

Cluster probabilities and outcomes

Authors:

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

In this script we get the outcomes and we examine how the clusters we have validated across cohorts are associated with prevalent diseases at the time of clustering, and also assess whether they add significant information for prediction of MACE events and diabetes progression on top of commonly used risk stratification tools.

Libraries and functions

```
[In [1]] library(readr)
library(dplyr, warn.conflicts = FALSE)
library(tidyverse)
library(purrr)
library(R.utils)

And the functions we have prepared to facilitate some steps:

[In [2]] source("cross_sectional_FX2.R")
```

Loading data needed

Initial input table of biomarkers and basic covariates

The input table is the same table of 17 traits we had prior to run UNMAP. Here is a description of this table:

System targeted	Biomarker	Units	Column name
Individual ID	-	-	id
Blood pressure	Systolic blood pressure	millimeters of mercury (mmHg)	sbp
	Diastolic blood pressure	millimeters of mercury (mmHg)	dbp
Lipid fractions	High-density lipoprotein	mmol/L	hdl
	Low-density lipoprotein	mmol/L	ldl
	Triglycerides	mmol/L	tg
Glycemia	Fasting glucose	mmol/L	fg
Liver metabolism	Aspartate transaminase	U/L	alt
Fat distribution	Waist-to-hip ratio	cm/cm	whr
Kidney function	Serum creatinine	umol/L	scr
Inflammation	C-reactive protein	mg/L	crp
Basic covariates	Current smoking status	1 if yes, 0 if not	smoking
	Sex	String ('Female' or 'Male')	sex
	Age	Years	age

Important note: All columns should be there in the units required, and the names should match, so that the functions we have prepared for the analyses work properly. This is true for this and all the following tables we require for our analysis.

This input table has been preprocessed by:

1. Filtering out values that are possible errors in measurement (± 5 SD away from the mean in continuous variables).
2. Only including complete cases.
3. Stratifying by sex.

Here is how the input table should look like - a list of two data frames, one for each sex:

```
[In [3]] load(female, male, start.dat, end.dat)

[In [4]] head(female, head)

#Female
# A tibble: 6 x 15
  id 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>
1 1000117 47 Female 23.8408 0.7254902 147.5 84.0 14.07 61.0 0.24 ... 0.991 2.352 4.395 0
2 1000120 47 Female 38.0558 0.8603501 137.0 100.5 18.89 60.5 0.31 1.296 2.007 3.088 5.214 0
3 1000126 69 Female 38.1271 0.8897638 137.5 93.5 36.39 66.9 1.601 1.969 4.561 4.266 0
4 1000223 63 Female 25.4603 0.7708474 163.0 84.0 6.10 67.1 1.29 1.453 2.629 3.491 5.876 0
5 1000282 42 Female 25.4297 0.7708333 135.5 89.0 9.63 46.2 0.16 2.165 0.722 3.084 5.212 0
6 1000387 42 Female 19.3230 0.6777778 107.0 72.5 9.34 57.1 0.89 2.346 0.395 3.072 4.646 0

#Male
# A tibble: 6 x 15
  id 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>
1 1000059 67 Male 36.6669 0.9911504 124.5 64.5 34.97 92.0 3.60 1.158 2.800 3.956 5.427 0
2 1000071 44 Male 36.4807 0.8887143 174.5 103.0 46.74 68.7 9.41 1.372 1.127 2.311 7.079 0
3 1000088 40 Male 24.7786 0.6762945 107.0 86.0 12.14 80.2 1.80 0.683 1.966 4.263 5.481 0
4 1000096 41 Male 26.5744 0.9261629 143.0 96.0 30.32 80.1 1.541 2.541 2.713 4.629 4.629 0
5 1000109 62 Male 33.8719 1.0316182 156.5 104.5 36.26 89.3 14.42 0.890 2.437 5.525 6.000 0
6 1000125 66 Male 36.1100 1.0620000 155.0 102.5 25.59 88.7 1.91 1.061 1.320 2.538 4.531 1
```

Table of validated clusters

The second thing needed is the clusters we have validated. We have put this in an R file called `validclusters`:

```
[In [5]] load(female, male, validclusters, start.dat)
print(validclusters)

# A tibble: 2 x 2
  sex outcome
  <chr> <list>
1 Female <list1 [77,287 x 2]>
2 Male <list1 [67,984 x 2]>

This object contains, for each sex:
• residded: The model to obtain residuals for each variable, i.e., the variability beyond what is explained by BMI, adjusting for age and smoking.
• clustered: The clustering model to apply to the residuals.
```

Table of pre-existing conditions and medications

The third thing we need is a table of pre-existing conditions and medications participants are currently taking.

```
[In [6]] covar.dat <- read_csv("../data/covar.dat.csv", show_col_types = FALSE)
head(covar.dat)

# A tibble: 6 x 10
  HT CKD CKD_Sbdo PAD CKD CKD_LiverFailure RA T2D T2Dage Insulin AntiDM AntiHT LipidLower
  <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 1000027 0 0 0 0 0 0 0 0 0 0 0 0 0
2 1000059 0 0 0 0 0 0 0 0 0 0 0 0 0
3 1000040 1 0 0 0 0 0 0 0 0 0 0 0 0
4 1000035 0 0 0 0 0 0 0 0 0 0 0 0 0
5 1000064 1 0 0 0 0 0 0 1 0 49.5 0 1 1
6 1000071 1 0 0 0 0 0 0 1 0 65.5 0 0 1 1
```

All the columns in this table are coded 1 or 0 representing current diagnosis of a disease or whether the person is taking the medications specified. The exception is `T2Dage`, which is the age of onset of T2D. This is what each column represent

Group	Cluster name	Meaning
Diagnoses	HT	Hypertension
	CKD	Chronic heart disease
	Stroke	Stroke
	PAD	Peripheral artery disease
	CKD	Chronic kidney disease
	LiveFailure	Liver failure
	RA	Rheumatoid arthritis
	T2D	Type 2 diabetes
	T1D	Type 1 diabetes
Age at onset	T2Dage	Age at onset of T2D - If not a T2D, is 0. Needed in SCORE2.
Medication	Insulin	Taking insulin
	AntiDM	Taking medication for diabetes other than insulin
	AntiHT	Taking medication for hypertension
	LipidLower	Taking lipid lowering medication

If any of the columns in this table are missing in your data, one option is to assume that none in your population had the disease, i.e., you should have a column with 0 for all individuals.

Survival data

Lastly, we need survival data for MACE and diabetes progression. They should look like this:

```
[In [7]] survivedat <- read_csv("../data/survivedat.csv", show_col_types = FALSE)
head(survivedat)

# A tibble: 6 x 3
  eid outcome_value outcome_timeys
  <dbl> <dbl> <dbl>
1 1000071 0 10.00369
2 1000120 1 6.91612
3 1000124 1 3.761897
4 1000283 1 3.761897
5 1001175 1 4.53667
6 1001382 1 9.18489
```

```
[In [8]] survivedat <- read_csv("../data/survivedat.csv", show_col_types = FALSE)
head(survivedat)

# A tibble: 6 x 3
  eid outcome_value outcome_timeys
  <dbl> <dbl> <dbl>
1 1000109 1 3.90274
2 1000120 1 2.34660
3 1004387 1 4.950309
4 1006281 1 1.907563
5 1007454 0 6.423662
6 1009195 1 6.802641
```

These two tables include individuals followed up to **age 30 years**. This means that any outcome after 10 years is censored. `outcome_value` is 1 if the person experienced the event during the follow-up time and 0 if not. `outcome_timeys` is the time of follow-up in years, up to the first event or 30 years.

It is important that these tables **do not include** individuals who already experienced the events we will study. In any case, we will make sure of this in the next step, when we combine all the data. For example, any individual in the `survivedat` table with a value of 1 in the columns `Stroke` or `PAD` or `CKD` or `CKD_Sbdo` will be excluded from the analysis.

In case your cohort does not have survival data, then follow this guideline until the section below entitled "Prevalent diseases and medication".

Calculation of cluster probabilities

With the data needed in place, we can start by calculating cluster allocation probabilities given the biomarker data. For that we will first add a new column called `data` to the `validclusters` table where we will put the biomarker data for each sex:

```
[In [9]] clusterf <- clusterprobcalc(clusterf ~ validclusters, stratdat ~ strat.dat)

print(clusterf)

# A tibble: 2 x 2
  sex outcome
  <chr> <list>
1 Female <list1 [77,287 x 2]>
2 Male <list1 [67,984 x 2]>

Checking that the probabilities were calculated for each sex:
```

```
[In [10]] head(clusterfdata[[1]])

# A tibble: 6 x 21
  eid age sex bmi whr sbp dbp alt scr crp ... CKD LiveFailure RA T2D T2Dage Insulin AntiDM AntiHT LipidLower
  <dbl> <dbl> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 1000117 47 Female 23.8408 0.7254902 147.5 84.0 14.07 61.0 0.24 ... 0.991 2.352 4.395 0 0.009059 0.9821796-01 0.146269e-08 0.000624965 0.000023962 0.816396e-05
2 1000120 47 Female 38.0558 0.8603501 137.0 100.5 18.89 60.5 0.31 ... 2.037 2.052 4.216 0 0.3956035 0.842613e-01 0.669902e-03 0.0004612784 0.000098692 0.106767e-03
3 1000126 69 Female 38.1271 0.8897638 137.5 93.5 36.39 66.9 1.601 ... 1.988 4.561 4.266 0 0.9885208 0.907371e-06 5.304406e-05 0.012028762 0.000241394 0.725176e-05
4 1000223 63 Female 25.4603 0.7708474 163.0 84.0 6.10 67.1 1.29 ... 2.629 3.491 5.876 0 0.9391540 1.237050e-06 5.224040e-02 0.000243480 0.000154409 1.786232e-03
5 1000282 42 Female 25.4297 0.7708333 135.5 89.0 9.63 46.2 0.16 ... 0.703085 0.07705e-01 0.60971e-07 0.001804800 0.001070480 3.600200e-03
6 1000387 42 Female 19.3230 0.6777778 107.0 72.5 9.34 57.1 0.89 ... 0.395 3.072 4.646 0 0.9827697 2.826543e-03 2.959585e-07 0.003310679 0.000235946 1.772521e-03
```

Descriptive statistics

At this point we will check some of the characteristics of the clusters as we did in our previous script, weighting calculations by cluster probabilities.

The distribution of biomarkers per cluster:

```
[In [12]] markerdistrib <- markerdistrib(clusterf)
head(markerdistrib)

# A tibble: 6 x 7
  sex Variable Cluster Type N N_weighted Summary%
  <chr> <chr> <chr> <chr> <dbl> <dbl>
1 Female whr BC Numeric 5870.879 0.82 (0.7) 0.81 (0.7 - 0.76 - 0.86)
2 Female whr DBP Numeric 7483.777 0.79 (0.00) 0.78 (0.69 - 0.75 - 0.82 - 0.9)
3 Female whr DAL Numeric 2952.388 0.87 (0.00) 0.86 (0.76 - 0.83 - 0.91 - 0.99)
4 Female whr CLF Numeric 2055.284 0.84 (0.07) 0.84 (0.71 - 0.79 - 0.89 - 0.86)
5 Female whr DSS Numeric 2795.474 0.84 (0.07) 0.83 (0.71 - 0.79 - 0.89 - 0.9)
6 Female whr DMG Numeric 1477.887 0.85 (0.00) 0.85 (0.71 - 0.79 - 0.91 - 1.02)
```

The effect of BMI on biomarkers specifically within each cluster, adjusted for age and smoking:

```
[In [14]] bsafffkerf <- bsafffkerf(clusterf)
head(bsafffkerf)

# A tibble: 6 x 5
  sex Variable Cluster Type estimate se
  <chr> <chr> <chr> <chr> <dbl> <dbl>
1 Female whr BC (residual) 0.57875260 1.917426e-03
2 Female whr BC age 0.00228926 2.786274e-05
3 Female whr BC smoking 0.00211460 7.707979e-05
4 Female whr BC bmi 0.00209209 4.309794e-05
5 Female whr DBP (residual) 0.56864206 1.566971e-03
6 Female whr DBP age 0.00188705 2.086438e-05
```

Prevalent diseases and medication

To add covariate data to the `allData` table we will do the following:

```
[In [16]] clusterf <- addcovdat(x = clusterf, covar.dat ~ covar.dat)

print(clusterf)

# A tibble: 2 x 2
  sex outcome
  <chr> <list>
1 Female <list1 [73,378 x 3]>
2 Male <list1 [60,348 x 3]>
3 Female DM <list1 [34,581 x 3]>
4 Male DM <list1 [29,988 x 3]>

Checking again if the columns were added as expected:
```

```
[In [18]] head(clusterfdata[[1]])

# A tibble: 6 x 35
  eid age sex bmi whr sbp dbp alt scr crp ... CKD LiveFailure RA T2D T2Dage Insulin AntiDM AntiHT LipidLower
  <dbl> <dbl> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 1000117 47 Female 23.8408 0.7254902 147.5 84.0 14.07 61.0 0.24 ... 0 0 0 0 0 0 0 0 0 0 0
2 1000120 47 Female 38.0558 0.8603501 137.0 100.5 18.89 60.5 0.31 ... 0 0 0 0 0 0 0 0 0 0 1
3 1000126 69 Female 38.1271 0.8897638 137.5 93.5 36.39 66.9 1.601 ... 0 0 0 0 0 0 0 0 0 0 1
4 1000223 63 Female 25.4603 0.7708474 163.0 84.0 6.10 67.1 1.29 ... 0 0 0 0 0 0 0 0 0 0 1
5 1000282 42 Female 25.4297 0.7708333 135.5 89.0 9.63 46.2 0.16 ... 0 0 0 0 0 0 0 0 0 0 0
6 1000387 42 Female 19.3230 0.6777778 107.0 72.5 9.34 57.1 0.89 ... 0 0 0 0 0 0 0 0 0 0 0
```

We will first count the number of individuals with disease in each cluster. Here we will also count the number of individuals taking each class of medication in each cluster.

```
[In [19]] countcovarsf <- countcovarsf(clusterf)
head(countcovarsf)

# A tibble: 6 x 6
  sex Cluster Covariate Ncases Nnoncases
  <chr> <chr> <chr> <dbl> <dbl>
1 Female whr BC Numeric 5870.879 492056.11
2 Female whr CKD Numeric 58673.33 3659.0550 9704.93
3 Female whr CKD_Sbdo Numeric 132.00 12.00000
4 Female whr CKD_LiverFailure Numeric 635.00 67.00000
5 Female whr CKD_LiveFailure Numeric 10480.00 120.00000
6 Female whr CKD_PAD Numeric 8573.00 120.00000
```

We will use this table to calculate prevalences and compare prevalences across clusters.

We are also interesting in looking at the proportion of individuals receiving medications in each cluster, stratified by each condition. This is obtained with the following function:

```
[In [21]] countmedsf <- countmedsf(clusterf)
head(countmedsf)

# A tibble: 6 x 6
  sex Dx Cluster Med Ncases Nnoncases
  <chr> <chr> <chr> <chr> <dbl> <dbl>
1 Female CKD proBc Insulin 4350.640 200.00007
2 Female CKD proBc Insulin 132.00 12.00000
3 Female CKD proBc AntiDM 635.00 67.00000
4 Female CKD proBc AntiHT 10480.00 120.00000
5 Female CKD proBc LipidLower 8573.00 120.00000
6 Female CKD proBc NMAid 2795.48 9.48602
```

We will also formally test the association between cluster allocation and diseases using logistic regressions where the outcome is each disease and the predictors are the cluster allocations. We will have two models for each disease, one with only clusters, and a second one adjusting for medication.

```
[In [23]] assocdf <- assocdf(clusterf)
print(assocdf)

# A tibble: 36 x 5
  sex Dx cluster model estimates varcovmat Means AFit
  <chr> <chr> <chr> <list> <list> <list>
1 Female HT <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
2 Female HT <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
3 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
4 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
5 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
6 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
7 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
8 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
9 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
10 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
11 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
12 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
13 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
14 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
15 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
16 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
17 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
18 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
19 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
20 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
21 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
22 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
23 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
24 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
25 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
26 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
27 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
28 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
29 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
30 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
31 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
32 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
33 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
34 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
35 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
36 Female CKD <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]> <list1 [73,378 x 3]>
```

Adding survival data

As explained before, since we want to be careful when adding survival data for analysis, we have prepared a function separately for both outcomes, and making sure we exclude individuals who already experienced the events under study:

```
[In [25]] clusterSurvsf <- addcovdat(x = clusterf, survMACEsf ~ survMACEsf, survDMsf ~ survDMsf)

print(clusterSurvsf)

# A tibble: 4 x 3
  sex outcome data
  <chr> <chr> <list>
1 Female MACE <list1 [73,378 x 3]>
2 Male MACE <list1 [60,348 x 3]>
3 Female DM <list1 [34,581 x 3]>
4 Male DM <list1 [29,988 x 3]>

data no contain the data necessary to run survival analysis.
```

Rates of outcomes by cluster

Similar to what was done in the cross sections setting, we will calculate the number of cases and the total follow-up in each cluster using the weighted approach:

```
[In [27]] ratesbyclus <- ratesbyclusf(clusterSurvsf)
head(ratesbyclus)

# A tibble: 6 x 5
  sex outcome Cluster Ncases TPT
  <chr> <chr> <chr> <dbl> &lt
```