + 1
}

*Primary care, death & qualifications
foreach var of varlist death qual_0 qual_1-qual_3 {

    qui sum `var`
    local N = r(mean)*r(N)
    local N: dis %9.0fc `N`
    local N = strtrim("`N`")
    local percent = r(mean)*100
    local percent: dis %9.2f `percent`
    local percent = strtrim("`percent`")
    local x = "`N` (`percent`)"
    qui replace All = "`x`" in `i`

    local N2 = r(N)
    local N2: dis %9.0fc `N2`
    local N2 = strtrim("`N2`")
    qui replace N_All = "`N2`" in `i`

    qui sum `var` if sex == 1
    local N = r(mean)*r(N)
    local N: dis %9.0fc `N`
    local N = strtrim("`N`")
    local percent = r(mean)*100

```

```

local percent: dis %9.2f `percent´
local percent = strtrim("`percent´")
local x = "`N´ (`percent´)"
qui replace Men = "`x´" in `i´

local N2 = r(N)
local N2: dis %9.0fc `N2´
local N2 = strtrim("`N2´")
qui replace N_Men = "`N2´" in `i´

qui sum `var´ if sex == 0
local N = r(mean)*r(N)
local N: dis %9.0fc `N´
local N = strtrim("`N´")
local percent = r(mean)*100
local percent: dis %9.2f `percent´
local percent = strtrim("`percent´")
local x = "`N´ (`percent´)"
qui replace Women = "`x´" in `i´

local N2 = r(N)
local N2: dis %9.0fc `N2´
local N2 = strtrim("`N2´")
qui replace N_Women = "`N2´" in `i´

local i = `i´+1
}

*HES QALYs
foreach var of varlist qaly_hes {
  *All

  sum `var´, d
  local median = r(p50)
  local median: dis %9.2f `median´
  local median = strtrim("`median´")
  local p25 = r(p25)
  local p25: dis %9.2f `p25´
  local p25 = strtrim("`p25´")
  local p75 = r(p75)
  local p75: dis %9.2f `p75´
  local p75 = strtrim("`p75´")
  local x = "`median´ (`p25´ to `p75´)"
  qui replace All = "`x´" in `i´

  *Men
  su `var´ if sex == 1, d
  local median = r(p50)
  local median: dis %9.2f `median´
  local median = strtrim("`median´")
  local p25 = r(p25)
  local p25: dis %9.2f `p25´
  local p25 = strtrim("`p25´")
  local p75 = r(p75)
  local p75: dis %9.2f `p75´
  local p75 = strtrim("`p75´")
  local x = "`median´ (`p25´ to `p75´)"
  qui replace Men = "`x´" in `i´

  *Women
  su `var´ if sex == 0, d
  local median = r(p50)
  local median: dis %9.2f `median´
  local median = strtrim("`median´")
  local p25 = r(p25)
  local p25: dis %9.2f `p25´
  local p25 = strtrim("`p25´")
  local p75 = r(p75)
  local p75: dis %9.2f `p75´

```

```

        local p75 = strtrim("`p75`")
        local x = "`median` (`p25` to `p75`)"
        qui replace Women = "`x`" in `i`

        local i = `i` + 1
    }

    keep Var-N_Women
    drop N*
    drop if Variable == ""
    save "$stata_sbp_result\Table_1.dta", replace
}

```

4.7.1.2 Table 2: Main analysis

```
use "$stata_sbp_result\Result_sbp_table.dta", clear
replace beta = 10*beta
replace se = 10*se
*Change to percentages
replace beta = 100*beta
replace se = 100*se
gen lower = beta-1.96*se
gen upper = beta+1.96*se
format beta se lower upper %9.2f
*Effect estimate variable - sort of decimal places
foreach var of varlist beta se lower upper {
    tostring `var', gen(`var'_x) force

    replace `var'_x = substr(`var'_x,1,4) if `var' < 0 & `var' > -1
    replace `var'_x = substr(`var'_x,1,5) if `var' <= -1 & `var' > -10
    replace `var'_x = substr(`var'_x,1,6) if `var' <= -10 & `var' > -100
    replace `var'_x = substr(`var'_x,1,7) if `var' <= -100 & `var' > -1000
    replace `var'_x = substr(`var'_x,1,3) if `var' > 0 & `var' < 1
    replace `var'_x = substr(`var'_x,1,4) if `var' >= 1 & `var' < 10
    replace `var'_x = substr(`var'_x,1,5) if `var' >= 10 & `var' < 100
    replace `var'_x = substr(`var'_x,1,6) if `var' >= 100 & `var' < 1000

    replace `var'_x = substr(`var'_x,1,1) if substr(`var'_x,1,1) == "." | substr(`var'_
> x,1,2) == "-."
    replace `var'_x = `var'_x + "%"

    *replace `var'_x = "f"+`var'_x if strpos(outcome,"QALYs") == 0
}

gen effect = beta_x + " (" + lower_x + " to " + upper_x + ")"
keep outcome type n effect se_x p p_endog f_stat
rename se_x se
order outcome type n effect se p p_endog f_stat
save "$stata_sbp_result\Table_2.dta", replace
```


4.7.1.3 Table 3: Sensitivity analyses - No history of antihypertensive medication cohort

```

use "$stata_sbp_result\Result_sbp_exclusive_table_no_medication.dta", clear
replace beta = 10*beta
replace se = 10*se

*Change to percentages
replace beta = 100*beta
replace se = 100*se

gen lower = beta-1.96*se
gen upper = beta+1.96*se

format beta se lower upper %9.2f

*Effect estimate variable - sort of decimal places
foreach var of varlist beta se lower upper {
    tostring `var', gen(`var'_x) force

    replace `var'_x = substr(`var'_x,1,4) if `var' < 0 & `var' > -1
    replace `var'_x = substr(`var'_x,1,5) if `var' <= -1 & `var' > -10
    replace `var'_x = substr(`var'_x,1,6) if `var' <= -10 & `var' > -100
    replace `var'_x = substr(`var'_x,1,7) if `var' <= -100 & `var' > -1000
    replace `var'_x = substr(`var'_x,1,3) if `var' > 0 & `var' < 1
    replace `var'_x = substr(`var'_x,1,4) if `var' >= 1 & `var' < 10
    replace `var'_x = substr(`var'_x,1,5) if `var' >= 10 & `var' < 100
    replace `var'_x = substr(`var'_x,1,6) if `var' >= 100 & `var' < 1000

    replace `var'_x = subinstr(`var'_x,".", "0.",.) if substr(`var'_x,1,1) == "." | substr(`var'_
> x,1,2) == "-."
    replace `var'_x = `var'_x + "%"

    *replace `var'_x = "f"+`var'_x if strpos(outcome,"QALYs") == 0
}

gen effect = beta_x + " (" + lower_x + " to " + upper_x + ")"
keep outcome type n effect se_x p p_endog f_stat
rename se_x se
order outcome type n effect se p p_endog f_stat
save "$stata_sbp_result\Table_3.dta", replace

```

4.7.1.4 Table 4: Sensitivity analyses - Two-sample MR

```

use "$stata_sbp_result\Results_table_sensitivity_sbp.dta", clear

gen ll_ivw = ivw - ivw_se*1.96
gen ul_ivw = ivw + ivw_se*1.96
gen ll_egger_slope = egger_slope - egger_slope_se*1.96
gen ul_egger_slope = egger_slope + egger_slope_se*1.96
gen ll_egger_cons = egger_cons - egger_cons_se*1.96
gen ul_egger_cons = egger_cons + egger_cons_se*1.96
gen ll_median = median - median_se*1.96
gen ul_median = median + median_se*1.96
gen ll_mode = mode - mode_se*1.96
gen ul_mode = mode + mode_se*1.96

foreach var of varlist ivw ivw_se ll_ivw ul_ivw egger_slope egger_slope_se ll_egger_slope ul_egger_s
> lope egger_cons egger_cons_se ll_egger_cons ul_egger_cons median median_se ll_median ul_median mod
> e mode_se ll_mode ul_mode {
    replace `var' = `var'*10
    replace `var' = `var'*100
    tostring `var', gen(`var'_x) force

        replace `var'_x = substr(`var'_x,1,4) if `var' < 0 & `var' > -1
    replace `var'_x = substr(`var'_x,1,5) if `var' <= -1 & `var' > -10
    replace `var'_x = substr(`var'_x,1,6) if `var' <= -10 & `var' > -100
    replace `var'_x = substr(`var'_x,1,7) if `var' <= -100 & `var' > -1000
    replace `var'_x = substr(`var'_x,1,3) if `var' > 0 & `var' < 1
    replace `var'_x = substr(`var'_x,1,4) if `var' >= 1 & `var' < 10
    replace `var'_x = substr(`var'_x,1,5) if `var' >= 10 & `var' < 100
    replace `var'_x = substr(`var'_x,1,6) if `var' >= 100 & `var' < 1000

    replace `var'_x = substr(`var'_x,1,10) if substr(`var'_x,1,1) == "." | substr(`var'_
> x,1,2) == "-."
    replace `var'_x = `var'_x + "%"
}

drop ivw ivw_se egger_slope egger_slope_se egger_cons egger_cons_se median median_se mode mode_se ll
> _ivw ul_ivw ll_egger_slope ul_egger_slope ll_egger_cons ul_egger_cons ll_median ul_median ll_mode
> ul_mode

rename (ivw_x ivw_se_x egger_slope_x egger_slope_se_x egger_cons_x egger_cons_se_x median_x median_s
> e_x mode_x mode_se_x ll_ivw_x ul_ivw_x ll_egger_slope_x ul_egger_slope_x ll_egger_cons_x ul_egger_
> cons_x ll_median_x ul_median_x ll_mode_x ul_mode_x) (ivw ivw_se egger_slope egger_slope_se egger_c
> ons egger_cons_se median median_se mode mode_se ll_ivw ul_ivw ll_egger_slope ul_egger_slope ll_egg
> er_cons ul_egger_cons ll_median ul_median ll_mode ul_mode)

gen ivw_effect = ivw + " (" + ll_ivw + " to " + ul_ivw + ")"
gen egger_slope_effect = egger_slope + " (" + ll_egger_slope + " to " + ul_egger_slope + ")"
gen egger_cons_effect = egger_cons + " (" + ll_egger_cons + " to " + ul_egger_cons + ")"
gen median_effect = median + " (" + ll_median + " to " + ul_median + ")"
gen mode_effect = mode + " (" + ll_mode + " to " + ul_mode + ")"

keep outcome ivw_effect ivw_se ivw_p egger_slope_effect egger_slope_se egger_slope_p egger_cons_eff
> ct egger_cons_se egger_cons_p median_effect median_se median_p mode_effect mode_se mode_p
order outcome ivw_effect ivw_se ivw_p egger_slope_effect egger_slope_se egger_slope_p egger_cons_eff
> ct egger_cons_se egger_cons_p median_effect median_se median_p mode_effect mode_se mode_p

save "$stata_sbp_result\Table_4.dta",replace

```

4.7.1.5 Table 5: Sensitivity analyses - EQ-5D-index from UK Biobank survey

```

use "$stata_sbp_result\Result_sbp_table_ukb_utility.dta", clear
replace beta = 10*beta
replace se = 10*se
*Change to percentages
replace beta = 100*beta
replace se = 100*se
gen lower = beta-1.96*se
gen upper = beta+1.96*se
format beta se lower upper %9.2f
*Effect estimate variable - sort of decimal places
foreach var of varlist beta se lower upper {
    tostring `var', gen(`var'_x) force

    replace `var'_x = substr(`var'_x,1,4) if `var' < 0 & `var' > -1
    replace `var'_x = substr(`var'_x,1,5) if `var' <= -1 & `var' > -10
    replace `var'_x = substr(`var'_x,1,6) if `var' <= -10 & `var' > -100
    replace `var'_x = substr(`var'_x,1,7) if `var' <= -100 & `var' > -1000
    replace `var'_x = substr(`var'_x,1,3) if `var' > 0 & `var' < 1
    replace `var'_x = substr(`var'_x,1,4) if `var' >= 1 & `var' < 10
    replace `var'_x = substr(`var'_x,1,5) if `var' >= 10 & `var' < 100
    replace `var'_x = substr(`var'_x,1,6) if `var' >= 100 & `var' < 1000

    replace `var'_x = substr(`var'_x,1,1) if substr(`var'_x,1,1) == "." | substr(`var'_
> x,1,2) == "-."
    replace `var'_x = `var'_x + "%"

    *replace `var'_x = "f"+`var'_x if strpos(outcome,"QALYs") == 0
}

gen effect = beta_x + " (" + lower_x + " to " + upper_x + ")"
keep outcome type n effect se_x p p_endog f_stat
rename se_x se
order outcome type n effect se p p_endog f_stat
save "$stata_sbp_result\Table_5.dta", replace

```

4.7.1.6 Table 6: Secondary analysis

```
use "$stata_sbp_result\Result_mvmmr_sbp_table.dta", clear
replace beta = 10*beta
replace se = 10*se
*Change to percentages
replace beta = 100*beta
replace se = 100*se
gen lower = beta-1.96*se
gen upper = beta+1.96*se
format beta se lower upper %9.2f
*Effect estimate variable - sort of decimal places
foreach var of varlist beta se lower upper {
    tostring `var', gen(`var'_x) force

    replace `var'_x = substr(`var'_x,1,4) if `var' < 0 & `var' > -1
    replace `var'_x = substr(`var'_x,1,5) if `var' <= -1 & `var' > -10
    replace `var'_x = substr(`var'_x,1,6) if `var' <= -10 & `var' > -100
    replace `var'_x = substr(`var'_x,1,7) if `var' <= -100 & `var' > -1000
    replace `var'_x = substr(`var'_x,1,3) if `var' > 0 & `var' < 1
    replace `var'_x = substr(`var'_x,1,4) if `var' >= 1 & `var' < 10
    replace `var'_x = substr(`var'_x,1,5) if `var' >= 10 & `var' < 100
    replace `var'_x = substr(`var'_x,1,6) if `var' >= 100 & `var' < 1000

    replace `var'_x = substr(`var'_x,1,1) if substr(`var'_x,1,1) == "." | substr(`var'_
> x,1,2) == "-."
    replace `var'_x = `var'_x + "%"
}

gen effect = beta_x + " (" + lower_x + " to " + upper_x + ")"
keep outcome n type effect se_x p f_stat
rename se_x se
order outcome n type effect se p f_stat
save "$stata_sbp_result\Table_6.dta", replace
```

4.7.2 Figures

4.7.2.1 Figure 1

Flow chart showing selection criteria for genetic instruments primarily associated with systolic blood pressure

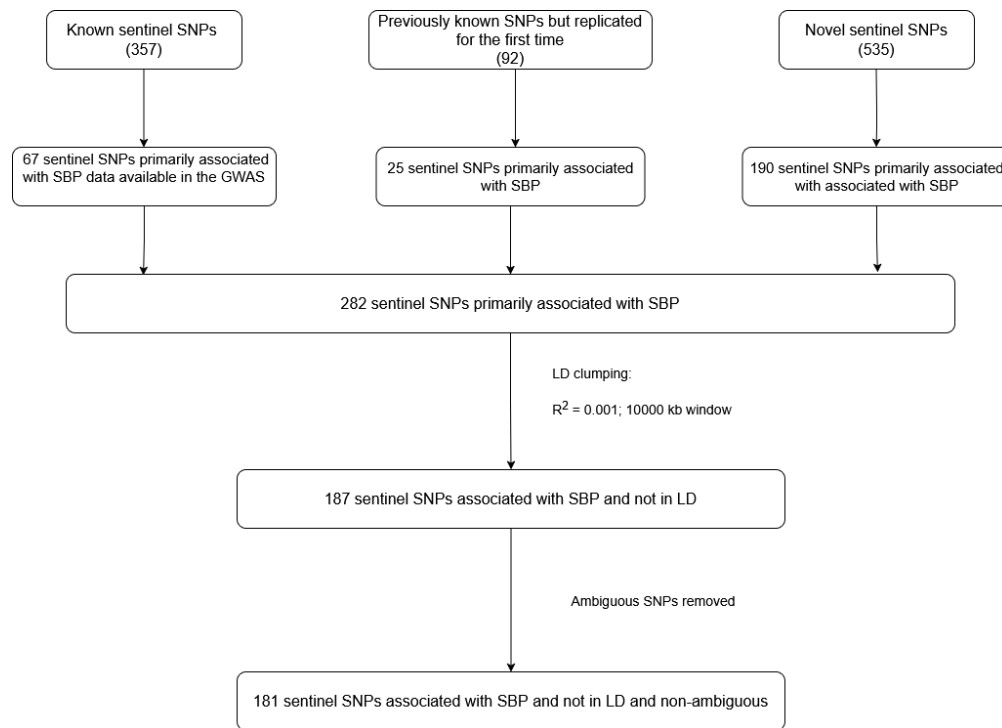


Figure 1: Flow chart showing selection criteria for genetic instruments primarily associated with systolic blood pressure

4.7.2.2 Figure 2

Flow chart showing selection criteria for genetic instruments primarily associated with diastolic blood pressure

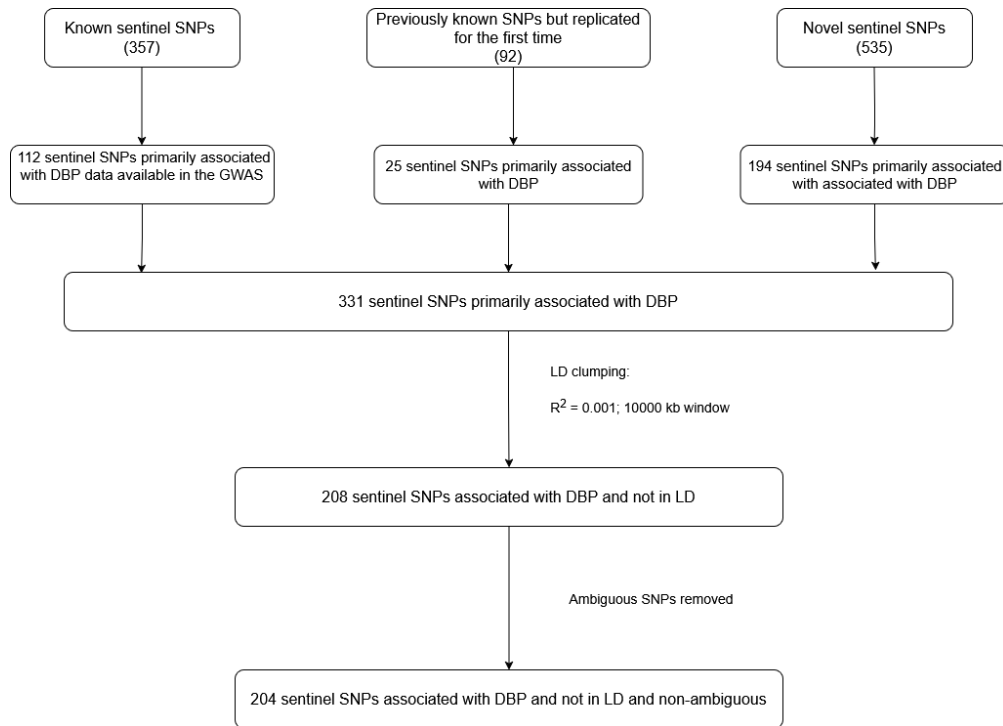


Figure 2: Flow chart showing selection criteria for genetic instruments primarily associated with diastolic blood pressure

4.7.2.3 Figure 3

Figure showing main analysis, and sub-group analyses by age, sex and PRS-free SBP.

```
1
2 #R plots
3
4 getwd()
5
6 #Folder paths
7 wd_path = "C:/Users/tabe0010/OneDrive - Monash University/MR_backup_file/
8 Articles/Evangelou/new_sbp_snps/stata/r_codes/r_data"
9 data_path = "C:/Users/tabe0010/OneDrive - Monash University/MR_backup_file
10 /Articles/Evangelou/new_sbp_snps/stata/r_codes/r_data"
11 png_plot_path = "C:/Users/tabe0010/OneDrive - Monash University/MR_backup_
12 file/Articles/Evangelou/new_sbp_snps/stata/stata_sbp_plot/png"
13 pdf_plot_path = "C:/Users/tabe0010/OneDrive - Monash University/MR_backup_
14 file/Articles/Evangelou/new_sbp_snps/stata/stata_sbp_plot/pdf"
15
16 #Create MR graphs from main analysis
17 setwd(wd_path)
18
19 #install.packages("TwoSampleMR", repos = c("https://mrcieu.r-universe.dev
20 ", "https://cloud.r-project.org"))
21 #install.packages("devtools")
22 #devtools::install_github("MRCIEU/MRInstruments")
23 #install.packages("data.table")
24
25 library(TwoSampleMR)
26 library(MRInstruments)
27 library(data.table)
28 library("ggplot2")
29 library(plyr); library(dplyr)
30 #install.packages("forestplot")
31 library(forestplot)
32
33 #####
34
35 #while (dev.cur() > 1) dev.off()
36 master_data = read.csv("metan_sbp_r.csv", stringsAsFactors = FALSE)
37 qaly_hes<-master_data%>%filter(outcome=="qaly_hes")
38
39 x_list_qaly_1 = unique(qaly_hes$outcome)
40 qaly_hes = qaly_hes[order(qaly_hes$outcome, qaly_hes$type),]
41 tabletext_1 = cbind(unique(qaly_hes$label), (paste(qaly_hes$effect[qaly_hes
42 $type=="Main Analysis MR"], "\n", "\n", qaly_hes$effect[qaly_hes$type=="
43 Multivariable Adjusted"], sep="")))
44 colnames(tabletext_1)<-c("variable", "effect")
45 hrzl_lines = list("1"=gpar(lty=0), "2"=gpar(lty=1), "3"=gpar(lty=1), "4"=
46 gpar(lty=2), "5"=gpar(lty=1), "6"=gpar(lty=2),
47 "7"=gpar(lty=2), "8"=gpar(lty=2), "9"=gpar(lty=2), "10" =
48 gpar(lty=1), "11" = gpar(lty=2), "12"=gpar(lty=2))
```

```

40
41
42
43 #png("qaly_hes.png", width = 1000, height = 1000)
44 #pdf("qaly_hes.pdf", width = 12, height = 12)
45 png(file.path(png_plot_path, "qaly_hes.png"), width = 1000, height = 1000)
46 #pdf(file.path(pdf_plot_path, "qaly_hes.pdf"), width = 12, height = 12)
47
48 forestplot(tabletext_1,
49             legend = c("Main Analysis MR", "Multivariable Adjusted"),
50             title = "QALYs per year",
51             mean = cbind(qaly_hes$beta[qaly_hes$type == "Main Analysis MR"
52                             ], qaly_hes$beta[qaly_hes$type == "Multivariable Adjusted"])
53             ,
54             lower = cbind(qaly_hes$lower[qaly_hes$type == "Main Analysis MR"
55                             ], qaly_hes$lower[qaly_hes$type == "Multivariable Adjusted"
56                             ]),
57             upper = cbind(qaly_hes$upper[qaly_hes$type == "Main Analysis MR"
58                             ], qaly_hes$upper[qaly_hes$type == "Multivariable Adjusted"
59                             ]),
60             col=fpColors(box=c("blue", "darkred"),
61                             zero=c("darkblue")),
62             boxsize = 0.1,
63             line.margin = 0.2,
64             #xticks = c(-5,-4,-3,-2,-1,0,1,2,3),
65             xticks = c(-3,-2,-1,0,1,2,3),
66             grid = FALSE,
67             hrzl_lines=hrzl_lines,
68             txt_gp = fpTxtGp(xlab=gpar(cex=1.2),
69                             ticks = gpar(cex=1),
70                             summary = gpar(cex=1.5),
71                             title = gpar(cex=1.75),
72                             legend = gpar(cex=1.3),
73                             label = list(gpar(cex=1.2),gpar(1.5),gpar(cex
74                                     =0.8))),
75             xlab = "Percentage change in QALYs for 10mmHg increase in SBP"
76
77 ) %>%
78   fp_add_header(effect = "Estimate (95% CI)" # Add custom header
79   dev.off()

```

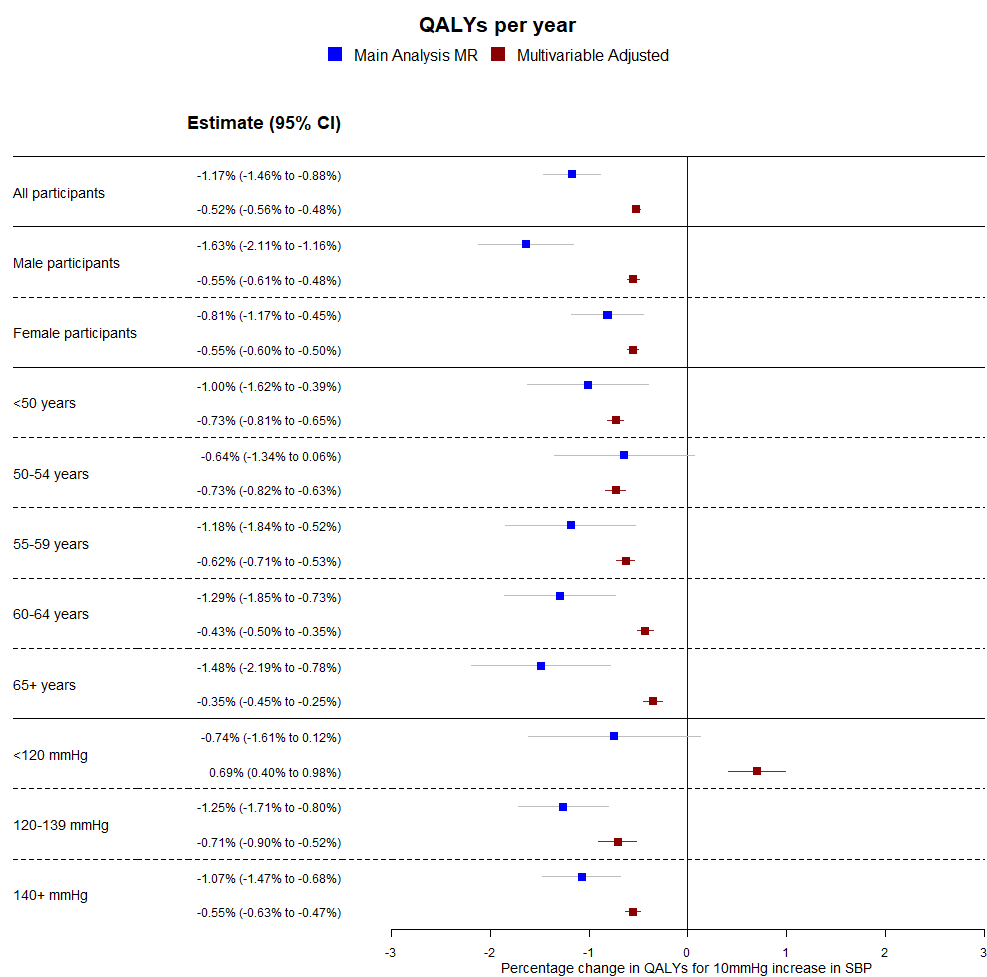



Figure 3: Mendelian randomisation estimates for QALYs on average year of follow-up

4.7.2.4 Figure 4

Associations between SNPs associated with systolic blood pressure and quality adjusted life years.

```
1
2 #R plots
3 # We will continue working on the previous R plot environment
4
5 dat <- read.csv("mr_analysis_sbp.csv", stringsAsFactors = FALSE)
6 dat <- rename(dat, beta.outcome = beta_outcome, se.outcome = se_outcome,
7               pval.outcome = pval_outcome, beta.exposure = beta_exposure,
8               se.exposure = se_exposure, id.exposure = id_exposure, id.
9               outcome = id_outcome)
10
11 dat$mr_keep <- TRUE
12 dat$exposure <- "SBP"
13
14 # MR analysis
15 setwd(paste(wd_path, sep=""))
16 res <- mr(dat)
17 p1 <- mr_scatter_plot(res, dat)
18
19 x = res[res$method == "Inverse variance weighted" | res$method == "Wald
20       ratio", c("outcome", "exposure")]
21
22 # Update axis labels, background theme, move legend to bottom, and remove
23   grid for each plot
24 for(i in 1:length(p1)){
25   outcome <- res[res$method == "Inverse variance weighted" | res$method ==
26     "Wald ratio", "outcome"][i]
27
28   # Customize the plot with new labels, white background, legend at the
29     bottom, and no grid
30   p1[[i]] <- p1[[i]] +
31     labs(x = "SNP effect on SBP", y = "SNP effect on QALYs") + # Change
32       these to your desired labels
33     theme_bw() + # Set background theme to white
34     theme(legend.position = "bottom", # Move legend to bottom
35           panel.grid = element_blank()) # Remove grid
36
37   # Save the modified plot
38   #ggsave(p1[[i]], file=paste("IVW - ", outcome, ".png", sep=""), width=7,
39     height=7)
40   ggsave(p1[[i]], filename = paste(png_plot_path, "/IVW - ", outcome, ".
41     png", sep=""), width=7, height=7)
42 }
```

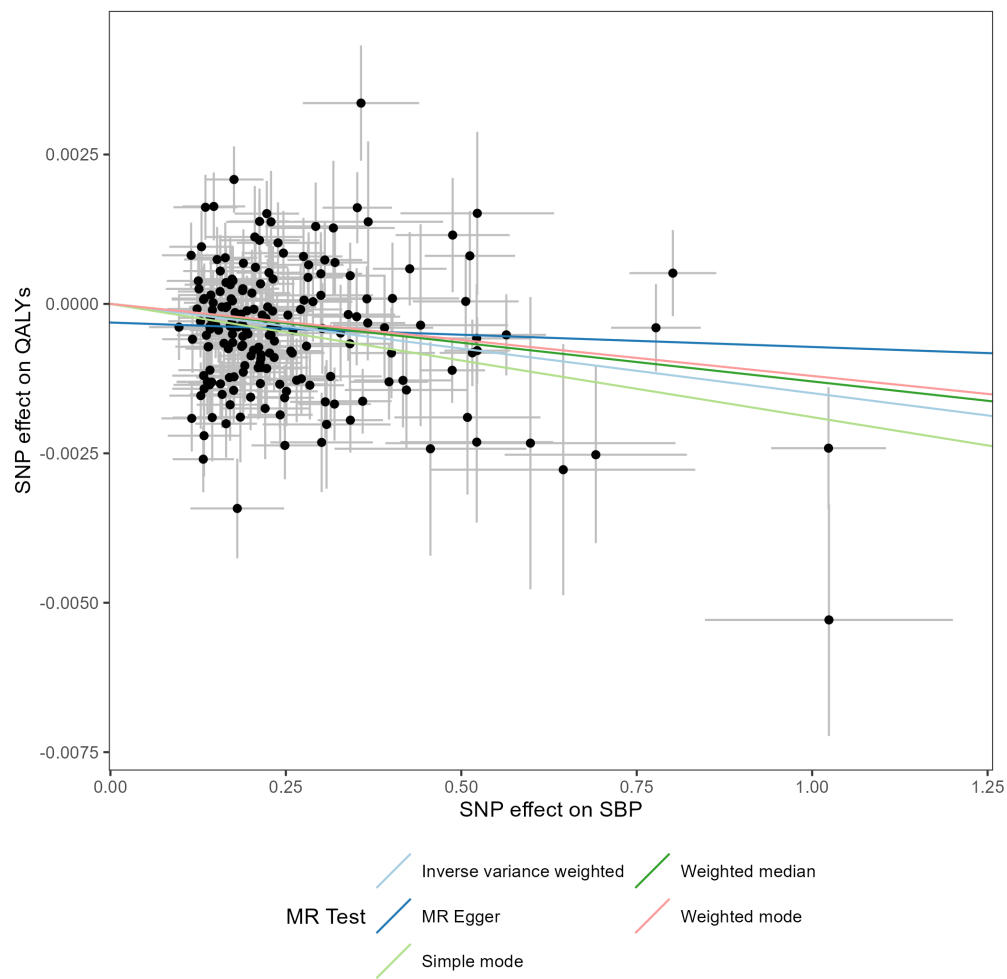


Figure 4: Associations between SNPs associated with systolic blood pressure and quality adjusted life years.

4.7.2.5 Figure 5

Non-linear MR showing the estimated effects of 1 mmHg increase in SBP on QALYs over an average year of follow-up, across SBP levels.

A positive value indicates an increase in SBP would increase in QALYs, and vice versa. There was little evidence of nonlinearity in the effect of SBP on QALYs. The SBP thresholds of 120 mmHg (for pre-hypertension SBP) and 140 mmHg (hypertension) are represented with dashed red lines. The green shaded area represents the 95

The code for this figure already done in **Section 4.5.4**.

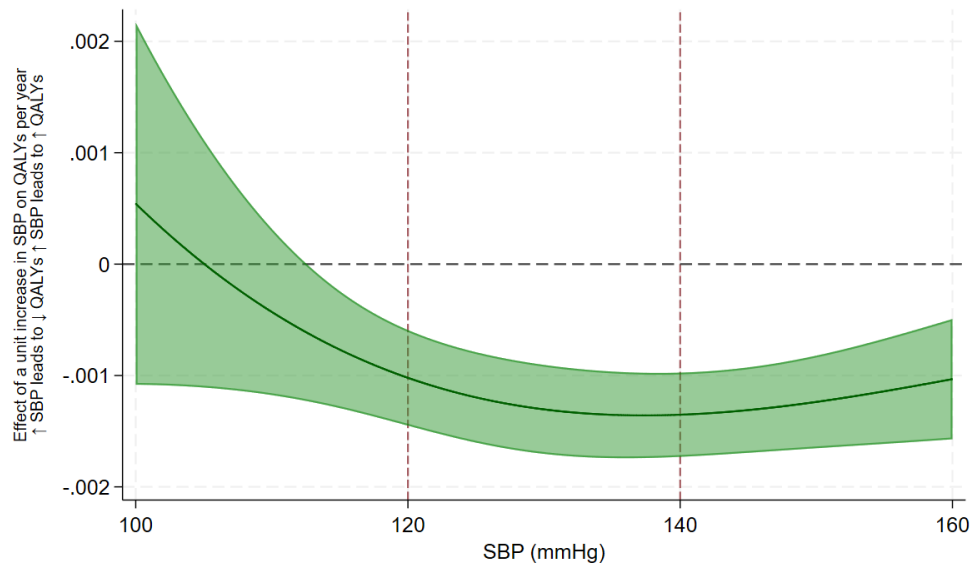


Figure 5: Estimated effects of 1 mmHg increase in SBP on QALYs over an average year of follow-up, across SBP levels.

References

- [1] Sean Harrison, Padraig Dixon, Hayley E Jones, Alisha R Davies, Laura D Howe, and Neil M Davies. Long-term cost-effectiveness of interventions for obesity: A mendelian randomisation study. *PLoS medicine*, 18(8):e1003725, 2021.
- [2] Jonathan C Brown, Thomas E Gerhardt, and Edward Kwon. Risk factors for coronary artery disease. <https://www.ncbi.nlm.nih.gov/books/NBK554410/>, 2023. Accessed 20-03-2025.
- [3] Seamus P Whelton, John W McEvoy, Leslee Shaw, Bruce M Psaty, Joao AC Lima, Matthew Budoff, Khurram Nasir, Moyses Szklo, Roger S Blumenthal, and Michael J Blaha. Association of normal systolic blood pressure level with cardiovascular disease in the absence of risk factors. *JAMA cardiology*, 5(9):1011–1018, 2020.
- [4] Brian A Ference, Deepak L Bhatt, Alberico L Catapano, Chris J Packard, Ian Graham, Stephen Kaptoge, Thatcher B Ference, Qi Guo, Ulrich Laufs, Christian T Ruff, et al. Association of genetic variants related to combined exposure to lower low-density lipoproteins and lower systolic blood pressure with lifetime risk of cardiovascular disease. *Jama*, 322(14):1381–1391, 2019.
- [5] Eric Yuk Fai Wan, Wing Tung Fung, C Mary Schooling, Shiu Lun Au Yeung, Man Ki Kwok, Esther Yee Tak Yu, Yuan Wang, Esther Wai Yin Chan, Ian Chi Kei Wong, and Cindy Lo Kuen Lam. Blood pressure and risk of cardiovascular disease in uk biobank: a mendelian randomization study. *Hypertension*, 77(2):367–375, 2021.
- [6] Debbie A Lawlor, Roger M Harbord, Jonathan AC Sterne, Nic Timpson, and George Davey Smith. Mendelian randomization: using genes as instruments for making causal inferences in epidemiology. *Statistics in medicine*, 27(8):1133–1163, 2008.
- [7] Patrick W Sullivan, Julia F Slejko, Mark J Sculpher, and Vahram Ghushchyan. Catalogue of eq-5d scores for the united kingdom. *Medical Decision Making*, 31(6):800–804, 2011.
- [8] Cathie Sudlow, John Gallacher, Naomi Allen, Valerie Beral, Paul Burton, John Danesh, Paul Downey, Paul Elliott, Jane Green, Martin Landray, et al. Uk biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS medicine*, 12(3):e1001779, 2015.
- [9] Martin D Tobin, Nuala A Sheehan, Katrina J Scurrah, and Paul R Burton. Adjusting for treatment effects in studies of quantitative traits: antihypertensive therapy and systolic blood pressure. *Statistics in medicine*, 24(19):2911–2935, 2005.
- [10] Nancy J Devlin, Koonal K Shah, Yan Feng, Brendan Mulhern, and Ben Van Hout. Valuing health-related quality of life: An eq-5 d-5 l value set for e ngland. *Health economics*, 27(1):7–22, 2018.
- [11] Stephen Burgess and Simon G Thompson. *Mendelian randomization: methods for using genetic variants in causal estimation*. CRC Press, 2015.
- [12] Neil M Davies, Michael V Holmes, and George Davey Smith. Reading mendelian randomisation studies: a guide, glossary, and checklist for clinicians. *bmj*, 362, 2018.

- [13] Evangelos Evangelou, Helen R Warren, David Mosen-Ansorena, Borbala Mifsud, Raha Pazoki, He Gao, Georgios Ntritsos, Niki Dimou, Claudia P Cabrera, Ibrahim Karaman, et al. Genetic analysis of over 1 million people identifies 535 new loci associated with blood pressure traits. *Nature genetics*, 50(10):1412–1425, 2018.
- [14] George Davey Smith and Shah Ebrahim. ‘mendelian randomization’: can genetic epidemiology contribute to understanding environmental determinants of disease? *International journal of epidemiology*, 32(1):1–22, 2003.
- [15] James Durbin. Errors in variables. *Revue de l’institut International de Statistique*, pages 23–32, 1954.
- [16] De-Min Wu. Alternative tests of independence between stochastic regressors and disturbances. *Econometrica: journal of the Econometric Society*, pages 733–750, 1973.
- [17] Jerry A Hausman. Specification tests in econometrics. *Econometrica: Journal of the econometric society*, pages 1251–1271, 1978.
- [18] Clare Bycroft, Colin Freeman, Desislava Petkova, Gavin Band, Lloyd T Elliott, Kevin Sharp, Allan Motyer, Damjan Vukcevic, Olivier Delaneau, Jared O’Connell, et al. The uk biobank resource with deep phenotyping and genomic data. *Nature*, 562(7726):203–209, 2018.
- [19] William G Cochran. The combination of estimates from different experiments. *Biometrics*, 10(1):101–129, 1954.
- [20] UK Biobank. Pain web questionnaire version 2.1, 2022. Accessed: 30-Mar-2025.
- [21] UK Biobank. Mental well-being web questionnaire version 1.1, 2023. Accessed: 30-Mar-2025.
- [22] Stephen Burgess and Simon G Thompson. Multivariable mendelian randomization: the use of pleiotropic genetic variants to estimate causal effects. *American journal of epidemiology*, 181(4):251–260, 2015.
- [23] Eleanor Sanderson, George Davey Smith, Frank Windmeijer, and Jack Bowden. An examination of multivariable mendelian randomization in the single-sample and two-sample summary data settings. *International journal of epidemiology*, 48(3):713–727, 2019.
- [24] Eleanor Sanderson. Multivariable mendelian randomization and mediation. *Cold Spring Harbor perspectives in medicine*, 11(2):a038984, 2021.
- [25] Helen R Warren, Evangelos Evangelou, Claudia P Cabrera, He Gao, Meixia Ren, Borbala Mifsud, Ioanna Ntalla, Praveen Surendran, Chunyu Liu, James P Cook, et al. Genome-wide association analysis identifies novel blood pressure loci and offers biological insights into cardiovascular risk. *Nature genetics*, 49(3):403–415, 2017.