
Machine Learning for Glucose Prediction to Identify Diabetes-related Metabolic Pathways

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Diabetes is a major global health issue, with cases predicted to rise from 451
2 million today to 642 million by 2040. We use three machine learning methods
3 (k-NN, regression and neural networks) to predict the glucose levels of mice based
4 on genetic expression data across five tissues. Based on the best-performing neural
5 network model, we derive modules that correspond to metabolic pathways by
6 retraining networks on permuted feature clusters. The neural networks performed
7 the best of the three model families, achieving a mean absolute percentage error of
8 26.0% on the adipose tissue. From the neural network models, we produce a list of
9 9 modules that have high impact in their respective models.

10 1 Introduction

11 Diabetes is a significant health crisis worldwide, as the number of people with diabetes has risen
12 from 108 million in 1980 to 422 million in 2014. Determining causes and treatments for diabetes is
13 critical, as adverse health impacts of diabetes hinder daily life and increase mortality (1). However,
14 the causes of diabetes in humans, which include diet, weight, and genetics, are extremely complex.
15 By studying model organisms such as mice in a controlled laboratory setting, researchers can more
16 easily investigate causality between genetic features and diabetic outcomes. In this project, we use
17 supervised machine learning to predict glucose levels based on gene expression levels of obese mice,
18 with the goal of informing diabetes diagnoses in humans. Clusters of genes with high levels of
19 co-expression also investigated to determine the most important biological pathways for glucose level
20 prediction.

21 Genetics-based disease prediction models play important roles in clinical decisions. In recent years,
22 machine learning techniques are emerging as a key approach in the diagnosis of disease (3)(12).
23 Machine learning models are capable of learning complex functions based on large genetic and clinical
24 trait datasets. Various studies about applying machine learning algorithms to genetic sequence data
25 have been conducted. Liu (2019) constructed a web server, BioSeq-Analysis, which aggregates the
26 feature selection, predictor construction, and performance evaluation processes to help researchers
27 construct predictor models for DNA, RNA, and protein sequences (13). Alimadadi et al. (2020)
28 demonstrate how machine learning algorithms can be applied to RNA-Seq data from human left
29 ventricular heart tissue for classifying clinical cardiomyopathies types (2). Jiang et al.(2020) utilized a
30 Generative Adversarial Network to solve the small sample problem in brain-related disease prediction.
31 Machine learning methods have also been employed for diabetes prediction in particular (9). Carter
32 et al. (2016) used plasma metabolic profiling, genetics, and proteins to predict Type 2 diabetes using
33 Bayesian Networks (4).

34 This project builds on the existing literature by comparing the performance of several machine
35 learning techniques, determining the relative importance of co-expression clusters for glucose levels,
36 and mapping these clusters to biological pathways for future research targeted on specific processes.

37 **2 Dataset**

38 Our dataset was obtained from the Attie Lab Diabetes Dataset (<http://diabetes.wisc.edu/search.php>).
39 The dataset contains expression and clinical trait data from obese 10 week-old F2 offspring from
40 B6 and BTBR mice that segregate for diabetes; i.e., demonstrate a wide range of blood glucose
41 values. The dataset consists of 1) measurements of clinical traits such as blood glucose and insulin
42 values, and 2) Microarray-based (Agilent) measurements of gene expressions in six tissues; islet,
43 adipose, liver, gastrocnemius, kidney, and hypothalamus. All gene expression measurements used for
44 modeling were transformed to normal quantiles with respect to each tissue.

45 **2.0.1 Transcript Cleaning**

46 The Agilent microarray uses 60-mer nucleotide probes to measure the level of transcript abundance for
47 all genes. The array consists of 40,000 60-mers that are associated with 23156 genes. Therefore, our
48 first step was to identify probes that harbor genetic differences (i.e. single nucleotide polymorphisms
49 (SNPs), insertion/deletions (indels) and other complex structural variants) between B6 and BTBR
50 mice and remove them from further consideration. The presence of these genetic variants between
51 B6 and BTBR within a 60-mer sequence can cause false differences expression levels.

52 First, we performed a batch query using the Basic Local Alignment Search Tool (BLAST) to identify
53 the genomic location for each of the 40,000 60-mer sequences. This procedure excluded 5.5%
54 and 8% of the 60-mers due to lack of alignment to the mouse genome, or alignment with less than
55 58 nucleotides, respectively. This left us with robust alignment of 35,000 60-mers to the mouse
56 genome. We chose a 58/60 nucleotide match as our benchmark to ensure rigorous alignment of
57 probe-to-genome.

58 The second step was to identify all 60-mers that contain any genetic variant between BTBR and B6,
59 and their genomic locations. We first downloaded all BTBR-B6 variants from the Sanger Institute's
60 Mouse Genomes Project (19) and used it to identify all 60-mers that contain one or more of these
61 variants. This analysis resulted in 6.8% (2,677) of the remaining probes being excluded from our
62 analysis due to genetic differences between B6 and BTBR. Thus, we are left with approximately
63 32,600 60-mers that we have confidently mapped to the mouse genome, and are free of genetic
64 variants between B6 and BTBR, consisting 82% of the original dataset.

65 The final step was to link the cleaned transcripts to specific genes. Using the Ensembl Release 102
66 database, we mapped each 60-mer to it's nearest corresponding gene using the chromosome range
67 found using BLAST. This resulted in 87% (28362) of the clean 60-mers exactly matching the gene,
68 which resulted in 23156 unique genes total.

69 **2.0.2 Missing Value Imputation**

70 All tissues contained a significant number of missing gene expression values which hindered accurate
71 glucose predictions. To fill missing expression values for a given gene and tissue, a correlation matrix
72 was constructed to determine the within-tissue gene with the strongest correlation to the gene to fill.
73 A linear regression model was then trained with the gene with missing values as a function of the
74 highest-correlation gene. The trained regression model was then used to fill the missing values.

75 All 10-week glucose measurements were fully available. However, there were missing 4, 6, and
76 8-week glucose measurements. Missing glucose level values for a given mouse was filled by fitting
77 a univariate least squares regression function to the available values and linearly interpolating the
78 missing glucose measurement.

79 2.1 Exploratory Analysis

80 Having linked each transcript to its nearest gene, we also wanted to find how groups of transcripts
81 might correspond to specific biological pathways (modules). We used weighted gene correlation
82 network analysis (WGCNA) to identify modules within the cleaned transcripts (11)(20). WGCNA
83 takes in expression values from the 32,600 transcripts and forms a pairwise adjacency matrix
84 according to the formula $a_{i,j} = abs(0.5 + 0.5 * corr(x_i, x_j))^{\beta}$, where $\beta = 12$. Using this matrix,
85 we form a scale-free and “signed” co-expression network that is used to compute modules of highly
86 correlated transcripts. We did this for each tissue, yielding 29 (islet), 23 (liver), 37 (adipose), 24
87 (gastroc), and 33 (kidney) modules. Owing to their highly coordinated regulation, gene modules
88 often contain transcripts highly enriched for physiological pathways.

89 We used the R package *allez*, a pathway enrichment algorithm, to enrich modules for biological
90 pathways. *Allez* uses random-set scoring to find components in an enrichment signal (14). The
91 program ultimately produces a list of pathways for each module. To determine if a module was
92 “significantly enriched”, we evaluated each pathway using two criteria: a minimum z-score of 5 and
93 a minimum number of genes of 5. Overall, the enrichment was effective an average of 65.8% of
94 modules across the five tissues being significantly enriched.

95 Finally, we used the modules in each dataset to determine which sets of transcripts (and thereby
96 genes) greatly affected model performance. To find which modules were significant, we randomized
97 each module using the process described in Section 4.5.

98 3 Methods

99 Predictive models were constructed to predict mouse glucose levels using their gene expression data.
100 We used a train-test split ratio of 7:1, and stratified sampling on mouse sex was used to minimize
101 the impact of sex on predictions. Input data had dimensions (Number of mice for tissue x 31463)
102 with the number of mice ranging from 473 for kidney tissue data to 490 for gastrocnemius tissue data.
103 Output targets had dimensions (Number of mice for tissue x 4), where 4 represents the 4, 6, 8, and
104 10-week glucose measurements for each mouse. A combined dataset of all tissues was also used for
105 the baseline regression models. This dataset had dimensions (427 x 157315), corresponding to the
106 427 mice with expression data for all tissues and 31463 x 5 total genes.

107 Google Colab was used for preliminary code development, as well as initial model testing on smaller
108 datasets. All regression results were obtained from Amazon’s AWS Cloud Computing services.
109 Finally, all results from k-Nearest Neighbors and Neural Network models were computed using
110 the University of Wisconsin Center for High Throughput Computing’s (UW CHTC) distributed
111 computing system. The UW CHTC computing system is managed by HTCondor (a modified version
112 of Condor for UW purposes) which was used to allocate necessary computing resources (processors,
113 gpus, RAM, and disk memory).

114 3.1 k-Nearest Neighbors

115 k-Nearest Neighbors, a non-parametric statistical model, was used to obtain baseline results to
116 compare with neural network results. To predict glucose level for a given mouse, the kNN model
117 determines the k-nearest mice based on their gene expression data using the Euclidean distance
118 metric. The glucose level prediction is the average glucose level of these k-nearest mice. For each
119 tissue, k from 1 through 20 was used to determine the optimal number of nearest neighbors. As the
120 k-NN model is non-parametric, a train-test split has no purpose. To maintain comparability, k-NN
121 models were evaluated on the same testing data as the regression and neural network models.

122 3.2 Regression Analysis

123 Initially, regression modeling was conducted on the combined dataset of all tissues. Five models
124 were employed in order to test a wide range of models and determine the optimal algorithm: Linear

Model	Adipose	Gastrocnemius	Islet	Liver	Kidney
k-Nearest Neighbors	33.5	35.9	30.0	39.6	43.2
Regression	35.9	35.6	34.4	35.9	32.4
Neural Network	29.6	33.1	32.1	28.2	28.3

Table 1: Mean Absolute Percentage Error (MAPE) for each model family and tissue type

125 Regression, SGD Regression, Kernel Ridge Regression, SVM Regression, and Elastic Net Regression.
126 L1 and L2 regularization was used for Elastic Net Regression. In the Simple approach, the model was
127 trained using the gene expression data to predict each of the four weeks separately. In the iterative
128 approach, we used the same data to train the model four times and predict the result for each trial
129 separately. In the Reinforced approach, the outcome of the previous glucose measurement was
130 included as an input feature to predict the outcome of the next glucose measurement. As detailed
131 in Section 4.2, the highest performing model (Reinforced Elastic Net Regression) was used to form
132 models trained separately on each of the five tissues.

133 3.3 Neural Networks

134 Experiments were conducted to determine optimal artificial neural network structures for glucose
135 level prediction. Artificial neural networks (ANN) consist of fully connected layers composed of
136 nodes, with weights connecting nodes between layers. The gene expression data is passed into the
137 input layer of the ANN, and weights are trained in an iterative process by optimizing a loss function
138 to produce glucose level predictions. Regularization techniques such as L2, Dropout, and Early
139 Stopping were implemented. 7-fold cross validation was used to improve the generalization ability
140 of the model. Randomized Grid Search was conducted on the following hyperparameters to find an
141 optimal configuration: L2 regularization parameter, batch size, learning rate, optimizer, number of
142 layers, layer sizes, and dropout rates. Mean Squared Error (MSE) was selected as the loss function
143 for all neural networks, as this is a standard function used for neural network training.

144 3.3.1 Permutation Testing

145 After obtaining reliable glucose level predictions from the full set of features, permutation tests were
146 conducted to determine the gene clusters most impactful for glucose prediction. To achieve this,
147 neural network models were re-trained on the test dataset to obtain a baseline testing performance.
148 Then, each gene cluster obtained from the WGCNA was randomly shuffled before reevaluating the
149 testing performance. Calculating the difference in MAPE and MSE between the baseline data and
150 the permuted data enabled evaluation of the contribution of each cluster of genes. These steps of
151 permuting, reevaluating and comparing were repeated 100 times for each module. Each permutation
152 test has two meaningful outcomes: 1) Model performance improves after permutation. This suggests
153 permuting the features worsens the model performance. 2) Model performance decreases after
154 permutation. This suggests that permuting the features benefits the model performance.

155 3.4 Performance Metrics

156 We primarily used two metrics to evaluate our model performance: mean squared error (MSE) and
157 mean absolute percentage error (MAPE). Given that $e_t = (\hat{y}_t) - y_t$ is the model error on instance t ,
158 MAPE is calculated using the equation $MAPE = \frac{1}{n} \sum_{t=1}^n \left| \frac{e_t}{y_t} \right|$. Similarly, MSE is calculated using
159 the equation $MSE = \frac{1}{n} \sum_{t=1}^n e_t^2$. These two metrics provide both a raw evaluation (MSE) and a
160 scale-free evaluation (MAPE).

161 **4 Results**

162 **4.1 k-Nearest Neighbor Results**

163 The highest performing k-nearest neighbors model was trained on the islet expression data and
164 resulted in an MAPE of 30.0%. kNN models trained on adipose and gastrocnemius expression
165 data had slightly worse but comparable performance to the islet model, while liver and kidney-
166 based models resulted in less accurate glucose prediction performance. The islet data-base kNN
167 had both the most accurate predictions and the simplest model structure, as 13 neighbors yielded
168 optimal performance. All other tissues required at least 16 neighbors to achieve optimal performance.
169 This result shows that the islet gene expression data provides more useful information for glucose
170 prediction than other tissues.

171 **4.2 Regression Analysis Results**

172 As detailed in Section 3.2, Reinforced Elastic Net Regression model was selected to be trained on each
173 of the five tissues. This model yielded consistent results across all tissue types, ranging from a mean
174 absolute percentage error (MAPE) of 32.38% for kidney to 35.94% for liver. The regression model
175 had greater glucose prediction accuracy for liver and kidney tissue than the kNN model, comparable
176 adipose and gastrocnemius-based performance, and significantly worse islet-based performance.
177 Overall, the regression model had less variability in performance across tissues than the kNN model.

178 **4.3 Neural Network Results**

179 The neural network models were tuned and optimized on each tissue individually. Models across all
180 tissues shared the following optimal hyperparameters: L2 regularization rate of 0.3, learning rate of
181 0.0001, Adam optimizer, swish activation, 4 dense/dropout layers and batch size of 100. Adipose and
182 gastroc performed optimally with dense layers with 30000, 15000, 1000, and 200 nodes respectively.
183 The other three tissues performed best with dense layers with 25000, 150000, 5000 and 200 nodes
184 respectively. All five tissues used dropout rates of 0.25, 0.25, 0.25 and 0.1. Overall, the neural
185 network trained on the adipose tissue had the best performance, with a mean absolute percentage
186 error (MAPE) of 26.0%. In contrast, gastrocnemius had the worst performance with a MAPE of
187 29.5%. For all five tissues, neural networks performed best.

188 **4.4 Co-expression Network Results**

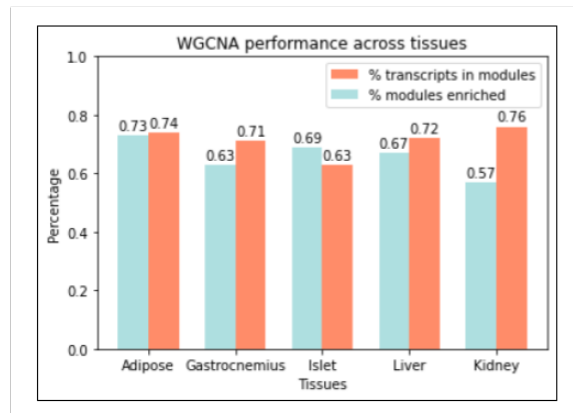


Figure 1: Percentage of transcripts in modules and percentage of modules enriched for each tissue.

189 The WGCNA script yielded 146 total modules across the five tissues we examined (islet, liver, adipose,
 190 gastroc, kidney) in the F2 mice. Adipose had the most modules at 37, while liver had the least at 23.
 191 The remaining tissues had 24 (gastroc), 29 (islet) and 33 (kidney). In each tissue, the majority of
 192 transcripts ($\geq 63\%$) uniquely belonged to a module. Furthermore, a large portion ($\geq 57\%$) of the
 193 modules were significantly enriched (z -score ≥ 5) for gene ontology (GO) pathways, as shown in
 194 Figure 1. These results demonstrate that the co-expression network successfully grouped transcripts
 195 into biologically meaningful modules. In tandem with the models we trained, these modules pave the
 196 way for us to better understand the underlying biology in the following randomization tests.

197 4.5 Permutation Testing Results

198 We ran permutation tests on each of the modules in each tissue to determine the significance of each
 199 module in determining obesity in the F2 mice. Permutation testing was conducted only on neural
 200 networks due to their greater prediction accuracy on the unpermuted base dataset. The overall results
 201 are displayed in Figure 2 with more prominent results tabulated in Table 2, which are also discussed
 202 in more detail in the following section.

Tissue	Module Color	Pathway	MSE Diff.	MAPE Diff (%)
Adipose	Red	cilium organization	+171.6	+0.23
Gastroc	Turquoise	RN-protein complex biogenesis	+806.8	+1.10
Gastroc	Blue	tricarboxylic acid cycle	+463.7	+0.54
Islet	Turquoise	natural killer cell activation	+1410.7	+1.51
Islet	Brown	pos. regulation of triglyceride	+911.7	+0.77
Liver	Turquoise	ncRNA metabolic process	+323.7	+0.30
Liver	Red	proteasome-med. process	+171.5	+0.18
Kidney	Turquoise	actin filament bundle assembly	+500.9	+0.70
Kidney	Light Cyan	isoprenoid biosynthetic process	+369.3	+0.83

Table 2: Most important modules from each tissue based on average MSE differential, average MAPE differential, and significant enrichment. Color corresponds with the point in Figure 2.

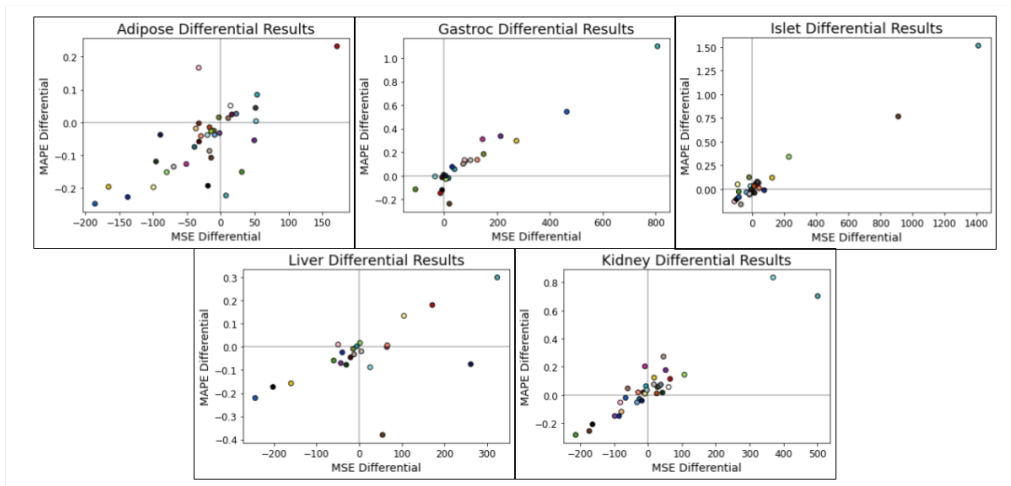


Figure 2: Performance differentials of modules for each tissue. Origin represents differential of 0 from original results in both metrics. Modules closer to the top-right corner are deemed important.

203 The adipose results were fairly ambiguous overall. The most significant module was red, which
 204 enriches for cilium organization. Ciliopathies, or diseases related to cilia dysfunction, have been

205 shown to be associated with type 2 diabetes (18). Although this module did not significantly differ
206 from optimal test results, these results may serve to reinforce the existing literature suggesting that
207 cilia function may play a role in diabetes.

208 In the gastrocnemius tissue, we found two modules (turquoise and blue) that seem exceptionally
209 important compared to the rest. The turquoise module, which is associated with ribonucleoprotein
210 complex biogenesis, showed a strong performance. However, there is a lack of significant literature
211 linking this pathway with diabetes, with the most closely related literature discussing RNA binding
212 in diabetes (15). This is a field of interest and further research. The blue module enriches for the
213 tricarboxylic acid cycle (Kreb cycle), which is one of the most important metabolic cycles for ATP
214 production. There is an enormous amount of work focusing on the relationship between this pathway
215 and diabetes (17)(7) and the importance of this module supports this literature.

216 The islet results delivered two promising modules (turquoise and brown), which both enrich for highly
217 relevant pathways. The turquoise module enriches for natural killer (NK) cell activation involved
218 in immune response. NK cell activity has been shown to be an indicator of both type 1 and type 2
219 diabetes (10)(6). The strong performance of this module continues to suggest that NK cell activation
220 is closely tied to diabetes. The brown module enriches for the positive regulation of sequestering of
221 triglyceride. Triglyceride accumulation is closely related to glucolipotoxicity, the detrimental effects
222 of high glucose and fat levels on pancreatic B-cell function, which has been shown to contribute to
223 type 2 diabetes (16)(8). The performance of this module supports these analyses.

224 The liver produced two interesting modules (turquoise and red). The turquoise module enriches
225 for ncRNA (non-coding RNA) metabolic processes. Non-coding RNA regions are essential for the
226 regulation of genes throughout the genome. This is a broad topic that encompasses many different
227 pathways and thus it is unclear how to directly relate it to diabetes. This permutation test result
228 suggests that there may be some ncRNA regions that play a more important role in diabetes, but
229 further research is required. The red module enriches for proteasome-mediated ubiquitin-dependent
230 protein catabolic process, which is not related to diabetes in existing literature.

231 The kidney has two clear significant modules (turquoise and light cyan). The turquoise had relatively
232 strong results, with an average MSE differential of +500.9 and average MAPE differential of +0.70%.
233 This module enriches for actin filament bundle assembly, which has no diabetes-related literature
234 associated with it. The strong performance of this module suggests that this pathway could potentially
235 be related to diabetes and is a viable future research topic. The light cyan module enriches for the
236 isoprenoid biosynthetic process. This inhibition of this pathway has been associated with increased
237 insulin resistance and likeliness of type 2 diabetes (5). Thus, our model supports the idea that this
238 pathway could be impactful to diabetes.

239 **5 Conclusion**

240 This work provides three main contributions: 1) a novel application of machine learning and neural
241 networks to determine gene significance for the task of predicting mouse diabetes, 2) a reproducible
242 and generalizable process of cleaning and associating genetic expression data from probes, 3) a series
243 of biological pathways that the our testing has deemed significant. Among the biological pathways
244 we deemed important, several have an existing literature base that our results serve to reinforce
245 and support; this includes the tricarboxylic acid cycle in the gastrocnemius, the natural killer cell
246 activation and positive regulation of triglyceride in the islet, and more. However, there were also
247 several pathways that have very little relevant literature or are too general of a pathway, despite
248 having strong results in the permutation tests. These pathways are promising starting points for future
249 research.

250 In our future work, we could expand our work to other datasets and investigate if results match the list
251 of key pathways identified in this work. The current dataset contains a set of parental data that could
252 be used to reinforce our existing results. Furthermore, we could dive into some specific pathways and
253 tissues to determine causality between each pathway and diabetes.

254 **References**

- 255 [1] Diabetes.
- 256 [2] Ahmad Alimadadi, Ishan Manandhar, Sachin Aryal, Patricia B. Munroe, Bina Joe, and Xi Cheng.
257 Machine learning-based classification and diagnosis of clinical cardiomyopathies. *Physiological*
258 *Genomics*, 52(9):391–400, September 2020.
- 259 [3] Danilo Bzdok, Naomi Altman, and Martin Krzywinski. Statistics versus machine learning.
260 *Nature Methods*, 15(4):233–234, April 2018.
- 261 [4] Tonia C. Carter, Dietrich Rein, Inken Padberg, Erik Peter, Ulrike Rennefahrt, Donna E. David,
262 Valerie McManus, Elisha Stefanski, Silke Martin, Philipp Schatz, and Steven J. Schrodi. Vali-
263 dation of a metabolite panel for early diagnosis of type 2 diabetes. *Metabolism: Clinical and*
264 *Experimental*, 65(9):1399–1408, September 2016.
- 265 [5] Luke H Chamberlain. Inhibition of isoprenoid biosynthesis causes insulin resistance in 3T3-L1
266 adipocytes. *FEBS Letters*, 507(3):357–361, November 2001.
- 267 [6] Chris Fraker and Allison L. Bayer. The Expanding Role of Natural Killer Cells in Type 1
268 Diabetes and Immunotherapy. *Current Diabetes Reports*, 16(11):109, November 2016.
- 269 [7] Marta Guasch-Ferré, José L Santos, Miguel A Martínez-González, Clary B Clish, Cristina
270 Razquin, Dong Wang, Liming Liang, Jun Li, Courtney Dennis, Dolores Corella, Carlos Muñoz-
271 Bravo, Dora Romaguera, Ramón Estruch, José Manuel Santos-Lozano, Olga Castañer, Angel
272 Alonso-Gómez, Luis Serra-Majem, Emilio Ros, Sílvia Canudas, Eva M Asensio, Montserrat
273 Fitó, Kerry Pierce, J Alfredo Martínez, Jordi Salas-Salvadó, Estefanía Toledo, Frank B Hu,
274 and Miguel Ruiz-Canela. Glycolysis/gluconeogenesis- and tricarboxylic acid cycle-related
275 metabolites, Mediterranean diet, and type 2 diabetes. *The American Journal of Clinical Nutrition*,
276 111(4):835–844, April 2020.
- 277 [8] Jung-Hee Hong, Dae-Hee Kim, and Moon-Kyu Lee. Glucolipotoxicity and GLP-1 secretion.
278 *BMJ Open Diabetes Research & Care*, 9(1):e001905, February 2021.
- 279 [9] Xue Jiang, Jingjing Zhao, Wei Qian, Weichen Song, and G. Lin. A Generative Adversarial
280 Network Model for Disease Gene Prediction With RNA-seq Data. *IEEE Access*, 2020.
- 281 [10] Jung Hye Kim, Kahui Park, Sang Bae Lee, Shinae Kang, Jong Suk Park, Chul Woo Ahn, and
282 Ji Sun Nam. Relationship between natural killer cell activity and glucose control in patients
283 with type 2 diabetes and prediabetes. *Journal of Diabetes Investigation*, 10(5):1223–1228,
284 September 2019.
- 285 [11] Peter Langfelder and Steve Horvath. WGCNA: an R package for weighted correlation network
286 analysis. *BMC bioinformatics*, 9:559, December 2008.
- 287 [12] Zhijun Liao, Dapeng Li, Xinrui Wang, Lisheng Li, and Quan Zou. Cancer Diagnosis Through
288 IsomiR Expression with Machine Learning Method. *Current Bioinformatics*, 13(1):57–63.
- 289 [13] Bin Liu. BioSeq-Analysis: a platform for DNA, RNA and protein sequence analysis based on
290 machine learning approaches. *Briefings in Bioinformatics*, 20(4):1280–1294, July 2019.
- 291 [14] Michael A. Newton, Fernando A. Quintana, Johan A. den Boon, Srikumar Sengupta, and Paul
292 Ahlquist. Random-set methods identify distinct aspects of the enrichment signal in gene-set
293 analysis. *The Annals of Applied Statistics*, 1(1):85–106, June 2007. Publisher: Institute of
294 Mathematical Statistics.
- 295 [15] Curtis A. Nutter and Muge N. Kuyumcu-Martinez. Emerging roles of RNA-binding proteins
296 in diabetes and their therapeutic potential in diabetic complications. *Wiley interdisciplinary*
297 *reviews. RNA*, 9(2), March 2018.

- 298 [16] Vincent Poitout, Julie Amyot, Meriem Semache, Bader Zarrouki, Derek Hagman, and Ghis-
299 laine Fontés. Glucolipototoxicity of the pancreatic beta cell. *Biochimica Et Biophysica Acta*,
300 1801(3):289–298, March 2010.
- 301 [17] P. Schrauwen and M. K. C. Hesselink. Reduced tricarboxylic acid cycle flux in type 2 diabetes
302 mellitus? *Diabetologia*, 51(9):1694–1697, 2008.
- 303 [18] Francesco Volta and Jantje M. Gerdes. The role of primary cilia in obesity and diabetes. *Annals*
304 *of the New York Academy of Sciences*, 1391(1):71–84, March 2017.
- 305 [19] Binnaz Yalcin, Kim Wong, Avigail Agam, Martin Goodson, Thomas M. Keane, Xiangchao
306 Gan, Christoffer Nellåker, Leo Goodstadt, Jérôme Nicod, Amarjit Bhomra, Polinka Hernandez-
307 Pliego, Helen Whitley, James Cleak, Rebekah Dutton, Deborah Janowitz, Richard Mott, David J.
308 Adams, and Jonathan Flint. Sequence-based characterization of structural variation in the mouse
309 genome. *Nature*, 477(7364):326–329, September 2011.
- 310 [20] Bin Zhang and Steve Horvath. A general framework for weighted gene co-expression network
311 analysis. *Statistical Applications in Genetics and Molecular Biology*, 4:Article17, 2005.