

# Deep learning for predicting disease status using genomic data

**Qianfan Wu**<sup>1</sup>, **Adel Boueiz**<sup>2,3</sup>, **Alican Bozkurt**<sup>4</sup>, **Arya Masoomi**<sup>4</sup>, **Allan Wang**<sup>5</sup>, **Dawn L DeMeo**<sup>2</sup>, **Scott T Weiss**<sup>2</sup>, **Weiliang Qiu**<sup>Corresp. 2</sup>

<sup>1</sup> Questrom School of Business, Boston University, Boston, USA

<sup>2</sup> Brigham and Women's Hospital/Harvard Medical School, Boston, USA

<sup>3</sup> Pulmonary and Critical Care Division, Brigham and Women's Hospital/Harvard Medical School, Boston, USA

<sup>4</sup> Department of Computer Science, Northeastern University, Boston, USA

<sup>5</sup> Belmont High School, Boston, USA

Corresponding Author: Weiliang Qiu

Email address: stwxq@channing.harvard.edu

Predicting disease status for a complex human disease using genomic data is an important, yet challenging, step in personalized medicine. Among many challenges, the so-called curse of dimensionality problem results in unsatisfied performances of many state-of-art machine learning algorithms. A major recent advance in machine learning is the rapid development of deep learning algorithms that can efficiently extract meaningful features from high-dimensional and complex datasets through a stacked and hierarchical learning process. Deep learning has shown breakthrough performance in several areas including image recognition, natural language processing, and speech recognition. However, the performance of deep learning in predicting disease status using genomic datasets is still not well studied. In this article, we performed a review on the four relevant articles that we found through our thorough literature review. All four articles used auto-encoders to project high-dimensional genomic data to a low dimensional space and then applied the state-of-the-art machine learning algorithms to predict disease status based on the low-dimensional representations. This deep learning approach outperformed existing prediction approaches, such as prediction based on probe-wise screening and prediction based on principal component analysis. The limitations of the current deep learning approach and possible improvements were also discussed.

1                   **Deep learning for predicting disease status using genomic data**

2   Qianfan Wu<sup>1</sup>, Adel Boueiz<sup>2,3</sup>, Alican Bozkurt<sup>4</sup>, Arya Masoomi<sup>4</sup>, Allan Wang<sup>5</sup>, Dawn L. DeMeo<sup>2</sup>,  
3                   Scott T. Weiss<sup>2</sup>, Weiliang Qiu<sup>2\*</sup>

4           <sup>1</sup> Questrom School of Business, Boston University, 595 Commonwealth Avenue, Boston,  
5           MA, 02215

6           <sup>2</sup> Channing Division of Network Medicine, Brigham and Women’s Hospital/Harvard  
7           Medical School, 181 Longwood Avenue, Boston MA 02115

8           <sup>3</sup> Pulmonary and Critical Care Division, Department of Medicine, Brigham and Women’s  
9           Hospital, Harvard Medical School, Boston, MA

10          <sup>4</sup> Department of Computer Science, Northeastern University, Boston, MA

11          <sup>5</sup> Belmont High School, Boston, MA

12   \*: Corresponding author’s email address: stwxq@channing.harvard.edu

13

14   **Abstract**

15

16   Predicting disease status for a complex human disease using genomic data is an important, yet

17   challenging, step in personalized medicine. Among many challenges, the so-called curse of

18   dimensionality problem results in unsatisfied performances of many state-of-art machine

19   learning algorithms. A major recent advance in machine learning is the rapid development of

20   deep learning algorithms that can efficiently extract meaningful features from high-dimensional

21   and complex datasets through a stacked and hierarchical learning process. Deep learning has

22   shown breakthrough performance in several areas including image recognition, natural language

23   processing, and speech recognition. However, the performance of deep learning in predicting

24 disease status using genomic datasets is still not well studied. In this article, we performed a  
25 review on the four relevant articles that we found through our thorough literature review. All  
26 four articles used auto-encoders to project high-dimensional genomic data to a low dimensional  
27 space and then applied the state-of-the-art machine learning algorithms to predict disease status  
28 based on the low-dimensional representations. This deep learning approach outperformed  
29 existing prediction approaches, such as prediction based on probe-wise screening and prediction  
30 based on principal component analysis. The limitations of the current deep learning approach and  
31 possible improvements were also discussed.

32

### 33 **1. Introduction**

34

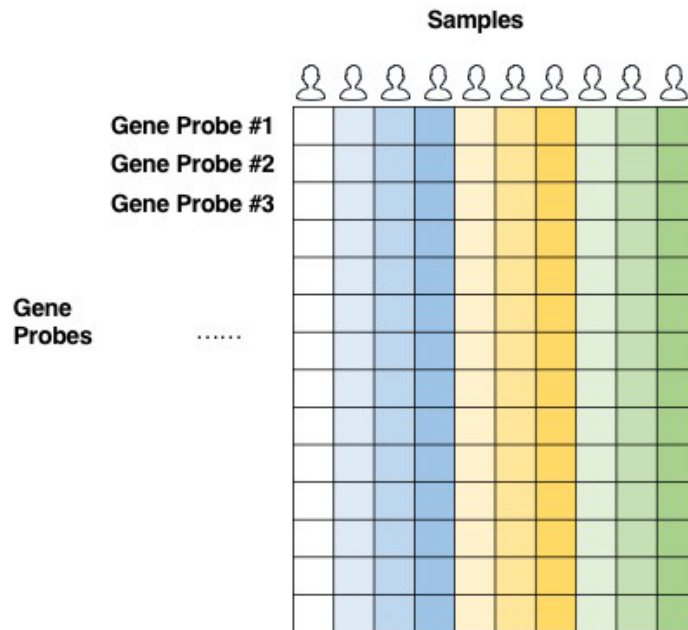
35 Complex human diseases, such as cancers, cardiovascular diseases, and respiratory diseases,  
36 have caused huge public health concerns and economic burdens [1, 2]. It is believed that both  
37 environmental factors (e.g., smoking exposure, nutrient intake, physical exercise, etc.) and  
38 genomic factors contribute to the development of complex human diseases[3]. We refer genomic  
39 factors to any molecular factors related to genes, such as genotype, gene expression, DNA  
40 methylation, microRNA expression, metabolites, proteins, etc. Cutting-edge technologies, e.g.,  
41 genotyping and next-generation whole genome sequencing, greatly facilitate the investigations of  
42 the associations of genomic factors to complex human diseases so that researchers can  
43 unbiasedly detect disease-associated factors. In addition to uncovering the underlying molecular  
44 mechanisms, researchers expect that the disease-associated genomic factors could also help  
45 diagnose disease, personalize treatment, and develop new medicines[4].

46

47 Several machine learning methods, such as support vector machine[5] (SVM), Random  
48 Forest[6], and K-Nearest Neighbors[7] have been successfully applied in disease prediction  
49 based on clinical data[8-10]. For genomic data generated by high-throughput technologies  
50 (Figure 1), the major challenge in disease prediction is the “curse of dimensionality”[11-13] (i.e.,  
51 the number of genomic factors is far larger than the number of samples), resulting in model over-  
52 fitting and computational inefficiency.

53

54 A reasonable approach[14, 15] to handle the curse of dimensionality is to first apply feature  
55 selection techniques to select key features relevant to the disease of interest, and then to predict  
56 the disease status based on these key features (Figure 2). In genomic data analysis, a feature can  
57 be a gene probe/transcript or a (non)linear combination of several gene probes/transcripts.  
58 Traditional feature selection techniques (e.g., forward variable selection, backward variable  
59 deletion, stepwise variable selection, probe-wise tests, or principal component analysis) have  
60 limited performance in genomic data analysis. Forward variable selection, backward variable  
61 deletion, and stepwise variable selection are time-consuming. Hence they are not suitable for  
62 whole genome-wide analysis. Probe-wise tests ignore the fact that many omics variables are  
63 correlated and therefore carry redundant information regarding prediction. Ignoring the  
64 redundancy would result in the selected probes are non-reproducible in independent cohorts [13,  
65 16, 17]. In addition, contributions of different genomic risk factors might be different, however,  
66 probe-wise tests implicitly assign equal weights to all selected probes. Principal component  
67 analysis (PCA) explicitly assigns different weights to different probes. However, PCA produces  
68 linear combination of probes and ignores the possible non-linear relationship between probes.

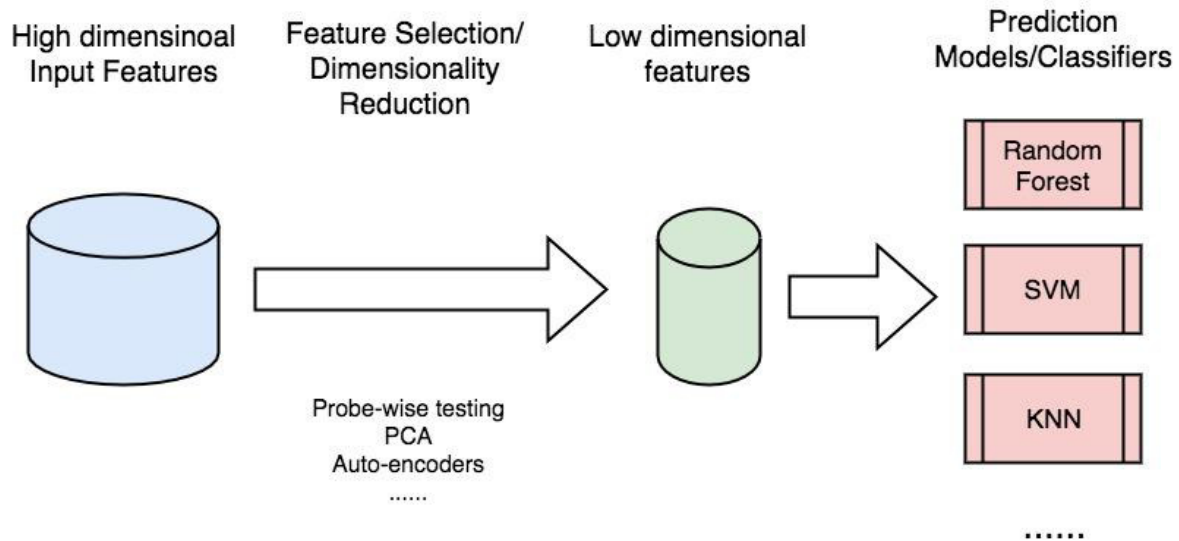


69

70 *Figure 1: An illustration of gene expression data. In the above figure, each row represents 1 gene probe and each*  
71 *column represents one sample (one person). The (i,j) cell records the expression level of the i-th gene probe for the*  
72 *j-th sample. Gene expression data typically have high dimensionality (20,000 – 50,000 gene probes) and small*  
73 *sample size (<1000), resulting in the “curse of dimensionality problem.”*

74

75



76

77 *Figure 2: An illustration of building prediction models using genomic datasets. The idea is to first reduce the*  
 78 *dimensionality of the input features and then feed the low dimensionality features into prediction model/classifiers.*  
 79 *Dimensionality reduction techniques typically include probe-wise testing, principal component analysis (PCA), and*  
 80 *auto-encoders.*

81

82 Recently, deep learning methods have made breakthrough progress in image/video  
 83 recognition[18], natural language processing[19], and robotics[20, 21]. Through a stacked and  
 84 hierarchical learning system, deep learning methods could efficiently capture complex  
 85 relationships between high-dimensional features, either spatial or consequential[22].

86

87 In bioinformatics, deep learning methods have fruitful and innovative applications in medical  
 88 image classification[23, 24], predicting DNA- and RNA-binding proteins sequences[25], and  
 89 DNA sequence noncoding variants effects predicting[26]. However, using deep learning methods  
 90 to predict disease status is not a well-researched area.

91

92 Most investigators in genomic data analysis fields might hear about deep learning and would like  
93 to learn more about deep learning and how deep learning could be used to predict disease status  
94 based on genomic data. In this review, we will first introduce the main components of deep  
95 learning and the most frequently used deep learning feature extraction methods in genomic data  
96 analysis. We will then review the papers that used deep learning to predict complex human  
97 diseases based on genomic data. The limitations of the current deep learning approach and  
98 possible improvements will also be discussed.

99

## 100 **2. Survey Methodology**

101 To thoroughly search recent literatures on deep learning applications in disease prediction, we  
102 carefully reviewed previous works, searched popular sites: Google Scholar, PubMed, IEEE  
103 Xplore, and PMC, and examined related online blogs and tutorials, such as GitHub  
104 (<http://github.com/>), Kaggle (<http://www.kaggle.com/>), and Cross Validated  
105 (<https://stats.stackexchange.com/>). We identified four papers[13, 31, 46, 47] published between  
106 January 2013 and December 2017, which applied deep learning methods in disease prediction  
107 using genomic data.

108

109 Before we review the details of the four studies, we first introduce in the following sections the  
110 main components of deep learning and the most frequently used deep learning feature extraction  
111 methods in genomic data analysis.

112

## 113 **3. Artificial Neural Networks (ANNs) and Deep Learning Methods in Predicting Disease**

114

115 The main component of all deep learning algorithms is Artificial Neural Networks (ANNs).  
116 Understanding how ANNs are constructed and trained is the first step to understand deep  
117 learning methods.

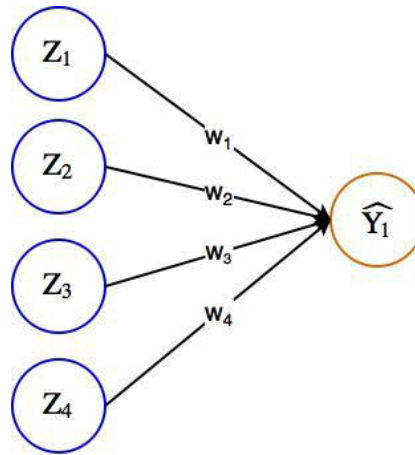
118

### 119 **Artificial Neural Networks (ANNs):**

120

121 Artificial Neural Networks are computing systems that are inspired by the biological neural  
122 networks constituting brains. Typically, an ANN is a network of nodes with multilayers (one  
123 input layer, one output layer, and several hidden internal layers). Within a layer, nodes are not  
124 connected, while between the layers nodes are fully connected (*Figure 3, Figure 4*). Each node  
125 can store a value (e.g.,  $Z_i$  for the  $i$ -th node in a given layer) and each edge can have a weight  
126 (e.g.,  $w_{ji}$  for the edge connecting node  $i$  in the given layer with node  $j$  in the previous layer). The  
127 value of a node on a given layer, except for the first layer (i.e., the input layer), is a function of a  
128 bias (i.e., threshold; e.g.,  $b_i$  for the  $i$ -th node) and the weighted average values of all nodes on the  
129 previous layer. The function is called activation function. For instance,  $\hat{Y}_1 = 1$  if  
130  $(b_i + w_{1i} * Z_1 + \dots + w_{ni} * Z_n) > 0$  and  $\hat{Y}_1 = 0$  otherwise, where  $n$  is the number of nodes in the  
131 previous layer and  $Z_j$  is the value for the  $j$ -th node in the previous layer. Usually, activation  
132 functions, such as Sigmoid, Rectified Linear Unit (ReLU)[27], and Hyperbolic Tangent (Tanh),  
133 are non-linear.





134

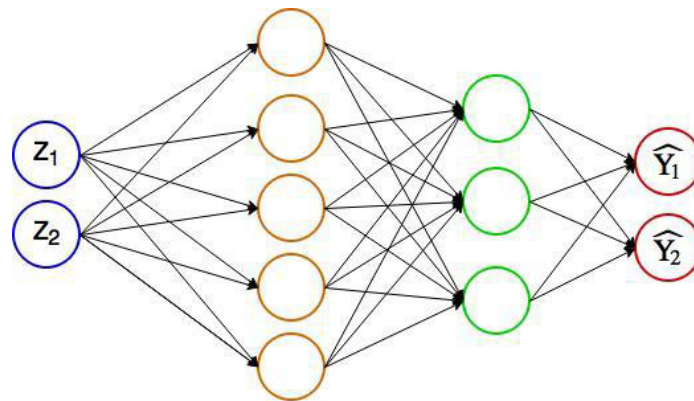
135

136 *Figure 3: An illustration of a simple ANN: This simple feed-forward ANN has four input nodes and one output node.*

137 *On the edges,  $w_1 - w_4$  represent the weights of the input nodes. The value  $Y_1$  for the output node is computed as  $\hat{Y}_1$*

138 *=  $f(b + Z_1 * w_1 + Z_2 * w_2 + Z_3 * w_3 + Z_4 * w_4)$ , where  $b$  is the bias term, and  $f$  is the activation function.*

139



140

141 *Figure 4: An illustration of a multiple-layer ANN. This multiple layer ANN has one input layer, two hidden layers,*

142 *and one output layer, with each layer connected to the previous layer. The activation function  $f$  is applied to each*

143 *node on the hidden layer and the output layer.*

144

145 **Training ANNs:**

146

147 A training data set and a validation set, in which the values of the nodes in the output layers are  
148 known (e.g., 1 for a positive outcome and 0 for a negative outcome), are needed to estimate the  
149 optimal values of the biases and edge weights (i.e., to train the ANN). The idea is to find a set of  
150 biases and edge weights (parameters) that minimize the difference between the true values and  
151 predicted values of nodes in the output layer. The difference is a function of the biases and edge  
152 weights and is usually called loss function.

153

154 Gradient descent is an optimization method for updating the parameters of a neural network to  
155 minimize the loss function (Figure 5). It uses the fact that optimal parameters are achieved when  
156 gradient of the loss function with respect to the parameters are zero. However, finding  
157 parameters that are the solution to zero gradient equation is a nontrivial task for complex  
158 networks with large number of parameters. An alternative method to solving the gradient  
159 equation is, starting with an initial point, to iteratively update each parameter proportional to the  
160 negative of the gradient of the loss function with respect to the parameter, and continue this  
161 procedure until amount of change of parameters are below a predefined threshold. An important  
162 part of this method is to calculate the gradient of loss function with respect to every parameter in  
163 the network. Backpropagation is an algorithm for efficiently calculating the gradient for each  
164 parameter, using chain rule: For the simple network in Figure 3,  $\frac{\partial Loss(w)}{\partial w_1} = \frac{\partial Loss(w)}{\partial \hat{Y}_1} \frac{\partial \hat{Y}_1}{\partial w_1}$ , where  
165  $Loss(w)$  is the loss function. This implies that once we know the gradients at some layer, we can  
166 easily calculate the gradients for the layer before it.

167

168

169

170

171

172

173

174

175

176

177

178

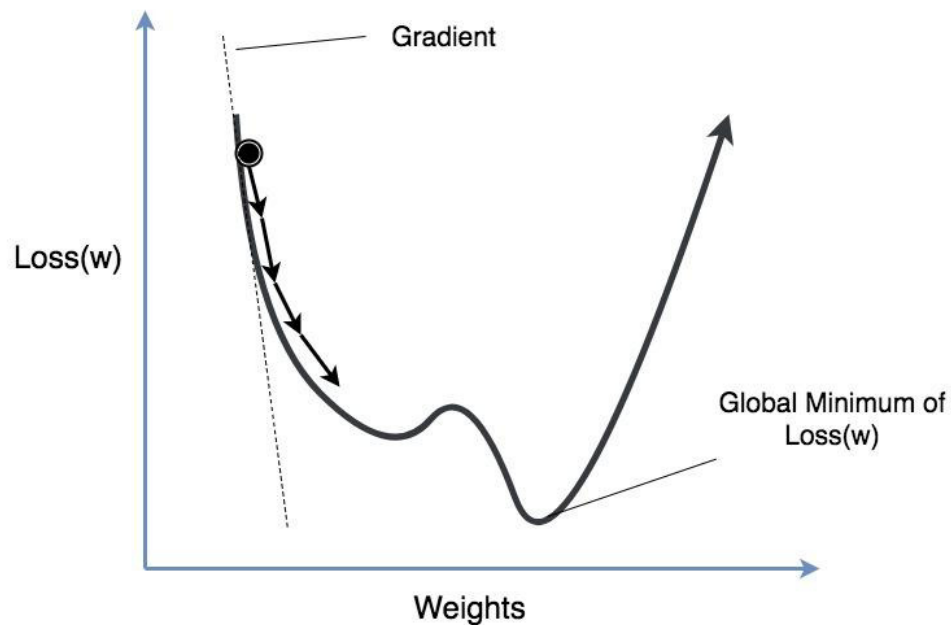
179

180

181

182

183



184 *Figure 5: Gradient Descent Training. The x-axis is the weight  $w$  and the y-axis is the loss function  $Loss(w)$ . In*

185 *Gradient Descent optimization, learning rate represents how much the edge weights are adjusted in each step before*

186 *the global minimum is achieved. Learning rate could also be seen as the “step size” in the learning process. With a*

187 *higher learning rate, the gradients are adjusted by a greater amount each step. With a lower learning rate, the*

188 *gradients are adjusted by a smaller amount each step.*

189

190 **Deep Learning and Deep Neural Networks (DNNs):**

191

192 ANNs with only one or two hidden layers have a shallow architecture, which contains only two  
193 levels of data-dependent computational elements and can be very inefficient regarding the  
194 number of computational units (e.g., hidden nodes), and in terms of required training  
195 examples[11]. In contrast, ANNs with more than two hidden layers (i.e., deep neural networks)  
196 have a deep architecture, which can compactly represent a large number of computational  
197 elements via the composition of many nonlinearities[11]. Deep learning methods are defined as  
198 computational models that are composed of multiple processing layers to learn representations of  
199 data with multiple levels of abstraction[22].

200

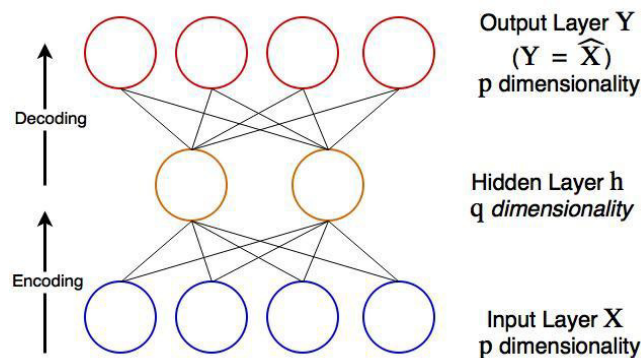
201 The performance of deep learning relies on the methods to train the parameters in DNNs.  
202 Intuitively, we can train the parameters by minimizing the prediction error rates (the loss  
203 function) through applying gradient descent. However, empirical experiments showed that this  
204 supervised approach has poor performance for DNNs[11, 28], in the regime where number of  
205 input features are comparable to (or even far larger than) number of training samples, which is  
206 the case in genomic datasets. In contrast, unsupervised learning at each stage of a deep network  
207 proposed by the seminal works of Hinton et al. (2006)[29] and Hinton and Salakhutdinov  
208 (2006)[30] pretrains each hidden layer as the encoder of an auto-encoder trying to reconstruct the  
209 output of the previous layer. . Hence, combining unsupervised approach with the supervised  
210 approach, such as combining an auto-encoder with a supervised fine-tuning phase (i.e., fine-tune  
211 all the parameters of the ANN using backpropagation and gradient descent on a global  
212 supervised cost function), can significantly improve the performance of deep learning methods  
213 for data-sparse datasets[11, 28].

214

215 **Auto-encoder (AE):**

216 An auto-encoder is a type of ANN that aims to find a new representation of input nodes (e.g.,  
 217 gene probes in genomic data analysis) in an *unsupervised* manner, from which the input can be  
 218 reconstructed without too much loss of information[28]. Like ANN, an auto-encoder has one  
 219 input layer, one output layer, and one or multiple hidden layers (*Figure 6*). Suppose  $X$  is the  
 220 original data in a  $p$ -dimensional space. An auto-encoder would first project  $X$  to a  $q$ -dimensional  
 221 space  $Y=g_1(X)$ , where  $g_1$  is a non-linear projection function. Then it transforms  $Y$  back to the  $p$ -  
 222 dimensional space  $Z=g_2(Y)$ , where  $g_2$  is also a non-linear projection function. The optimal  
 223 projection  $Y^*$  minimizes the loss function  $loss[X, g_2(Y)]$  that measures the differences between  $X$   
 224 and  $Z=g_2(Y)$ . Note that since  $q$  is different from  $p$ , both the projection function  $g_1$  and the  
 225 projection function  $g_2$  are not one-to-one mapping functions. Hence, the inverse functions  $g_1^{-1}$   
 226 and  $g_2^{-1}$  do not exist.

227



228

229

230 *Figure 6: Illustration of a basic auto-encoder. This auto-encoder has 2 hidden units.  $X$  is the inputs,  $Y=\hat{X}$  is the*  
 231 *reconstructed inputs in the output layer,  $h$  is the hidden layer. The dimension of the original input data is reduced*  
 232 *from  $p=4$  to  $q=2$ . The optimal representation in the  $q$ -dimensional space is obtained by minimizing the difference*  
 233 *between the inputs  $X$  and the reconstructed inputs  $Y$ .*

234

235

236 Similar to training ANNs, backpropagation and gradient descent can be applied to train an auto-  
237 encoder, in which the output layer has the dimension as the original data  $\mathbf{Z} = g_2(\mathbf{Y}) = g_2(g_1(\mathbf{X}))$ .

238

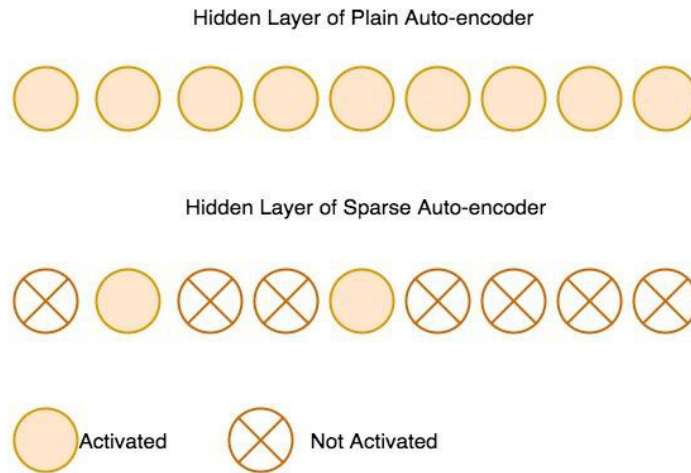
239 The nodes  $\mathbf{Y} = g_1(\mathbf{X})$  within the hidden layer are the representations of original features. The  
240 hidden layer is “under-complete” if the number ( $q$ ) of nodes in the hidden layer is smaller than  
241 that ( $p$ ) in the input layer ( $q < p$ ). In most cases, auto-encoder outperforms Principal Component  
242 Analysis in processing high dimensional complex datasets because auto-encoder performs both  
243 linear and non-linear projections, while PCA performs only linear projection. Auto-encoders  
244 have been successfully used to efficiently extract meaningful features in disease diagnosis based  
245 on high-throughput genomic data[31, 32].

246

247

#### 248 **Sparse auto-encoder (Sparse AE):**

249 Performing backpropagation and gradient descent could be inefficient if there are too many free  
250 nodes with complex dependencies in each layer[33, 34]. Sparse auto-encoder is developed to  
251 restrict the number of hidden nodes to be activated by introducing sparsity-constraints on the  
252 hidden units (*Figure 7*). Sparse auto-encoder have been proved to have favorable performance in  
253 image recognition[35] and speech emotion recognition[36], due to its efficiency in extracting  
254 meaningful features from high-dimensional data.



255

256 *Figure 7: Illustration of a sparse auto-encoder: A sparse auto-encoder restricts the number of hidden layers*  
257 *activated by adding a sparsity term to the loss function. The sparsity term set the expected activation value of the*  
258 *hidden nodes to a small constant so that most of the hidden nodes' activations are near zero. Hence, very few hidden*  
259 *nodes are activated in a sparse auto-encoder.*

260

### 261 **Stacked auto-encoder (Stacked AE):**

262

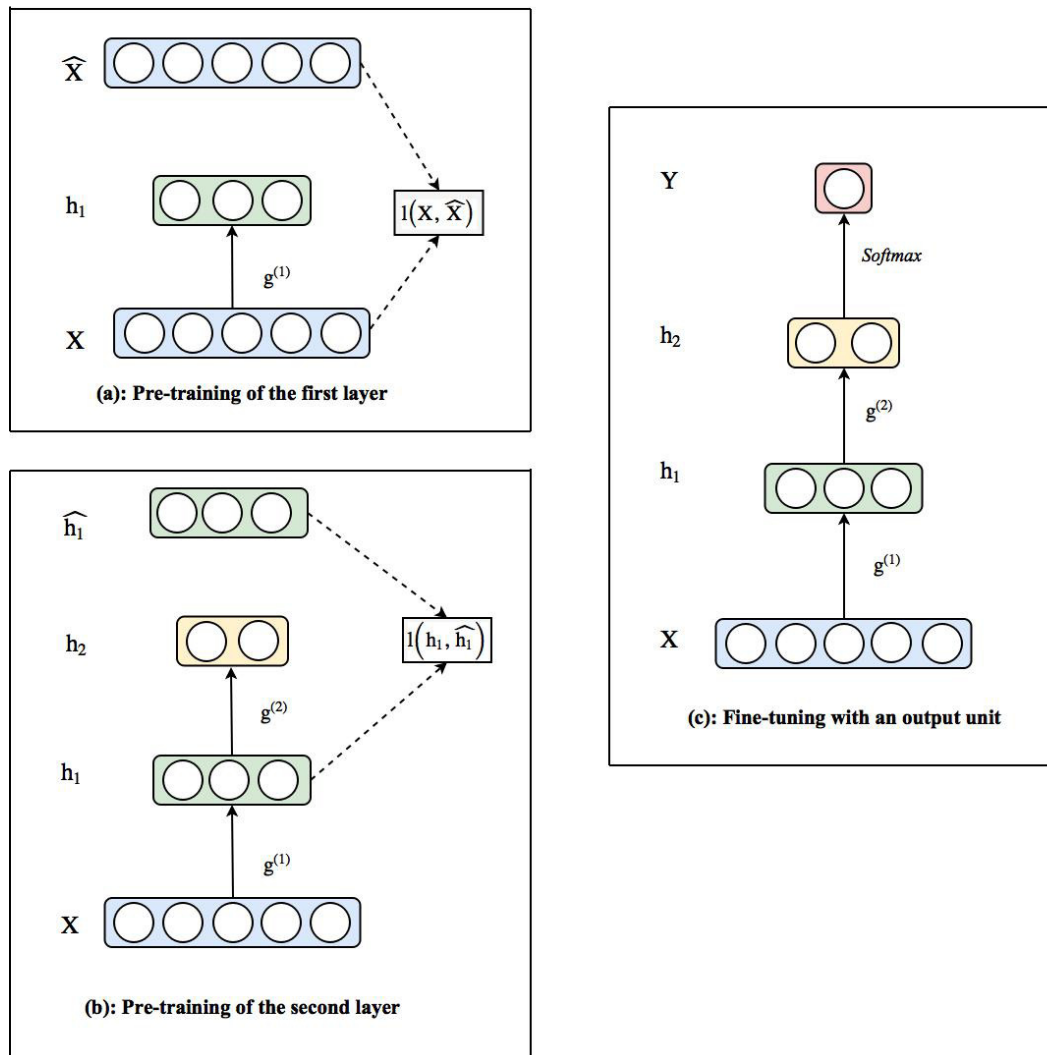
263 A stacked auto-encoder[11, 37, 38] is a multi-layer auto-encoder, each hidden layer of which is a  
264 representation of previous layer obtained by an auto-encoder with one hidden layer (Figure 8).

265 The training of stacked auto-encoders is often completed by applying the greedy layer-wise pre-

266 training approach[11]. Given extremely high-dimensional input data, a stacked auto-encoder

267 could extract features layer by layer and finally forms a better representation to be passed into

268 classifiers.



269

270 Figure 8: Illustration of stacked auto-encoder and greedy layer-wise pre-training: The stacked auto-encoder has 2  
 271 hidden layers  $h_1$  and  $h_2$ . Under the greedy layer-wise pre-training, hidden layer  $h_1$  is first trained under the same  
 272 way as training a simple 1-layer auto-encoder by minimizing  $l(X, \hat{X})$ . The function  $g^{(1)}$  that maps  $X$  to  $h_1$  is learned  
 273 from the first layer training, which is shown in (a). Then nodes values on  $h_1$  are passed to the second layer  $h_2$  to  
 274 train the function  $g^{(2)}$  that maps  $h_1$  to  $h_2$  by minimizing  $l(h_1, \hat{h}_1)$ , which is shown in (b). After pre-training all  
 275 hidden layers, an output unit  $Y$ , which serves as a classifier, could be wired on top of the hidden layers to make  
 276 predictions. The whole architecture could be fine-tuned together using backpropagation and labeled data, which is  
 277 shown in (c).

278

279



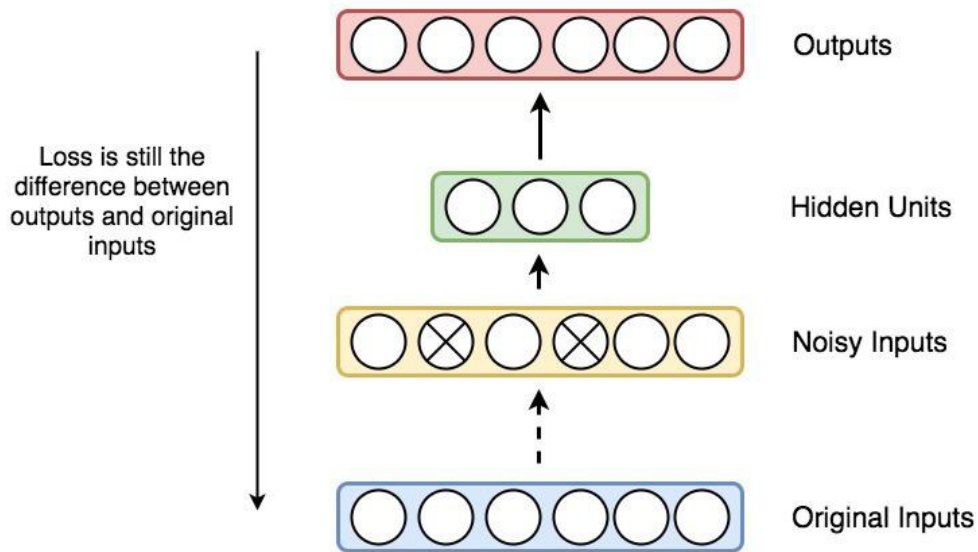
280

281 **Denoising auto-encoder (DA):**

282

283 A basic auto-encoder could successfully retain much of the information from the inputs in new  
284 features within the hidden layer. However, Vincent et al. (2010)[38] demonstrated that simply  
285 retaining information from the inputs does not guarantee that the extracted features are “good  
286 features”, which could achieve high-performance in supervised learning tasks. Denoising auto-  
287 encoder has been proposed to overcome this challenge by generating a noisy representation  
288 based on the inputs, such as setting values to 0 for a small proportion of input nodes or adding a  
289 noise term with a Gaussian distribution, and then feeding the noisy term into the auto-encoder  
290 (*Figure 9*). With the introduction of the noise term to the original inputs, denoising auto-  
291 encoders construct more robust feature representations and thereby could generalize better to  
292 unseen examples and datasets.

293



294

295 *Figure 9: Illustration of a denoising auto-encoder. A denoising auto-encoder first transforms original inputs into*  
 296 *noisy inputs. However, the loss in each step of the training process is still computed by the difference between the*  
 297 *reconstructed representations in the output layer and the original inputs in the input layer.*

298

299

### 300 **Stacked denoising auto-encoder (SDAE):**

301 A SDAE is a multi-layer auto-encoder, each hidden layer of which is a representation of the  
 302 previous layer obtained by a denoising auto-encoder with one hidden layer. For example, when  
 303 pre-train the 2 hidden layers  $h_1$  and  $h_2$  in Figure 8, one could add a noise term to the pre-training  
 304 inputs  $X$  and  $h_1$  to construct SDAE. Vincent et al. (2010)[38] showed that the features extracted  
 305 by SDAE are stable and robust under noisy inputs, by achieving the best classification results  
 306 under 9 out of 10 image databases. These features could efficiently capture useful information in  
 307 the input distribution and have yield equivalent or better classification performance over most of  
 308 the image data processing benchmarks.

309 Table 1 summarizes the 5 auto-encoders described above.

310 **Table 1. A summary of different auto-encoders**

Method	Description
Regular auto-encoder (AE)	Find low-dimensional representation of input using an unsupervised approach (i.e., no outcome information is used)
Sparse AE	Restrict the number of hidden nodes to be activated to avoid too many free nodes with complex dependencies in each layer
Stacked AE	Each hidden layer is a low-dimensional representation of the previous layer obtained by AE
Denoising AE (DA)	Introduce noises to input to make AE more robust to noises
Stacked denoising AE (SDAE)	Combine stacked AE and DA (i.e., introduce noises to input in a stacked AE)

311

312

313

314

#### 315 **4. Deep Learning Applications in Disease Prediction**

316

317 **Previous Works of Disease Prediction in Genomic Data Analysis using non-deep learning**

318 **approach:**

319

320 Plenty of methods have been proposed in disease prediction using genomic data (e.g., [39-44]).  
321 Due to the large number of predictors (i.e., gene probes), the main approach in disease  
322 detection/prediction is to first obtain a subset of gene probes (e.g., a few top gene probes in  
323 probe-wise tests) or a subset of representations of gene probes (e.g., a few top principal  
324 components), and then to predict disease status based on the selected probes or representations  
325 using machine learning algorithms.

326

327 Furey et al. (2000)[39] used SVMs to classify cancer tissue samples using gene expression  
328 datasets. The study showed that SVMs are able to classify tissue and cell types based gene  
329 expression data and have similar performances to other machine learning methods. Khan et al.  
330 (2001)[40] was among the first to adopt basic ANNs (ANNs without hidden layers) to classify  
331 cancer samples and to identify relevant genes. In their study, the 10 top PCA components were  
332 used as inputs to the ANN to classify the small, round blue-cell tumors (SRBCT) to four distinct  
333 diagnostic categories. All 63 samples in the training set and all 25 samples in the independent  
334 testing set were correctly classified based on the 96 selected genes. Pal et al. (2007)[41]  
335 proposed to combine modified perceptron network and relational fuzzy clustering algorithms[45]  
336 to select a gene subset used for cancer subgroup classification. They applied their method to the  
337 SRBCT dataset analyzed by Khan et al. (2001)[40] and identified 7 genes that can accurately  
338 classify the samples in both training set and testing set. Chang et al. (2011)[42] used an ANN  
339 with one hidden layer coupled with an additive step-wise approach for predicting colorectal  
340 cancer (CRC) using microRNAs (miRNAs). Three miRNAs were identified with median  
341 accuracy 100% by using an extensive Monte Carlo cross-validation strategy. Sharma et al.  
342 (2012)[15] proposed a top-r feature selection technique that repeatedly divides and merge gene

343 expression data to select the gene subset minimizing the loss of information. The selected genes  
344 are then tested on three tumor datasets and achieved higher accuracies than other feature  
345 selection methods, such as probe-wise tests. Nanni et al. (2012)[43] examined the SVM  
346 classification performance using multiple feature reduction and data transformation approaches,  
347 including neighborhood preserving embedding, orthogonal wavelet coefficients, and texture  
348 descriptors. The study showed that a combination of different feature extraction methods could  
349 enhance genomic classification performance. For instance, the two combined methods achieved  
350 the highest average area under ROC curves (AUC = 92.18% for the WF method and 92.07% for  
351 the FUS method), while the AUC values for the 8 individual feature extraction methods were  
352 ranged from 79.24% to 91.85%. Jordan and Do (2018)[44] reviewed the studies that predict  
353 disease using full genomic information. Their review focused on polygenic risk scores (PRS),  
354 which is the most common method of integrating information from across the genome into a  
355 single estimate of genetic risk. A PRS is a weighted average of the genetic status at each  
356 associated risk locus. The weighting of each locus is usually the regression coefficient of GWAS  
357 association for the locus. Jordan and Do (2018) mentioned that the power of most PRSs to  
358 predict disease risk has been very low due to several reasons, such as small sample size, genetic  
359 ancestry, heterogeneity of risk factors and causation.

360

361 The main limitations of these previous works[13] include (1) ignoring potential non-linear  
362 relationships among the features; (2) ignoring the contribution of features with weak signals to  
363 distinguish diseases; and (3) over-simplifying the complex prediction problem, such as using  
364 single-layer ANNs.

365

## 366 Deep Learning Applications in Disease Prediction

367

368 Through thorough literature search, we identified four papers[13, 31, 46, 47] published between  
369 January 2013 and December 2017, which applied deep learning methods in disease prediction  
370 using genomic data (Table 2). The details of the four studies will be discussed below.

371

372 Fakoor et al. (2013)[13] is among the first to apply deep learning methods to extract key features  
373 from gene microarray data in predicting cancers. Fakoor et al. (2013)[13] first applied PCA to  
374 eliminate the effects of redundant and noisy dimensions, then applied three auto-encoders  
375 methods (a sparse auto-encoder with one hidden layer, a stacked auto-encoder with 2 hidden  
376 layers, and a stacked auto-encoder with fine tuning) to further extract non-linearly-correlated  
377 discriminating features based on the top principal components combined with some randomly  
378 selected original features, and finally used softmax regression to do classification based on the  
379 low-dimensional representations. Thirteen gene microarray datasets (the range of sample sizes is  
380 20 – 1,047; the range of the numbers of features is 2,000 – 54,613) were used to compare the  
381 performances of deep learning methods and two traditional prediction methods: Softmax based  
382 on the top principal components (PCA+Softmax) or SVM with Gaussian kernel based on the top  
383 principal components (PCA+SVM). Ten-fold cross-validation was applied to estimate the  
384 average and standard deviation of the prediction accuracies and compared the average ACCuracy  
385 (ACC) of the three deep-learning methods with the maximum of the accuracy of the two  
386 traditional methods. For 8 of the 13 genomic datasets, at least one of the three deep learning  
387 methods has significantly higher average accuracy than the maximum accuracy of PCA+Softmax  
388 and PCA+SVM. The median [min, max] increase of average ACC is 1.5% [0.7%, 8.3%]. The

389 sample sizes of the 8 datasets range from 20 to 1,047. However, stacked auto-encoder without  
390 fine-tuning usually had much worse accuracy than the traditional methods. The stacked auto-  
391 encoder with fine-tuning achieved the best accuracy in six datasets with ACC ranging from  
392 76.67% to 95.15%, while the single-layer sparse auto-encoder perform the best in 5 datasets with  
393 ACC ranging from 46.76% to 91.50%.

394

395 Tan et al. (2015)[31] used denoising auto-encoders to learn compact and efficient representations  
396 in predicting disease status. Tan et al. (2015)[31] used the Molecular Taxonomy of Breast  
397 Cancer International Consortium (METABRIC) cohort as the training set (1,424 samples) and  
398 the testing set (712 samples) and the cohort from The Cancer Genome Atlas (TCGA) as the  
399 independent evaluation set (547 samples). The DA used in Tan et al. (2015)[31] has four layers:  
400 an input layer, a corrupted input layer (i.e., noisy input layer), a hidden layer, and a reconstructed  
401 input layer. Each node in the hidden layer was used to predict disease status (e.g., tumor vs. non-  
402 tumor, or ER+ vs. ER-) depending on whether the node value for a sample in the evaluation set is  
403 greater than the optimal threshold that was obtained based on the discovery set and testing set.  
404 Tan et al. (2015)[31] showed that each of the top three hidden nodes in the discovery set could  
405 also have high prediction accuracy ( $> 0.9$ ) in the evaluation set when they used their method to  
406 predict tumor status (tumor sample vs. non-tumor sample).

407

408 Danaee et al. (2016)[46] used SDAE to transform high dimensional, noisy RNA-seq gene  
409 expression data to lower dimensional, meaningful representations, based on which they applied  
410 different machine learning methods to classify breast cancer samples from the healthy control  
411 samples. They also identified a set of “Deeply Connected Genes” (DCGs) that have strongly

412 propagated influence on the reduced-dimension SDAE-encoding. Inspired by the classic study  
413 that applies SDAE to extract features in image data[38], Danaee et al. (2016)[46] built a SDAE  
414 model with four stacked layers of dimensions of 15,000, 10,000, 2,000, and 500, to obtain  
415 representations of genomic features to be fed into classifiers. A RNA-seq dataset (1,210 samples:  
416 1,097 breast cancer samples and 113 healthy samples) from TCGA is used to train and validate  
417 the model in the study. Danaee et al. (2016)[46] compared their prediction method with  
418 prediction methods based on PCA, Kernel PCA (KPCA, a non-linear PCA), the 206  
419 differentially expressed genes (DIFFEXP0.05) that were significant at an FDR of 0.05 in gene-  
420 wised tests, and top 500 most significant differentially expressed genes (DIFFEXP500). Three  
421 classifiers, including a single-layer ANN, SVM, and SVM-RBF (SVM with a radial basis  
422 function kernel), were used to do the prediction based on extracted features. Like Tan et al.  
423 (2015)[31], Danaee et al. (2016)[46] used a training set and a testing set to train classifiers, and  
424 used a validation set to evaluate the performance of the prediction methods. The classification  
425 result shows that the low-dimensional representations by SDAE outperformed other four sets of  
426 extracted features. For example, SDAE+SVM-RBF had accuracy (98.26%), sensitivity  
427 (97.61%), specificity (99.11%), precision (99.17%), and F-score [48] (0.983). Furthermore,  
428 Danaee et al. (2016)[46] showed that DCGs had slightly lower prediction accuracy than SDAE-  
429 extracted low-dimensional representations, but much higher prediction accuracy than the other  
430 methods.

431

432 Singh et al. (2016)[47] applied a stacked sparse auto-encoder (SSAE) to extract features to  
433 predict disease status for each of 36 datasets from the Gene Expression Machine Learning  
434 Repository (GEMLeR)[49]. The SSAE used by Singh et al. (2016)[47] has three hidden layers.



435 The input layer contains top 800 features selected based on Individual Training Error Reduction  
436 (ITER) ranking. The three hidden layers have 700, 600, and 500 nodes, respectively. The three  
437 classifiers, Softmax Regression, kernel SVM, and Random Forest, were applied to the 500  
438 extracted features to perform binary classification. Singh et al., (2016)[47] applied 10-cross-  
439 validation to estimate the classification accuracy and area under the ROC curve (AUC).  
440 Compared with the benchmark classification results taken from the GEMLeR website[49], the  
441 deep learning approach achieved slightly higher performance: ACC > 90.8% for 35 datasets  
442 (ACC>83.7% for all 36 datasets), and AUC>90.2% for 34 datasets (AUC >79.6 for all 36  
443 datasets).

444

#### 445 **Software packages for deep-learning-based feature extraction**

446 Since deep learning algorithms usually are complicated, it is important to have open-source  
447 software packages available so that investigators can directly use these packages to their  
448 genomic data analysis. Both Tan et al. (2015)[31] and Danaee et al. (2016)[46] used *Theano*  
449 software that provides the implementation of auto-encoder algorithms. Fakoor et al. (2013)[13]  
450 and Singh et al. (2016)[47] did not mention the software packages that they used for auto-  
451 encoding.

452

453 Several software packages/libraries are available to build auto-encoder models and fine-tune  
454 model parameters, such as Python packages (*Scikit\_learn*, *Theano*, *Keras*, and *TensorFlow*) and  
455 R packages (*h2o*, *kerasR*, and *autoencoder*). Wikipedia provides a table of deep learning  
456 software ([https://en.wikipedia.org/wiki/Comparison\\_of\\_deep\\_learning\\_software](https://en.wikipedia.org/wiki/Comparison_of_deep_learning_software)).

457

458 **5. Discussion**

459

460 In this article, we aimed to review all papers that applied the deep learning approach to predict  
461 disease status based on genomic data, which first obtains low-dimensional representations of  
462 high-dimensional genomic features, and then inputs these representations to the state-of-art  
463 classifiers that have excellent performance in low-dimensional classification problems. We  
464 found only 4 such papers, indicating that it is still in its infancy to predict disease status using  
465 deep learning on genomic data. However, the results of these 4 papers showed that the deep  
466 learning approach could extract useful genomic features from high-throughput whole genome  
467 data for prediction purpose with high accuracy.

468

469 Compared with commonly-used dimension-reduction methods (such as PCA and probe-wise  
470 testing), the deep learning approach could have better performance in terms of a variety of  
471 accuracy measurements: accuracy, AUC, sensitivity, specificity, precision, and F-score.  
472 Especially, it is impressive that probe-wise testing, which is currently the most popular approach  
473 to identify disease-associated probes, performed poorly compared with PCA or auto-encoders  
474 [46]. However, whether the performance of the deep learning approach is significantly better  
475 than the commonly used approaches was not investigated in the 4 papers, among which only  
476 Fakoor et al. (2013)[13] provided standard errors for the estimated ACC. However, Fakoor et al.  
477 (2013)[13] neither provided some key details (e.g., the number of principal components used and  
478 the number of randomly selected raw features), nor provided p-values for testing if the mean  
479 ACC obtained using a deep learning approach is significantly better than that by using the PCA  
480 approach. Moreover, Fakoor et al. (2013)[13] showed that not all auto-encoders could

481 outperform PCA. For example, Table 2 of Fakoor et al. (2013)[13] showed that for the first  
482 dataset, mean ACC (standard error) is 74.36% (0.062%) by using PCA+sparse auto-encoder,  
483 51.35% (0.019%) by using PCA+stacked auto-encoder, while PCA approach had mean ACC  
484 94.04% (SE 0.03%), although PCA+stacked auto-encoder with fine tuning (95.15% (0.047%))  
485 performed better than PCA.

486

487 Different auto-encoders were used in the 4 papers, such as sparse auto-encoder, stacked auto-  
488 encoder, stacked auto-encoder with fine-tuning, denoising auto-encoder, stacked denoising auto-  
489 encoder, and stacked sparse auto-encoder. Except Fakoor et al. (2013)[13], the other three papers  
490 did not compare the auto-encoders used in the paper with other auto-encoders. Table 2 of Fakoor  
491 et al. (2013)[13] showed that PCA+stacked auto-encoder performed worse than PCA+sparse  
492 auto-encoder and PCA+stacked auto-encoder with fine-tuning in 12 of the 13 datasets. However,  
493 neither PCA+sparse auto-encoder nor PCA+stacked auto-encoder with fine-tuning could  
494 outperform each other in all 13 datasets. For a fair comparison, it could be beneficial for future  
495 studies to compare the deep learning methods mentioned above using the same datasets.

496

497 All four papers mentioned the number of hidden layers and the number of nodes in each hidden  
498 layer used for the auto-encoders. However, no justifications and guidance were given on why  
499 choosing those specific numbers of hidden layers and those specific numbers of nodes in each  
500 hidden layer. This is probably one of the main reasons why deep learning has not been widely  
501 used in the genomic research area. There are some existing methods to choose the number of  
502 layers and nodes, such as (1) starting from a small neural network and adding layers and nodes  
503 until the error stops decreasing, and (2) starting from a big neural network and remove layer and

504 nodes until the error increases significantly[50]. Optimization methods such as grid search and  
505 random search are also proposed and discussed[51] to optimize the parameters in model training.  
506 However, these methods are still not well studied in genomic data analysis and could not  
507 eliminate the risks of over-fitting and under-fitting. Future research is still needed in choosing  
508 and optimizing deep learning parameters, especially in genomic data analysis.

509

510 Another possible reason why deep learning has not been widely used in the genomic research  
511 area is the lack of software packages that implement deep learning algorithms for genomic data  
512 analysis. Many investigators in genomic research area use the R language and use packages in  
513 Bioconductor (a repository of R packages specifically for genomic data analysis). Although there  
514 are a couple of R packages, such as *keras* and *kerasR*, connecting R to the Keras deep learning  
515 library, there is lack of examples and tutorials on how to use them to analyze genomic data and  
516 to visualize the low-dimensional representations that are obtained by auto-encoders.

517

518 It is a non-trivial task to interpret the low-dimensional representations (features) of the original  
519 expression data obtained by auto-encoders because the representations are non-linear functions  
520 of gene probes and the hidden layers in deep learning algorithms are like “black box”[52].

521 Among the 4 papers that we reviewed, Tan et al. (2015)[31] and Danaee et al. (2016)  
522 [46] suggested interpreting the representations based on the probes having strongly propagated  
523 influence on the reduced-dimension auto-encoding. However, no details were given on how to  
524 select these probes, except that these probes have high edge weights.

525

526 To evaluate classification performance, several measurements were used in the four papers that  
527 we reviewed, including accuracy (ACC), area under the ROC curve (AUC), sensitivity,  
528 specificity, precision, and F-measure. When the dataset is imbalanced (i.e., number of  
529 cases/positive samples is much different from that of controls/negative samples), using ACC  
530 could be biased. For example, given a dataset with 99% true negative samples and 1% true  
531 positive samples, a classifier could achieve 99% ACC even if it wrongly classifies all the true  
532 positive samples to the negative group. Fakoor et al. (2013)[13] only used ACC as the  
533 performance metric, while several genomic datasets analyzed in Fakoor et al. (2013)[13] are  
534 imbalanced. Tan et al. (2015)[31] also only used ACC to evaluate the performances of different  
535 prediction methods, while both the training and testing datasets are highly imbalanced. For  
536 imbalanced data, other performance metrics can be used, such as AUC, F-measure, and G-  
537 measure[48, 53], which are less sensitive to the case/control imbalance.

538

539 Over-fitting is a big issue in prediction. Using the same data set to both train the prediction  
540 model and evaluate the performance of the prediction model usually causes over-estimation of  
541 the prediction accuracy. Ideally, a testing set from a population independent of the training  
542 population is required in evaluating prediction accuracy. However, genomic data are usually  
543 expensive to collect. Hence, it is usually hard to obtain independent testing set in genomic  
544 research. Thanks to the policy of the National Institute of Health of the United States, numerous  
545 genomic datasets are now publicly available in the Gene Expression Omnibus  
546 (<https://www.ncbi.nlm.nih.gov/geo/>), an online repository of genomic datasets. Other public  
547 genomic repositories are also available, such as TCGA (<https://cancergenome.nih.gov>) and  
548 GTEx (<https://www.gtexportal.org/home/>). Hence, nowadays it is relatively easy to obtain an

549 independent testing set for most of complex human diseases. Among the 4 papers that we  
550 reviewed, only Tan et al. (2015)[31] used an independent testing set. The other 3 papers used K-  
551 fold cross-validation technique to alleviate the over-fitting issue.

552

553 Genomic data usually contain technical noises, such as batch effects (large samples have to be  
554 handled in multiple batches due to capacity limits of machines). Several methods, such as  
555 ComBat[54], have been proposed to remove the effects of technical batches before downstream  
556 data analysis. We can apply ComBat to the training set and the testing set, separately. Suppose  
557 after removing technical noises we build and validate a prediction model based on the training  
558 set and the testing set, with excellent prediction accuracy. Now a new subject's genomic data are  
559 obtained. Can we apply the prediction model to this new subject? The answer probably is “no”,  
560 since we do not know how to remove technical noises for only one new sample. One possible  
561 solution is to collect genomic data for a batch of subjects together. Then we can apply the  
562 prediction model to subjects in this batch after removing possible batch effects. A possibly  
563 better solution is to improve the technology to reduce technical noises. With the advancements in  
564 sequencing technology and a rapid decline in sequencing costs, DNA sequencing has gained  
565 remarkable popularity among biomedical researchers. Compared to microarrays, DNA  
566 sequencing data is believed to deliver faster, more complete, and more scientifically accurate  
567 genomic analysis[55].

568

569 The four deep-learning papers identified in this review compared the performances of deep  
570 learning approaches with PCA approach and probe-wise test approach. There are many more

571 advanced feature selection methods in the literature, such as the stable feature selection method  
 572 [16] and the Boruta algorithm [17]. More comprehensive comparisons are warranted.

573

574 Recently, semi-supervised learning and reinforcement learning are receiving a lot of attention in  
 575 image recognition, gaming, and robotics[56-58]. How to apply the frontier deep learning  
 576 innovations to genomic data analysis could be an interesting future research topic[59].

577

## 578 6. Conclusion.

579 In summary, this review showed that applying deep learning to find a low-dimensional  
 580 representation for high-throughput genomic data is a promising future trend in disease prediction  
 581 based on high-dimensional genomic data. The low-dimensional representation obtained by deep  
 582 learning could capture both linear and non-linear relationship among the probes. Deep learning is  
 583 a new technology for most scientists in genetics. Scientists in genetics should collaborate to  
 584 understand how deep learning could help predict disease status using genomic data, hence to  
 585 move this field forward.

586

587

588 **Table 2. Summary of the four studies that applied deep learning to predict disease status in the genomic research**

Author/Year	Datasets	Total Number of Samples	Feature Extraction Method	Classifier	Using cross-validation or not	Performance based on deep learning	Traditional Methods compared with
<b>Fakoor:2013</b>	gene expression data from 13 datasets	Various number of samples (20 - 1047) and	PCA+Sparse Auto-encoder PCA+Stacked	Softmax regression	10-fold cross-validation to evaluate	ACC+/- standard error: (33.7%+/-	Dimensional Reduction: PCA

		various number of features (2000 – 54675) across 13 datasets	Auto-encoder PCA+Stacked Auto-encoder with Fine-tuning	SVM with Gaussian kernel	classification performance	0.038% – (97.5%+/- 0.079%)	Classifier: Softmax regression SVM
<b>Tan:2015</b>	<p>Training: A gene expression dataset from METABRIC (Illumina HT-12 v3 platform)</p> <p>Independent Testing: A gene expression dataset from TCGA</p>	<p>METABRIC: 2136 samples (1992 breast cancer specimens and 144 tumor-adjacent normal tissues); 2520 genes after data cleaning</p> <p>TCGA: 547 samples (522 primary tumors, 3 metastatic tumors, and 22 tumor-adjacent normal samples); 2520 genes after data cleaning</p>	Denoising Auto-encoder (DA)	Sigmoid Activation	10-fold cross-validation to determine the appropriate parameter setting for the training set	ACC in testing set: 75.0%-99.6%	N/A



<b>Danaee:2016</b>	An RNA-seq expression dataset from TCGA	1210 samples (1097 breast cancer samples and 113 healthy samples)	Stacked Denoising Auto-encoder(SDAE)	ANN SVM SVM-RBF	5-fold cross-validation to evaluate classification performance	ACC: 96.95%- 98.26% Sensitivity: 97.21%- 98.73% Specificity: 95.29%- 99.11% Precision: 95.42%- 99.17%  F-measure 0.970-0.983	PCA KPCA Differentially Expressed Genes
<b>Singh: 2016</b>	36 gene microarray datasets from GEMLeR (Affymetrix GeneChip U133 Plus 2.0 arrays)	1545 samples (9 cancers, no control samples); 54676 features	Stacked Sparse Auto-encoder(SSAE)	Softmax Regression Random Forest Linear SVM RBF SVM	10-fold cross-validation to evaluate classification performance	AUC: 80%- 100%  ACC: 76%- 100%	KNN; SVM-RFE

589

590

591

592

593 **References:**

594

595

- 596 1. Ferlay, J., et al., *GLOBOCAN 2012 v1.0, Cancer Incidence and Mortality Worldwide: IARC CancerBase*. 2013, International Agency for Research on Cancer: Lyon, France.
- 597 2. Stewart, B. and C. Wild, *World cancer report*. 2014, International Agency for Research
- 598 on Cancer: Lyon, France.
- 599 3. World Health Organization. *Genes and noncommunicable diseases*. Genes and human
- 600 disease [cited 2018 April, 15]; Available from:
- 601 <http://www.who.int/genomics/public/geneticdiseases/en/index3.html>.
- 602 4. Whitcomb, D.C., *What is personalized medicine and what should it replace?* Nat Rev
- 603 Gastroenterol Hepatol, 2012. **9**(7): p. 418-24.
- 604 5. Cortes, C. and V. Vapnik, *Support-vector networks*. Machine learning, 1995. **20**(3): p.
- 605 273-297.
- 606 6. Breiman, L., *Random forests*. Machine learning, 2001. **45**(1): p. 5-32.
- 607 7. Altman, N., *An introduction to kernel and nearest-neighbor nonparametric regression*.
- 608 The American Statistician, 1992. **46**(3): p. 175-185.
- 609 8. Mastoi, Q.U., et al., *Automated Diagnosis of Coronary Artery Disease: A Review and*
- 610 *Workflow*. Cardiol Res Pract, 2018. **2018**: p. 2016282.
- 611 9. Zhang, Y.D., et al., *An imaging-based approach predicts clinical outcomes in prostate*
- 612 *cancer through a novel support vector machine classification*. Oncotarget, 2016. **7**(47): p.
- 613 78140-78151.
- 614 10. C, G., et al., *Model-based and Model-free Machine Learning Techniques for Diagnostic*
- 615 *Prediction and Classification of Clinical Outcomes in Parkinson's Disease* SCIENTIFIC
- 616 RePoRTS, 2018. **8**.
- 617 11. Bengio, Y., et al. *Greedy layer-wise training of deep networks*. in *Advances in Neural*
- 618 *Information Processing Systems 19*. 2007. MIT Press.
- 619 12. Catchpole, D.R., et al., *The curse of dimensionality: a blessing to personalized*
- 620 *medicine*. J Clin Oncol, 2010. **28**(34): p. e723-4; author reply e725.
- 621 13. Fakoor, R., et al. *Using deep learning to enhance cancer diagnosis and classification*. in
- 622 *The 30th International Conference on Machine Learning (ICML 2013)*. 2013.
- 623 14. Aliferis, C.F., et al. *Machine Learning Models for Classification of Lung Cancer and*
- 624 *Selection of Genomic Markers Using Array Gene Expression Data*. in *Proceedings of the*
- 625 *Sixteenth International Florida Artificial Intelligence Research Society Conference*.
- 626 2003. St. Augustine, Florida, USA.
- 627 15. Sharma, A., S. Imoto, and S. Miyano, *A top-r feature selection algorithm for microarray*
- 628 *gene expression data*. IEEE/ACM Trans Comput Biol Bioinform, 2012. **9**(3): p. 754-64.
- 629 16. Moon, M. and K. Nakai, *Stable feature selection based on the ensemble L*. BMC
- 630 Genomics, 2016. **17**(Suppl 13): p. 1026.
- 631 17. Degenhardt, F., S. Seifert, and S. Szymczak, *Evaluation of variable selection methods for*
- 632 *random forests and omics data sets*. Brief Bioinform, 2017.
- 633 18. Higham, C.F., et al., *Deep learning for real-time single-pixel video*. Sci Rep, 2018. **8**(1):
- 634 p. 2369.
- 635 19. Collobert, R., et al., *Natural Language Processing (Almost) from Scratch*. Journal of
- 636 Machine Learning Research, 2011. **12**: p. 2493-2537.
- 637

- 638 20. Lenz, I., H. Lee, and A. Saxena, *Deep learning for detecting robotic grasps*. The  
639 International Journal of Robotics Research, 2015. **34**(4-5): p. 705-724.
- 640 21. Levine, S., et al., *Learning hand-eye coordination for robotic grasping with deep*  
641 *learning and large-scale data collection*. The International Journal of Robotics Research,  
642 2017.
- 643 22. LeCun, Y., Y. Bengio, and G. Hinton, *Deep learning*. Nature, 2015. **521**(7553): p. 436-  
644 44.
- 645 23. Wang, X., et al., *Searching for prostate cancer by fully automated magnetic resonance*  
646 *imaging classification: deep learning versus non-deep learning*. Sci Rep, 2017. **7**(1): p.  
647 15415.
- 648 24. Bychkov, D., et al., *Deep learning based tissue analysis predicts outcome in colorectal*  
649 *cancer*. Sci Rep, 2018. **8**(1): p. 3395.
- 650 25. Alipanahi, B., et al., *Predicting the sequence specificities of DNA- and RNA-binding*  
651 *proteins by deep learning*. Nat Biotechnol, 2015. **33**(8): p. 831-8.
- 652 26. Zhou, J. and O.G. Troyanskaya, *Predicting effects of noncoding variants with deep*  
653 *learning-based sequence model*. Nat Methods, 2015. **12**(10): p. 931-4.
- 654 27. Nair, V. and G.E. Hinton. *Rectified linear units improve restricted boltzmann machines*.  
655 *in the 27th International Conference on Machine Learning*. 2010. Haifa, Israel.
- 656 28. Larochelle, H., et al., *Exploring Strategies for Training Deep Neural Networks*. Journal  
657 of Machine Learning Research, 2009. **1**: p. 1-40.
- 658 29. Hinton, G.E., S. Osindero, and Y.W. Teh, *A fast learning algorithm for deep belief nets*.  
659 Neural Comput, 2006. **18**(7): p. 1527-54.
- 660 30. Hinton, G.E. and R.R. Salakhutdinov, *Reducing the dimensionality of data with neural*  
661 *networks*. Science, 2006. **313**(5786): p. 504-7.
- 662 31. Tan, J., et al., *Unsupervised feature construction and knowledge extraction from genome-*  
663 *wide assays of breast cancer with denoising autoencoders*. Pac Symp Biocomput, 2015:  
664 p. 132-43.
- 665 32. Gupta, A., H. Wang, and M. Ganapathiraju. *Learning structure in gene expression data*  
666 *using deep architectures, with an application to gene clustering*. in *2015 IEEE*  
667 *International Conference on Bioinformatics and Biomedicine (BIBM)*. 2015. Washington,  
668 DC, USA: IEEE.
- 669 33. Nair, V. and G. Hinton, *3D Object Recognition with Deep Belief Nets*, in *Advances in*  
670 *Neural Information Processing Systems 22*, Y. Bengio, et al., Editors. 2009, Curran  
671 Associates, Inc. p. 1339-1347.
- 672 34. Lee, H., C. Ekanadham, and A. Ng, *Sparse deep belief net model for visual area V2*, in  
673 *Advances in Neural Information Processing Systems 20*. 2008, Curran Associates, Inc. p.  
674 873-880.
- 675 35. Le, Q., et al. *Building High-level Features Using Large Scale Unsupervised Learning*  
676 *in the 29th International Conference on International Conference on Machine Learning*.  
677 2012. Edinburgh, Scotland: Omnipress, USA.
- 678 36. Deng, J., et al. *Sparse autoencoder-based feature transfer learning for speech emotion*  
679 *recognition*. in *2013 Humaine Association Conference on Affective Computing and*  
680 *Intelligent Interaction*. 2013. IEEE.
- 681 37. Larochelle, H., et al. *An Empirical Evaluation of Deep Architectures on Problems*  
682 *with Many Factors of Variation*. in *the 24th International Conference on Machine*  
683 *Learning*. 2007. Corvallis, OR: ACM New York, NY , USA.

- 684 38. Vincent, P., et al., *Stacked Denoising Autoencoders: Learning Useful Representations*  
685 *in a Deep Network with a Local Denoising Criterion*. Journal of Machine Learning  
686 Research 2010. **11**: p. 3371-3408.
- 687 39. Furey, T.S., et al., *Support vector machine classification and validation of cancer tissue*  
688 *samples using microarray expression data*. Bioinformatics, 2000. **16**(10): p. 906-14.
- 689 40. Khan, J., et al., *Classification and diagnostic prediction of cancers using gene expression*  
690 *profiling and artificial neural networks*. Nat Med, 2001. **7**(6): p. 673-9.
- 691 41. Pal, N.R., et al., *Discovering biomarkers from gene expression data for predicting cancer*  
692 *subgroups using neural networks and relational fuzzy clustering*. BMC Bioinformatics,  
693 2007. **8**: p. 5.
- 694 42. Chang, K.H., et al., *MicroRNA signature analysis in colorectal cancer: identification of*  
695 *expression profiles in stage II tumors associated with aggressive disease*. Int J Colorectal  
696 Dis, 2011. **26**(11): p. 1415-22.
- 697 43. Nanni, L., S. Brahnam, and A. Lumini, *Combining multiple approaches for gene*  
698 *microarray classification*. Bioinformatics, 2012. **28**(8): p. 1151-7.
- 699 44. Jordan, D.M. and R. Do, *Using Full Genomic Information to Predict Disease: Breaking*  
700 *Down the Barriers Between Complex and Mendelian Diseases*. Annu Rev Genomics  
701 Hum Genet, 2018.
- 702 45. Hathaway, R. and J. Bezdek, *NERF c-Means: Non-Euclidean relational fuzzy clustering*.  
703 Pattern Recognition, 1994. **27**: p. 429-437.
- 704 46. Danaee, P., R. Ghaeini, and D.A. Hendrix, *A DEEP LEARNING APPROACH FOR*  
705 *CANCER DETECTION AND RELEVANT GENE IDENTIFICATION*. Pac Symp  
706 Biocomput, 2016. **22**: p. 219-229.
- 707 47. Singh, V., et al. *Layerwise feature selection in Stacked Sparse Auto-Encoder for tumor*  
708 *type prediction*. in *2016 IEEE International Conference on Bioinformatics and*  
709 *Biomedicine (BIBM)*. 2016.
- 710 48. Powers, D.M.W., *Evaluation: from precision, recall and F-measure to ROC,*  
711 *informedness, markedness and correlation*. Journal of Machine Learning Technologies,  
712 2011. **2**(1): p. 37-63.
- 713 49. Stiglic, G. and P. Kokol, *Stability of ranked gene lists in large microarray analysis*  
714 *studies*. J Biomed Biotechnol, 2010. **2010**: p. 616358.
- 715 50. Hagan, M., et al., *Generalization*, in *Neural Network Design*. 2015. p. 468-519.
- 716 51. Bergstra, J. and Y. Bengio, *Random Search for Hyper-Parameter Optimization*. Journal  
717 of Machine Learning Research, 2012. **13**: p. 281-305.
- 718 52. Snoek, J., H. Larochelle, and R. Adams. *Practical Bayesian Optimization of Machine*  
719 *Learning Algorithms*. in *the 25th International Conference on Neural Information*  
720 *Processing Systems*. 2012. Lake Tahoe, Nevada: Curran Associates Inc., USA.
- 721 53. He, H. and E.A. Garcia. *Learning from Imbalanced Data*. in *IEEE Transactions on*  
722 *Knowledge and Data Engineering*. 2009. IEEE.
- 723 54. Johnson, W.E., C. Li, and A. Rabinovic, *Adjusting batch effects in microarray expression*  
724 *data using empirical Bayes methods*. Biostatistics, 2007. **8**(1): p. 118-27.
- 725 55. Shendure, J. and H. Ji, *Next-generation DNA sequencing*. Nat Biotechnol, 2008. **26**(10):  
726 p. 1135-45.
- 727 56. Rasmus, A., et al. *Semi-supervised learning with Ladder networks*. in *the 28th*  
728 *International Conference on Neural Information Processing Systems 2015*. Montreal,  
729 Canada: MIT Press Cambridge, MA, USA.

- 730 57. Cutler, M. and J.P. How. *Efficient Reinforcement Learning for Robots Using Informative*  
731 *Simulated Priors*. in *2015 IEEE International Conference on Robotics and Automation*  
732 *(ICRA)*. 2015. Seattle, WA, USA: IEEE.
- 733 58. Mnih, V., et al., *Human-level control through deep reinforcement learning*. *Nature*, 2015.  
734 **518**(7540): p. 529-33.
- 735 59. Ching, T., et al., *Opportunities And Obstacles For Deep Learning In Biology And*  
736 *Medicine*. 2017: bioRxiv.  
737