

Heterogeneity in the relationship between BMI and risk biomarkers

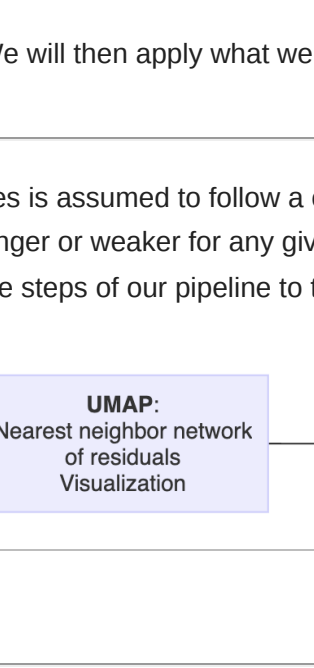
Authors:

- Daniel E. Coral
- Femke Smit
- Elena Santos
- Ali Farzaeehi

Introduction

This is a guideline to be followed by analysts in SOPHIA who are participating in the cross-sectional clustering project in the general population in Working Group 1. The idea is to standardize every step of the analysis across cohorts.

Not everyone has to follow all the steps, as we require different things from every cohort. We have divided the participating cohorts into 4 groups:



This guideline is designed to be applied in discovery and validation cohorts. We will then apply what we learn in these cohorts to the mental health and intervention cohorts.

As a background, generally the relationship between BMI and multiple diseases is assumed to follow a continuum – the higher the BMI, the higher the risk. However, it has also been found that in certain groups of people this relationship is disproportionately stronger or weaker for any given BMI. Our objective is to test the hypothesis that clustering-based approaches can be used to better capture these subgroups. An overview of the steps of our pipeline to test this hypothesis is shown in Figure 1.

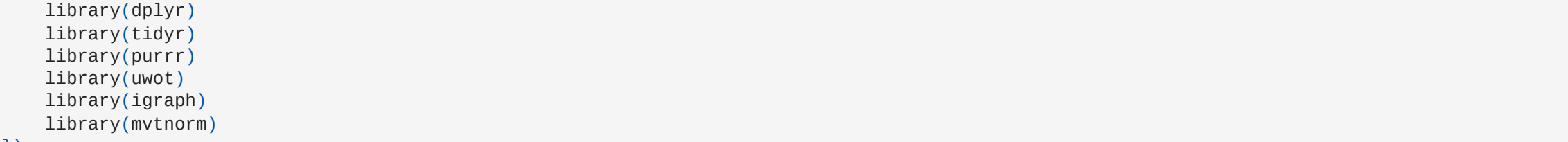


Figure 1 Pipeline overview

Software

All steps are intended to be followed in the R environment. To facilitate the analyses we have put together a list of functions that can be used to run every step of this guideline. They are located in the accompanying file `cross_sectional_FX.R`, which you can find like this:

```
In [1]: source("cross_sectional_FX.R")
```

These functions have dependencies on the following packages:

```
In [2]: suppressMessages({
  library(tibble)
  library(readr)
  library(dplyr)
  library(tidy)
  library(purrr)
  library(wget)
  library(sigraph)
  library(mvtnorm)
})
```

The environment where discovery analysis in UK Biobank was executed is then the following:

```
In [3]: sessionInfo()

R version 4.1.2 (2021-11-01)
Platform: x86_64-conda-linux-gnu (64-bit)
Running under: RHEL

Matrix products: default
BLAS/LAPACK: /gpfs/gpfs0/home/daniel_c/miniconda3/envs/Next/1/libopenblas-r0.3.18.so

locale:
 [1] C

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  methods   base

other attached packages:
[1] mvtnorm_1.1-3      sigraph_1.3.0      readr_1.3.0      dplyr_1.1.2      tibble_3.2.1      wget_0.1.1.9000    Matrix_1.4-1
[8] purrr_3.0.1        tidyrr_1.3.0      lwol_0.1.1.9000 readr_1.1.2      readr_2.1.4

loaded via a namespace (and not attached):
 [1] Rcpp_1.0.11      pillar_1.0.9      compiler_4.1.2    base64enc_0.1-3
 [5] tools_4.1.2      digest_0.6.33     uuid_1.0-4        jsonlite_1.8.7
 [9] evaluate_0.21    lifecycle_1.0.3   lattice_0.20-45   pkgconfig_2.0.3
[13] rlang_1.1.1      RMarkdown_2.11    cli_3.0.6         R.methodsS3_1.3
[17] fastmap_1.1.1    rproj4_1.1.4       generics_0.1.3    vctrs_0.6.0
[21] hms_1.1.3        grid_4.1.2        tidyrselect_1.2.0 glue_1.6.2
[25] nlme_2.5.1        fast1_2.0.4        plot2mg_0.2-7     t2d3_0.4-9
[29] magrittr_2.0.3   httools_0.5.5     utf8_1.2.3        crayon_1.5.2
```

Initial input

We have selected 10 traits, based on biological systems that are commonly affected by obesity:

- Blood pressure: SBP and DBP
 - Lipids: HDL, LDL, TG
 - Fat distribution: WHR
 - Glycemia: Fasting glucose
 - Liver metabolism: ALT
 - Kidney function: Creatinine
 - Inflammation: CRP
- The covariates that will be needed are:
- Sex
 - Age
 - Current smoking status, coded as 1 if current smoker, 0 otherwise.

The initial input table should look like the following:

```
In [4]: recoded_dat <- read_tsv("../data/recoded_dat.tsv", show_col_types = FALSE)
head(recoded_dat)
```

A tibble: 6 × 15														
eid	age	sex	bmi	whr	sbp	dbp	alt	scr	crp	hdl	tg	ldl	fg	smoking
<dbl>	<dbl>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1000039	44	Male	36.6959	0.9911504	124.5	64.5	34.87	93.0	3.60	1.158	2.800	3.956	5.427	0
1000071	67	Male	39.4807	0.8857143	179.5	103.0	46.74	68.7	9.41	1.372	1.127	2.311	7.079	0
1000088	60	Male	24.2786	0.8761905	152.0	89.0	13.14	80.6	1.20	0.983	1.590	4.200	5.401	0
1000096	41	Male	26.5744	0.9587629	143.0	90.0	30.32	80.1	6.13	1.041	2.713	4.029	4.239	0
1000109	62	Male	33.8719	1.0818182	156.5	104.5	16.26	89.3	14.42	0.890	2.437	3.525	6.100	0
1000117	47	Female	23.8408	0.7254902	147.5	84.0	14.07	61.0	0.24	1.972	0.591	2.252	4.395	0

The units of the continuous variables we expect to run the analysis:

- Age in years
- Sex as a string of either "Female" or "Male".
- BMI in kg/m²
- WHR is unitless, calculated by dividing waist and hip circumferences measured in cm
- SBP and DBP in mmHg
- ALT in U/L
- SCR in umol/L
- CRP in mg/L
- HDL in mmol/L
- TG in mmol/L
- LDL in mmol/L
- FG in mmol/L
- Smoking as a dummy variable: 1 if currently smoking, 0 otherwise.

Note

For the functions included in the `cross_sectional_FX.R` to work, the input table should be **exactly** as shown above.

Missing data

Since our clustering method ignores individuals with missing values for any biomarker, the input data should only contain individuals who have all biomarker values. From the previous analyses, we know that:

- Only including complete cases without losing too much data is possible in UK Biobank, Maastricht, GHS and ABOS.
- In Rotterdam the initial input table contains some values have been imputed using a random forest algorithm.
- In Grona only a small subset of individuals have CRP values.
- In SCALE there are values for waist but not for hip circumference, so it is not possible to calculate WHR.

Based on these observations we have waited for the following decisions on how to deal with missing values:

How to handle missing values	Cohorts
Only include complete cases.	UK Biobank
	Maastricht
	GHS
	ABOS
Use data that has been imputed	Rotterdam
	Grona
	SCALE

For the cohorts in the last group, to be able to apply our method, we will assume that BMI explains the variability in the biomarkers that are missing. This assumption is based on what we have observed in the other cohorts. In practice, this means that the clustering method will focus on the biomarkers that are available to group individuals into clusters. The input table should still have the same columns so that the functions in `cross_sectional_FX.R` work properly.

Remove possible errors in measurement

In discovery and validation cohorts we will exclude biomarker measurements that are 5 SD away from the mean, under the assumption that these are most likely measurement errors. This can be done using the `remove_outliers` function that we have provided, which replaces outliers with `NA` values. Then we again make sure to have only complete cases:

```
In [5]: recoded_dat <- mutate(recoded_dat,
  across(c(bmi, whr, scr, dbp, alt, scr, crp, hdl, tg, ldl, fg),
    ~remove_outliers(x, sdunits = 5)))
recoded_dat <- recoded_dat[complete.cases(recoded_dat),]
```

Stratify by sex

All the pipeline is applied separately in each sex group. The functions we have designed work on a list containing two dataframes for each sex group, which we can obtain like this:

```
In [6]: strat_dat <- split(recoded_dat, ~sex)

To see the first lines of the two elements in the list:
```

```
In [7]: lapply(strat_dat, head)
```

A tibble: 6 × 15														
eid	age	sex	bmi	whr	sbp	dbp	alt	scr	crp	hdl	tg	ldl	fg	smoking
<dbl>	<dbl>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1000117	47	Female	23.8408	0.7254902	147.5	84.0	14.07	61.0	0.24	1.972	0.591	2.252	4.395	0
1000132	43	Female	35.0559	0.8403361	137.0	100.5	18.89	60.5	4.31	1.236	2.037	3.686	5.214	0
1000176	69	Female	38.1271	0.8897638	137.5	93.5	36.39	68.9	3.69	1.601	1.988	4.551	4.266	0
1000223	63	Female	25.4603	0.7789474	163.0	94.0	6.10	67.1	1.29	1.453	2.629	3.491	5.876	0
1000282	48	Female	25.4297	0.7786333	135.5	89.0	9.63	46.2	0.16	2.185	0.722	3.584	5.212	0
1000387	42	Female	19.3280	0.6777778	107.0	72.5	9.34	57.1	0.69	2.346	0.395	3.072	4.649	0

A tibble: 6 × 15														
eid	age	sex	bmi	whr	sbp	dbp	alt	scr	crp	hdl	tg	ldl	fg	smoking
<dbl>	<dbl>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1000039	44	Male	36.6959	0.9911504	124.5	64.5	34.87	93.0	3.60	1.158	2.800	3.956	5.427	0
1000071	67	Male	39.4807	0.8857143	179.5	103.0	46.74	68.7	9.41	1.372	1.127	2.311	7.079	0
1000088	60	Male	24.2786	0.8761905	152.0	89.0	13.14	80.6	1.20	0.983	1.590	4.200	5.401	0
1000096	41	Male	26.5744	0.9587629	143.0	90.0	30.32	80.1	6.13	1.041	2.713	4.029	4.239	0
1000109	62	Male	33.8719	1.0818182	156.5	104.5	16.26	89.3	14.42	0.890	2.437	3.525	6.100	0
1000125	66	Male	36.1100	1.0625000	155.0	102.5	25.59	88.7	1.91	1.061	1.320	2.538	4.531	1

Summary of initial input

We need a table summarising the initial input, which can be generated like this:

```
In [8]: gendesc_tab <- get_general_descriptives(strat_dat)
gendesc_tab
```

A data frame: 28 × 7						
sex	Variable	Type	N	N miss	Summary1	Summary2
<chr>	<chr>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
Female	bmi	Numeric	77207	0	27.11 (5.07)	26.18 (19.68 - 23.52 - 29.78 - 39.67)
Female	age	Numeric	77207	0	57.45 (8.04)	59 (42 - 52 - 64 - 69)
Female	smoking	Categorical	77207	0	0	70355 (91.13%)
Female	whr	Numeric	77207	0	1	6852 (8.87%)
Female	sbp	Numeric	77207	0	0.82 (0.07)	0.81 (0.7 - 0.77 - 0.86 - 0.96)
Female	dbp	Numeric	77207	0	137.51 (14.47)	135.5 (105.5 - 123.5 - 150 - 180)
Female	alt	Numeric	77207	0	81.49 (9.92)	81 (63.5 - 74.5 - 88 - 102)
Female	scr	Numeric	77207	0	19.96 (9.71)	17.57 (8.97 - 13.96 - 23 - 46.41)
Female	crp	Numeric	77207	0	64.31 (10.74)	63.2 (46.9 - 57.2 - 70 - 88.1)
Female	hdl	Numeric	77207	0	2.52 (1.19)	1.42 (0.21 - 0.67 - 2.99 - 12.33)
Female	tg	Numeric	77207	0	1.61 (0.38)	1.57 (0.99 - 1.34 - 1.88 - 2.46)
Female	ldl	Numeric	77207	0	1.54 (0.8)	1.33 (0.6 - 0.96 - 1.84 - 3.54)
Female	fg	Numeric	77207	0	3.67 (0.87)	3.62 (2.14 - 3.05 - 4.28 - 6.76)
Female	tg	Numeric	67904	0	4.97 (0.6)	4.91 (4.06 - 4.63 - 5.22 - 6.33)
Male	bmi	Numeric	67904	0	27.9 (4.23)	27.37 (21.02 - 25.05 - 30.14 - 37.82)
Male	age	Numeric	67904	0	57.74 (8.02)	60 (41 - 52 - 64 - 69)
Male	smoking	Categorical	67904	0	0	59703 (87.92%)
Male	whr	Numeric	67904	0	1	8201 (12.08%)
Male	sbp	Numeric	67904	0	0.94 (0.06)	0.94 (0.81 - 0.9 - 0.98 - 1.07)
Male	dbp	Numeric	67904	0	142.31 (17.66)	141 (112 - 130 - 153 - 181)
Male	alt	Numeric	67904	0	84.94 (9.58)	84.5 (66 - 76 - 91.5 - 105.5)
Male	scr	Numeric	67904	0	26.59 (12.33)	23.64 (11.25 - 18.34 - 31.43 - 60.22)
Male	crp	Numeric	67904	0	81.48 (13.07)	80.2 (59.8 - 72.8 - 88.4 - 111.2)
Male	hdl	Numeric	67904	0	2.24 (2.07)	1.3 (0.22 - 0.68 - 2.56 - 11.13)
Male	tg	Numeric	67904	0	1.3 (0.32)	1.25 (0.81 - 1.08 - 1.47 - 2.54)
Male	ldl	Numeric	67904	0	1.89 (1.01)	1.65 (0.65 - 1.26 - 1.45 - 4.07)
Male	fg	Numeric	67904	0	3.5 (0.87)	3.48 (1.91 - 2.89 - 4.08 - 5.28)
Male	tg	Numeric	67904	0	5.03 (0.7)	4.94 (3.95 - 4.64 - 5.28 - 6.76)

For smoking:

- Summary1 contains the categories.
- Summary2 contains the proportion of each category.

For the rest (continuous) variables:

- Summary1 contains the mean and standard deviation.
- Summary2 contains the median and percentiles 2.5, 25, 75 and 97.5.

Estimates of BMI-biomarker associations

The first step of the pipeline is to generate sex-specific linear models of BMI for each variable, adjusting for age, and smoking. To do that we have the following function:

```
In [9]: mods <- get_bmi_mods(strat_dat)
```

The result is a table with a column that contains the models specific for each sex and biomarker:

```
In [10]: print(mods)
```

```
# A tibble: 28 × 3
  sex   biomarker and
  <chr> <chr>      <list>
1 Female whr    <lm>
2 Female sbp    <lm>
3 Female dbp    <lm>
4 Female alt    <lm>
5 Female scr    <lm>
6 Female crp    <lm>
7 Female hdl    <lm>
8 Female tg     <lm>
9 Female ldl    <lm>
10 Male whr     <lm>
11 Male sbp     <lm>
12 Male dbp     <lm>
13 Male alt     <lm>
14 Male scr     <lm>
15 Male hdl     <lm>
16 Male tg      <lm>
17 Male ldl     <lm>
18 Male fg      <lm>
19 Male ldl     <lm>
20 Male fg      <lm>
```

As an example, we can print the summary of the female model for CRP:

```
In [11]: summary(mods$mod[mods$sex == "Female" & mods$biomarker == "crp"][[1]])
```

```
Call:
lm(formula = reformulate(response = biomarker, termLabels = c(BODY$SIZEINDEX,
COVARIATES)), data = X, na.action = na.exclude)

Residuals:
    Min       1Q   Median       3Q      Max
-7.7411 -1.4397 -0.6721  0.3533 24.1163

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -5.494409    0.093488  -58.77   <2e-16 ***
bmi           0.251869    0.002697   121.88   <2e-16 ***
age           0.018570    0.001341   14.60   <2e-16 ***
smoking       0.735783    0.036928   19.93   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.966 on 77203 degrees of freedom
Multiple R-squared:  0.1677, Adjusted R-squared:  0.1677
F-statistic: 5187 on 3 and 77203 DF, p-value: < 2.2e-16

We then use this table to generate a table containing the estimates of the effect of BMI as well as the covariates on every biomarker:
```

```
In [12]: bmi_coeffs_tab <- get_bmi_coeffs(mods)
bmi_coeffs_tab
```

A tibble: 80 × 7						
sex	Biomarker	term	Estimate	SE	lowerCI	upperCI
<chr>	<chr>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
Female	whr	(Intercept)	0.57262	0.00195	0.56879	0.57645
Female	whr	bmi	0.00623	0.00004	0.00615	0.00632
Female	whr	age	0.00131	0.00003	0.00125	0.00136
Female	whr	smoking	0.02104	0.00077	0.01952	0.02255
Female	sbp	(Intercept)	69.54710	0.57139	68.42717	70.66703
Female	sbp	bmi	0.59574	0.01263	0.57099	0.62050
Female	sbp	age	0.90572	0.00619	0.88966	0.92178
Female	sbp	smoking	-2.46168	0.22570	-2.90406	-2.01930
Female	dbp	(Intercept)	63.15773	0.30613	62.55989	63.75578
Female	dbp	bmi	0.59594	0.00674	0.54532	0.57276
Female	dbp	age	0.05639	0.00438	0.04791	0.06497
Female	alt	(Intercept)	1.47995	0.30069	0.62039	1.79991
Female	alt	bmi	0.20959	0.00665	0.46296	0.48902
Female	alt	age	0.10352	0.00431	0.09507	0.11197
Female	alt	smoking	-0.96029	0.11878	-1.19309	-0.72749
Female	scr	(Intercept)	53.38579	0.34249	52.71451	54.05708
Female	scr	age	0.08609	0.00451	0.07646	0.09572
Female	scr	smoking	-1.28687	0.13629	-1.55203	-1.02170
Female	crp	(Intercept)	-5.49			