1    # Deep learning accurately predicts estrogen receptor status in breast

2    # cancer metabolomics data

3    Fadhl M Alakwaa[1], Kumardeep Chaudhary[1], Lana X Garmire[1,2*]

4    [1]Epidemiology Program, University of Hawaii Cancer Center, Honolulu, HI 96813, USA.

5    [2]Molecular Biosciences and Bioengineering Graduate Program, University of Hawaii at Manoa,

6    Honolulu, HI 96822, USA.

7    *To whom the correspondence should be addressed:

8    Lana X Garmire, PhD

9    Associate Professor

10    Email address: lgarmire@cc.hawaii.edu

11    Phone: +1 (808) 441-8193

12

13

14

15

16      **ABSTRACT**

17      Metabolomics holds the promise as a new technology to diagnose highly heterogeneous diseases.

18      Conventionally, metabolomics data analysis for diagnosis is done using various statistical and machine

19      learning based classification methods. However, it remains unknown if deep neural network, a class of

20      increasingly popular machine learning methods, is suitable to classify metabolomics data. Here we use a

21      cohort of 271 breast cancer tissues, 204 positive estrogen receptor (ER+) and 67 negative estrogen receptor

22      (ER-), to test the accuracies of autoencoder, a deep learning (DL) framework, as well as six widely used

23      machine learning models, namely Random Forest (RF), Support Vector Machines (SVM), Recursive

24      Partitioning and Regression Trees (RPART), Linear Discriminant Analysis (LDA), Prediction Analysis for

25      Microarrays (PAM), and Generalized Boosted Models (GBM). DL framework has the highest area under

26      the curve (AUC) of 0.93 in classifying ER+/ER- patients, compared to the other six machine learning

27      algorithms. Furthermore, the biological interpretation of the first hidden layer reveals eight commonly

28      enriched significant metabolomics pathways (adjusted P-value<0.05) that cannot be discovered by other

29      machine learning methods. Among them, protein digestion & absorption and ATP-binding cassette (ABC)

30      transporters pathways are also confirmed in integrated analysis between metabolomics and gene expression

31      data in these samples. In summary, deep learning method shows advantages for metabolomics based breast

32      cancer ER status classification, with both the highest prediction accurcy (AUC=0.93) and better revelation

33      of disease biology. We encourage the adoption of autoencoder based deep learning method in the

34      metabolomics research community for classification.

35      **KEYWORDS:** Breast cancer, metabolomics, estrogen receptor, deep learning, bioinformatics

36

37

38

1

39

## Introduction

41      According to Global Health Estimates (WHO 2013), more than half million women died due of breast

42      cancer worldwide[1]. Breast cancer is the second leading cause of cancer-related deaths among women in the

43      United States[2]. Based on human epidermal growth factor receptor 2 (Her2), progesteron receptor (PR) and

44      estrogen receptor (ER), breast cancer can be categorized into four molecular subtypes[3]: Luminal A (ER+,

45      PR+/- and Her2-), Luminal B (ER+, PR+/- and Her2+/-), Her2-enriched (ER-, PR- and Her2+), and triple

46      negative (ER-, PR- and Her2)[4]. The survival outcomes differ significantly among these subtypes. Luminal

47      A and B subtypes have a relatively good prognosis, however triple negative tumors and Her2 tumors have

48      very poor prognosis[5]. Identification of molecular subtypes is crucial in determining cancer prognosis and

49      therapeutic selection. Recently, many studies used metabolomics data to segregate molecular subtypes,

50      given that breast cancer is manifested as a metabolic disease[6, 7]. For example, glutamate-to-glutamine ratio

51      and aerobic glycolysis were proposed as biomarkers of ER and Her2 status, respectively[8, 9].

52      Metabolomics studies are usually done by three major platforms: gas chromatography-mass spectrometry

53      (GC-MS), liquid chromatography (LC-MS), and nuclear magnetic resonance (NMR). The parallel use of

54      these instruments allows detecting more metabolites for the same sample. Coupling with the development

55      in the instrumentations, state-of-the-art data analysis tools are much needed to handle the large amount of

56      metabolite data generated. For problems of metabolomics data classification and regression, machine

57      learning algorithms have been applied[10]. For example, Random Forest (RF) is a widely used machine

58      learning algorithm based on decision tree theory. It works with high-dimensional data and can deal with

59      unbalanced and missing values in the data[11]. Support Vector Machine (SVM) is another machine learning

60      algorithm that separates the metabolites data with N data points into (N-1) dimensional hyperplane[12]. SVM

61      was used to classify healthy and pneumonia patients based on nuclear magnetic resonance (NMR)

62      metabolomics data[12].

63 DL or deep neural network, is a new class of machine learning methods that have been successfully applied

64 to various areas of genomics research[13, 14], including predicting the intrinsic molecular subtypes of breast

65 cancer[15], inferring expression profiles of genes[16] and predicting the functional activity of genomic

66 sequence[17]. In a recent study, denoising autoencoder (DAs), a type of DL algorithm, was applied to gene

67 expression data of the breast cancer[15]. It successfully extracted features that stratify normal/tumor samples,

68 ER+/ER- status, and intrinsic molecular subtypes. In another study based on gene expression data, DL

69 outperformed linear regression in inference of the expression of target genes from the expression of

70 landmark genes[16]. Moreover, an open source conventional neural networks (CNNs) package "Basset" was

71 developed to learn the functional activity of 164 cell types DNA sequences from genomics data, and to

72 annotate the non-coding genome[17]. Compared to the flourishing applications of DL in genomics, it remains

73 unknown if deep neural network is suitable to classify metabolomics data, esp. when the samples are of

74 medium size (i.e. several hundred).

75 Here we applied feed-forward networks, a type of DL framework, as an alternative to the machine learning

76 methods such as those listed earlier, to classify metabolomics data. We examined the predictive accuracy

77 of the DL and other machine learning algorithms to predict ER status from a public metabolomics dataset[18].

78 We demonstrated this DL method performs better than a wide cluster of machine learning methods,

79 including Random Forest (RF), Support Vector Machines (SVM), Recursive Partitioning and Regression

80 Trees (RPART), Linear Discriminant Analysis (LDA), Prediction Analysis for Microarrays (PAM), and

81 Generalized Boosted Models (GBM). Furthermore, the biological interpretation of the hidden layers reveals

82 eight breast cancer related pathways such as central carbon metabolism in cancer and glutathione

83 metabolism. Moreover, we further analyzed the extracted features from our DL model, by mapping the

84 biosynthetic enzymes involved in the metabolomics pathways.

85 **Materials and Methods**

86 **Data set**

87    The metabolomics data used in this study consists of 271 breast cancer samples (204 ER+ and 67 ER-)

88    collected from a biobank at the Pathology Department of Charité Hospital, Berlin, Germany[18].

89    Metabolomics profiles of these BC patients can be downloaded from the supporting material of this study[19].

90    A total of 162 metabolites with known chemical structure were measured using gas chromatography

91    followed by time of flight mass spectroscopy (GC-TOFMS) for all tissues samples. A detailed description

92    of the protocols and the platforms used in this study were described in [18]. For validation, we downloaded

93    gene expression dataset GSE59198[20] from the Gene Expression Omnibus (GEO) database, which is

94    composed of 154 samples, a subset of the 271 samples. In this data set, the gene expression profiles of BC

95    tumor tissues (122 ER+ and 32 ER-) were analyzed using the cDNA-mediated Annealing, Selection,

96    Extension and Ligation (DASL) assay. A total of 15,927 genes were detected (p<0.01) in at least 10% of

97    the samples after applying spline normalization. Data can be downloaded from GEO repository

98    http://www.ncbi.nlm.nih.gov/geo.

99    **Data Preprocessing**

100   We used K-Nearest Neighbors (KNN) method to impute missing metabolomics data[21]. To adjust for the

101   offset between high and low-intensity features, and to reduce the heteroscedasticity, the logged value of

102   each metabolite was centered by its mean ($\bar{x}$) and autoscaled by its standard deviation ($s$) as described in

103   Equation 1[22]. We used quantile normalization to reduce sample-to-sample variation[23].

104   $$\hat{x}_{ij} = (\frac{log_2(x_{ij}) - \bar{x}_i}{s}) \qquad (1)$$

105   **Deep Learning**

106   DL refers to deep neural network framework, which is widely applied in pattern recognition, image

107   processing, computer vision, and recently in bioinformatics[13, 24, 25]. Similar to other feed-forward artificial

108   neural networks (ANNs), DL employs more than one hidden layer ($y$) that connects the input ($x$) and output

109   layer ($z$) via a weight ($W$) matrix as shown in equation (2). Here we use sigmoid function as the transitioning

110   function.

111   $$y = sigmoid(Wx + b) \qquad (2)$$

4

112     Activation value of the hidden layer ($y$) can be calculated by sigmoid of the multiplication of the input

113     sample $x$ with the weight matrix $W$ and bias $b$. The transpose of the weight matrix W and the bias $b$ can

114     then be used to construct the output ($z$) layer, as described in equation (3).

$$z = sigmoid(W'y + b') \qquad (3)$$

116     The best set of the weight matrix $W$ and bias $b$ are expected to minimize the difference between the input

117     layer ($x$) and the output layer ($z$). The objective function is called cross-entropy in equation (4) below, in

118     which the optimal parameters are obtained by stochastic gradient descent searching.

$$L_H(x,z) = -\sum_{k=1}^{d}[x_k \, logz_k + (1 - x_k)\log(1 - z_k)] \qquad (4)$$

120     To train the model, we first supplied sample input ($x$) to the first layer and obtained the best parameters ($W$,

121     $b$) and the activation of the first hidden layer ($y$), and then used y to learn the second layer. We repeated

122     this process in subsequent layers, updating the weights and bias in each epoch. We then used back-

123     propagation to tune the parameters of all layers.  Finally, we fed the output of the last hidden layer to a

124     softmax classifier which assigned new labels to the samples[26]. We used *h2o* R package to tune the

125     parameters of the DL model[27].

126     **Other machine learning algorithms**

127     We selected a representative set of six machine-learning algorithms that are highly recommended by the

128     metabolomics community and applied widely in the literature reports: Random Forest (RF), Support Vector

129     Machines (SVM), Recursive Partitioning and Regression Trees (RPART), Linear Discriminant Analysis

130     (LDA), Prediction Analysis for Microarrays (PAM), and Generalized Boosted Models (GBM). To get the

131     optimal predictions, we used the *caret* R package[28] to tune the parameters in the models.

132     **Modeling and evaluation**

133     We randomly split metabolomics samples into 80% training set and 20% testing set. The 80/20 split is a

134     common practice of splitting ratio for samples of a moderate size in the machine learning applications. We

135     chose this ratio in order to having enough training samples to build a good model and sufficient testing

136     samples to evaluate the model. We performed 10-fold cross-validation on the 80% training data during the

137     model construction process, and tested the model on the hold out 20% of data. We used pROC R package[29]

5

138    to compute area under the curve (AUC) of a receiver-operating characteristic (ROC) curve to assess the

139    overall performance of the models. To avoid sampling bias, we repeated the above splitting process ten

140    times and calculated the average AUC on the hold out 10 test samples. To control overfitting, we used two

141    regularization parameters: L1, which increases model stability and causes many weights to become 0 and

142    L2, which prevents weights enlargement.

143    We tuned DL model and other machine learning algorithms, on the following parameters: DL model:

144    Epochs (number of passes of the full training set), $l1$ (penalty to converge many weights to 0) and $l2$ (penalty

145    to prevent weights enlargement), and input dropout ratio (ratio of ignored neurons in the input layer during

146    training), number of hidden layers; RPART model: complexity parameters (cost of adding node to the tree);

147    GBM model: number of trees and interaction depths; SVM model: cost of classification; RF model: number

148    of trees to fit; PAM model: threshold amount by for each of the class's centroid shrinking towards the all

149    classes' centroid.

**Feature importance**

151    Features importance was estimated based on model based approach[28]. In other words, a feature is

152    considered important if it contributes to the model performance[30]. We used the variable importance

153    functions varimp in *h2o* and varImp in *caret* R packages, to evaluate the top 20 features.

**Identifiers standardization and differentially expressed genes**

155    We used the PubChem Identifier exchange service[31] to convert metabolites into their corresponding KEGG

156    compound IDs; we then used KEGG API[32] to get the compound pathways and enzyme IDs. We used *limma*

157    R package[33] to find enzymes with high fold changes as well as significant adjusted p-values between ER+

158    and ER- samples.

**Metabolomics enzymes network reconstruction and visualization**

160    We used MetaScape[34] v3.1.3, a Cytoscape plug-in to generate gene-metabolite network which integrates

161    reaction and pathway information from KEGG and Edinburgh human metabolic network (EHMN)

162    databases. To build enzyme-metabolite network, we selected a pathway based network from Metsacpe

6

163    analysis options. The input of this step were two files. The first file included the compound KEGG IDs, p-

164    value and the fold change values of the top 20 metabolites extracted from the DL model.  The second file

165    included the enzyme KEGG IDs, p-value and the fold change values of the 898 genes whose expression

166    values were statistically significantly different between ER- and ER+ samples.

167    **Metabolites enzymes correlation**

168    We calculated the correlations between the intensity levels of the metabolites and enzymes using

169    Spearman's Correlation Coefficient in R. We plot the Circos plot of the strongest correlation using Circlize

170    R package v0.4.0.

171    **Joint significant pathway analysis**

172    To perform joint significant pathway analysis on metabolomics and gene expression data from the same

173    samples, we considered a comprehensive list of pathways from Reactome, EHMN, and KEGG databases,

174    using online web tool IMPaLA[35], and calculated hypergeometric p-values of genes ($P_G$) and metabolites

175    ($P_M$). The joint P-value ($P_j$) between metabolites and genes for pathway $i$ was calculated as $P_{ji} = P_{Gi} P_{Mi}$[36].

176    This value was adjusted to control for multiple testing with the False Discovery Rate method.

177    **Code availability**

178    We include all preprocessing and the learning steps of the DL method as an R script in the supplementary

179    file 1.

180    # Results

181    **Workflow of autoencoder based classification**

182    We aim to assess the predictive ability of the DL framework to separate breast cancer patients based on

183    their ER status, using metabolomics data.  Towards this goal, we implemented the workflow of DL

184    framework as in **Figure 1**. We applied preprocessing steps (log transformation, centering, autoscaling , and

185    quantile normalization) before constructing the DL model, as recommended by others[18, 22]. Before training

186    the model, we pre-trained the model using autoencoder and the whole data without labels. This step

187    improves the model performance, avoids random initialization of the weights, and selects the best model

188    architecture[37]. Then we trained the DL model using a wide range of parameters and selected the best model

189    with the minimum mean square error (see Materials and Methods).

190    **Performance of the autoencoder based deep learning classification**

191    We compared DL with six other machine-learning methods commonly used in the community: Random

192    Forest (RF), Support Vector Machines (SVM), Recursive Partitioning and Regression Trees (RPART),

193    Linear Discriminant Analysis (LDA), Prediction Analysis for Microarrays (PAM), and Generalized

194    Boosted Models (GBM).  To assess the predictive power of the models, we partitioned the data into 80%

195    training and 20% testing subsets. We performed 10-fold cross-validation on the 80% training data, and

196    tested the model on the hold out 20% of data. To avoid sampling bias, we performed 10 independent

197    splitting of training and testing subsets. We reported the averaged AUCs calculated on the hold out test

198    sets. As shown in **Figure 2A**, the average AUC of DL yields the best AUC of 0.93, compared to other six

199    classification methods. The superiority of DL accuracy is statistically significant (Wilcoxon signed-rank

200    test P<0.05) than other methods, except RF and GBM. LDA and RPAT had the worst accuracy, likely due

201    to their sensitivity to overfitting and unfit to the non-linear problems[38].

202    DL as other machine learning algorithm needs more samples to achieve high accuracy[39]. To assess the

203    effect of sample size on various models, we randomly removed ¼, ½, and ¾ of the data sets (**Figure S1**).

204    As expected, decreasing in sample size decreases the averaged AUCs of all classification methods in

205    general except LDA on ¼ samples, due to overfitting. Notably, the reduction of average AUC in DL is most

206    pronounced among all methods, from the full to ¾ data set (**Figure S1**). While DL loses the best average

207    AUC status when the sample size is around 255, GBM, SVM and RF have the highest AUC for small

208    sample sizes of 203, 136 and 68, respectively. Similarly, we also experimented the effect of metabolite size

209    on various models (**Figure S2**). We randomly removed ⅛, ¼, and ½ of the 162 metabolites. Even with

210    reduced numbers of metabolites, deep learning and the robust machine learning method SVM still have

211    fairly good predictions, compared to other algorithms tested. This suggests that, due to colinearality, much

8

212    of information still exist in the remaining metabolites. Together, the drop-out experiments (**Figures S1 and**

213    **S2**) demonstrate that DL method is sensitive to sample size, but much less sensitive to metabolite size.


214    **Important features from DL**

215    To relate the importance of metabolites to ER status directly, we ranked the metabolites extracted from DL

216    model based on their functional contributions to the outputs. In this approach, features that provide unique

217    information to the trained network are ranked more importantly than those giving redundant information[40].

218    We listed the top 20 metabolites from DL in Table S1, and presented their heatmap and boxplots in **Figure**

219    **S3**. Note the choice of 20 metabolite is guided by the original study, in which 19 out of 162 metabolites

220    were claimed to change significantly among training and validation samples[19]. The original author divided

221    the 271 samples into two parts, the training (2/3) and the validation (1/3) set. Among the training set, 65

222    metabolites were different in ER- and ER+ and only 19 metabolites were validated in the validation set.


223    Among the 20 features, the top five features are beta-alanine, xanthine, isoleucine, glutamate, and taurine.

224    These five metabolites have been either proposed as breast cancer biomarkers or associated with breast

225    cancers in the original metabolomics report[19] and/or other studies[6, 8, 41-43]. For instance, Budczies et al. [19]

226    found that beta-alanine had the most significant and largest fold changes between ER-(n=67) and ER+

227    (n=204) tumor tissues. In another study, Glutamate was suggested as markers to segregate ER- from ER+

228    in the training (n=186) as well as validation dataset (n=88)[8]. Glutamate to glutamine ratio (GGR) was

229    significantly increased in the ER- tumors as compared to ER+. Overall survival analyses suggested GGR

230    as a positive prognostic marker for BC[8]. In another study, Fan et al. classified BC plasma samples into

231    subtypes i.e. ER+ vs ER- and HER2+ vs HER2-, based on a training set (n=51) and another test set (n=45)[6].

232    They found isoleucine had significant differential level between ER+ (lower) and ER- (higher) samples.

233    Similarly, a study among female breast cancer patients (n=50) suggested serum taurine as an early marker,

234    where its level was significantly lower than the normal (n=20) and high risk samples (n=15)[42]. In a cell line

235    based study, xanthine was suggested as potential biomarker of breast cancer metastasis[43], as it had the

9

236     highest variable influence on projection (VIP) in the three pair-wise comparisons among MCF-7/MCF-

237     10A, MDA-MB-231/MCF-10A and MDA-MB-231/MCF-7[43].

238     Further, we compared DL top 20 features with the same number of top features from all other methods in

239     a bipartite graph (**Figure 2B**). Twelve metabolites are shared between DL and one or more algorithms.

240     Among them, 1 (xanthine) is shared by six methods, 2 ( glyceric acid and citrulline) are shared by five

241     methods, 4 (glutamine, taurine, glutamine acid, and beta-alanine) are shared by four methods, 1 (2-

242     aminoadipic acid) is shared by three methods, 2 (nicotinamide acid and trehalose) are shared by two

243     methods, and two (linoleic acid and hypoxanthine) are shared by one method  (Table S1). Additionally, DL

244     has identified 8 unique metabolites: isoleucine, putrescine, glycerol, 5'-deoxy-5'-methylthioadenosine,

245     ornithine, tocopherol beta, phenylalanine, and arachidonic acid,

246     **The biological relevance of the hidden layers**

247     To understand the high performance of the DL model, we probed into the hidden layer and analyzed the 25

248     activation nodes from the first hidden layer. Among the top 12 nodes with the variances $> 0.1$, node 8, 22

249     and 25 are significantly correlated with the samples' ER- status (P=1.14e-12), whereas all other top 9 nodes

250     are associated with the ER+ status (**Figure 3A**). These results confirm that the nodes in DL have significant

251     biological meaning.

252     We identified a total of 129 metabolites which contribute most to the activation values of the top 12 nodes.

253     Their relationships between the 129 metabolites and 12 nodes are shown in **Figure S4**. We define that

254     metabolite $x$ contributes to the activation value ($y$) of node $n$, if the aboslute value of the weight connecting

255     metabolite $x$ and node $n$ is greater that 0.1. Beta-alanine and xanthine are the most common metabolites

256     from all top 12 nodes. Among nodes 8, 22, and 25 which are highly correlated with ER- (**Figure 3A**), four

257     common metabolites are shared: inositol, glutamate, xanthine, and uracil. Xanthine was among the panel

258     of prognostic markers of breast cancer metastasis based on the metabolic profiling of the three breast cancer

259     cell lines[43]. Glutamate have been reported as biomarkers to segregate ER- from ER+ in the training as well

260    as validation dataset, as described earlier[8]. Inositol phosphate metabolism pathway was previously reported

261    to be associated with breast cancer, but not between ER+ and ER- cancers[44]. Uracil is, however, a potencial

262    new marker for ER- breast cancer that was not reported previously, according to our knowledge.


263    To link the metabolites in **Figure S4** with biological functions, we conducted pathways enrichment analysis

264    using online web tool IMPALA[35.] The pathways are taken from Reactome, EHMN, and KEGG databases.

265    Eight significant breast cancer related pathways (Figure 3B) are enriched in all nodes: protein digestion and

266    absorption, central carbon metabolism in cancer, neuroactive ligand receptor interaction, ABC transporters,

267    mineral absorption, inositol phosphate metabolism, glutathione metabolism, and cysteine and methionine

268    metabolism. Albeit the name of "Neuroactive ligand-receptor interaction", this pathway is significantly

269    enriched (q-value=0.001) and it was shown changed in breast cancer cell lines [45] and naked mole rat [46].

270    Aspartate, glucine, taurine and glutamate are metabolites associated with this pathway in the metabolic

271    dataset. Another interesting pathway with the name "mineral absorption" also shows significance (q-

272    value=7.51E-06), attributed by five metabolites tryptophan, alanine, glycine, phosphoric acid, glutamine.

273    All these five metabolites were found related with breast cancer previously[47-49] .


**Integration of DL metabolites and enzymes**

275    We further aimed to validate the important metabolite features of DL model, by integrating metabolomics

276    and gene expression data from the same patients. Towards this, we first conducted a joint pathway analysis

277    between 20 metabolites extracted from DL model and 898 significantly differentiated enzymes between

278    ER+ and ER- samples, using IMPALA (Figure 4). Most of the top significant pathways are related to

279    metabolism of amino acids or protein digestion and absorption. Two pathways remain significant in joint

280    pathway analysis, by comparing to metabolomics based pathway analysis in Figure 3B:  protein digestion

281    & absorption and ABC transporters, with 6 and 9 metabolites over-represented respectively. Specifically,

282    urea, inositol allo-, phosphoric acid, glucose, glutamine, Isoleucine , and glutathione are the associated

283    metabolites in ABC transporters. For protein digestion, glutamine, lysine, isoleucine, and beta-alanine are

284    associated metabolites. Some literature evidence shows that protein digestion and ABC transporters are

11

285    related to breast cancer. For example, humans have 49 members of the ATP-binding cassette (ABC)

286    membrane proteins[50]. Several of them such as ABCB1 and ABCC1 have developed a resistance to drug

287    "multidrug resistance" (MDR) in breast cancer, when they are over-expressed over a period of time[51].

288    To gain insights at individual metabolite/enzyme level, we then calculated Spearman correlations between

289    the intensity levels of the top 20 metabolites and enzymes whose gene expression levels are significantly

290    different between ER+/ER- for the same patients[20]. The Circos plot in **Figure 5** shows the names of

291    metabolomics and enzymes that have correlations ($|r| > 0.35$). Impressively, beta-alanine, the top ranked

292    metabolite in DL, is the single most connected metabolite, correlated to more than 100 significantly

293    differentially expressed enzymes. Pathway analysis of these enzymes correlated with beta-alanine shows

294    strikingly significant enrichment (adjusted p-value =3.84e-05) with FOXM1 transcription factor network

295    pathway. FOXM1 is highly expressed in ER- samples and with a correlation coefficient r=0.5 with beta-

296    alanine.

297    Complementary to the correlation based analysis, we also used Metscape (Cytoscape plug-in) for gene-

298    metabolite network analysis, by combining the ER+/ER- metabolomics data[18] and gene expression (from

299    GSE59198)[20] for the same patients.   ABAT, the enzyme that catalyze beta-alanine to malonate

300    semialdehyde (Figure 6B), is highly correlated with beta-alanine (r=-0.62, Figure 6A).  To understand better

301    the connection between beta-alanine and FOX genes family, we performed motif enrichment analysis for

302    the enzymes interacted with beta-alanine in Figure 6B using PASTAA tool[52]. Interestingly, FOXO1 was

303    one of most significant transcription factors (p= 5.89e-04) that targeted the promoters regions of beta-

304    alanine interacted enzymes.

305

## Discussion

307    Metabolomics has become a new platform for biomarker discovery. Accompanying this technology, robust

308    and accurate classification methods to predict sample labels are in critical need. Recently, DL methods have

309    gained much attention in domains such as genomics and imaging analysis. However, there has not been any

12

310      systematic investigation of DL methods in the metabolomics space. In this report, we aimed to fill this void

311      and assessed the performance of feed-forward network, a widely used DL framework, on classifying

312      ER+/ER- breast cancer metabolomics data.

313      There are many advantages of DL over shallow machine learning algorithms, which are beyond the scope

314      of this study. The conventional machine learning algorithms require engineering domain knowledge to

315      create features from raw data, whereas DL automatically extracts simple features from the input data using

316      general purpose learning procedure. These simple features are mapped into outputs using a complex

317      architecture composed of a series of non-linear functions "hierarchical representations," to maximize the

318      predictive accuracy of the model optimally. By increasing number of layers and neurons per layers, robust

319      features may be constructed, and error signals can be diminished as they pass through multiple layers[13].

320      Therefore, DL succeeds to construct high-level transformed features from input data, making it more

321      desirable than shallow machine learning algorithms in this respect[14].

322      We demonstrated that DL has a higher predictive accuracy over the other six popular machine learning

323      methods, in detecting ER status from metabolomics data. DL exploits the idea that the higher "succeeding"

324      layer is learned from the lower "preceding" layer and selects the essential metabolites from DL model.

325      These metabolites are useful for the learning process and explain the high predictability of DL compared

326      to conventional machine learning algorithms. DL extracted features that could be considered as novel

327      biomarkers, such as uracil, which were not previously reported as breast cancer. Also, unlike other machine

328      learning methods, DL method offers additional insights on eight KEGG pathway being significantly

329      different due to ER status. All these new observations warrant further investigation.

330      An interesting new link we discover lies between FOXM1 family and beta-alanine. A recent study showed

331      FOXM1 to be a major cause for resistance to various chemotherapeutics[53], and reduction of FOXM1 levels

332      induced apoptosis of breast cancer cells[54]. The motif enrichment analysis of the beta-alanine interacted

333      enzymes indicates that the transcription factor FOXO1 targeted the promoter regions of these enzymes.

334      Thus the relationships among beta-alanine, FOXM1 and FOXO1 is worth further investigation. In addition,

335      we found many interesting involvement of DL unique metabolites in breast cancer diagnosis and treatment.

336     For example, phenylalanine is found significantly elevated in the advanced metastatic breast cancer[55] and

337     linoleic acid has been used to lower the risk of breast cancer[56]. Also, Putrescine has been known to play a

338     critical role in many metabolomics processes in breast cancer, such as apoptosis, and proliferation[57]. The

339     knock-down experiments on ornithine decarboxylase (ODC), an enzyme which converts ornithine to

340     putrescin, showed the growth inhibition in the ERα+ MCF7 and T47D and ERα- MDA-MB-231 breast

341     cancer cells[58]. Arachidonic acid was previously shown to be integral part of the new signaling for the cell

342     migrations in the MDA-MB-231 breast cancer cells[59].

343     Despite the outstanding performance of DL methods, one should be mindful of several caveats in its

344     application in metabolomics research. DL is time-consuming computation (Table S2), relative to some other

345     machine learning methods[40]. Also, metabolomics data sets are generally small, in comparison to imaging

346     data. Thus very small data sets may not be suitable for DL. We experimented with the effects of reducing

347     sample size and metabolite size on the seven methods in comparison, and found that DL is indeed sensitive

348     to the sample size of the study. On the contrary, due to colinearality among metabolites, autoencoder has

349     fairly robust predictions even when the number of metabolites are reduced. Another point of consideration

350     is the reproducibility of the technology itself. A platform with better reproducibility is expected to yield

351     biomarker models that predict more accurately in validation datasets (less overfitting). We thus speculate

352     that DL models based on NMR metabolomics data (more metabolites and better reproducibility) will be

353     more accurate than DL models based on LC-MS data, when other conditions are the same.

354     Lastly, in this report we compared the ML vs DL under the topic of classification of metabolomics data.

355     The advantages of DL on other non-classification problems in metabolomics research are yet to be explored.

356     For example, unsupervised machine learning algorithms such as PCA and hierarchical clustering were

357     applied to the metabolomics data[60], and our group is currently exploring using autoencoders for

358     unsupervised learning in metabolomics data. As another example, we have also worked on prognosis

359     prediction using shallow and deep neural network models in the genomics space [61, 62]. We successfully

360     used autoencoder to integrate multiple omics datasets (RNA-Seq, microRNA-Seq and DNA methylation)

14

361    to predict patient survival robustly, exemplified by liver cancer [2]. Compared to genomics data,

362    metabolomics data have higher multicolinearity and noise levels. Also the number of identifiable

363    metabolites are lower than the identifiable genes in genomics assays. These issues pose potential challenges

364    when extending genomics tools for metabolomics research. Nevertheless, it will be very interesting to test

365    these DL and neural network models on appropriate metabolomics data sets alone, or in combination with

366    coupled genomics data.

## Conclusions

368    We show evidence that DL outperforms other machine learning algorithms for ER status classification in

369    breast cancer metabolomics data. The biological interpretation of the hidden layer of the DL model also

370    reveals eight significant breast cancer related pathways, which are not able to obtain from the other machine

371    learning algorithms in comparison.

372

## Author Contributions

374    LXG and FMA envisioned the project and designed the work. FMA coded the project and conducted the

375    analysis. KC mapped metabolites and enzymes into KEGG pathway. FMA wrote the manuscript with help

376    from LXG and KC. LXG, FMA and KC have read, revised and approved the final manuscript.

377

## Competing financial interests

379    The author(s) declare no competing financial interests.

## Acknowledgements

384    NICHD to LX Garmire. We thank all members in Garmire group for reviewing and commenting the

385    manuscript.

386

387    **References**

388    (1)    Organization, W. H. Breast cancer: prevention and control.
389    http://www.who.int/cancer/detection/breastcancer/en/index1.html (October 10, 2017)
390    (2)    Society, A. C. About Breast Cancer. https://www.cancer.org/cancer/breast-
391    cancer/about/how-common-is-breast-cancer.html (September 21, 2017)
392    (3)    Carey, L. A.; Perou, C. M.; Livasy, C. A.; Dressler, L. G.; Cowan, D.; Conway, K.;
393    Karaca, G.; Troester, M. A.; Tse, C. K.; Edmiston, S.; Deming, S. L.; Geradts, J.; Cheang, M. C.;
394    Nielsen, T. O.; Moorman, P. G.; Earp, H. S.; Millikan, R. C. Race, breast cancer subtypes, and
395    survival in the Carolina Breast Cancer Study. *JAMA* **2006,** *295* (21), 2492-2502.
396    (4)    O'Brien, K. M.; Cole, S. R.; Tse, C. K.; Perou, C. M.; Carey, L. A.; Foulkes, W. D.;
397    Dressler, L. G.; Geradts, J.; Millikan, R. C. Intrinsic breast tumor subtypes, race, and long-term
398    survival in the Carolina Breast Cancer Study. *Clin Cancer Res* **2010,** *16* (24), 6100-6110.
399    (5)    Haque, R.; Ahmed, S. A.; Inzhakova, G.; Shi, J.; Avila, C.; Polikoff, J.; Bernstein, L.;
400    Enger, S. M.; Press, M. F. Impact of breast cancer subtypes and treatment on survival: an
401    analysis spanning two decades. *Cancer Epidemiol Biomarkers Prev* **2012,** *21* (10), 1848-1855.
402    (6)    Fan, Y.; Zhou, X.; Xia, T. S.; Chen, Z.; Li, J.; Liu, Q.; Alolga, R. N.; Chen, Y.; Lai, M.
403    D.; Li, P.; Zhu, W.; Qi, L. W. Human plasma metabolomics for identifying differential
404    metabolites and predicting molecular subtypes of breast cancer. *Oncotarget* **2016,** *7* (9), 9925-
405    9938.
406    (7)    Tang, X.; Lin, C. C.; Spasojevic, I.; Iversen, E. S.; Chi, J. T.; Marks, J. R. A joint analysis
407    of metabolomics and genetics of breast cancer. *Breast Cancer Res* **2014,** *16* (4), 415.
408    (8)    Budczies, J.; Pfitzner, B. M.; Gyorffy, B.; Winzer, K. J.; Radke, C.; Dietel, M.; Fiehn, O.;
409    Denkert, C. Glutamate enrichment as new diagnostic opportunity in breast cancer. *Int J Cancer*
410    **2015,** *136* (7), 1619-1628.
411    (9)    Lien, E. C.; Lyssiotis, C. A.; Juvekar, A.; Hu, H.; Asara, J. M.; Cantley, L. C.; Toker, A.
412    Glutathione biosynthesis is a metabolic vulnerability in PI(3)K/Akt-driven breast cancer. *Nat*
413    *Cell Biol* **2016,** *18* (5), 572-578.
414    (10)    Truong, Y.; Lin, X.; Beecher, C. In *Learning a complex metabolomic dataset using*
415    *random forests and support vector machines*, Proceedings of the tenth ACM SIGKDD
416    international conference on Knowledge discovery and data mining, 2004; ACM: 2004; pp 835-
417    840.
418    (11)    Huang, J.-H.; Yan, J.; Wu, Q.-H.; Duarte Ferro, M.; Yi, L.-Z.; Lu, H.-M.; Xu, Q.-S.;
419    Liang, Y.-Z. Selective of informative metabolites using random forests based on model
420    population analysis. *Talanta* **2013,** *117* (Supplement C), 549-555.
421    (12)    Mahadevan, S.; Shah, S. L.; Marrie, T. J.; Slupsky, C. M. Analysis of Metabolomic Data
422    Using Support Vector Machines. *Analytical Chemistry* **2008,** *80* (19), 7562-7570.
423    (13)    Min, S.; Lee, B.; Yoon, S. Deep learning in bioinformatics. *Brief Bioinform* **2016**.

424 (14) Angermueller, C.; Parnamaa, T.; Parts, L.; Stegle, O. Deep learning for computational
425 biology. *Mol Syst Biol* **2016,** *12* (7), 878.
426 (15) Tan, J.; Ung, M.; Cheng, C.; Greene, C. S. Unsupervised feature construction and
427 knowledge extraction from genome-wide assays of breast cancer with denoising autoencoders.
428 *Pac Symp Biocomput* **2015**, 132-143.
429 (16) Chen, Y.; Li, Y.; Narayan, R.; Subramanian, A.; Xie, X. Gene expression inference with
430 deep learning. *Bioinformatics* **2016,** *32* (12), 1832-1839.
431 (17) Kelley, D. R.; Snoek, J.; Rinn, J. L. Basset: learning the regulatory code of the accessible
432 genome with deep convolutional neural networks. *Genome Res* **2016,** *26* (7), 990-999.
433 (18) Budczies, J.; Denkert, C.; Muller, B. M.; Brockmoller, S. F.; Klauschen, F.; Gyorffy, B.;
434 Dietel, M.; Richter-Ehrenstein, C.; Marten, U.; Salek, R. M.; Griffin, J. L.; Hilvo, M.; Oresic,
435 M.; Wohlgemuth, G.; Fiehn, O. Remodeling of central metabolism in invasive breast cancer
436 compared to normal breast tissue - a GC-TOFMS based metabolomics study. *BMC Genomics*
437 **2012,** *13*, 334.
438 (19) Budczies, J.; Brockmoller, S. F.; Muller, B. M.; Barupal, D. K.; Richter-Ehrenstein, C.;
439 Kleine-Tebbe, A.; Griffin, J. L.; Oresic, M.; Dietel, M.; Denkert, C.; Fiehn, O. Comparative
440 metabolomics of estrogen receptor positive and estrogen receptor negative breast cancer:
441 alterations in glutamine and beta-alanine metabolism. *J Proteomics* **2013,** *94*, 279-288.
442 (20) Edgar, R.; Domrachev, M.; Lash, A. E. Gene Expression Omnibus: NCBI gene
443 expression and hybridization array data repository. *Nucleic Acids Res* **2002,** *30* (1), 207-210.
444 (21) Shao, J. C. a. J. Nearest Neighbor Imputation for Survey Data. *Journal of Official*
445 *Statistics* **2000,** *16* (2), 113–131.
446 (22) van den Berg, R. A.; Hoefsloot, H. C.; Westerhuis, J. A.; Smilde, A. K.; van der Werf, M.
447 J. Centering, scaling, and transformations: improving the biological information content of
448 metabolomics data. *BMC Genomics* **2006,** *7*, 142.
449 (23) Jauhiainen, A.; Madhu, B.; Narita, M.; Narita, M.; Griffiths, J.; Tavare, S. Normalization
450 of metabolomics data with applications to correlation maps. *Bioinformatics* **2014,** *30* (15), 2155-
451 2161.
452 (24) Li, H. Deep learning for image denoising. *International Journal of Signal Processing,*
453 *Image Processing and Pattern Recognition* **2014,** *7* (3), 171-180.
454 (25) LeCun, Y.; Kavukcuoglu, K.; Farabet, C. In *Convolutional networks and applications in*
455 *vision*, Circuits and Systems (ISCAS), Proceedings of 2010 IEEE International Symposium on,
456 2010; IEEE: 2010; pp 253-256.
457 (26) Lee, H. In *Tutorial on deep learning and applications*, NIPS 2010 Workshop on Deep
458 Learning and Unsupervised Feature Learning, 2010; 2010.
459 (27) Candel, A.; Parmar, V.; LeDell, E.; Arora, A., Deep learning with h2o. In H2O: 2015.
460 (28) Kuhn, M. Caret package. *Journal of Statistical Software* **2008,** *28* (5), 1-26.
461 (29) Robin, X.; Turck, N.; Hainard, A.; Tiberti, N.; Lisacek, F.; Sanchez, J. C.; Muller, M.
462 pROC: an open-source package for R and S+ to analyze and compare ROC curves. *BMC*
463 *Bioinformatics* **2011,** *12*, 77.
464 (30) Gedeon, T. D. Data mining of inputs: analysing magnitude and functional measures.
465 *International Journal of Neural Systems* **1997,** *8* (02), 209-218.
466 (31) Kim, S.; Thiessen, P. A.; Bolton, E. E.; Chen, J.; Fu, G.; Gindulyte, A.; Han, L.; He, J.;
467 He, S.; Shoemaker, B. A.; Wang, J.; Yu, B.; Zhang, J.; Bryant, S. H. PubChem Substance and
468 Compound databases. *Nucleic Acids Res* **2016,** *44* (D1), D1202-1213.

469 (32)    Kanehisa, M.; Goto, S.; Sato, Y.; Furumichi, M.; Tanabe, M. KEGG for integration and
470 interpretation of large-scale molecular data sets. *Nucleic Acids Res* **2012,** *40* (Database issue),
471 D109-114.
472 (33)    Smyth, G. K., Limma: linear models for microarray data. In *Bioinformatics and*
473 *computational biology solutions using R and Bioconductor*, Springer: 2005; pp 397-420.
474 (34)    Karnovsky, A.; Weymouth, T.; Hull, T.; Tarcea, V. G.; Scardoni, G.; Laudanna, C.;
475 Sartor, M. A.; Stringer, K. A.; Jagadish, H. V.; Burant, C.; Athey, B.; Omenn, G. S. Metscape 2
476 bioinformatics tool for the analysis and visualization of metabolomics and gene expression data.
477 *Bioinformatics* **2012,** *28* (3), 373-380.
478 (35)    Kamburov, A.; Cavill, R.; Ebbels, T. M.; Herwig, R.; Keun, H. C. Integrated pathway-
479 level analysis of transcriptomics and metabolomics data with IMPaLA. *Bioinformatics* **2011,** *27*
480 (20), 2917-2918.
481 (36)    Cavill, R.; Kamburov, A.; Ellis, J. K.; Athersuch, T. J.; Blagrove, M. S.; Herwig, R.;
482 Ebbels, T. M.; Keun, H. C. Consensus-phenotype integration of transcriptomic and metabolomic
483 data implies a role for metabolism in the chemosensitivity of tumour cells. *PLoS Comput Biol*
484 **2011,** *7* (3), e1001113.
485 (37)    Pasa, L.; Sperduti, A. In *Pre-training of recurrent neural networks via linear*
486 *autoencoders*, Advances in Neural Information Processing Systems, 2014; 2014; pp 3572-3580.
487 (38)    Lee, C.; Nkounkou, B.; Huang, C. H. Comparison of LDA and SPRT on Clinical Dataset
488 Classifications. *Biomed Inform Insights* **2011,** *4*, 1-7.
489 (39)    Cho, J.; Lee, K.; Shin, E.; Choy, G.; Do, S. How much data is needed to train a medical
490 image deep learning system to achieve necessary high accuracy? *arXiv preprint*
491 *arXiv:1511.06348* **2015**.
492 (40)    Goodfellow, I.; Bengio, Y.; Courville, A., *Deep learning*. MIT press: 2016.
493 (41)    Fini, M. A.; Monks, J.; Farabaugh, S. M.; Wright, R. M. Contribution of xanthine
494 oxidoreductase to mammary epithelial and breast cancer cell differentiation in part modulates
495 inhibitor of differentiation-1. *Mol Cancer Res* **2011,** *9* (9), 1242-1254.
496 (42)    El Agouza, I. M.; Eissa, S. S.; El Houseini, M. M.; El-Nashar, D. E.; Abd El Hameed, O.
497 M. Taurine: a novel tumor marker for enhanced detection of breast cancer among female
498 patients. *Angiogenesis* **2011,** *14* (3), 321-330.
499 (43)    Kim, H. Y.; Lee, K. M.; Kim, S. H.; Kwon, Y. J.; Chun, Y. J.; Choi, H. K. Comparative
500 metabolic and lipidomic profiling of human breast cancer cells with different metastatic
501 potentials. *Oncotarget* **2016,** *7* (41), 67111-67128.
502 (44)    Tan, J.; Yu, C. Y.; Wang, Z. H.; Chen, H. Y.; Guan, J.; Chen, Y. X.; Fang, J. Y. Genetic
503 variants in the inositol phosphate metabolism pathway and risk of different types of cancer. *Sci*
504 *Rep* **2015,** *5*, 8473.
505 (45)    Huan, J.; Wang, L.; Xing, L.; Qin, X.; Feng, L.; Pan, X.; Zhu, L. Insights into significant
506 pathways and gene interaction networks underlying breast cancer cell line MCF-7 treated with
507 17β-Estradiol (E2). *Gene* **2014,** *533* (1), 346-355.
508 (46)    Chen, Z. Y. a. Y. Z. a. L. In *in silico identification of novel cancer-related genes by*
509 *comparative genomics of naked mole rat and rat*, 2012 IEEE 6th International Conference on
510 Systems Biology (ISB), 2012; 2012; pp 285-290.
511 (47)    Hensley, C. T.; Wasti, A. T.; DeBerardinis, R. J. Glutamine and cancer: cell biology,
512 physiology, and clinical opportunities. *J Clin Invest* **2013,** *123* (9), 3678-3684.
513 (48)    Amelio, I.; Cutruzzola, F.; Antonov, A.; Agostini, M.; Melino, G. Serine and glycine
514 metabolism in cancer. *Trends Biochem Sci* **2014,** *39* (4), 191-198.

515 (49) Lyon, D. E.; Walter, J. M.; Starkweather, A. R.; Schubert, C. M.; McCain, N. L.
516 Tryptophan degradation in women with breast cancer: a pilot study. *BMC Res Notes* **2011,** *4,*
517 156.
518 (50) Sun, Y. L.; Patel, A.; Kumar, P.; Chen, Z. S. Role of ABC transporters in cancer
519 chemotherapy. *Chin J Cancer* **2012,** *31* (2), 51-57.
520 (51) Liu, Y.; Peng, H.; Zhang, J. T. Expression profiling of ABC transporters in a drug-
521 resistant breast cancer cell line using AmpArray. *Mol Pharmacol* **2005,** *68* (2), 430-438.
522 (52) Thomas-Chollier, M.; Hufton, A.; Heinig, M.; O'Keeffe, S.; Masri, N. E.; Roider, H. G.;
523 Manke, T.; Vingron, M. Transcription factor binding predictions using TRAP for the analysis of
524 ChIP-seq data and regulatory SNPs. *Nat Protoc* **2011,** *6* (12), 1860-1869.
525 (53) Park, Y. Y.; Jung, S. Y.; Jennings, N. B.; Rodriguez-Aguayo, C.; Peng, G.; Lee, S. R.;
526 Kim, S. B.; Kim, K.; Leem, S. H.; Lin, S. Y.; Lopez-Berestein, G.; Sood, A. K.; Lee, J. S.
527 FOXM1 mediates Dox resistance in breast cancer by enhancing DNA repair. *Carcinogenesis*
528 **2012,** *33* (10), 1843-1853.
529 (54) Bergamaschi, A.; Madak-Erdogan, Z.; Kim, Y. J.; Choi, Y. L.; Lu, H.; Katzenellenbogen,
530 B. S. The forkhead transcription factor FOXM1 promotes endocrine resistance and invasiveness
531 in estrogen receptor-positive breast cancer by expansion of stem-like cancer cells. *Breast Cancer*
532 *Res* **2014,** *16* (5), 436.
533 (55) Jobard, E.; Pontoizeau, C.; Blaise, B. J.; Bachelot, T.; Elena-Herrmann, B.; Tredan, O. A
534 serum nuclear magnetic resonance-based metabolomic signature of advanced metastatic human
535 breast cancer. *Cancer Lett* **2014,** *343* (1), 33-41.
536 (56) Arab, A.; Akbarian, S. A.; Ghiyasvand, R.; Miraghajani, M. The effects of conjugated
537 linoleic acids on breast cancer: A systematic review. *Adv Biomed Res* **2016,** *5,* 115.
538 (57) Lessard, M.; Zhao, C.; Singh, S. M.; Poulin, R. Hormonal and feedback regulation of
539 putrescine and spermidine transport in human breast cancer cells. *J Biol Chem* **1995,** *270* (4),
540 1685-1694.
541 (58) Zhu, Q.; Jin, L.; Casero, R. A.; Davidson, N. E.; Huang, Y. Role of ornithine
542 decarboxylase in regulation of estrogen receptor alpha expression and growth in human breast
543 cancer cells. *Breast Cancer Res Treat* **2012,** *136* (1), 57-66.
544 (59) Navarro-Tito, N.; Soto-Guzman, A.; Castro-Sanchez, L.; Martinez-Orozco, R.; Salazar,
545 E. P. Oleic acid promotes migration on MDA-MB-231 breast cancer cells through an arachidonic
546 acid-dependent pathway. *Int J Biochem Cell Biol* **2010,** *42* (2), 306-317.
547 (60) Xia, J. a. W., D.S Using MetaboAnalyst 3.0 for comprehensive metabolomics data
548 analysis. *Curr. Protoc. Bioinform* **2016,** *55* (14).
549 (61) Ching, T.; Zhu, X.; Garmire, L. Cox-nnet: an artificial neural network Cox regression for
550 prognosis prediction. *bioRxiv* **2016**.
551 (62) Chaudhary, K.; Poirion, O. B.; Lu, L.; Garmire, L. X. Deep Learning based multi-omics
552 integration robustly predicts survival in liver cancer. *Clinical Cancer Research* **2017**.
553

## Figure Legends

555 **Figure 1**: Block diagram of the proposed system. The first step is the preprocessing (log transformation,

556 centering, autoscaling and quantile normalization). We used Autoencoder pretraining (unsupervised step)

19

557    to initial model weights and select model architecture. Model used the 80% of data split to train the model

558    and the remaining 20% to measure model performance. The data was split 10 times to avoid the bias of

559    data sampling, and the average AUC was calculated on the 10 holds out test samples.

560    **Figure 2**: **A:** The average AUC on 10 hold out test samples of the DL framework against six machine

561    learning algorithms for prediction of ER status from metabolomics data: Recursive Partitioning and

562    Regression Trees (RPART) (0.83), Linear Discriminant Analysis (LDA) (0.74), Support Vector Machine

563    (SVM)(0.89), DeepLearning (DL)(0.93), Random Forest (RF)(0.89), Generalized Boosted Models

564    (GBM)(0.89), and Prediction Analysis for Microarrays (PAM)(0.88). The above algorithms were run 10

565    times on different train/test splits. We used pairwise Wilcoxon signed-rank test to estimate the statistical

566    significance of the difference in performance between DL and other methods (** $p<0.01$, * $p<0.1$). **B:**

567    Bipartite graph of the top 20 important metabolites extracted from DL model and other machine learning

568    algorithms. Large nodes represent the models and small nodes are metabolites. A connection between

569    metabolite and the model means this metabolite is one of the top 20 high importance metabolites extracted

570    by this model.

571    **Figure 3:** Biological relevance of the DL hidden layers. (A) Activation levels of the high variance nodes

572    extracted from the layer 1 of the DL model. Columns are samples and rows are the top 12 nodes with high

573    variance > 0.5. (B) Bipartite graph of enriched significant metabolomics pathways and top hidden nodes.

574    The nodes represent enriched pathways common to all top 12 nodes (green color) in the 1st hidden layer of

575    DL in KEGG pathway enrichment analysis (FDR< 0.05).

576    **Figure 4:** The joint pathway analysis between the top 20 DL metabolites and the high differentiated

577    enzymes. Only significant pathways with at least 5 overlapping metabolites are shown. X-axis shows the

578    number of overlapped metabolites with the number of genes (number in parentheses) involved in the same

579    pathway, y axis shows the adjusted joint *P*-value calculated from IMPALA tool[42]. The size of the nodes

580    represents the size of metabolomic pathway (number of metabolites involved in that pathway). The color

581    of the nodes represents the database source of these pathways.

20

582    **Figure 5:** Circos plot of Spearman correlation values between 20 top DL metabolites and high differentiated

583    enzymes with cut-off=|0.35|.

584    **Figure 6**: Beta-alanine and ABAT interaction network. (A) Metabolite level of beta-alanine and expression

585    of ABAT. (B) Beta-alanine-ABAT interaction network in ER– breast cancer tissues compared to ER+

586    breast cancer tissues.  Metscape, a Cytoscape plug-in, was used to integrate ER+/ER- metabolomics and

587    gene expression data (GSE59198) of the same patients.  Fold change of metabolites (hexagon nodes) or

588    enzymes (circle nodes) are represented by the size of the nodes. The input of Metascape are the top 20

589    metabolites from the DL model and the 898 genes whose expression values are statistically significantly

590    different between ER- and ER+ samples. Enzymes and metabolites of significant difference are marked by

591    green line(s) on the shapes.

592    **Supplementary Materials**

593    **Figure S1:** (A) The effect of sample size on the performance of the DL and other machine learning

594    algorithms.

595    **Figure S2:** The effect of metabolite size on the performance of the DL and other machine learning

596    algorithms.

597    **Figure S3:** DL 20 top important metabolites. **A.** Heatmap and **B.** Box plot of the 20 top important

598    metabolites extracted from the DL model.

599    **Figure S4:** Heatmap of the metabolites (columns) which most contribute to the activation value of the top

600    hidden nodes (rows).

601    **Table S1:** The list of the top 20 important features

602    **Table S2:** Running time of the seven algorithms on the metabolomics dataset

603    **Supplementary file 1**: R code of the preprocessing, models training and testing

604

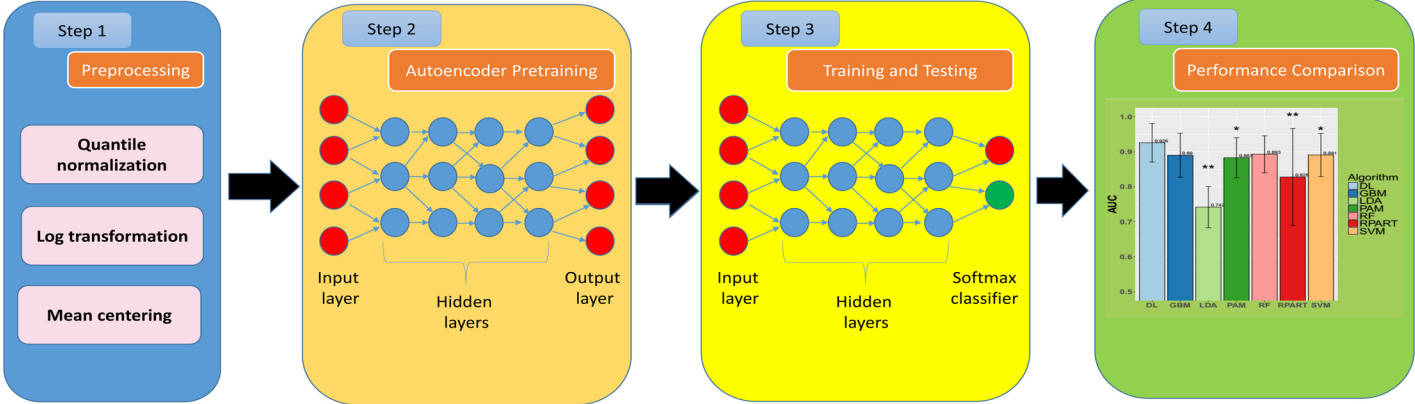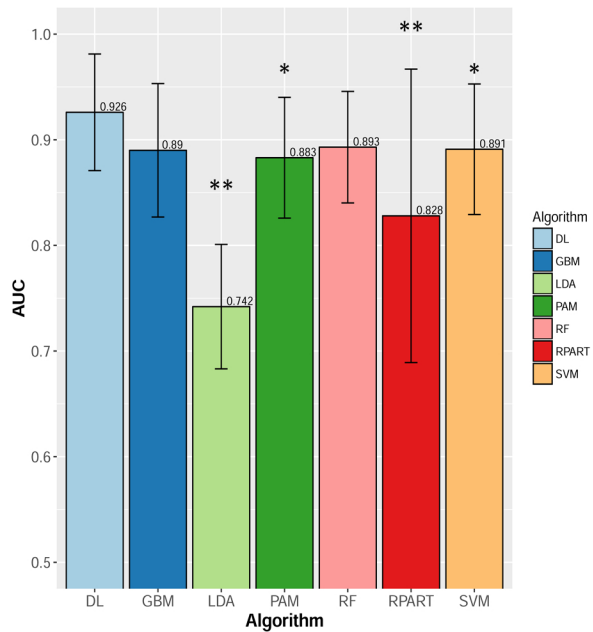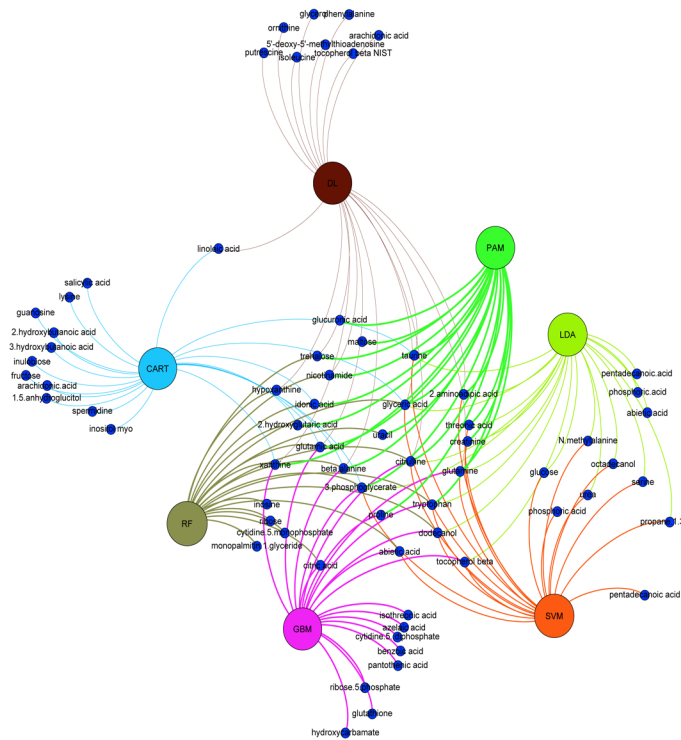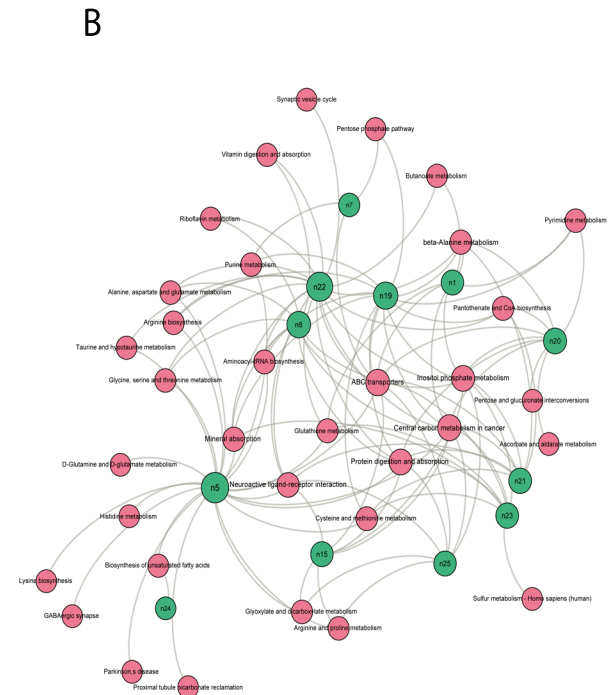For TOC only



605

606

Figure 1
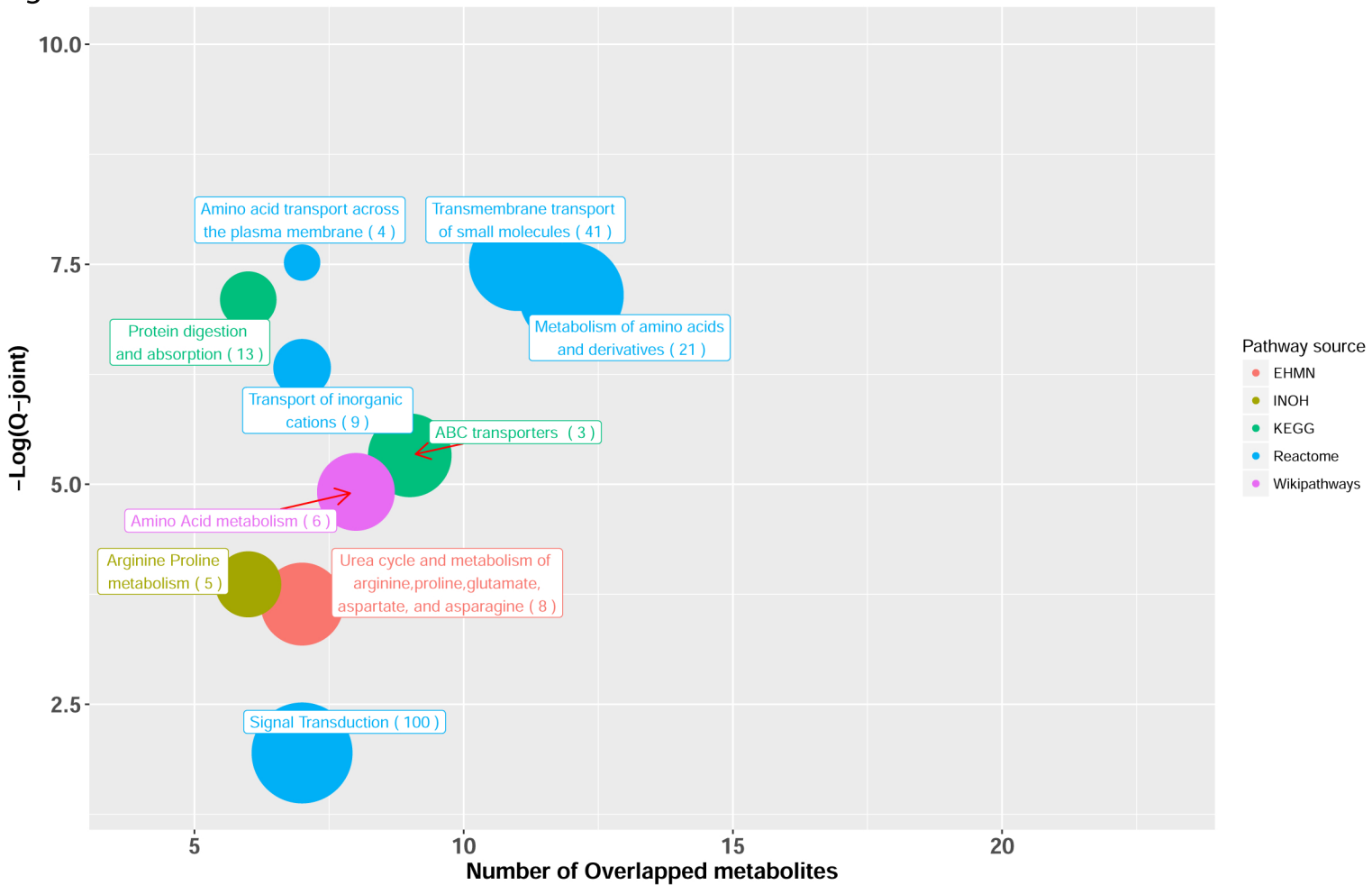
Figure 2

# Figure 3

# Figure 4



Metabolites–genes overlapping pathways

# Figure 5

# Figure 6