

EpiMetal: An open-source graphical web browser tool for easy statistical analyses in epidemiology and metabolomics

Journal:	<i>International Journal of Epidemiology</i>
Manuscript ID	IJE-2019-02-0272.R1
Manuscript Type:	Software Application Profile
Date Submitted by the Author:	n/a
Complete List of Authors:	Ekholm, Jussi; Oulun Yliopisto Laaketieteellisen tiedekunta Ohukainen, Pauli; Oulun Yliopisto Laaketieteellisen tiedekunta Kangas, Antti; University of Oulu, Institute of Clinical Medicine Kettunen, Johannes ; Oulun Yliopisto Laaketieteellisen tiedekunta Wang, Qin; Oulun Yliopisto Laaketieteellisen tiedekunta, Karsikas, Mari; Oulun Yliopisto Laaketieteellisen tiedekunta Khan, Anmar; Baker Heart and Diabetes Institute Kingwell, Bronwyn; Baker Heart and Diabetes Institute Kähönen, Mika; University of Tampere,, Department of Clinical Physiology, Tampere University Hospital and Medical School, Lehtimäki, Terho; Tampere University Hospital and School of Medicine, University of Tampere, Department of Clinical Chemistry, Fimlab Laboratories Raitakari, OT; University of Turku, Department of Clinical Physiology Jarvelin, M-R; Imperial College London Meikle, Peter; Baker Heart and Diabetes Institute Ala-Korpela, Mika; Baker Heart and Diabetes Institute, Systems Epidemiology
Key Words:	Software, epidemiology, metabolomics, statistics, self-organising map, data visualisation

Software Application Profile

EpiMetal: An open-source graphical web browser tool for easy statistical analyses in epidemiology and metabolomics

Jussi Ekholm,^{1,2,3} Pauli Ohukainen,^{1,2,3} Antti J. Kangas,⁴ Johannes Kettunen,^{1,2,3,5} Qin Wang,^{1,2,3,6} Mari Karsikas,^{1,2,3,7} Anmar A. Khan,^{8,9} Bronwyn A. Kingwell,¹⁰ Mika Kähönen,¹¹ Terho Lehtimäki,¹² Olli T. Raitakari,^{13,14} Marjo-Riitta Järvelin,^{2,3,15,16,17} Peter J. Meikle,⁸ and Mika Ala-Korpela^{1,2,3,6,18,19,20,21,*}

1. Computational Medicine, Faculty of Medicine, University of Oulu, Oulu, Finland
2. Biocenter Oulu, Oulu, Finland
3. Center for Life Course Health Research, Faculty of Medicine, University of Oulu, Oulu, Finland
4. Nightingale Health Ltd., Helsinki, Finland
5. THL: National Institute for Health and Welfare, Helsinki, Finland
6. Systems Epidemiology, Baker Heart and Diabetes Institute, Melbourne, VIC, Australia
7. Solita Ltd., Tampere, Finland
8. Metabolomics, Baker Heart and Diabetes Institute, Melbourne, VIC, Australia
9. Laboratory Medicine Department, Faculty of Applied Medical Sciences, Umm Al-Qura University, Kingdom of Saudi Arabia
10. Metabolic and Vascular Physiology, Baker Heart and Diabetes Institute, Melbourne, VIC, Australia
11. Department of Clinical Physiology, Tampere University Hospital and Finnish Cardiovascular Research Center - Tampere, Faculty of Medicine and Health Technology, Tampere University, Tampere, Finland
12. Department of Clinical Chemistry, Fimlab Laboratories and Finnish Cardiovascular Research Center - Tampere, Faculty of Medicine and Health Technology, Tampere University, Tampere, Finland
13. Research Centre of Applied and Preventive Cardiovascular Medicine, University of Turku, Turku, Finland
14. Department of Clinical Physiology and Nuclear Medicine, Turku University Hospital, Turku, Finland
15. Unit of Primary Health Care, Oulu University Hospital, OYS, Oulu, Finland
16. Department of Epidemiology and Biostatistics, MRC-PHE Centre for Environment and Health, School of Public Health, Imperial College London, London, UK
17. Department of Life Sciences, College of Health and Life Sciences, Brunel University London, UK
18. NMR Metabolomics Laboratory, School of Pharmacy, University of Eastern Finland, Kuopio, Finland

1
2
3 19. Medical Research Council Integrative Epidemiology Unit at the University of Bristol,
4 Bristol, UK

5
6 20. Population Health Science, Bristol Medical School, University of Bristol, Bristol, UK

7 21. Department of Epidemiology and Preventive Medicine, School of Public Health and
8 Preventive Medicine, Faculty of Medicine, Nursing and Health Sciences, The Alfred
9 Hospital, Monash University, Melbourne, VIC, Australia
10
11
12

13
14 ***Correspondence:**

15 Professor Mika Ala-Korpela

16 Computational Medicine, Faculty of Medicine

17 University of Oulu

18 Oulu, Finland
19
20

21
22 E-mail: mika.ala-korpela@oulu.fi
23
24

25 Mobile: +358 40 1977 657
26
27
28
29

- 30
31
- 32 • Abstract 144 words
 - 33 • Main text 2,052 words
 - 34 • 2 Figures
 - 35 • Supplementary material
 - 36 • EpiMetal software documentation online (including, e.g., a user guide and installation
37 instructions)
 - 38 • EpiMetal software usage exemplar open access online together with clinical, NMR
39 metabolomics and mass spectrometry lipidomics data
 - 40 • MIT licenced source code at the Github repository
- 41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Abstract

Motivation: An intuitive graphical interface that allows statistical analyses and visualisations of extensive data without any knowledge of dedicated statistical software or programming.

Implementation: EpiMetal is a single-page web application written in JavaScript, to be used via a modern desktop web browser.

General features: Standard epidemiological analyses and self-organising maps for data-driven metabolic profiling are included. Multiple extensive data sets with an arbitrary number of continuous and category variables can be integrated with the software. Any snapshot of the analyses can be saved and shared with others via a www-link. We demonstrate the usage of EpiMetal using pilot data with over 500 quantitative molecular measures for each sample as well as in two large-scale epidemiological cohorts (N>10,000).

Availability: The software usage exemplar and the pilot data are open access online at <http://EpiMetal.computationalmedicine.fi>. MIT licenced source code is available at the Github repository at <https://github.com/amerigin/epimetal>.

Keywords: statistical analyses, epidemiology, software, open-source, metabolomics, NMR, lipidomics

Note for the reviewers: The Github repository has been published. The samples.tsv containing the NMR metabolomics and MS lipidomics data for the 190 samples will be released when the paper is officially accepted.

Introduction

We are living in a multi-omics era of systems epidemiology.^{1,2} Quantitative high-throughput metabolomics³⁻⁶ and lipidomics^{7,8} have resulted in hundreds of molecular measures for up to tens of thousands of people in multiple cohorts and biobanks. Extensive and complex data create significant challenges for statistical analyses. It would therefore be beneficial, not only for omics beginners, but for all epidemiologists to have a simple visual tool for rapid exploratory analyses of these kinds of modern data sets without the immediate need of bioinformaticians fluent with currently available professional statistical analysis tools as, for example, the R software.⁹

To this end, we developed a web browser-based graphical software – EpiMetal – for standard statistical epidemiological analyses as well as for multivariate self-organising maps (SOMs) for data-driven analyses, metabolic profiling and potentially for clinical subgrouping.¹⁰⁻¹⁵ EpiMetal is versatile and any data set with an arbitrary number of continuous and categorical variables can be easily integrated with the software. Data from multiple cohorts can be imported for comparative analyses. The original data sets can be divided into subgroups via multiple ways, for example, based on SOMs, histograms or scatterplots; the created subgroups can be saved and analysed separately or in comparison to any other data set. Regression analyses with covariate adjustments are available with graphical visualisation of the results. Publication quality visualisations can be made and exported. Any snapshot of the analyses pipeline can be saved and shared with others via a www-link. Though it might not be an optimal choice to use EpiMetal for final publishable results, an additional good usage might be to utilise it as a benchmarking tool for newly written scripts and functions in another software.

As a usage exemplar, we present explorative analyses in a pilot cohort of 190 samples¹⁶⁻¹⁸ for which serum nuclear magnetic resonance (NMR) metabolomics³⁻⁶ and mass spectrometry (MS) lipidomics^{8,19} data are available. The data include over 500 quantitative molecular measures for each individual from these complementary methodologies that are getting increasingly popular in epidemiological applications. This is apparently the first time these comprehensive data are combined in an epidemiological setting. The data are made public along with the software. The exemplar demonstrates how the graphical interface of EpiMetal can be used to visualise extensive data, select subgroups and ultimately gain

1
2
3 epidemiological insights via a combination of various statistical analysis options. In the
4 supplement we also illustrate comparative analyses for two large-scale epidemiological
5 cohorts.
6
7
8
9

10 **Implementation**

11 EpiMetal consists of three major components: (i) the database (MongoDB) is the long-term
12 store for data set samples, computational results and stored sessions; (ii) a single-page web
13 application written in JavaScript (JS) that is accessed by users via a web browser; and (iii) a
14 back end software written in Python that serves as an intermediary between the web
15 application and the database to retrieve data and record user sessions. The application
16 utilises third-party open-source libraries (versions and licencing information is available at
17 Github). The software is encapsulated inside Docker containers to facilitate easy deployment
18 across server platforms and to isolate the software from host system. Several Plotter
19 instances can be run in parallel with differing configurations and data. The overall
20 architecture of the software is presented in Supplementary Figure 1. Key data handling,
21 visualisation and statistical analyses features are summarised in Figure 1.
22
23
24
25
26
27
28
29
30
31
32
33

34 **Installation**

35 A step-by-step installation guide is provided in the EpiMetal software user guide. The source
36 data are imported from a machine-readable format file (usually a tab-separated .tsv- or
37 comma-separated values .csv-files) where each row corresponds to an individual sample.
38 These samples have a unique identifier and usually belong to a single data set. An import
39 configuration file defines the column names and the column separator. A metadata file is
40 needed for the source data to indicate variable name, free description, unit of measurement
41 (e.g. mmol/L) and the group name for each variable. For variables that follow a common
42 pattern, a regular expression pattern can be employed. Imported variable types are either
43 numerical or categorical. Example files for the sample data are provided in the Github
44 repository. While there is no hard-limit for the number of samples or variables that can be
45 imported to the software, the larger the data set, the longer the download and processing
46 times. The software has been tested with over 500 variables and some 30 000 samples,
47 which present a realistic upper-bound for current usage.
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 During the installation phase, a Docker container is set up that initialises the database from
4 the source data file, along with the metadata descriptions. A second container is utilised for
5 compiling the front-end application from the JS source files, forming a bundle. This container
6 contains an http server to serve the bundle to user's web browser. A third container runs the
7 back end and its application programmatic interface (API). Using the Docker system, the
8 EpiMetal software can be set up to serve an individual user locally, or to allow the software
9 to be accessed by the public. An example configuration on how to limit the access to the
10 software instance with a username and a password is provided.
11
12
13
14
15
16
17
18
19

20 Architecture

21 The separation of concerns in EpiMetal is achieved by the common distinction between the
22 presentation layer (front end) and the data access layer (back end). Modern desktop
23 computers have considerable computing resources available, and thanks to the
24 developments in web browser JS engines and web technologies, those resources can be fully
25 appreciated in web applications. The philosophy of EpiMetal is to perform these calculations
26 on the client side and store them to the database for later retrieval. This is achieved by
27 asynchronously downloading chunks of the sample data and then performing the
28 computations as requested by the user. Computationally heavy actions are processed in
29 parallel using Web Workers, if supported by the browser.
30
31
32
33
34
35
36
37

38 AngularJS was chosen as the front end web framework as it was popular at the time of
39 starting the project, it had an active user-base and several useful libraries, and its two-way
40 data binding feature was appropriate considering the interactive nature of the application.
41
42
43
44

45 Back and front ends

46 The back end is a Python script developed with a Flask framework utilising Mongo Engine for
47 object data mapping and served with Gunicorn HTTP server. The back end defines actions
48 for retrieving the settings for the software, the metadata for variables, and samples for the
49 requested variables. In addition, the back end is called to request previously stored SOM
50 computations and SOM planes, and to store new ones.
51
52
53
54
55

56 The front end is a single-page application written in JS utilising AngularJS framework.
57 Several open-source auxiliary JS libraries are employed, most notably DC.js for interactive
58 charting throughout the application and Data-Driven Documents (D3.js) as a dependency for
59
60

1
2
3 DC.js and for SOM planes and other chart types. Visual appearance and user interface (UI)
4 stylings depend on Bootstrap framework and Angular-strap library. The UI allows the user to
5 freely create, resize, move and close window-like objects containing figures. The front end
6 fetches necessary data samples by querying the back-end asynchronously as user navigates
7 on the page.
8
9

10
11
12 Figures produced with EpiMetal can be exported either in SVG or in PNG format. A
13 particular state of the application can always be saved by creating a link to it and sharing the
14 link with collaborators.
15
16
17
18

19 **Usage exemplar: Combined comprehensive metabolomics and lipidomics data**

20
21 We present here explorative analyses in a unique pilot cohort of 190 blood samples¹⁶⁻¹⁸ for
22 which serum NMR metabolomics³⁻⁶ and MS lipidomics^{8,19} data are available. All these
23 molecular data, including basic clinical characteristics (age, sex, systolic and diastolic blood
24 pressure, body mass index and height) have been made open access along with the EpiMetal
25 software. The NMR metabolomics data comprise of over 200 metabolic measures, including
26 standard lipids, lipoprotein subclass and composition data, fatty acids, amino acids, ketones,
27 glycolysis and gluconeogenesis-related substrates and an inflammatory marker, glycoprotein
28 acetyls.³⁻⁶ The MS lipidomics data consist of over 350 individual lipid concentrations in 20
29 lipid classes including, e.g., ceramides, sphingomyelins, phosphatidylcholines,
30 phosphatidylinositols, cholesteryl esters, and triacylglycerols.^{7,8,19} We used EpiMetal to
31 conduct a multifaceted exploratory analysis of this pilot data set. We sought to demonstrate
32 some commonly known epidemiological and molecular features of these types of data.
33
34
35
36
37
38
39
40
41
42

43 First, we plotted the distribution of high-density lipoprotein cholesterol (HDL-C) and
44 the correlation of HDL-C with triglycerides (TG) in the entire data set (Figure 2B). As
45 expected, the histogram follows roughly a normal distribution. The scatterplot for HDL-C and
46 TG association reveals the well-known negative population-level correlation.²⁰
47
48
49

50 We then applied the SOM analysis to organise the samples in the data set via their
51 systemic metabolic profiles. Readers interested in the comparison of the SOM methodology
52 with other subgrouping methods in epidemiology are referred to a recent Software
53 Application Profile in the *Journal*.¹⁴ Additional details of the statistical issues in SOM analyses
54 can be found in references.¹⁰⁻¹³ We based the SOM profiling on 26 metabolic measures
55 including multiple amino acids, 14 lipoprotein subclasses, standard cholesterol measures,
56
57
58
59
60

glycoprotein acetyls and glycolysis-related measures. It should be noted that users could freely modify the initial SOM training data according to their preferences and data characteristics. The SOM planes for low-density lipoprotein cholesterol (LDL-C), HDL-C, and TG are shown in Figure 2C. These planes reveal, on average, that people with high circulating HDL-C (Circle A in the SOM; subgroup marked lowest TG) are indeed those that have low TG and vice versa (area marked B in the SOM; subgroup marked lowest HDL) as expected by the previous scatter plot. The SOM analysis also reveals that circulating LDL-C concentrations are rather indifferent regarding HDL-C and TG in this pilot data set; this is also emphasized in Figure 2D by the box plot for LDL-C. These associations can also be illustrated via formal regression analyses; we considered HDL-C as an outcome variable and LDL-C or TG as an exposure with age and sex as covariates. The results are given in Figure 2E for the entire cohort and the abovementioned SOM-derived subgroups. The negative association between HDL-C and TG, depicted in Figure 2B, is well replicated in the formal regression analysis. These demonstrations indicate the internal consistency of various software functions and illustrate that the pilot cohort represents well-known features of lipoprotein metabolism with respect to lipoprotein lipid measures.

Exploration of associations between the NMR metabolomics and MS lipidomics data can be found in Supplementary Figure 2, which shows a heatmap of Spearman's rank correlation co-efficients between selected lipoprotein (NMR) and lipid variables (MS). Overall the correlations are very well in line with the known molecular characteristics of lipoprotein subclasses and their lipid compositions²¹ and demonstrate robust agreement between the NMR metabolomics and MS lipidomics platforms.

To additionally demonstrate the properties of EpiMetal, we performed an additional set of analyses using data from two large-scale population-based epidemiological cohorts including over 10,000 individuals (see the Supplement).

Conclusion

The new EpiMetal software is used via a modern web browser and it provides an intuitive easy-to-use graphical interphase for multiple statistical methods relevant in epidemiological analyses. It easily handles data for tens of thousands of people and for hundreds of measures – numbers that are a reality nowadays in many metabolomics applications. It provides instant data visualisations and allows convenient sharing of results and data via

1
2
3 data captures accessible via an automatically created www-link. The data sets can be fully
4 customised by the users. We illustrated the usage and opportunities of EpiMetal in real
5 large-scale epidemiological data sets (Figure 1 and the Supplement). In addition, we provide
6 an open access usage exemplar of EpiMetal for a pilot cohort in which over 500 quantitative
7 molecular measures are available from each sample.
8
9

10
11
12 With increasing amounts of complex molecular data in epidemiology, sophisticated
13 software is required for both convenient data handling and statistical analyses. Without
14 statistical or programming expertise, the learning curve to conveniently use typical modern
15 data analysis software, for example R,⁹ can be steep. From the epidemiology perspective,
16 extensive molecular data may challenge traditional hypothesis-driven data analyses. These
17 are common situations in which the EpiMetal software can help researchers. Firstly, by
18 enabling instant graphical exploration and analyses of a (new) data set without the hurdles
19 of programming-based data analyses, and, secondly, by also allowing data-driven options to
20 find unknown relations in the data without pre-existing hypotheses. As far as we are aware,
21 the EpiMetal software is first-of-a-kind versatile tool for both traditional and data-driven
22 analyses of extensive large-scale epidemiological data sets.
23
24
25
26
27
28
29
30
31
32
33
34
35

36 **Funding**

37
38 MAK is supported by a Senior Research Fellowship from the National Health and Medical
39 Research Council (NHMRC) of Australia (APP1158958). He also works in a unit that is
40 supported by the University of Bristol and UK Medical Research Council (MC_UU_12013/1).
41 QW is supported by the Novo Nordisk Foundation (NNF17OC0027034). BAK is supported by
42 a Senior Principal Research Fellowship from the NHMRC of Australia (APP1154331). The
43 Sigrid Juselius Foundation and the Academy of Finland have also funded this work. The
44 Northern Finland Birth Cohort 1966 and the Young Finns Study have been financially
45 supported by multiple funding bodies (see the Supplement). The Baker Institute is supported
46 in part by the Victorian Government's Operational Infrastructure Support Program.
47
48
49
50
51
52
53
54
55
56
57
58
59
60

References

1. Sudlow C, Gallacher J, Allen N, *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 Med* 2015;**12**:e1001779.
2. Ala-Korpela M, Davey Smith G. Metabolic profiling-multitude of technologies with great research potential, but (when) will translation emerge? *Int J Epidemiol* 2016;**45**:1311–8.
3. Soininen P, Kangas AJ, Würtz P, Suna T, Ala-Korpela M. Quantitative serum nuclear magnetic resonance metabolomics in cardiovascular epidemiology and genetics. *Circ Cardiovasc Genet* 2015;**8**:192–206.
4. Würtz P, Kangas AJ, Soininen P, Lawlor DA, Davey Smith G, Ala-Korpela M. Quantitative Serum Nuclear Magnetic Resonance Metabolomics in Large-Scale Epidemiology: A Primer on -Omic Technologies. *Am J Epidemiol* 2017;**186**:1084–96.
5. Würtz P, Wang Q, Soininen P, *et al.* Metabolomic Profiling of Statin Use and Genetic Inhibition of HMG-CoA Reductase. *J Am Coll Cardiol* 2016;**67**:1200–10.
6. Sliz E, Kettunen J, Holmes M V., *et al.* Metabolomic Consequences of Genetic Inhibition of PCSK9 Compared With Statin Treatment. *Circulation* 2018;**138**:2499–512.
7. Mundra PA, Shaw JE, Meikle PJ. Lipidomic analyses in epidemiology. *Int J Epidemiol* 2016;**45**:1329–38.
8. Huynh K, Barlow CK, Jayawardana KS, *et al.* High-Throughput Plasma Lipidomics: Detailed Mapping of the Associations with Cardiometabolic Risk Factors. *Cell Chem Biol* 2019;**26**:71–84.
9. R Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria; 2014.
10. Mäkinen V-P, Forsblom C, Thorn LM, *et al.* Metabolic phenotypes, vascular complications, and premature deaths in a population of 4,197 patients with type 1 diabetes. *Diabetes* 2008;**57**:2480–7.
11. Kumpula LS, Mäkelä SM, Mäkinen V-P, *et al.* Characterization of metabolic interrelationships and in silico phenotyping of lipoprotein particles using self-organizing maps. *J Lipid Res* 2010;**51**:431–9.
12. Mäkinen V-P, Tynkkynen T, Soininen P, *et al.* Metabolic diversity of progressive kidney disease in 325 patients with type 1 diabetes (the FinnDiane Study). *J Proteome Res*

- 1
2
3 2012;**11**:1782–90.
4
5 13. Lithovius R, Toppila I, Harjutsalo V, *et al.* Data-driven metabolic subtypes predict
6 future adverse events in individuals with type 1 diabetes. *Diabetologia* 2017;**60**:1234–
7 43.
8
9
10 14. Gao S, Mutter S, Casey A, Mäkinen V-P. Numero: a statistical framework to define
11 multivariable subgroups in complex population-based datasets. *Int J Epidemiol*
12 2019;**48**:369–74.
13
14
15 15. Mäkinen V-P, Kangas AJ, Soininen P, Würtz P, Groop P-H, Ala-Korpela M. Metabolic
16 phenotyping of diabetic nephropathy. *Clin Pharmacol Ther* 2013;**94**:566–9.
17
18 16. Khan AA, Mundra PA, Straznicky NE, *et al.* Weight Loss and Exercise Alter the High-
19 Density Lipoprotein Lipidome and Improve High-Density Lipoprotein Functionality in
20 Metabolic Syndrome. *Arterioscler Thromb Vasc Biol* 2018;**38**:438–47.
21
22
23 17. Straznicky NE, Lambert EA, Nestel PJ, *et al.* Sympathetic Neural Adaptation to
24 Hypocaloric Diet With or Without Exercise Training in Obese Metabolic Syndrome
25 Subjects. *Diabetes* 2010;**59**:71–9.
26
27
28 18. Straznicky NE, Grima MT, Sari CI, *et al.* A Randomized Controlled Trial of the Effects of
29 Pioglitazone Treatment on Sympathetic Nervous System Activity and Cardiovascular
30 Function in Obese Subjects With Metabolic Syndrome. *J Clin Endocrinol Metab*
31 2014;**99**:E1701–7.
32
33
34 19. Weir JM, Wong G, Barlow CK, *et al.* Plasma lipid profiling in a large population-based
35 cohort. *J Lipid Res* 2013;**54**:2898–908.
36
37
38 20. Schaefer EJ, Levy RI, Anderson DW, Danner RN, Brewer HB, Blackwelder WC. Plasma-
39 triglycerides in regulation of H.D.L.-cholesterol levels. *Lancet* 1978;**2**:391–3.
40
41
42 21. Kumpula LS, Kumpula JM, Taskinen M-R, Jauhiainen M, Kaski K, Ala-Korpela M.
43 Reconsideration of hydrophobic lipid distributions in lipoprotein particles. *Chem Phys*
44 *Lipids* 2008;**155**:57–62.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure captions

Figure 1. Key data handling, visualisation and statistical analyses features of the EpiMetal software illustrated using real epidemiological data (Northern Finland Birth Cohort 1966; N=5,713). A generalized flow of analysis begins by choosing a dataset(s) from user-uploaded options. This can be, e.g., one population cohort but also a combination of many. Main analysis options are located in the top of the graphical interface and divided into three categories; “Explore and filter”, “Regression analysis” and “SOM”. Under “Explore and filter”, the user can quickly generate basic plots to gain an overview of the data structure. Variables can be plotted and compared using histograms, scatterplots and boxplots. Heatmaps can also be created for an overall visualization of variable (Pearson’s rank) correlations. Active filters can also be applied to select subsets of the data. For example, one can choose to analyse only individuals with HDL-C < 1.0 mmol/L in a given population cohort. The main category “Regression analysis” allows the user to choose an outcome and exposure variables with an optional number of covariates and generate a forest plot displaying the point estimate and 95% confidence intervals. Under “SOM”, the user can calculate a self-organizing map trained according to selected variables. Map can then be used to choose a subset of the entire data set on the basis of this metabolic profiling. It should be noted that the analyses made in the “Explore and filter” and “SOM” sections are fully compatible with each other enabling, for example, the SOM-based subgroups to be analysed via histograms and vice versa.

Figure 2. Explorative analysis of a cohort of 190 samples with serum NMR metabolomics and mass spectrometry lipidomics measures available. A: The control panel of EpiMetal that contains clickable buttons for generating graphs, selecting, naming and generating subgroups. Colours indicate the entire cohort (cyan) and selected subgroups based on the self-organizing map (SOM) analysis. B: The histograms of HDL-C in the entire cohort and in the subgroups and the scatterplot of HDL-C vs. triglycerides. C: The SOM component planes for serum triglycerides, HDL-C and LDL-C (note that the individuals in the entire cohort are identically distributed in each plane). Colours indicate high (red) and low (blue) concentration values of the variable in each plane. Individuals with similar metabolic profiles cluster close to each other in the SOM component planes. The user can specify and select

1
2
3 different subgroups via the circular selection tools. D: A box plot for LDL-C in the entire
4 cohort and in the two subgroups. E: Regression analyses with a forest plot showing standard
5 deviation (SD) increment in the outcome variable as a function of 1-SD increment in selected
6 exposure variables. Point estimates are indicated by a dot surrounded by 95% CI. Plotting
7 HDL-C as the outcome and triglycerides as an exposure illustrates the same negative
8 association as already indicated via the scatterplot in B.
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Review Only

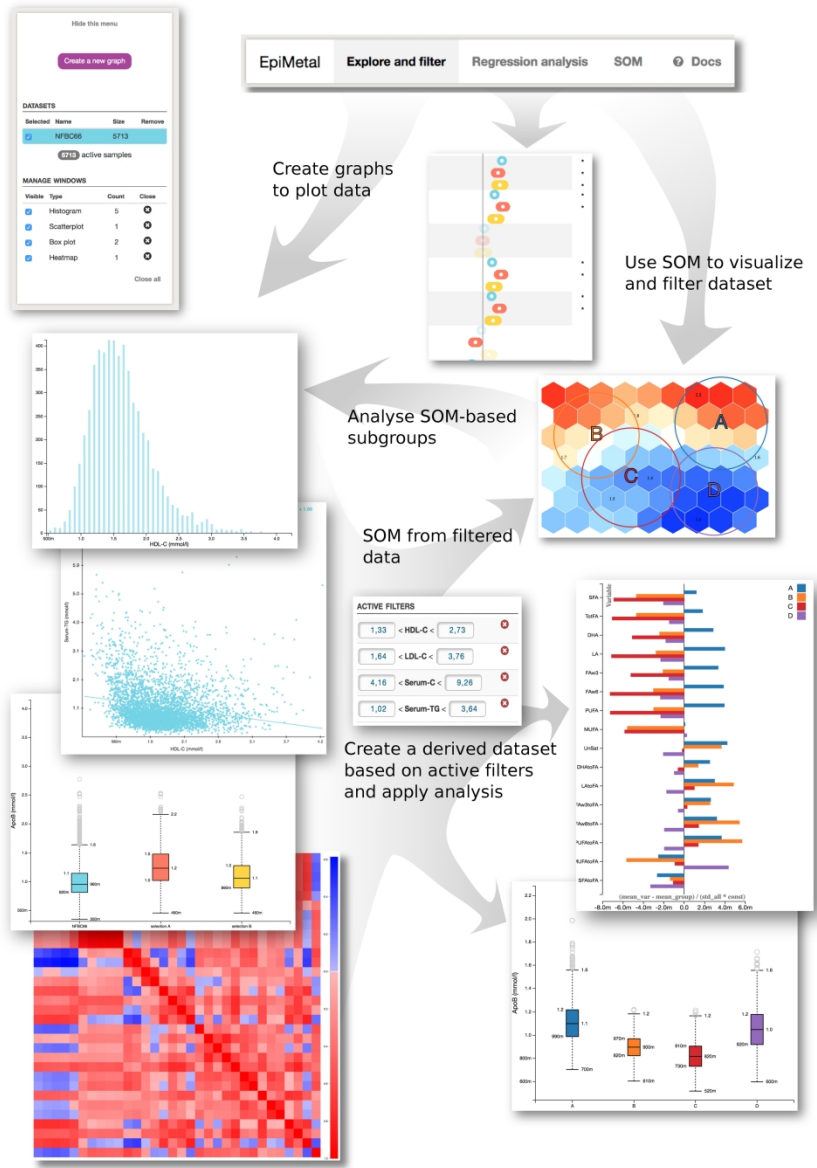


Figure 1

209x297mm (300 x 300 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

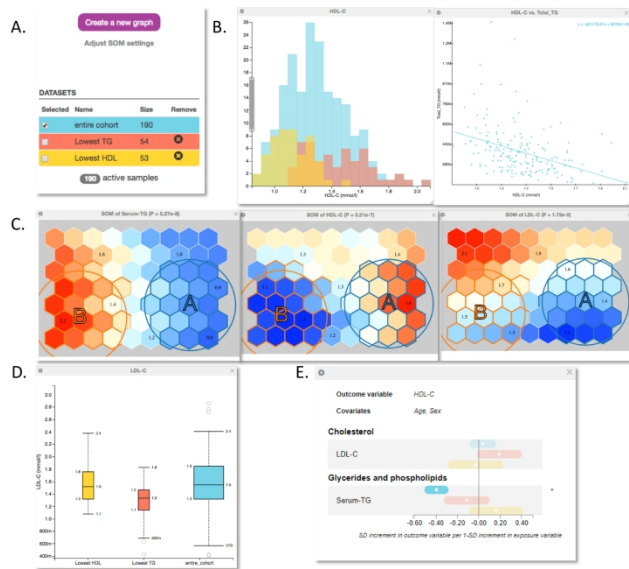


Figure 2

266x150mm (300 x 300 DPI)