# Screening for inborn errors of metabolism using untargeted metabolomics and out-of-batch controls

**Michiel Bongaerts[1*], Ramon Bonte[1], Serwet Demirdas[1], Ed H. Jacobs[1], E. Oussoren[2], Ans T. van der Ploeg[2], Margreet A.E.M. Wagenmakers[3], Robert M.W. Hofstra[1], Henk J. Blom[1], Marcel J.T. Reinders[4] and George J. G. Ruijter[1*]**

[1] Department of Clinical Genetics, Erasmus Medical Centre, Dr. Molewaterplein 40, 3015 GD Rotterdam, The Netherlands; r.bonte@erasmusmc.nl; s.demirdas@erasmusmc.nl; h.j.blom@erasmusmc.nl; e.jacobs@erasmusmc.nl; r.hofstra@erasmusmc.nl

[2] Erasmus University Medical Center, Department of Pediatrics, Center for Lysosomal and Metabolic Diseases, Dr. Molewaterplein 40, 3015 GD Rotterdam, The Netherlands; e.oussoren@erasmusmc.nl; a.vanderploeg@erasmusmc.nl

[3] Erasmus University Medical Center, Center for Lysosomal and Metabolic Diseases, Department of Internal Medicine, Dr. Molewaterplein 40, 3015 GD Rotterdam, The Netherlands; m.wagenmakers@erasmusmc.nl

[4] Faculty of Electrical Engineering, Mathematics and Computer Science, TU Delft, Van Mourik Broekmanweg 6, 2628 XE, Delft, The Netherlands; M.J.T.Reinders@tudelft.nl

[*] Correspondence: m.bongaerts@erasmusmc.nl; g.ruijter@erasmusmc.nl

**Motivation**: Untargeted metabolomics is an emerging technology in the laboratory diagnosis of inborn errors of metabolism (IEM). In order to judge if metabolite levels are abnormal, analysis of a large number of reference samples is crucial to correct for variations in metabolite concentrations resulting from factors such as diet, age and gender. However, a large number of controls requires the use of out-of-batch controls, which is hampered by the semi-quantitative nature of untargeted metabolomics data, i.e. technical variations between batches. Methods to merge and accurately normalize data from multiple batches are urgently needed.

**Methods & results**: Based on six metrics, we compared existing normalization methods on their ability to reduce batch effects from eight independently processed batches. Many of those showed marginal performances, which motivated us to develop *Metchalizer*, a normalization method which uses 17 stable isotope-labeled internal standards and a mixed effect model. In addition, we propose a regression model with age- and sex as covariates fitted on control samples obtained from all eight batches. *Metchalizer* applied on log-transformed data showed the most promising performance on batch effect removal as well as in the detection of 178 known biomarkers across 45 IEM patient samples and performed at least similar

36  to an approach using 15 within-batch controls. Furthermore, our regression model indicates that 10-24% of

37  the considered features showed significant age-dependent variations.

38  **Conclusions**: Our comprehensive comparison of normalization methods showed that our *Log-Metchalizer*

39  approach enables the use out-of-batch controls to establish clinically-relevant reference values for

40  metabolite concentrations. These findings opens possibilities to use large scale out-of-batch control samples

41  in a clinical setting, increasing throughput and detection accuracy.

42  **Availability**: *Metchalizer* is available at https://github.com/mbongaerts/Metchalizer/

43

## Introduction

45  Screening of patients suspected for inborn errors of metabolism (IEM) is currently based on measuring

46  panels of specific groups of metabolites like amino acids or organic acids using a number of different tests

47  and techniques such as ion-exchange chromatography, LC-MS/MS and GS-MS. This targeted approach

48  with several different tests is time consuming and limited in the number of metabolites being analyzed.

49  Untargeted metabolomics using High Resolution Accurate Mass Liquid Chromatography Mass

50  Spectrometry (HRAM LC-MS) can detect hundreds to thousands of metabolites within one test, and, as a

51  consequence, receives increasing interest to be used in IEM screening (Miller, et al., 2015) (Coene, et al.,

52  2018) (Körver-Keularts, et al., 2018) (Haijes, et al., 2019) (Bonte, et al., 2019). Moreover, untargeted

53  metabolomics can also reveal new biomarkers or increase our understanding of disease mechanism when

54  exploited in epidemiological studies (Glinton, et al., 2019).

55

56  In traditional targeted diagnostic laboratory tests hundreds of reference samples are required to establish

57  robust reference intervals. When using untargeted metabolomics the establishment of reference values is

58  complicated due to the semi-quantitative nature of the data owing to several sources of variation like

59  injection volume, retention time, temperature, or ionization efficiency in the mass spectrometer that cannot

60  easily be amended. Moreover, these variations are even larger between different measurement runs in which

61  a batch of samples is being measured simultaneously, hampering the resemblance between different

62  batches. As a result, within-batch variation is smaller than between-batch variation. Therefore, to conquer

63  these batch effects, current approaches include reference samples in each single batch of measurements

64  (Miller, et al., 2015) (Coene, et al., 2018) (Haijes, et al., 2019) (Körver-Keularts, et al., 2018) (Bonte, et al.,

65  2019) to improve detection sensitivity (due to tighter reference values as a result of lower variation in the

66  in-batch reference samples).

67

68    Clearly, this reduces the throughput efficiency of IEM screening as the number of patient samples that can

69    be included in a batch is considerably lower when the reference samples need to be measured as well. But,

70    more importantly, the number of reference samples in one batch might fall short in the establishment of

71    adequate reference ranges as variations in certain metabolites are not captured well enough in the relatively

72    small reference panel. For example, factors like age, sex and BMI can affect abundancies of metabolites,

73    and, to establish reliable reference ranges, one thus needs to correct for these factors by using a large number

74    of reference samples (Chaleckis, et al., 2016) (Rist, et al., 2017) (Yu, et al., 2012). Consequently, for reliable

75    untargeted metabolomics in clinical testing, a large set of reference samples is needed, while for throughput

76    efficiency a small set is preferred. Altogether, this calls for an approach that can establish reference values

77    based on reference samples being measured in several batches (out-of-batch controls).

78

79    When relying on reference samples from different batches, one needs to correct for the batch effects to

80    obtain reliable estimates for the reference ranges. This is generally solved by normalization methods and

81    some have already been proposed within the context of untargeted metabolomics and mass spectrometry

82    (Veselkov, et al., 2011) (Li, et al., 2017) (Välikangas, et al., 2016). Only a few groups have used out-of-

83    batch controls to determine the reference values and used relatively simple normalization techniques like

84    median scaling (Miller, et al., 2015), using a reference internal standard per metabolite (Körver-Keularts,

85    et al., 2018) or using anchor samples (Glinton, et al., 2019). However, there has not been an extensive

86    exploration of normalization techniques within the context of diagnostic testing for IEM's.

87

88    We explore several known normalization methods on their ability to remove batch effects and to detect

89    biomarkers from patients with known IEM. Furthermore, we introduce a new normalization method, which

90    we called *Metchalizer,* which uses internal standards and a mixed effect model to remove batch effects. As

91    this allows for a large set of (out-of-batch) reference samples, we also explore a regression model that uses

92    age and sex as covariates to correct for potential age and sex effects on the reference values. Using the

93    regression model combined with the *Metchalizer* normalization, we achieve similar performances in

94    biomarker detection compared to the use of within-batch controls. Hence, this opens the possibility to

95    increase the throughput of untargeted metabolomics in IEM screening as well as including more complex

96    confounder strategies.

97

98

99

置

## Materials and methods

**Untargeted metabolomics datasets**

Human plasma samples of 260 control samples and 53 IEM patients were measured over eight batches over the period 10-12-2018 to 03-05-2019 (Bonte, et al., 2019) having in total 33 unique IEMs. For every patient a technical triplicate was included. A QC (Quality Control) sample was included in all eight batches and more than four technical replicates were present in every batch. Since the QC sample was a commercial sample, the sample differed in concentration of several metabolites when compared to the (average) concentrations of the human plasma samples analyzed in these datasets. Features were annotated as described in Bonte et al. (Bonte, et al., 2019). Note that within each batch about 30 normal controls have been measured, which allows us to establish reference values based on within-batch controls, whereas the controls being measured for the other (seven) batches can be used for out-of-batch strategies. In this study we will refer to 'feature' as being either a single m/z-value (with unique retention time) or a merge of multiple features, where the adduct type and/or isotope was determined with corresponding neutral mass and consequently merged to a single feature.

The following internal standards have been added to each batch to facilitate normalization based on these internal standards: 1,3-$^{15}$N uracil (+/-), 5-bromotryptophan (+/-), D$_{10}$-isoleucine (+/-), D$_3$-carnitine (+/-), D$_4$-tyrosine (+/-), D$_5$- phenylalanine (+/-), D$_6$-ornithine (+), dimethyl-3,3-glutaric acid (+/-), $^{13}$C-thymidine (+/-), D$_4$-glycochenodeoxycholic acid (-), where + indicates positive ion mode, and – indicates the negative ion mode.

**Data processing**

Previous pre-processing steps (alignment, peak picking etc.) were performed per batch using Progenesis QI v2.4 (Newcastle-upon-Tyne, UK) (Bonte, et al., 2019). In-house software was developed to match features from each batch to a reference batch which in this case was the fifth batch when sorting on chronologically order. Chromatograms between batches were initially aligned to the reference batch by using lowess regression where features were matched based on retention time difference, m/z-value and median abundancy difference similar to the criteria described below.

Matching features was performed based on several criteria:

131    1)    When features were annotated in reference batch and the batch being merged, these features were
132        pooled to the merged dataset.
133    2)    When MS/MS spectra were present for a potential matching pair of features, the cosine similarity
134        metric was calculated and had to be > 0.8.
135    3)    Retention time difference in percentage was calculated between potential matches, and had to be <
136        2.5%.
137    4)    Progenesis QI determined per feature an isotope distribution and we required sufficient overlap of
138        these distributions between potential matching pairs. This was determined by calculating a difference
139        in percentage between each bin of this distribution. The maximum difference of these bins had to be <
140        50%.
141    5)    As we expect matching features to have similar within-batch median abundancies (despite of batch
142        effects), we calculated the differences between these medians in percentages, which had to be < 300%.
143    6)    Neutral masses were known for the matching pair but not the MS/MS spectra, the ppm-error had to be
144        < 1.
145    7)    m/z-values were known for the matching pair but not the MS/MS spectra and neutral masses, the ppm-
146        error of between the m/z-values had to be < 1.

147

148 Features matching multiple other features in the reference batch were discarded (and vice versa). The
149 resulting merged dataset contained only features which were matched across all eight batches.

150

151 **Quantitative evaluation set**
152 For the evaluation of the normalization methods, the following 16 metabolites were quantitatively (µmol/L)
153 measured in two separate assays: leucine (+/-), C0 | L-carnitine (+/-), methionine (+/-), C2 | acetylcarnitine
154 (+), 5-aminolevulinic acid/4-hydroxyproline (+), serine (+/-), citrulline (+/-), aspartic acid (+), glutamine
155 (+/-), (allo)isoleucine (+/-), proline (+/-), tyrosine (+), phenylalanine (+/-), taurine (+/-), asparagine (+/-),
156 arginine (+/-). Amino acids were determined by ion-exchange chromatography according to protocols
157 described by the manufacturer (Biochrom). Free carnitine and acylcarnitines analysis was performed as
158 described by Vreken et al. (Vreken, et al., 2002).

159

160

161

162 **Normalization methods**

163

164 **Initial transformations**

165 Prior to normalization raw abundancies were for some methods transformed using a log-transform or Box-

166 Cox transformation. The latter was given by:

167

$$\hat{y} = \frac{(y + \lambda_2)^{\lambda_1} - 1}{\lambda_1}$$

168 (1)

169 with $\lambda_1 = 0.5$ and $\lambda_2 = 1$. If an initial transformation was applied this was indicated in the name of the

170 (normalization) method, where 'BC-' refers to the Box-Cox transformation and 'Log-' to the log

171 transformation. When no transformation was performed this was indicated with 'None-'.

172

173 **Normalization by Metchalizer**

174 *Metchalizer* assumes a linear mixed effect relationship between the abundancies of the internal standards

175 and the feature of interest. Since the internal standards were expected to be correlated, we represented them

176 by an orthogonal set of covariates. These covariates are obtained as the Latent Variables (LV) from the

177 Partial Least Squares (PLS) of the set of internal standard abundancies (represented in matrix $\mathbf{X}$) and the

178 (categorical) information about which sample belonged to which batch (represented by matrix $\mathbf{Y}$). The

179 number of LV's were chosen from the metric *I(K)*:

180

$$I(K) = \sum_{k=1}^{K} \sum_{b,i} \left( x_{ib}^{\text{LV}_k} - \bar{x}_{.b}^{\text{LV}_k} \right)^2$$

181 (2)

182

183 where $\bar{x}_{.b}^{\text{LV}_k}$ is the center of batch $b$ in the direction of $\text{LV}_k$. We selected that $K$ for which *I(K)* reached 75

184 % of its maximum value.

185

186 The mixed effect model then considers the LV's as fixed effects and all variations not explained by the LV's

187 is considered as (random) batch effects:

188

$$\hat{y}_{ijb} = \beta_j^0 + \sum_k \beta_j^k x_i^{\text{LV}_k} + \gamma_{jb} + \epsilon_{ijb}$$

189 (3)

6

190 with $x_i^{\mathrm{LV}_k}$ indicating the covariate (score) of the $k^{\mathrm{th}}$ Latent Variable (LV) of sample $i$. $\gamma_{jb}$ is the (random)

191 batch intercept for feature $j$. Note, that when the LV's are sufficient in explaining $y_{ijb}$ the random intercept

192 $\gamma_{jb}$ will not contribute much. Before fitting the model, we remove outlier samples per batch $b$ and feature

193 $j$ based on their within-batch Z-score ($|Z| > 2$) determined from all samples in that batch. These Z-scores

194 were different than the Z-scores defined in other parts of this study.

195

196 The batch corrected abundancy then becomes:

197

198
$$y_{ijb}^{\mathrm{batch\ corrected}} = y_{ijb} - \hat{y}_{ijb} + \mathrm{Median}(\hat{y}_{.jb})\tag{4}$$

199

200

201 **Normalization by Best Correlated Internal Standard**

202 The internal standard, $m$, that best correlates with a feature $j$ is being used to normalize the abundances of

203 feature $j$. The correlation is measured within each batch using the spearman correlation between feature $j$

204 and each internal standard individually across all samples and subsequently averaged across all eight

205 batches. The internal standard which (positively) correlated the best was used for normalization according:

206

207
$$\hat{y}_{ij} = \frac{y_{ij}}{y_{im}} \mathrm{Median}(y_{.m})\tag{5}$$

208 with $m$ being the best correlated internal standard.

209

210 **Normalization methods from literature**

211 We compared *Metchalizer* with a number of different normalization methods. For a description we refer to

212 the original articles, here we only specify our settings:

213 **Anchor** (Glinton, et al., 2019)**:** *Anchor* assumes a linear response between the features in the anchor

214 samples and samples in the batch. An anchor sample is a fixed sample which is analyzed in all eight batches,

215 and was included more than four times in each batch . Normalization was performed per batch by dividing

216 each feature by the median of the anchor samples for that same feature per batch [1]. In this study we used

217 our QC samples as the anchor samples.

218 **CRMN** (Redestig, et al., 2009) : We used function normFit from the *crmn* R package with input argument

219 "crmn" and ncomp=3. As a design matrix we chose QC samples versus human plasma's.

220    **EigenMS** (Karpievitch, et al., 2015) : QC samples and plasma samples were treated as two different groups.

221    We chose three 'eigentrends'.

222    **Fast Cyclic Loess** (Ballman, et al., 2004) : We used the *normalizeCyclicLoess* function from the *limma* R

223    package using the method "`fast`" and `span=0.7`.

224    **NOMIS** (Sysi-Aho, et al., 2007) : We used the function *normFit* from the *crmn* R package with input

225    argument "`nomis`".

226    **PQN** (Filzmoser & Walczak, 2014) : *PQN* was implemented as described by Filzmose et al. The reference

227    spectrum was given by the median of every feature *j*.

228    **RUV** (Livera, et al., 2015) : We used the function *RUVRand* with `k=8` from the *MetNorm* R package.

229    **VSN** (Huber, et al., 2002) : We used the *vsn* R package using the *vsn2* function.

230

### Evaluation of normalization methods

232    Six metrics were used to evaluate the performance of normalization methods.

233    **WTR$_j$ score**: The WTR score (**W**ithin variance **T**otal variance **R**atio) calculates the ratio between the

234    'overall' within-batch variance and the total variance from the QC samples:

$$\text{WTR}_j = \frac{\sigma^2_{j,\text{within}}}{\sigma^2_{j,\text{tot}}} = \frac{\sigma^2_{j,\text{tot}} - \sigma^2_{j,\text{between}}}{\sigma^2_{j,\text{tot}}}$$

235    (6)

236    where $\sigma_{j,\text{between}}$ is the variance of all eight batch averages for metabolite *j* in the QC samples, and $\sigma^2_{j,tot}$

237    the 'overall' variance based on all QC samples. The WTR score is between 0 and 1. As we would like batch

238    averages to be similar for the QC samples (resulting in $\sigma_{j,\text{between}}$ approaching zero), we are interested in

239    WTR scores close to one.

240

241    $\Delta R$ **score**: Since normalization might also lead to the removal of variations of interest (for example

242    biological variations), we tested whether the ranks of the features ordered by their abundancies within the

243    QC samples were preserved after normalization. Per feature *j,* we determined the average rank the feature

244    is assigned across all QC samples (across all batches) for both the raw abundancies ($\bar{R}^{\text{raw}}_j$) as well as the

245    normalized abundancies ($\bar{R}^{\text{normalized}}_j$). The $\Delta R_j$ score then looks at the difference in rank positions due to

246    normalization per feature *j*:

247

248
$$\Delta R_j = \left| \bar{R}_j^{\text{raw}} - \bar{R}_j^{\text{normalized}} \right|$$
(7)

249 $\Delta R_j \in [0, p]$, with $p$ the number of features. Lower $\Delta R_j$ values indicate a better preservation of the ranks

250 of the normalization method.

251 **Spearman score**: For the set of 16 quantitatively measured metabolites, we calculated the Spearman

252 correlation between their quantitative measurements and the normalized abundancies. Overall

253 normalization performance could be judged based on the median Spearman score of these 16 scores, having

254 scores $\in [-1, 1]$. Higher values indicate better resemblance with the quantitative measurements.

255 **$R^2$ score**: The $R^2$ between the quantitative measurements and the normalized abundancies of the 16

256 quantitatively measured metabolites. Overall performance could be judged from the median $R^2$ score, with

257 scores $\in [0, 1]$. Higher values indicate better (linear) fits with the quantitative measurements.

258 **QC prediction score:** Since the QC samples were different from the human plasma samples in terms of

259 concentrations for several metabolites/features, we expect this difference to be observed in the first few

260 principal components (PCs) of a Principal Component Analysis (PCA) analysis applied to all features (excl.

261 standards). We fitted a logistic function using the first four PC's as covariates and with class labels: 'human

262 plasma' and 'QC'. The fitted model returns per sample a probability of belonging either to the class 'human

263 plasma' or 'QC'. The probabilities for all samples are averaged into the *QC prediction score* $\in [0, 1]$

264 Increasing normalization performances should result in higher scores, as QC - and human plasma samples

265 should be nicely separated. We used *LogisticRegression* from the Python package *scikitlearn* with

266 parameters `penalty='l1'`, `solver='saga'`, `multi_class='auto'`, `max_iter=10000`

267 (Pedregosa, et al., 2011).

268 **Batch prediction score:** Increasing normalization performances should result in less batch clustering when

269 examining the first few PC's of the PCA analysis (see *QC prediction score*). We fitted a logistic function

270 for each batch versus all other seven batches using the first four PC's as covariates and obtained the

271 probability scores for all human plasma's having the correct batch label. These scores were than averaged

272 for all human plasma samples into a *batch prediction scores* $\in [0, 1]$. Scores closer to 1 indicate decreased

273 normalization performances since batch separation is (still) present.

274

**Methods to determine aberrated metabolic abundancies**

Reference values for metabolites were determined by using a Z-score methodology: a set of reference values was Z-transformed (corrected for mean and divided by the standard deviation) which was then assumed to be normally distributed. Aberrations can then be called by considering significant Z-scores using a chosen cutoff level. We use four different methods to determine the Z-scores.

**Method *15in*: best matching controls within batch**: Z-scores were calculated by selecting 15 control samples originating from the same batch as the patient based on age and sex as described in Bonte et al. (Bonte, et al., 2019).

**Method *15out*: best matching controls from other batches**: Z-scores were calculated similarly as in *method 15in* using explicitly 15 out-of-batch controls. Note, that since there a more out-of-batch controls than within-batch controls that age and sex matching can be done more accurately.

**Method *All controls*:** This method used all available control samples from all eight batches, including within-batch controls, for Z-score calculation.

**Method *Regression*:** We fitted a linear model on all 260 available controls excluding outliers which were first removed based on their within-batch |Z-score| > 3, this Z-score is different from other Z-scores mentioned in this study, and only used to remove outliers. The regression model is given by:

$$
\begin{aligned}
\hat{y}_i &= \hat{\beta}^{\text{Intercept}} + \hat{\beta}^{\text{Sex}} x_i^{\text{Sex}} + \hat{\beta}^{\text{Sex,Age}} x_i^{\text{Sex}} x_i^{\text{Age}} \\
&\quad + \sum_{p=1}^{P} \hat{\beta}_p^{\text{Age}} (x_i^{\text{Age}})^p + \hat{\epsilon}_i
\end{aligned}
\tag{8}
$$

$$
\hat{y}_i = \vec{x}_i^T \vec{\hat{\beta}} + \hat{\epsilon}_i
\tag{9}
$$

where $\hat{y}_i$ is the predicted (normalized) abundancy of feature $j$ for sample $I$, $\hat{\beta}^{\text{Intercept}}$ is an intercept. $\hat{\beta}^{\text{Sex}}$, $\hat{\beta}^{\text{Sex,Age}}$ (interaction) and $\hat{\beta}_p^{\text{Age}}$ indicate slopes. $P$ is the degree of the polynomial used for regression on age and set to $P=3$ in this study. $x_i^{\text{Sex}}$ is 1 for women and 0 for men. $\hat{\epsilon}_i$ is the estimated error. The latter expression is the model in vector notation with $\vec{x}_i^T = [1, x_i^{\text{Sex}}, ..., (x_i^{\text{Age}})^P]$.

10

305 The coefficients were determined from the OLS estimator:

306

$$\vec{\hat{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\vec{y} \tag{10}$$

308

309 where the rows of $\mathbf{X}$ are given by $\vec{x}_i^T$ and the variance in $\hat{y}_i$ is determined by the variance in $\vec{\hat{\beta}}$ and the

310 variance in $\hat{\epsilon}_i$:

311

$$
\begin{aligned}
\mathrm{Var}[\hat{y}_i] &= \mathrm{Var}\left[\vec{x}_i^T\vec{\hat{\beta}}\right] + \mathrm{Var}[\hat{\epsilon}_i] \\
&= \vec{x}_i^T\,\mathrm{Cov}[\vec{\hat{\beta}}\,]\vec{x}_i + \hat{\sigma}_i^2
\end{aligned}
\tag{11}
$$

313

314 The covariance matrix of $\vec{\hat{\beta}}$ is given by:

315

$$
\begin{aligned}
\mathrm{Cov}[\vec{\hat{\beta}}] &= \mathrm{Cov}[\beta + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\vec{\epsilon}] \\
&= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\,\mathrm{E}[\vec{\epsilon}\,\vec{\epsilon}^T]\,\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}
\end{aligned}
$$

317 $$\tag{12}$$

318

319 with $\mathrm{E}[\vec{\epsilon}\,\vec{\epsilon}^T]$ estimated according:

320

$$
\mathrm{E}[\vec{\epsilon}\,\vec{\epsilon}^T] =
\begin{bmatrix}
\hat{\sigma}_1^2 & 0 & \cdots & 0 \\
0 & \hat{\sigma}_2^2 & \cdots & 0 \\
\vdots & \vdots & \ddots & \vdots \\
0 & 0 & \cdots & \hat{\sigma}_N^2
\end{bmatrix}
\tag{13}
$$

322

323 Since we expected $\sigma_i^2$ to be dependent on age (neglecting sex), we do estimate $\hat{\sigma}_i^2$ differently from a

324 weighted mean on the squared residuals:

325

$$
\begin{aligned}
\hat{\sigma}_i^2 &= \sum_{k=1}^{N} \frac{w_k(x_i^{\mathrm{Age}})}{\sum_{k'=1}^{N} w_{k'}(x_i^{\mathrm{Age}})}(y_k - \hat{y}_k)^2 \\
w_k(x_i^{\mathrm{Age}}) &= \exp\left(-\frac{|x_i^{\mathrm{Age}} - x_k^{\mathrm{Age}}|}{a + bx_i^{\mathrm{Age}}}\right)
\end{aligned}
\tag{14}
$$

11

327     where *a* and *b* determine how the weights decay (*a*) or increase (*b*) over age (we set *a, b = 1* years). Z-

328     scores were obtained by subtracting the predicted average $\hat{y}_i$ and dividing by the variance $\mathrm{Var}[\hat{y}_i]$

329     (Equation 11).

330

331     **Significance of regression coefficients:** Significance of the regression coefficients (Equation 8, 9) was

332     obtained by considering the statistic:

$$\frac{(\hat{\beta}_i - \beta_i)}{\sqrt{\mathrm{Var}[\hat{\beta}_i]}} \sim \mathcal{N}(0, 1)$$

333                                                                                                            (15)

334     The variances of the coefficients were found in the diagonal elements of $\mathrm{Cov}[\vec{\hat{\beta}}]$ (Equation 13). We tested

335     the hypotheses that $\beta_i = 0$ with a two-tailed test. A robust p-value was obtained from a bootstrap procedure

336     by taking the median p-value from a series of p-values obtained from 50 bootstraps on the above test

337     statistics taking 95 % of the data each bootstrap.

338

339     **Final Z-scores**

340     Since the patient samples were measured in triplicate, we determined the final Z-scores from the average

341     of these three Z-scores (Bonte, et al., 2019). These average  Z-score were determined for all Z-score

342     methods i.e. *15in, 15out, All controls, Regression* and IEM patient.

343

344     **P-values from Welch's t-test**

345     As an alternative to using the (average) Z-scores we also considered the p-values obtained from the Welch's

346     t-test to be informative, as it indicates whether the mean of triplicates differs significantly from the

347     population average. Note that the triplicate was expected to have only technical variance whereas the

348     reference population has variance consisting of technical- plus biological variance. For every Z-score

349     method (*15in, 15out, All controls, Regression*) these p-values were obtained per feature (and patient).

350

351     When using the regression model, we used an adjusted Welch's t-test assuming that variance in the estimate

352     of the average of the population (which is Z=0) was negligible :

$$t_j = \frac{\text{Mean}(Z_{j.})}{\sqrt{\frac{s_j^2}{3}}}$$

353

(16)

354

355  where $s_j$ is the sample standard deviation of the triplicate Z-scores, $\text{Mean}(Z_{j.})$ indicates the average of the

356  triplicate for feature $j$.

357

**Detection of the expected IEM biomarkers**

359  To explore how normalization and the method of determining these Z-scores (*15in, 15out, All controls* and

360  *Regression*) affected the detection of biomarkers, we plotted the number of detected biomarker of the

361  known IEM patients against the average number of detected features per patients for various (final) Z-score

362  and p-value cutoff levels, similar to a ROC curve. Improved biomarker detection was believed to increase

363  the area under the ROC(-like) curve (AUC).

364

365  Establishing this ROC curve was done by assigning a status for every biomarker (if present and annotated

366  in the MS-data). A database was established containing the expected biomarkers for each IEM including

367  the expected Z-score sign (up or down regulated) as can be found in supplement S5 Table 5. For every IEM

368  patient, we assigned for all expected biomarkers the status 'positive' or 'negative'. The status 'positive'

369  was assigned when 1) |Z-score| > $Z_{abnormal}$ , and 2) the sign of the Z-score corresponded with the expected

370  sign for that biomarker in the IEM patient. Criteria 1 and 2 were also used for the ROC-curve created by

371  the p-values. When a biomarker was found in both positive and negative ion mode, the Z-score(s) from the

372  mode having the largest population average abundancy was taken. The average number of detected features

373  (per patient) was obtained by considering features from both ion modes.

374

375  Some of the expected biomarkers were not matched across all eight batches and therefore were absent in

376  the merged dataset and analysis in this study. In the merged dataset, we obtained 178 patient-biomarker

377  combinations (one patient could have multiple biomarkers) associated with 45 patients (hence, for 8 IEM

378  patients no biomarkers were found in the merged dataset).
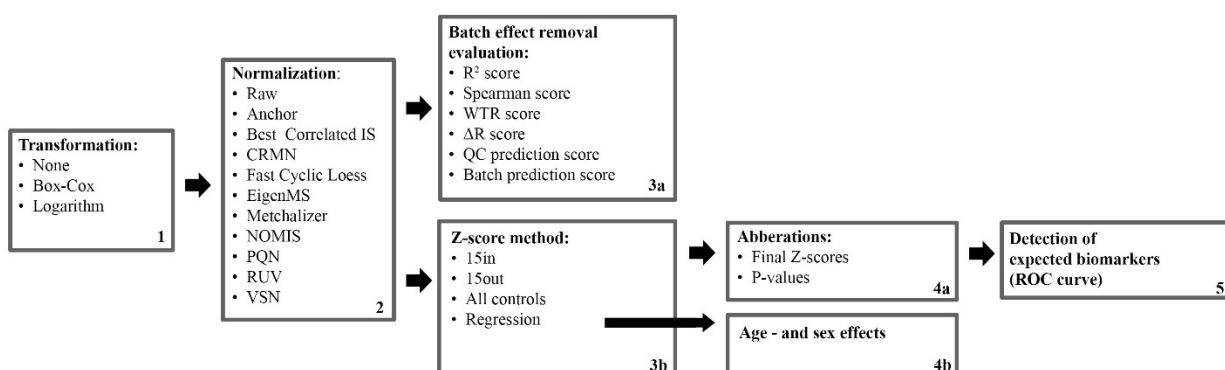
379

380

381

382

383

13

# Results



Figure 1. Flow diagram of different methods used in this study. 1) An initial transformation was applied. 2) A normalization method was applied. 3a) Multiple metrics were calculated to investigate batch effect removal. 3b) Normalized data was used to determine Z-scores for IEM patients using different (control) reference methods. 4a) Final Z-scores were calculated together with p-values. 4b) Regression analysis on all features/biomarkers was used to explore age- and sex dependency of abundancies. 5) Detection of the expected biomarkers was investigated using a ROC-like curve for Z-scores and p-values

## Batch characteristics

Eight untargeted metabolomics runs/batches were merged containing 260 control samples and 53 IEM patients, together having 33 unique IEMs. After merging, 773 positively ionized features were obtained, among which 121 were annotated, and 598 negatively ionized features were attained with 106 annotated features. We only included features which were merged across all eight batches to ensure consistency among the findings. Intra-batch coefficients of variation (CV) on 17 (internal and external) standards were smaller (median CV=14%) than inter-batch CV's (median CV=27%) indicating that batch effects were present (for more details see S1). Principle Component Analysis (PCA) further elucidated the presence of batch effects as shown in Figure 2A, showing the first three PC's for the log-transformed raw abundancies (*Log-Raw*).

## Comparing normalization methods

We investigated the performance of several normalization methods on batch effect removal by evaluating multiple metrics based on quantitative measurements, the Quality Control (QC) samples and PCA analysis

14

409    (see Methods and Figure 1). Some normalization methods were excluded from the following analysis

410    because of their marginal performance on the considered metrics (as evaluated in supplement S2).

411

412    *Reduced batch effects*: We visually observe in the PCA plots that most normalization methods reduced

413    batch effects since batch clustering seemed to be reduced after normalization (Figure 2), and is confirmed

414    when looking at the *batch prediction score* (Figure 3A) showing lower scores for normalized abundancies

415    when compared with the raw data (*None-Raw* or *Log-Raw*). *BC-Metchalizer*, *Log-Metchalizer* and *None-*

416    *Anchor* had the lowest *batch prediction scores*, with a median score of 0.13 (0.13) , 0.14 (0.14), 0.17 (0.16)

417    for positive (negative) ion mode respectively.

418

419    *Improved separation of QC samples*: QC samples (squares in Figure 2) were included in every batch and

420    were expected to separate from the human plasma samples (squares vs circles in Figure 2) in the first four

421    Principle Components (PC) due to overall abundancy differences for several metabolites. Normalization

422    should maintain this separation which was measured by the *QC prediction score* (Figure 3B). *Log-CRMN*

423    conserved QC/plasma separation, with a median *QC prediction score* of 1.00 (1.00) for positive (negative)

424    ion mode, but was less able to reduce batch effects since it had a median *batch prediction score* of 0.76

425    (0.39) for positive (negative) ion mode respectively. *Log-NOMIS* and *Log-RUV* were better in reducing

426    batch effects, with a median *batch prediction score* of 0.21 (0.21), 0.24 (0.19) for positive (negative) ion

427    mode respectively, but were less able to conserve the separation between QC and human plasma samples,

428    since the median *QC prediction scores* were 0.32 (0.88) and 0.39 (0.94) for positive (negative) ion mode

429    respectively. It is therefore likely that these two methods removed variations other than batch related

430    variation. QC samples were almost perfectly separated from the human plasma sample by *BC-Metchalizer*,

431    *Log-Metchalizer* and *None-Anchor.*

432

433    *Resemblance with quantitative measurements*: To further quantify batch effect removal, we calculated the

434    Spearman score and $R^2$ score between quantitative plasma concentrations (in µmol/L) and the normalized

435    abundancies of our evaluation set of amino acids and (acyl)carnitines (Methods). To ensure high signal-to-

436    noise ratio's in the quantitative measurements, we selected only metabolites having a population average

437    concentration above 1 µmol/L. Matching this evaluation set with the annotated features in the untargeted

438    metabolomics data resulted in 16 and 13 metabolites in positive - and negative ion mode, respectively.

439    Figure 3C and D shows both metrics for the investigated normalization methods. Again, for most

440    normalization methods both metrics improved when compared to the raw data (*None-Raw*). *BC-*

441    *Metchalizer*, *Log-Metchalizer* and *None-Anchor* appeared to perform the best on these metrics with median

15

442   $R^2$ scores of 0.56 (0.55), 0.57 (054), 0.57 (0.47), and median Spearman scores of 0.75 (0.74), 0.74 (0.79),

443   0.73 (0.71), respectively, for positive (negative) ion mode.

444

445   *Reduced between-batch variation in QC samples*: Next, we compared the within-batch variance of the QC

446   samples with respect to the total variance which is expressed by the WTR score (Methods) for each

447   normalization method. WTR scores close to 1 indicate the absence of batch effects. *None-Raw* and *Log-*

448   *Raw* had low WTR scores and after normalizing these scores increased (Figure 2E). *BC-Metchalizer* and

449   *Log-Metchalizer* scored among the highest on this WTR score. *None-Anchor* had high WTR scores, but

450   since *None-Anchor* uses the QC samples for normalization the WTR scores are biased towards higher

451   values.

452

453   *Preserved feature ranks in QC samples*: Removal of variation results in higher WTR scores but potentially

454   removes also variation(s) of interest. Therefore, we investigated whether the ranks of the abundancies of

455   the different features in the QC samples remained the same as in the raw data (expressed as the QC rank

456   differences, $\Delta R$, see Methods for details). A lower rank difference indicates that metabolic differences

457   present in the QC samples were conserved after normalization. Figure 2F shows the QC rank differences

458   for each normalization method. These results confirm the previous believe that *Log-NOMIS* and *Log-RUV*

459   also removed non-batch related variations (higher $\Delta R$), since they had relatively high $\Delta R$'s. *BC-*

460   *Metchalizer* and *Log-Metchalizer* showed rank differences but were lower than most other competing

461   methods. *None-Anchor* showed high QC rank differences, but this is again the result of the fact that *None-*

462   *Anchor* uses the QC samples for normalization.

463

464   Taken together, *BC-Metchalizer*, *Log-Metchalizer* and *None-Anchor* showed the most consistent

465   improvement across the evaluation metrics.

466

467

468

469

Figure 2. PCA plots for raw data and normalized data as indicated by the title of each panel. Each batch is indicated with a unique color. PCA was performed on 758 features (excluding the internal – and external standards) in positive ion mode. The squares indicate QC samples whereas the circles indicate patients and controls samples.
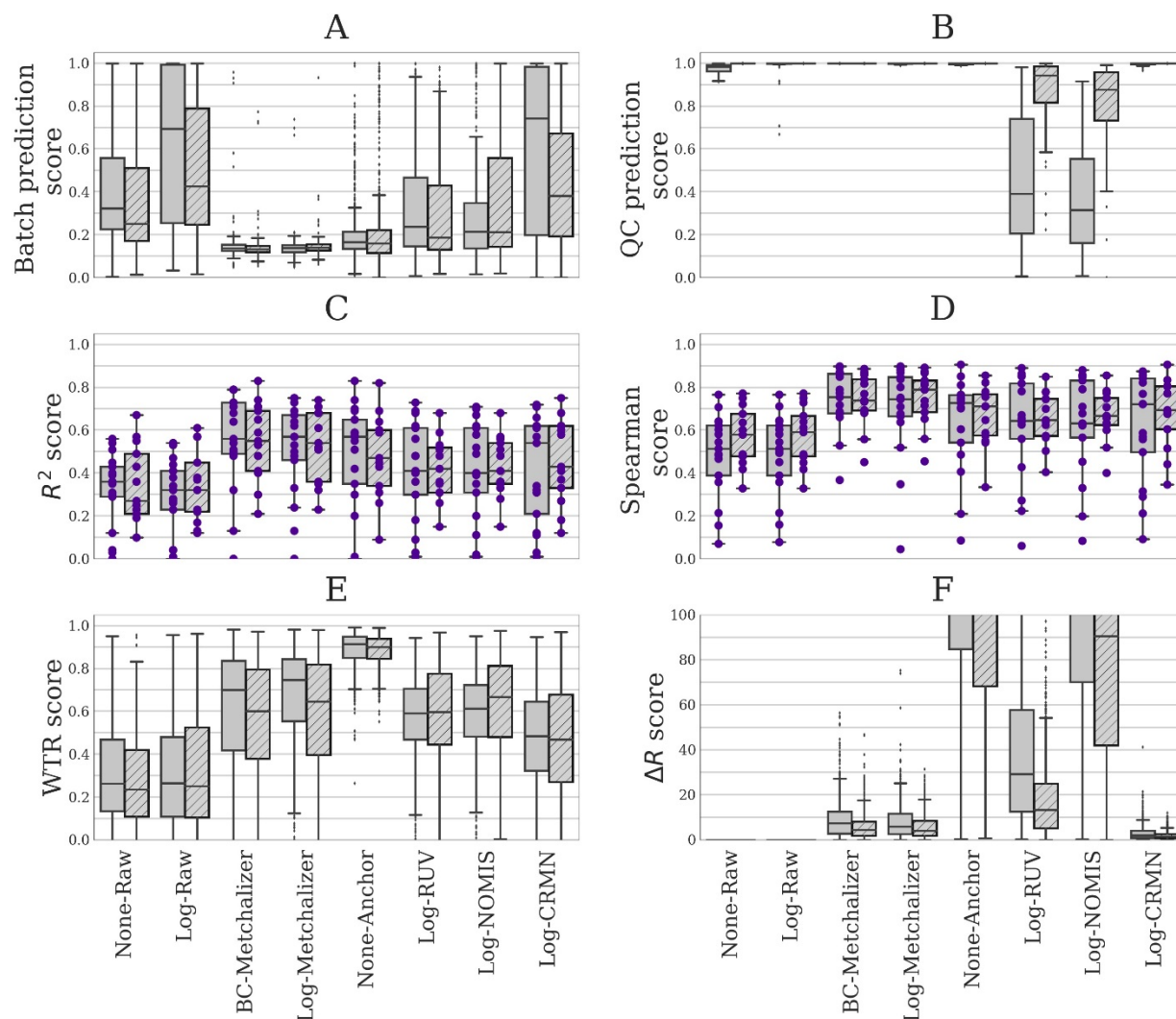
Figure 3. Six different performance metrics for batch effect removal (see Methods for more details). Data from positive – or negative ion mode is indicated by plain and stripped boxplots, respectively. A) *Batch prediction score* measures the presence of batch effects in the first four PC's from PCA analysis. B) *QC prediction score* measures how well QC samples are separated from human plasma sample in the first four PC's. C) $R^2$ score between (normalized) abundancies and quantitative measurements. D) Spearman score of (normalized) abundancies with quantitative measurements. E) The WTR score measuring the overall within batch variation with respect to the total variance using the QC samples. F) $\Delta R$ score measuring the preservation of the rank of features based on their abundancy in the QC samples before and after normalization.

**Confounder effects of age and sex**

To explore confounding effects of age and sex on metabolite abundancies, we developed a regression model with sex as linear covariates and age as a polynomial ($p=1,2,3$) covariate (see Methods). After normalization, we fitted the model parameters for every feature using all control samples present in the

487   eight batches and determined the significance of the coefficients in the regression model (see Methods).

488   Table 1 shows the percentages of (strong) significant coefficients ($\alpha = 2.7e^{-3}$) per ion mode and (selected)

489   normalization methods. Our findings suggest that 6-24% of all features showed age dependency when

490   looking at coefficient $\hat{\beta}_1^{\text{Age}}$ (i.e. the linear term in the model). It is noteworthy that more age-related features

491   were found in the negative ion mode.

492

493   Age-dependent metabolites (supplement S3 Table S3), when using normalization by *BC-Metchalizer*,

494   include known IEM biomarkers, such as: guanidinoacetic acid(+), homoarginine(-) and N-acetyltyrosine(-

495   ) , 2-ketoglutaric acid(-), citrulline(-) and ornithine(-). As an example, we plotted the regression model for

496   guanidinoacetic acid (Figure 3), illustrating that the Z-score for a fixed abundance depends on age (and

497   slightly on sex at later ages). This also shows a non-linear trend with age. Our analyses showed that more

498   metabolites have significant non-linear trends over age ($\hat{\beta}_2^{Age}$ and $\hat{\beta}_3^{Age}$ in Table 1). Moreover, age dependent

499   features have the tendency to increase/decrease in abundance faster for decreasing age, implying that a

500   matching reference population on younger ages seems to be more important (supplement S3 Figure 5).

501

502   Hardly any significant gender-related features were found (Table 1). When significance on $\hat{\beta}^{\text{Sex,Age}}$ was

503   relaxed ($\alpha = 0.05$), we found some biomarkers showing an interaction between age and sex, such as:

504   malonic acid(+/-), guanidinoacetic acid(+), homoarginine(-), ornithine(-), sebacic acid(+/-). See

505   supplement S3 Table 4 for a full list.

506

507

508   Table 1. Percentage of strongly significant ($\alpha = 2.7e^{-3}$) regression coefficients of the covariates age and sex when using

509   the regression model (Methods) predicting 758 positively - and 583 negatively ionized features, for the different

510   normalization method and ion modes.

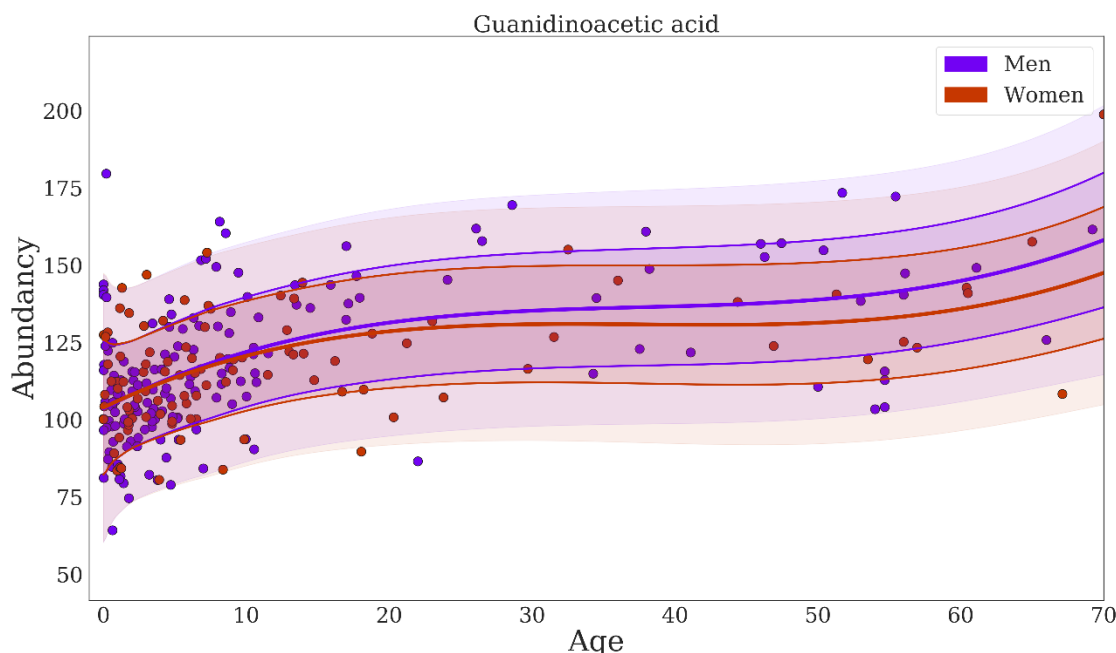| Normalization method | Ion mode | $\hat{\beta}^{\text{Intercept}}$ (%) | $\hat{\beta}_1^{\text{Age}}$ (%) | $\hat{\beta}_2^{\text{Age}}$ (%) | $\hat{\beta}_3^{\text{Age}}$ (%) | $\hat{\beta}^{\text{Sex}}$ (%) | $\hat{\beta}^{\text{Sex,Age}}$ (%) |
|---|---|---|---|---|---|---|---|
| *None-Anchor* | − | 97 | 15 | 6 | 4 | 0 | 1 |
| *None-Anchor* | + | 99 | 6 | 4 | 2 | 0 | 0 |
| *BC-Metchalizer* | - | 100 | 23 | 12 | 8 | 0 | 1 |
| *BC-Metchalizer* | + | 100 | 11 | 6 | 4 | 0 | 0 |
| *Log-Metchalizer* | - | 100 | 24 | 15 | 10 | 1 | 0 |
| *Log-Metchalizer* | + | 100 | 10 | 6 | 4 | 0 | 0 |

511

512

513

19

515 Figure 1. Regression of guanidinoacetic acid when using *BC-Metchalizer* normalized data. The different colors indicate the

516 sex as shown in the legend. The thick red/blue line indicates the average obtained from the fit on all controls for a given sex.

517 The first standard deviation is indicated by the thin(ner) line whereas the second standard deviation ends at the shaded region.

518

## Detection of the expected IEM biomarkers

519

520 Next we investigated the impact of normalization and using out-of-batch controls on expected biomarker

521 detection in the 45 IEM patients by plotting the number of detected expected biomarkers against the average

522 number of positives features per patient at various Z-score or p-value thresholds (Methods), similar to a

523 Receiver Operator Curve (ROC). Untargeted metabolomics did not allow us to make a distinction between

524 false positives (FP) and true positives (TP), due to unannotated features and even unknown disease related

525 features/biomarkers. Assuming that the majority of the positives per patient are false positives, we used

526 the average number of positives per patient as proxy for the false positives. Improved performance was

527 considered to increase the number of detected expected biomarkers (true positives of which we are certain)

528 while lowering the average number of positives per patient, thereby increasing the Area Under the Receiver

529 Operator Curve (AUC). We decided to take the method that uses 15 within-batch controls and raw

530 abundancies (*15in&None-Raw*) as the reference approach, where the performance was expressed as a

531 percentage of this reference AUC, named $\mathrm{AUC}^{x}_{15in\&None\text{-}Raw}$. Here X indicates if the AUC was created from

532 the average Z-scores or p-values. These p-values were obtained from the Welch's t-test which tests whether

20

533    the average Z-score of an expected biomarker or feature across the triplicate significantly differs from the

534    average Z-score of the reference population  (Methods).

535

536    *Log-transform improves biomarker detection for p-values*: Our first observation is that, when considering

537    the Z-scores,  the log-transformed raw abundancies (*15in&Log-Raw*) has an AUC approximately equal to

538    $AUC^Z_{15in\&None\text{-}Raw}$ (Figure 4), implying that this transformation hardly affected this performance metric.

539    However, when using the p-values, the log-transformation improved the detection of the expected

540    biomarkers, as $AUC^p_{15in\&Log\text{-}Raw}$ is 8% higher than the $AUC^p_{15in\&None\text{-}Raw}$ (Figure 4).

541

542    *Reduced performance with age/sex matched out-of-batch controls*: When comparing the performance of

543    using 15 out-of-batch controls (*15out&Raw*) to the *15in&Raw* reference, the performance for *15out* was

544    clearly reduced (Figure 4 A), achieving only 80% of the reference $AUC^Z_{15in\&None\text{-}Raw}$. This difference was

545    also present when looking at the p-values, resulting in a clear reduction of the $AUC^p_{15out\&None\text{-}Raw}$  (74%).

546    Hence, potential improved age/sex matching for *15out,* due to the increased number of available controls

547    (supplement S4 Figure 6), did not result in improved performance, most likely due to the dominance of

548    batch effects.

549

550    *Normalization improves performance of age/sex matched out-of-batch controls*: After normalizing with

551    *BC-Metchalizer*, *Log-Metchalizer* or *None-Anchor* and using 15 out-of-batch controls (*15out*), the

552    performance increased when compared to *15out&None-Raw* (Figure 4 A, B and C), and came close to the

553    $AUC^Z_{15in\&None\text{-}Raw}$*;* for *BC-Metchalizer* (94%) and *Log-Metchalizer* (96%), while *None-Anchor* stayed

554    behind (90%). Interestingly, when considering biomarker detection performance using the p-values, *BC-*

555    *Metchalizer* performed on par with  *15in&None-Raw* (99%), *Log-Metchalizer* improved over *15in&None-*

556    *Raw* (105%)  and *None-Anchor* stayed behind (90%). *Log-Metchalizer* performed similar to *15in&Log-*

557    *Raw* (105% and 108%, respectively), indicating that out-of-batch controls can be used instead of in-batch

558    controls to determine reference values.

559

560    *Regression model effectively models age and sex effects*: The regression model (*Regression*) slightly

561    improved $AUC^Z$ with respect to *15out* for *BC-Metchalizer* (+2%) and *None-Anchor* (+4%), but not for *Log-*

562    *Metchalizer* (-1%), see also Figure 4 A, B and C. When considering the p-values, $AUC^p$, only *BC-*

563    *Metchalizer* (+1%) improved but not *None-Anchor* (-2%) and *Log-Metchalizer* (-1%), although these

564    performance differences in all cases were small (Figure 4 D, E and F). Interestingly, when we took all

565    controls to determine the Z-scores (*All controls*, Methods), similar $AUC^Z$ performances were observed when

21

566 compared to *Regression,* i.e. -1% for *BC-Metchalizer* and +2% *Log-Metchalizer* and +1% for *None-Anchor.*

567 When considering the p-values the difference were larger, i.e. -5% for *BC-Metchalizer* and -1% *Log-*

568 *Metchalizer*, and -5% for *None-Anchor,* suggesting an influence of age- and sex effects on the detection
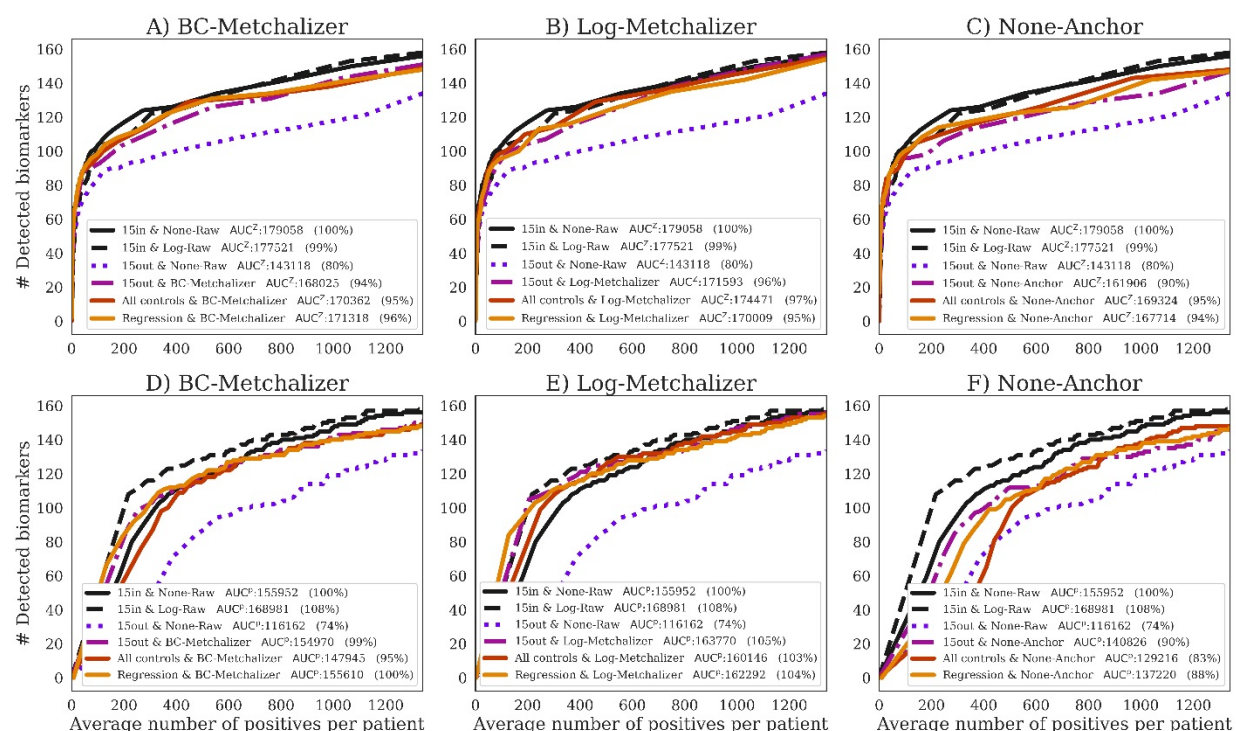
569 of biomarkers.

570



571

572

573 Figure 2. The number of detected expected biomarkers versus the average number of positives per patient. A curve in each

574 (sub)figure was formed by increasing the Z-score or p-values threshold ($Z_{abnormal}$, Methods). The legend indicates (per curve)

575 the methods used to determine Z-scores and how data was normalized, the AUC and AUC expressed as percentage of the

576 $\mathrm{AUC}^{Z}_{15in\&None\text{-}Raw}$. Performances using A) *BC-Metchalizer* using Z-scores, B) *Log-Metchalizer* using Z-scores, C) *None-*

577 *Anchor* using Z-scores, D) *BC-Metchalizer* using p-values, E) *Log-Metchalizer* using p-values, and F) *None-Anchor* using

578 p-values.

579

22

580

581 **Discussion**

582 Targeted measurements of metabolites in body fluids using various platforms such as HPLC, GC-MS and

583 LC-MS/MS are traditionally applied for laboratory diagnosis of IEM. For each individual metabolite, age-

584 and, sometimes, sex-dependent reference ranges are established using hundreds of reference samples.

585 Untargeted metabolomics is a promising alternative enabling the determination of many metabolites in one

586 analysis. This can speed up the diagnostic process and extends the number of different IEMs that can be

587 screened in a single run. A major drawback of current approaches is that reference samples need to be

588 included in the same experimental batch to ensure proper reference ranges (or Z-score transformations).

589 Some methods do exist that use reference samples measured in different batches (out-of-batch controls) to

590 determine age and sex corrected Z-scores, and they are based on normalizing methods that remove the batch

591 effects. There has not been a comprehensive comparison of the different normalization methods with

592 approaches that use out-of-batch controls, which we have set out in this work. Moreover, we developed a

593 new normalization method, *Metchalizer*, that makes use of isotope-stable internal standards, an approach

594 that has been shown to be useful when mapping specific metabolites to specific internal standards (Körver-

595 Keularts, et al., 2018) which we generalize to all features measured. Because more reference samples are

596 available when using the out-of-batch controls, we additionally propose a regression model that

597 incorporates sex and age effects as (non-linear) covariates. Alltogether, we have shown that our

598 methodology has biomarker detection performances at least similar to using 15 within-batch controls.

599

600 Typically, around 20,000 features in both negative and positive mode were detected per batch. When we

601 require a feature to have been measured (and matched) in all eight batches, we retained 598 positive and

602 773 negative ionized features, respectively. As some normalization methods use a statistical approach

603 (*PQN, Fast Cyclic Loess*), the reduction in features might explain the reduced performance of these

604 methods. In addition, the requirement of features being measured (and matched) across all eight batches

605 also resulted in the loss of some biomarkers, which hampered the performance of all out-of-batch methods

606 with respect to the within-batch methods. As an alternative, we could have made the inclusion of features

607 dependent on fewer batches (for example being present in >5 out of 8 batches). We decided not to do that

608 as this would have resulted in an unequal number of control samples for the different features, leading to

609 inconsistent results between the out-of-batch methods. The availability of more batches could have solved

610 this issue because an equal number of control samples could likely be obtained per feature, even when these

611 features were not present/matched in some batches. It is interesting to note that our proposed normalization

612 method (*Metchalizer*) showed consistent performances when data from various number of batches is being

613 used (supplement Figure S7). Some biomarkers, for example isobutyrylglycine, were only detected in the

23

614　　batches having patients with elevated levels of these specific metabolites. We anticipate that for this kind

615　　of biomarkers out-of-batch strategies are not useful since abundancies in (healthy) controls are (very) low,

616　　thereby making Z-score calculation unsuitable.

617

618　　*Anchor* uses an anchor (fixed) samples, measured in all batches, to normalize the features. *Anchor*

619　　normalization on none-transformed data performed well when compared to most of the other normalization

620　　methods explored, but slightly less than *BC-Metchalizer* and *Log-Metchalizer* when considering the

621　　performance metrics Spearman score, $R^2$ score, batch prediction score and performance on biomarker

622　　detection. We anticipate that the anchor samples may not correlate with all types of variation like, for

623　　example, injection volume which is a source of variation at the sample level, whereas the abundancy of the

624　　internal standards (used by *Metchalizer*) is directly linked to the injection volume. *Anchor* also assumes

625　　that metabolite levels remain constant over time in the anchor samples. As a consequence, if for example

626　　storage effects take place, *Anchor* is impeded. The use of *Anchor* may also be less practical because it

627　　requires the same anchor samples in every batch. The introduction of a new anchor sample requires an

628　　'overlapping batch' containing a set of both the former anchor sample together with the newly introduced

629　　anchor samples.

630

631　　*Metchalizer* exploits the linear relationship between the abundancy of a feature and those of the latent

632　　variables that are derived from the partial least squares between the internal standards and the features

633　　measured across all samples and capturing the covariance between the standards and the features (Methods).

634　　*Metchalizer* assumes that this relationship holds across batches and with that assumption determines (batch)

635　　intercepts that correct for the 'unexplained' batch/technical variations. Consequently, when such linear

636　　relationship between internal standards and features does not exist, the normalization would be fully based

637　　on the (batch) intercepts, deteriorating the power of this approach. Alternatively, when batch differences

638　　(represented by the intercepts) differ from each other due to biological variations between batches, this will

639　　be interpreted as 'unexplained' batch/technical variations, and consequently wrongly removed by

640　　*Metchalizer*. For this reason, it is important to use randomized control samples in every batch (in terms of

641　　age, sex etc) to minimize the possibility of biological variations between batches.

642

643　　*Log-Metchalizer* log transforms the abundancies before applying *Metchalizer*, whereas *BC-Metchalizer*

644　　uses a less strong Box-Cox transformation. The effect of this stronger transformation on most investigated

645　　metrics in this study was small, although we did observe that a stronger initial transformation led to

646　　improved biomarker detection performances when considering the p-values. *15in&None-Raw* had a lower

647　　$\text{AUC}^\text{p}$ than *15in&Log-Raw* and could therefore also explain the improved performance of *Log-Metchalizer*

24

648 over *BC-Metchalizer* on this metric. A simulation showed that log-transforming the raw abundancies indeed

649 caused differences in the obtained Z-scores and p-values when compared to the raw abundancies

650 (supplement S10). Positive Z-scores had relatively lower p-values (and vice versa) for log-transformed

651 abundancies and this could therefore explain the improved performance on biomarker detection, since most

652 of the considered biomarkers had positive Z-scores, thus biasing this performance metric. Increasing the

653 number of internal standards did not improve the normalization performance when considering metrics

654 based on the quantitative measurements, although we observed that certain combinations of internal

655 standards improved normalization of specific metabolites (supplement S6). This suggests that *Metchalizer*

656 might be improved by matching features/metabolites with a certain set of internal standards (for example

657 based on retention time).

658

659 We were a bit surprised that biomarker detection performance using the Z-scores ($AUC^Z$) for the regression

660 model was similar or slightly less than using all controls, as abundancies are known to be dependent on age

661 and sex. One explanation might be that only a subset of the considered (expected) biomarkers have an age

662 and/or sex dependency. Indeed, when we considered only these age-dependent biomarkers (19 biomarker-

663 patient combinations, supplement S3 Table 3), the performance of *Regression* was more improved than *All*

664 *controls* (supplement S8). However, this set was small, so substantial evidence to support this improvement

665 is lacking. Furthermore, our proposed performance metric assumed that the average number of positives

666 was a proxy for the average number of false positives. Using *Regression* resulted generally in more positives

667 (data not shown), but these were not necessarily merely false positives, which therefore could have affected

668 the performance of *Regression* negatively. Though, when judging biomarker detection using the p-values,

669 we did see that *Regression* slightly outperformed *All controls*.

670

671 In conclusion, out of all explored normalization methods, the removal of batch effects was best performed

672 by *Log-Metchalizer*. Fitting our regression model on the corresponding normalized data showed that 10-

673 24% (Table 1) of all considered features were depending on age, underlining the need for using age

674 corrected Z-scores. On average, biomarker detection performance using *Log-Metchalizer* using out-of-

675 batch controls was at least similar to the best performing *Log-Raw* approach when using the 15 within-

676 batch controls (*15in&Log-Raw*). We anticipate that the success of *Metchalizer* and age- and sex correcting

677 strategies such as our regression model depend on three factors: 1) a feature of interest being measured in

678 a number of other batches (not necessarily all), 2) batch effects containing (only) technical variations, and

679 3) abundancies being affected by age or other covariates (the presence of an effect-size). Together our

680 proposed approach opens new opportunities to improve abnormality detection, especially for age-dependent

681 features/biomarkers.

25

682

# References

684  Ballman, K. V., Grill, D. E., Oberg, A. L. & Therneau, T. M., 2004. Faster cyclic loess: normalizing RNA
685  arrays via linear models. *Bioinformatics,* 5, Volume 20, pp. 2778-2786.

686  Bonte, R. et al., 2019. Untargeted Metabolomics-Based Screening Method for Inborn Errors of
687  Metabolism using Semi-Automatic Sample Preparation with an UHPLC- Orbitrap-MS Platform.
688  *Metabolites,* 11, Volume 9, p. 289.

689  Chaleckis, R. et al., 2016. Individual variability in human blood metabolites identifies age-related
690  differences. *Proceedings of the National Academy of Sciences,* Volume 113, pp. 4252-4259.

691  Coene, K. L. M. et al., 2018. Next-generation metabolic screening: targeted and untargeted
692  metabolomics for the diagnosis of inborn errors of metabolism in individual patients. *Journal of*
693  *Inherited Metabolic Disease,* Volume 41, pp. 337-353.

694  Filzmoser, P. & Walczak, B., 2014. What can go wrong at the data normalization step for identification of
695  biomarkers?. *Journal of Chromatography A,* 10, Volume 1362, pp. 194-205.

696  Glinton, K. E. et al., 2019. Untargeted metabolomics identifies unique though benign biochemical
697  changes in patients with pathogenic variants in UROC1. *Molecular Genetics and Metabolism Reports,*
698  Volume 18, pp. 14-18.

699  Haijes, H. A. et al., 2019. Direct Infusion Based Metabolomics Identifies Metabolic Disease in Patients'
700  Dried Blood Spots and Plasma. *Metabolites,* Volume 9.

701  Huber, W. et al., 2002. Variance stabilization applied to microarray data calibration and to the
702  quantification of differential expression. *Bioinformatics,* 7, Volume 18, pp. S96--S104.

703  Karpievitch, Y. V. et al., 2015. Metabolomics Data Normalization with EigenMS. *PLOS ONE,* 12, Volume 9,
704  pp. 1-10.

705  Körver-Keularts, I. M. L. W. et al., 2018. Fast and accurate quantitative organic acid analysis with LC-
706  QTOF/MS facilitates screening of patients for inborn errors of metabolism. *Journal of Inherited*
707  *Metabolic Disease,* Volume 41, pp. 415-424.

708  Lawton, K. A. et al., 2008. Analysis of the adult human plasma metabolome. *Pharmacogenomics,* Volume
709  9, pp. 383-397.

710  Li, B. et al., 2017. NOREVA: normalization and evaluation of MS-based metabolomics data. *Nucleic Acids*
711  *Research,* 5, Volume 45, pp. W162-W170.

712  Livera, A. M. D. et al., 2015. Statistical Methods for Handling Unwanted Variation in Metabolomics Data.
713  *Analytical Chemistry,* 3, Volume 87, pp. 3606-3615.

714  Miller, M. J. et al., 2015. Untargeted metabolomic analysis for the clinical screening of inborn errors of
715  metabolism. *Journal of Inherited Metabolic Disease,* Volume 38, pp. 1029-1039.

716  Pedregosa, F. et al., 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning*
717  *Research,* Volume 12, pp. 2825-2830.

718  Redestig, H. et al., 2009. Compensation for Systematic Cross-Contribution Improves Normalization of
719  Mass Spectrometry Based Metabolomics Data. *Analytical Chemistry,* Volume 81, pp. 7974-7980.

720  Rist, M. J. et al., 2017. Metabolite patterns predicting sex and age in participants of the Karlsruhe
721  Metabolomics and Nutrition (KarMeN) study. *PLOS ONE,* 8, Volume 12, pp. 1-21.

722  Sysi-Aho, M., Katajamaa, M., Yetukuri, L. & Orešič, M., 2007. Normalization method for metabolomics
723  data using optimal selection of multiple internal standards. *BMC Bioinformatics,* Volume 8, p. 93.

724  Välikangas, T., Suomi, T. & Elo, L. L., 2016. A systematic evaluation of normalization methods in
725  quantitative label-free proteomics. *Briefings in Bioinformatics,* 10, Volume 19, pp. 1-11.

726  Veselkov, K. A. et al., 2011. Optimized Preprocessing of Ultra-Performance Liquid Chromatography/Mass
727  Spectrometry Urinary Metabolic Profiles for Improved Information Recovery. *Analytical Chemistry,*
728  Volume 83, pp. 5864-5872.

729  Vreken, P. et al., 2002. Rapid Diagnosis of Organic Acidemias and Fatty-acid Oxidation Defects by
730  Quantitative Electrospray Tandem-MS Acyl-Carnitine Analysis in Plasma. In: *Current Views of Fatty Acid*
731  *Oxidation and Ketogenesis.* sl:Springer US, pp. 327-337.

732  Yu, Z. et al., 2012. Human serum metabolic profiles are age dependent. *Aging Cell,* Volume 11, pp. 960-
733  967.

734

735

**Author contribution**

737  Ramon Bonte performed all the experimental work and developed the chromatographic- and mass
738  spectrometric method. Compound identification was also done by him. Michiel Bongaerts designed the
739  statistical models, the computational framework and analyzed the data. The manuscript was written by
740  Michiel Bongaerts, Henk Blom and George Ruijter. Serwet Demirdas and Ed Jacobs contributed in
741  establishing the IEM database used in this study, and actively contributed in giving feedback on the
742  methods. Marcel Reinders contributed to in-depth reviewing of the manuscript, all analytical methods and
743  suggested adjustments to initial work. Esmee Oussoren, Ans van der Ploeg, Margreet Wagenmakers and
744  Robert Hofstra provided resources. The research was under supervision of George Ruijter.

745

746  **Conflicts of Interest:** All authors state that they have no conflict of interest to declare. None of the authors
747  accepted any reimbursements, fees, or funds from any organization that may in any way gain or lose
748  financially from the results of this study. The authors have not been employed by such an organization. The
749  authors have not act as an expert witness on the subject of the study. The authors do not have any other
750  conflict of interest.

753

754     **Data and Code availability**

755     The regression model, *Best Correlated Internal Standard, PQN, Anchor* and *Metchalizer(Log)* were

756     developed in Python and are available at https://github.com/mbongaerts/Metchalizer. The code developed

757     for merging the batches can also be found here. The Progenesis QI processed data for all 8 batches is

758     available at https://github.com/mbongaerts/Metchalizer/Data. We removed the patient samples for privacy

759     reasons.

760

761
762