

Deep Learning Generates Custom-made Logistic Regression Models for Explaining how Breast Cancer Subtypes are Classified

Takuma Shibahara (✉ takuma.shibahara.nj@hitachi.com)

Hitachi Ltd., Japan

Chisa Wada

Daiichi Sankyo RD Novare Co., Ltd.

Yasuho Yamashita

Hitachi Ltd., Japan

Kazuhiro Fujita

Daiichi Sankyo RD Novare Co., Ltd.

Masamichi Sato

Daiichi Sankyo RD Novare Co., Ltd.

Atsushi Okamoto

Daiichi Sankyo RD Novare Co., Ltd.

Yoshimasa Ono

Daiichi Sankyo RD Novare Co., Ltd.

Research Article

Keywords: Breast cancer, Deep Learning, Logistic Regression Models, genetic analysis

Posted Date: July 6th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-598333/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Deep learning generates custom-made logistic regression models for explaining how breast cancer subtypes are classified

Takuma Shibahara^{1,*}, Chisa Wada², Yasuho Yamashita¹, Kazuhiro Fujita², Masamichi Sato², Atsushi Okamoto², and Yoshimasa Ono³

¹Research and Development Group, Hitachi Ltd., Tokyo, Japan

²Bioinformatics Group, Translational Research Department, Daiichi Sankyo RD Novare Co., Ltd., Tokyo, Japan

³Translational Research Department, Daiichi Sankyo RD Novare Co., Ltd., Tokyo, Japan

*Corresponding to: takuma.shibahara.nj@hitachi.com

ABSTRACT

Breast cancer is the most frequently found cancer in women and the one most often subjected to genetic analysis. Nonetheless, it has been causing the largest number of women's cancer-related deaths. PAM50, the intrinsic subtype assay for breast cancer, is beneficial for diagnosis but does not explain each subtypes mechanism. Deep learning can predict the subtypes from genetic information more accurately than conventional statistical methods. However, the previous studies did not directly use deep learning to examine which genes associate with the subtypes. To reveal the mechanisms embedded in the PAM50 subtypes, we developed an explainable deep learning model called a point-wise linear model, which uses meta-learning to generate a custom-made logistic regression for each sample. We developed an explainable deep learning model called a point-wise linear model, which uses meta-learning to generate a custom-made logistic regression for each sample. Logistic regression is familiar to physicians, and we can use it to analyze which genes are important for prediction. The custom-made logistic regression models generated by the point-wise linear model used the specific genes selected in other subtypes compared to the conventional logistic regression model: the overlap ratio is less than twenty percent. Analyzing the point-wise linear models inner state, we found that the point-wise linear model used genes relevant to the cell cycle-related pathways.

1 Introduction

2 Molecular subtyping of cancer is important for providing specified therapies to patients and for developing novel drugs
3 targeting different molecules for the different subtypes. Breast cancer was classically classified by protein expression of the
4 estrogen receptor (ER), progesterone receptor (PR) and the epidermal growth factor receptor ErbB2/Her2. And expressions of
5 those receptors were used as clinicopathological variables for treatment decisions¹. Since the beginning of this millennium,
6 with high-throughput genomics technologies, breast cancer has been considered to have five clinically relevant molecular
7 subtypes defined by intrinsic gene expression patterns of the cancer¹⁻⁵. These intrinsic subtypes are Luminal A, Luminal B,
8 Her2-enriched, basal-like and normal breast-like cancer. While the subtypes do not perfectly reflect the clinical features, most
9 breast cancers of the luminal subtypes are ER/PR-positive, most Her2-enriched ones have amplification of the Her2 gene and
10 most basal cancers are triple negative (ER/PR/Her2). In the original PAM50 study the classification of normal-like subtype was
11 trained with normal breast tissue⁴. Therefore, cancer samples classified to the normal-like subtype are often interpreted as low
12 tumor content sample^{4,6}.

13 PAM50 was originally developed as a predictor of the five intrinsic subtypes from the expression pattern of 50 genes
14 determined using a microarray³. Since the PAM50 subtyping was clinically relevant and added significant prognostic and
15 predictive information to diagnosis, different measurement platforms such as RNA-seq have been adapted to PAM50⁷. While
16 PAM50 subtyping is useful for diagnosis and stratified treatment, it remains elusive which genes contribute to the mechanisms
17 of action and/or mechanisms of resistance to treatment for each subtype. Therefore, we investigated whether clues to those
18 mechanisms were embedded in the PAM50-classified subtypes with a recent machine learning approach.

19 Seven decades after the birth of the learning machine⁸, we now have deep learning neural networks (NNs) that provide
20 various predictive analyses of cancer classification. As deep learning spreads into more and more applications, there is a growing
21 need to explain the reasons for its predictions. In fact, to make deep learning models explainable, a lot of methods measuring
22 the importance of individual features (i.e., how much each feature contributes to the output) have been devised. These methods
23 can be roughly classified into perturbation-based and saliency-based ones. Perturbation-based methods calculate an importance
24 score based on how the output behaves in relation to perturbed input⁹⁻¹¹. The problem is their high computational cost: these

25 methods require a large number of perturbations to be generated around each sample’s input and make all the perturbed input
 26 for all the samples propagate through the whole network. In saliency-based methods, on the other hand, the importance score
 27 depends on each feature’s saliency evaluated by the gradient of output with respect to the input¹²⁻¹⁵. Saliency-based methods
 28 are computationally inexpensive because this gradient can be calculated in a single backward pass.

29 In the work reported here we developed a point-wise linear model as an innately explainable deep learning model. In the
 30 conventional deep learning models, the network computes new feature vectors whose linear combination sufficiently expresses
 31 the objective model. The network of the point-wise linear model, in contrast, derives a weight function for each original feature
 32 vector as a function of the original feature vectors. The weighted sum of the original feature vectors with these weight functions
 33 gives the objective model. Since the weight functions depend on the original feature vector, the model, unlike a simple linear
 34 model, involves the nonlinear interactions between the original features. Because the model is a simple weighted sum of the
 35 original feature vectors, the importance of each feature can be evaluated by its weight function, as in a linear model.

36 The importance analysis using the point-wise linear model has two novel properties. One is that the importance is computed
 37 in the same way as in a linear model, a traditional model with which medical informatics experts are very familiar. This property
 38 is highly advantageous in medical applications because one can use the know-how of medical-data analysis accumulated
 39 throughout its long history. LIME (Local Interpretable Model-Agnostic Explanations)¹¹, a method that produces a linear model
 40 for each data point by sampling perturbed data points, also has this advantageous property but, like other perturbation-based
 41 methods, is exceptionally computationally expensive. In contrast, our model produces an explainable model and incurs no
 42 extra computational time to compute importance scores (other than the time needed for training the model); this is the other
 43 novel property of our model. The saliency-based methods are also free from computational-cost problems, but they do not have
 44 the familiarity of a method computing importance the way a linear model does. In Layerwise Relevance Propagation¹⁴ and
 45 DeepLIFT (Deep Learning Important FeaTures)¹⁵, for example, the importance score is derived from how much the output
 46 deviates when the input differs from some reference value. It is difficult to compare this score with the results of previous
 47 medical studies. In contrast, linear and piecewise linear models such as logistic regression, random forest, and gradient boosting
 48 can provide feature importances explaining what learning machines are doing. A deep learning model that accurately predicts
 49 the cancer subtypes might use unknown and nonstandard knowledge with gene expressions. In this regard, we have devised a
 50 way to extract novel knowledge from the deep learning model.

51 Results

52 Comparison of learning ability and explainability

53 To investigate the nonlinear prediction ability and explainability of the explainable point-wise linear model, we trained a logistic
 54 regression model and a self-normalizing neural networks (SNNs) model, as a state-of-the-art deep learning¹⁶, by using a simple
 55 dataset (*sklearn.datasets.make_circle*¹⁷), respectively. Figure 1 shows three architectures of the machine learning models: (a)
 56 logistic regression, (b) deep learning, and (c) point-wise linear. Let $x^{(n)} \in \mathcal{R}^D$ represent a feature vector with N denoting the
 57 sample size and \mathcal{R} indicating the real number set. Firstly, we define a logistic regression model (Fig. 1 (a)) as follows:

$$y^{(n)} = \sigma(w \cdot x^{(n)}), \quad (1)$$

58 where $w \in \mathcal{R}^D$ is a weight vector for $x^{(n)}$, σ is a sigmoid function and \cdot is the inner product. $y^{(n)}$ is a probability value such as
 59 one expressing the likelihood of tumor tissues or normal tissues. The weight vector w is bound to the feature vector $x^{(n)}$. We can
 60 understand the importance of each feature variable by analyzing the magnitude of the elements in w . However, the circle in a
 61 circle is not a linearly separable problem, so as shown in Fig. 2, the logistic regression model could not classify the two circles.

62 Next, we give a usual deep learning model like that shown in Fig. 1 (b). A new feature vector $\boldsymbol{\varphi}(x^{(n)}) \in \mathcal{R}^{D'}$ is nonlinearly
 63 generated from the original feature vector $x^{(n)}$ through the L -layer neural network ($L \geq 2$). Note that the notation $v(u)$ for
 64 arbitrary vectors v and u indicates that every element of v is a function of the elements of u . A deep-learning-based nonlinear
 65 classification function predicts probability $y^{(n)}$ (Fig. 1 (b)) as follows:

$$y^{(n)} = \sigma\left(w' \cdot \boldsymbol{\varphi}(x^{(n)})\right), \quad (2)$$

66 where $w' \in \mathcal{R}^{D'}$ is a universal weight vector for $\boldsymbol{\varphi}$. The magnitude of each w' element represents the contribution of the
 67 corresponding element of $\boldsymbol{\varphi}$ to the prediction as shown in Fig. 1 (b). The SNNs correctly classified the blue and orange dots
 68 as shown in Fig. 2 (c). However, one cannot “explain” the machine’s prediction by w' because one cannot understand the
 69 meanings of $\boldsymbol{\varphi}$ with which the machine makes its predictions.

70 In order to make a deep NN explainable, we considered a meta-learning approach to generate a logistic regression model
 71 defined (Fig. 1 (c)) as

$$y^{(n)} = \sigma \left(\xi(x^{(n)}) \cdot x^{(n)} \right), \quad (3)$$

72 where each element of $\xi \in \mathcal{R}^D$ is a function of $x^{(n)}$ that the NN determines. ξ behaves as the weight vector for the original
73 feature vector $x^{(n)}$. The magnitude of each element of ξ describes the importance of the corresponding feature variable. Here
74 we should notice that this weight vector is tailored to each sample because ξ depends on $x^{(n)}$. We called Eq. (3) a point-wise
75 linear model given by a straightforward method over the sample index (n). The architecture of the point-wise linear model
76 consists of two blocks as shown in Fig. 1 (c). Also, the ξ is called the point-wise weight. The upper block is a meta-learning
77 machine of generating logistic regression models. The lower block is the logistic regression models for the inference task.
78 However, the tailored weight vector ξ can easily lead to poor generalization (i.e., overfitting or underfitting). The point-wise
79 linear model (Eq. (3)) tried to learn the labels of all samples as it generated a weight vector optimized for each sample, which
80 led to underfitting in Fig. 2 (c).

81 We newly found an equation constructing a point-wise weight ξ without losing generalization ability as the following:

$$\xi(x^{(n)}) \equiv w \odot \eta(x^{(n)}), \quad (4)$$

82 where the reallocation vector $\eta \in \mathcal{R}^D$ is nonlinearly generated from the original feature vector $x^{(n)}$ through the L -layer NN
83 ($L \geq 2$). $w \in \mathcal{R}^D$ is a universal weight vector that is independent of $x^{(n)}$ and \odot is the Hadamard product. In contrast to the model
84 defined straightforwardly by Eq. (3), the model defined by Eq. (4) accurately predicts the classification boundary, as shown
85 in Fig. 2 (e). The weight vectors $\xi^{(n)}$ in Eq. (4) smoothly change for each data sample (Fig. 2 (g)). The reallocation-based
86 point-wise linear model thus enables generalization. Additionally, we call $\rho(x^{(n)}) \equiv \eta(x^{(n)}) \odot x^{(n)}$ a reallocated feature vector
87 in \mathcal{R}^d . NNs have a versatile ability to map a linear feature space to a nonlinear feature space. By utilizing this ability, $\eta^{(n)}$
88 reallocates the feature vector $x^{(n)}$ into the new vector ρ that is linearly separable by a single hyperplane drawn by w . The
89 mechanism of Eq. (4) is discussed in Supplementary Appendix 1.

90 Prediction performance

91 In advance of the breast cancer subtypes analysis, we evaluated three subtype prediction models: one built by using logistic
92 regression with regularization, one using SNNs (state-of-the-art feed-forward NNs), and the point-wise linear. The aim of this
93 paper is to obtain a more accurate explainable model using deep learning. Therefore we compared the explainability of the
94 point-wise linear model, our proposed model, with the explainability of representative logistic regression.

95 Table 1 shows the number of samples for each subtype. The number D of mRNAs types (i.e., feature dimension) was
96 17,837 due to genes with same gene symbols among RNA-seq data, and copy number alteration data were used. Table 2 shows
97 the training and test area under the curve (AUC) values of 10-fold double cross-validation (DCV) yielded from evaluation with
98 each technique (all data was stored in Supplementary Table 1–3). Values in the column labeled 'All' in Table 2 were calculated
99 by using Eq. (7). The AUC value of each subtype was averaged by the 10-fold DCV. The hyperparameter search ranges and
100 values of each prediction model were stored in Supplementary Tables 4–6 (ranges) and Supplementary Tables 7–9 (optimized
101 values), respectively. The point-wise linear model marked the best values of 'All' for training and test set (training: 0.859,
102 test: 0.862). The AUC values of the test were more increased than those of the training in the point-wise linear and logistic
103 regression models. The training and test AUC values of the SNNs model were lower than 0.8.

104 Subtype analysis

105 First, the top 500 genes in terms of the relative score calculated by the importance analysis (see Methods) were selected as
106 highly contributing feature variables to predict subtypes. As shown in Fig. 3 and Supplementary Table 12, both the point-wise
107 linear model and logistic regression model classified subtypes by specific genes not selected in other subtypes, accounting
108 for 87.695.2 %, 72.091.0 % in each 500 features, respectively. In particular, prediction of Luminal A and B was determined
109 by a higher number of specific genes in the point-wise linear model than in the logistic regression one. We then compared
110 the overlap of the specific genes between the two models. There was not much commonality in the specific genes that are
111 important for predicting the same class in both models, with 206 genes in the most overlapping Her2 subtype (Table 3). The
112 Her2 subtype is a class in which the amplification of the HER2 gene is enriched, and many genes in the vicinity of the HER2
113 gene were commonly contained (Supplementary Table 10 and Supplementary Fig. 1). Since the model was generated from
114 copy number data, it was reasonable that the two models had a relatively high degree of commonality. As shown in Table 3,
115 for the specific genes of Luminal A in both models, they had 43 genes in common. The specific genes of Luminal A in the
116 logistic regression model had 39 genes in common with those of Luminal B, the similar luminal type, in the pointwise linear
117 model. On the other hand, those of Luminal A in the point-wise linear model had nothing in common with those of Luminal B

118 in the logistic regression model. Moreover, they had more in common with those of the Normal-like and Basal subtypes of the
119 logistic regression model, with 21 and 60 genes, respectively.

120 As another approach to the subtype classification criteria of deep learning, we used a manifold learning technique called
121 UMAP¹⁸. Figure 4 shows the one-dimensional (1D) and 2D embeddings of the reallocation vector (RV η) and the reallocated
122 feature (RF ρ), and the final inner vector of SNNs (ϕ in Eq. (2)). The RV vector reflected the clinical features (ER/PR/Her2).
123 As shown in Fig. 4 (a), Luminal A and B stuck together. However, each peak of their 1D embeddings stood in a distinct
124 position. Her2-enriched samples are close to the cluster of Luminal A and B samples. The 2D embedding of RV η (Fig. 4
125 (a)), shows that basal-like samples stay away from other subtypes. Normal-like samples take a position under the Luminal
126 and Her2-enriched samples. In the case of RF ρ embeddings, all subtypes formed a single cluster, as shown in Fig. 4 (c). The
127 embeddings of SNNs were separated for each subtype, as shown in Fig. 4 (d).

128 These subtypes were originally grouped by PAM50 using mRNA expression level. For this reason, we investigated what
129 kinds of genes' mRNA expression levels are associated with the 1D embeddings of RV η , RF ρ , and the inner vector of SNNs
130 model in order to figure out the difference of three values. We obtained each of the top 500 of genes correlated with the mRNA
131 expression values and then examined their functionally enriched pathways using ingenuity pathway analysis (IPA, QIAGEN)¹⁹.
132 Notably, those genes derived from RV η were overlapped with the cell cycle-related pathways such as "G2/M DNA Damage
133 Checkpoint Regulation", "Estrogen-mediated S-phase Entry", "Mitotic Roles of Polo-Like Kinase", and "Cell Cycle Control
134 of Chromosomal Replication", whereas the ones derived from RF ρ showed little significant enrichment as summarized in
135 Table 4 (the full list is available in Supplementary Table 11). The gene from SNNs model that resulted in the most enriched
136 pathway was "ErbB Signaling", which according to IPA is related to the cell cycle during growth and development of a number
137 of tissues, however, its statistical significance was low. The cell cycle is an essential function of cell proliferation and affects
138 the characteristics and malignancy of cancer. The results that both the point-wise linear and SNNs models recognized those
139 pathways in the subtype classification were acceptable and thought-provoking. In terms of giving suggestions for further
140 analysis, the SNNs model was inadequate and the point-wise linear η model was preferred because it presented multiple related
141 pathways.

142 Discussion

143 This study has established the point-wise linear model as an innately explainable deep learning model that can evaluate breast
144 cancer subtypes' biological mechanism. The point-wise linear model generates a custom-made logistic regression model for
145 each sample, which allows us to analyze which genes are important for subtype prediction. We have also shown the new scoring
146 method for selecting genes for the subtype classification. Then, we demonstrated the point-wise linear model had the highest
147 prediction performance and explained the breast cancer subtypes.

148 The AUC values of the point-wise linear model were better than those of logistic regression and SNNs models throughout
149 the 10-fold DCV. The SNNs model had overfitting, as demonstrated by the higher AUC values in the training results and
150 the lower AUC values in the test results (Table 2). A possible reason for this overfitting is that the progression of learning is
151 confined to the upper-layer parameters: i.e., there is a lack of advanced learning in the lower layers. Klambauer et al. show
152 the SNNs could use 32 layers in the conventional machine learning data set¹⁶. In our task, the number of the SNNs model's
153 inner layers (average 13, minimum 10, and maximum 16 as summarized in Supplementary Table 8) was lower than that of the
154 point-wise linear model (average 19, minimum 12, and maximum 24 as summarized in Supplementary Table 9). We optimized
155 the number of both models' inner layers from 10 to 25 layers. Figure 1 shows that the point-wise linear model contains a deep
156 learning block to generate the custom-made logistic regression model. Thus, it model likely causes overfitting when the number
157 of layers is increased. This study used a unified architecture (see Supplementary Appendix 1) as a newly developed architecture
158 of deep learning characterized by the binding of each network layer neurons in a mesh-like form, as shown in Fig. A4 of
159 Supplementary Appendix 1 and Ref.²⁰. The unified architecture has horizontally shallow and vertically deep layers to prevent
160 gradient vanishing and explosion. No matter how many layers are stacked vertically, there are only two horizontal layers from
161 the data unit nodes to the output node, as shown in Fig. A4 of Supplementary Appendix 1.

162 The logistic regression model and the point-wise linear model had almost equivalent prediction performance. On the
163 other hand, Table 3 suggests that the genes important for predicting the breast cancer subtypes were different in the logistic
164 regression model and the point-wise linear model. The point-wise linear model tended to select specific genes for each subtype,
165 as shown in Fig. 3 and Supplementary Table 12. These differences are considered to result from the models' ability to treat the
166 nonlinear relationships between feature and target variables. The logistic regression expresses the target variables only as linear
167 combinations of the feature variables, while the point-wise linear can express the target variables nonlinearly as in Eqs. (3) and
168 (4). In addition, those equations (Eqs. (3) and (4)) can be modified as $w \cdot (\eta(x) \odot x)$. The RV η corrects the feature variables
169 x so that the universal weight w can linearly separate the feature vector x . The differences of important genes between the
170 point-wise linear method and the logistic regression method rely on the feature variables' correction mechanism. Namely, the
171 differences of important genes between both methods are due to the feature variables' correction mechanism of the point-wise

linear method. The mechanism might also help the point-wise linear model express the target variables without using some of the features not specific to the subtype. In contrast, the logistic regression model had to use such nonspecific features to express the target variables without the corrections. We expected the point-wise linear model's correction mechanism to be vital in the tasks that demand high nonlinearity and consequently result in low AUC values in logistic regression. One of the tasks was to predict subtypes Luminal A and B (See Table 2). Consistently, the difference of the number of specific genes in the two models was large in Luminal A and B, and the specific genes for Luminal A and B in the point-wise linear model had much in common with the specific genes of other subtypes in the logistic regression model.

As is discussed above, the results (Fig. 3, Table 2, Table 3, and Supplementary Table 1) are consistent with the interpretation of the point-wise linear model from the viewpoint of the correction to the features. This consistency suggests that our scoring method described in the Methods section works properly. Our scoring method was designed to extract the features that contribute to predicting a subtype, especially among the corresponding subtype samples, rather than features not found in the other samples. This scoring method helps examine which features contribute to the prediction result with the aid of corrections by RV η but cannot examine which features contribute to the corrections η . Therefore, the relative score is not the perfect measure to investigate the mechanism of the classification of breast cancer subtypes.

For further investigation, we analyzed our predictive model's internal state in detail and have shown that we could consider the gene sets' biological implications contributing to the classification. From the comparison of RV η , RF ρ as the internal structural data to be analyzed, we found that RV η was better suited to elicit candidate hypotheses. In this paper, we chose UMAP to analyze the internal state in detail and reduce the dimensionality and then combine it with mRNA expression levels, which derives the biological implication that η was involved in the cell cycle. The fact that genes involved in the cell cycle affect subtypes has also been well-known²¹, and such genes have been reported as promising drug targets²², indicating that our model internal state was worth analyzing. Notice that the UMAP embedding of RV η , rather than the ones of RF ρ and ϕ , contains richer information related to interpretable pathways. While RF ρ is created as new features with which the problem is linearly separable in the point-wise linear model as well as ϕ is in SNNs model, RV η is considered as the corrections to the features by which the features are transformed into RF ρ and is unique to the point-wise linear model. This result suggests that it is essential to analyze the corrections to the features for our task to investigate the breast cancer subtypes classification mechanism. Considering this case as an example, we think that biologists and informaticians can apply our analysis with RV η to other tasks, such as exploring new drug target molecules or investigating mechanisms.

The use of the point-wise linear method, fully explainable deep learning, could classify the breast cancer subtype and depict its genomic characteristics. The point-wise linear model performed better than other techniques predicting breast cancer subtypes, including state-of-the-art deep learning. The point-wise linear model used the specific genes not selected in other subtypes or used by the logistic regression model. Additionally, as the result of analyzing the deep learning models' inner state, we found that the point-wise linear model used genes relevant to the cell cycle-related pathways. The results of this study suggest the potential of our technique to play a vital role in cancer treatment.

Methods

Importance analysis

Using the point-wise linear model Eq. (3), we can calculate the importance of each feature variable, i.e., how much each feature contributes to the model's prediction. The point-wise weight vector ξ depends on $x^{(n)}$ and consequently describes each sample's own feature importance. Therefore, we invented a method to derive a feature importance for a sample group so as to reveal a group macroscopic property contained in the point-wise weight vector of the group's samples. The idea of this feature importance was also used in Ref.²³.

First, we calculated feature importance for each sample from the weight vector $\xi^{(n)}$ in Eq. (3). On the basis of the idea of Shapley value²⁴, we introduced a sample-wise importance score for a k -th feature x_k of a sample with index (n) as

$$s_k^{(n)} \equiv \xi_k^{(n)} x_k^{(n)} - \frac{1}{|U_{(n)}|} \sum_{i \in U_{(n)}} \left[\xi_k^{(i)} x_k^{(i)} \right] \quad (5)$$

where $U_{(n)}$ is the set of samples whose weights are close to that of sample (n). This sample-wise importance score expresses the contribution of a sample (n) to raising the output probability $y^{(n)}$ compared with the average contribution among a sample group whose members obey similar linear models. In this study, we defined $U_{(n)}$ as follows: a sample (i) is in the set $U_{(n)}$ if $|\xi^{(i)} - \xi^{(n)}|/|\xi^{(n)}|$ is smaller than $4|\sigma|/|\bar{\xi}|$, where σ and $\bar{\xi}$ are vectors whose elements are given as the standard deviation and the mean, respectively, of the corresponding element of ξ .

Next, we defined group-wise importance scores, importance scores for a group (e.g., subtype Her2 samples), by using the sample-wise importance score. We performed voting among the group in which each sample votes to its top 10% features

221 whose sample-wise importance scores are highest, i.e., the features that significantly raise each sample’s output probability. We
222 defined a group-wise importance score v for each feature as the rate of samples who vote for the feature in the above voting.

223 Finally, we introduced relative score to extract the features that characterize a subtype. We divided the samples into samples
224 of a target subtype and the others and evaluated the group-wise importance score for each group. We refer to the group-wise
225 importance score for the target subtype samples group and the one for the others as v_{target} and v_{others} , respectively. v_{target} is not
226 necessarily appropriate for extracting the features that characterize the target subtype because even when v_{target} is high for a
227 feature, if v_{others} is also high the feature might be important for all the subtype samples, not for only the target subtype samples.
228 Therefore we defined relative score v_{rel} so as to compute the feature importance for the subtype samples relative to the one for
229 the others:

$$v_{\text{rel}} \equiv (v_{\text{target}})^2 - (v_{\text{others}})^2. \quad (6)$$

230 Since the ranges of both v_{target} and v_{others} are $[0, 1]$, the range of v_{rel} becomes $[-1, 1]$. Supplementary Fig. 2 shows the
231 distribution of v_{rel} in $v_{\text{others}}-v_{\text{target}}$ space. If a relative score is large, it is expected that both the summation and difference
232 of the group-wise importance scores v_{target} and v_{others} are large. Namely, features with large relative scores are important to
233 a certain degree for all the samples, and simultaneously much more important for the target subtype samples than for the
234 others. We considered the features with high relative scores for each subtype to be the important features that characterize the
235 corresponding subtype.

236 Analysis using manifold learning

237 In addition to the importance scores, the point-wise linear model can present another approach to analyze the nature of the
238 subtypes of breast cancer. The outputs of NNs’ inner layers give us some hints as to what criteria NNs use to classify the
239 subtypes of breast cancer from the feature vector. The output of the final inner layer (a new feature vector ϕ) can be linearly
240 separated by a single hyperplane spanned by w' as shown in Eq. (2). When the prediction accuracy is high, the output of the
241 final inner layer provides a well-summarized representation of the feature vector x for the classification task. The final inner
242 layer of the point-wise linear model is the reallocation vector η as shown in Eq. (4). Here we analyzed the reallocation vector x
243 by using a manifold learning technique called UMAP to reveal the subtype classification criteria.

244 We analyzed the 1D values of the reallocation feature vector η projected by UMAP. Then we calculate the Spearman
245 correlation coefficient for the relation between the one-dimensional values of the point-wise weight vectors projected by UMAP
246 and RNA-seq $\log_2(\text{RSEM} + 1)$ values, after which we selected the top 500 genes for which the correlation coefficient was
247 positive and the top 500 genes for which the correlation coefficient was negative. Those genes were analyzed using IPA to
248 interpret the canonical pathways.

249 Feature vectors and target variables

250 We used breast cancer TCGA²⁵ dataset retrieved by UCSC public Xena hub²⁶ for gene expression RNA-seq dataset (dataset
251 ID: TCGA.BRCA.sampleMap\HiSeqV2_PANCAN), copy number alteration (gene-level) dataset (dataset ID: TCGA.BRCA-
252 .sampleMap\Gistic2_CopyNumber_Gistic2_all_thresholded.by_genes), and phenotype dataset (dataset ID: TCGA.BRCA-
253 .sampleMap\BRCA_clinicalMatrix). RNA-seq values were calculated by UCSC Xena as follows: $\log_2(x + 1)$ value were
254 mean-normalized per-gene across all TCGA samples (x is RSEM value²⁷). In the copy number dataset of UCSC Xena, GISTIC2
255 values were discretized to $-2, -1, 0, 1, 2$ by Broad Firehose. We used subtypes precalculated by PAM50 (Luminal A, Luminal
256 B, Basal-like, Her2-enriched and normal-like²) in the Xena dataset as the target variables of the prediction model⁵. Table 1
257 shows the number of samples for each subtype. The number D of mRNAs types (i.e., feature dimension) was 17,837 because
258 we adopted gene symbols that overlapped with both the RNA-seq data and the copy number alteration data.

259 Statistical evaluations

260 Prediction performance evaluation

261 The subtype prediction models were built by using logistic regression with regularization (implemented by scikit-learn v0.22,
262 Python 3.7.6), SNNs (implemented by PyTorch 1.5, Python 3.7.6), and point-wise linear (implemented by PyTorch 1.5, Python
263 3.7.6). The aim of this paper is to obtain a more accurate explainable model using deep learning. Therefore, we compared
264 the explainability of the point-wise linear method, our proposed method, with the explainability of representative logistic
265 regression.

266 If the hyperparameters of a prediction model are optimized by using all samples, we may overlook the hyperparameter
267 overfitting. Addressing the problem, the prediction model evaluation was carried out by a K -fold DCV²⁸. The K -fold DCV can
268 measure the prediction performance of the entire learning process including its hyperparameter optimization. The procedure of
269 K -fold DCV has internal (training) and outer (test) loops. In this study, each internal loop searches the best hyperparameters set

270 (i.e., combinations of the hyperparameters) of the prediction model L times by using the tree-structured parzen estimator^{29,30}; a
271 single nested loop in the inner loops uses M -fold CV to evaluate the prediction performance of the prediction model with a
272 hyperparameters set. Each inner loop trains the prediction models with different hyperparameters sets $L \times M$ times. Then each
273 outer loop tests the prediction model with the best hyperparameters set obtained by its internal loop. Thus the hyperparameter
274 optimization process is completely separated from the test data. The procedure of K -fold DCV trains the classifier $K \times L \times M$
275 times. In our experiment we set $M = K$ and $L = 100$ ($K \times L \times M = 10,000$).

276 In each iteration of the 10-fold DCV and its internal 10-fold CV, the mean AUC was calculated as seen in the following
277 equation:

$$\text{Mean AUC} = \frac{1}{KC} \sum_{k=1}^K \sum_{t \in \text{subtypes}} \text{AUC}(k, t), \quad (7)$$

278 where *subtypes* is a set of subtype categories ($C = 5$: Normal, Luminal A, Luminal B, Basal and Her2), and $\text{AUC}(k, t)$ is the
279 k -th and subtype t 's AUC.

280 **Subtype analysis evaluation**

281 According to IPA help and support pages, the p-value is calculated using the right-tailed Fisher Exact Test. Thus the p-value for
282 a pathway is calculated by considering:

- 283 1. the number of genes that participate in that pathway; and
- 284 2. the total number of genes in the QIAGEN Knowledge Base that are known to be associated with that pathway.

285 **Data availability**

286 The authors declare that the main data supporting the findings of this study are publicly available at <https://doi.org/10.5281/zenodo.4769028>. Extra data are available from the corresponding author upon request.

288 **References**

- 289 1. Reis-Filho, J. S. & Pusztai, L. Gene expression profiling in breast cancer: classification, prognostication, and prediction. *The Lancet* **378**, 1812–1823 (2011).
- 290 2. Perou, C. M. *et al.* Molecular portraits of human breast tumours. *Nature* **406**, 747–752 (2000).
- 291 3. Sørlie, T. *et al.* Repeated observation of breast tumor subtypes in independent gene expression data sets. *Proceedings of the National Academy of Sciences* **100**, 8418–8423 (2003).
- 292 4. Parker, J. S. *et al.* Supervised risk predictor of breast cancer based on intrinsic subtypes. *Journal of Clinical Oncology* **27**, 1160–1167 (2009).
- 293 5. Cancer Genome Atlas Network. Comprehensive molecular portraits of human breast tumours. *Nature* **490**, 61–70 (2012).
- 294 6. Weigelt, B. *et al.* Breast cancer molecular profiling with single sample predictors: a retrospective analysis. *The Lancet Oncology* **11**, 339–349 (2010).
- 295 7. Ciriello, G. *et al.* Comprehensive molecular portraits of invasive lobular breast cancer. *Cell* **163**, 506–519 (2015).
- 296 8. Turing, A. M. I.Computing machinery and intelligence. *Mind* **LIX**, 433–460 (1950).
- 297 9. Zeiler, M. D. & Fergus, R. Visualizing and understanding convolutional networks. *Proceedings of the European Conference on Computer Vision*, 818–833 (Springer, 2014).
- 298 10. Zintgraf, L. M., Cohen, T. S., Adel, T. & Welling, M. Visualizing deep neural network decisions: Prediction difference analysis. Preprint at <https://arxiv.org/abs/1702.04595> (2017).
- 299 11. Ribeiro, M. T., Singh, S. & Guestrin, C. Why should I trust you?: Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 1135–1144 (2016).
- 300 12. Simonyan, K., Vedaldi, A. & Zisserman, A. Deep inside convolutional networks: Visualising image classification models and saliency maps. Preprint at <https://arxiv.org/abs/1312.6034> (2013).
- 301 13. Sundararajan, M., Taly, A. & Yan, Q. Gradients of counterfactuals. Preprint at <https://arxiv.org/abs/1611.02639> (2016).
- 302 14. Bach, S. *et al.* On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS One* **10**, e0130140 (2015).
- 303
- 304
- 305
- 306
- 307
- 308
- 309
- 310
- 311

- 312 **15.** Shrikumar, A., Greenside, P. & Kundaje, A. Learning important features through propagating activation differences.
313 Preprint at <https://arxiv.org/abs/1704.02685> (2017).
- 314 **16.** Klambauer, G. *et al.* Self-Normalizing Neural Networks. *Proceedings of the Advances in Neural Information Processing*
315 *Systems* **30**, 972–981 (2017).
- 316 **17.** Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830
317 (2011).
- 318 **18.** McInnes, L., Healy, J. & Melville, J. UMAP: Uniform manifold approximation and projection for dimension reduction.
319 Preprint at <http://arxiv.org/abs/1802.03426> (2018).
- 320 **19.** Krämer, A., Green, J., Pollard Jr, J. & Tugendreich, S. Causal analysis approaches in ingenuity pathway analysis.
321 *Bioinformatics* **30**, 523–530 (2014).
- 322 **20.** Golas, S. B. *et al.* A machine learning model to predict the risk of 30-day readmissions in patients with heart failure: a
323 retrospective analysis of electronic medical records data. *BMC Medical Informatics and Decision Making* **18**, 44 (2018).
- 324 **21.** Breastcancer.org. Molecular subtypes of breast cancer.
325 <https://www.breastcancer.org/symptoms/types/molecular-subtypes> (2021).
- 326 **22.** Otto, T. & Sicinski, P. Cell cycle proteins as promising targets in cancer therapy. *Nature Reviews Cancer* **17**, 93–115
327 (2017).
- 328 **23.** Kumagai, S. *et al.* The PD-1 expression balance between effector and regulatory T cells predicts the clinical efficacy of
329 PD-1 blockade therapies. *Nature Immunology* **21**, 1346–1358 (2020).
- 330 **24.** Kuhn, H. W. & Tucker, A. W. *Contributions to the Theory of Games (AM-28), Volume II.* (Princeton Univ. Press, Princeton,
331 1953).
- 332 **25.** Cancer Genome Atlas Research Network. *et al.* The cancer genome atlas pan-cancer analysis project. *Nature Genetics* **45**,
333 1113–1120 (2013).
- 334 **26.** Goldman, M. J. *et al.* Visualizing and interpreting cancer genomics data via the Xena platform. *Nature Biotechnology* **38**,
335 675–678 (2020).
- 336 **27.** Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome.
337 *BMC Bioinformatics* **12**, 323 (2011).
- 338 **28.** Wang, L., Chu, F. & Xie, W. Accurate cancer classification using expressions of very few genes. *IEEE/ACM Transactions*
339 *on Computational Biology and Bioinformatics* **4**, 40–53 (2007).
- 340 **29.** Bergstra, J., Yamins, D. & Cox, D. Making a science of model search: Hyperparameter optimization in hundreds of
341 dimensions for vision architectures. *Proceedings of the 30th International Conference on International Conference on*
342 *Machine Learning* **28**, I-115–I-123 (2013).
- 343 **30.** Akiba, T., Sano, S., Yanase, T., Ohta, T. & Koyama, M. Optuna: A next-generation hyperparameter optimization framework.
344 *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery and data mining*, 26232631
345 (2019).

346 **Author contributions statement**

347 C. W., K. F., T. S. and Y. Y. have equal contribution. T.S. and Y.Y. designed the architecture of the statistical models. A. O.,
348 C. W., K. F., M. S. and Y. O. designed the biological analysis. C. W., K. F., T. S. and Y. Y. performed the experiments and
349 analyzed the data. All authors discussed the results and contributed to the final manuscript.

350 **Additional information**

351 **Competing interests:** T. S. has a pending US patent 16117260 belonging to Hitachi, Ltd. C. W., Y. Y., K. F., M. S., A. O., and
352 Y. O. declare no potential conflict of interest.

353 **Correspondence** and requests for materials should be addressed to T. S.

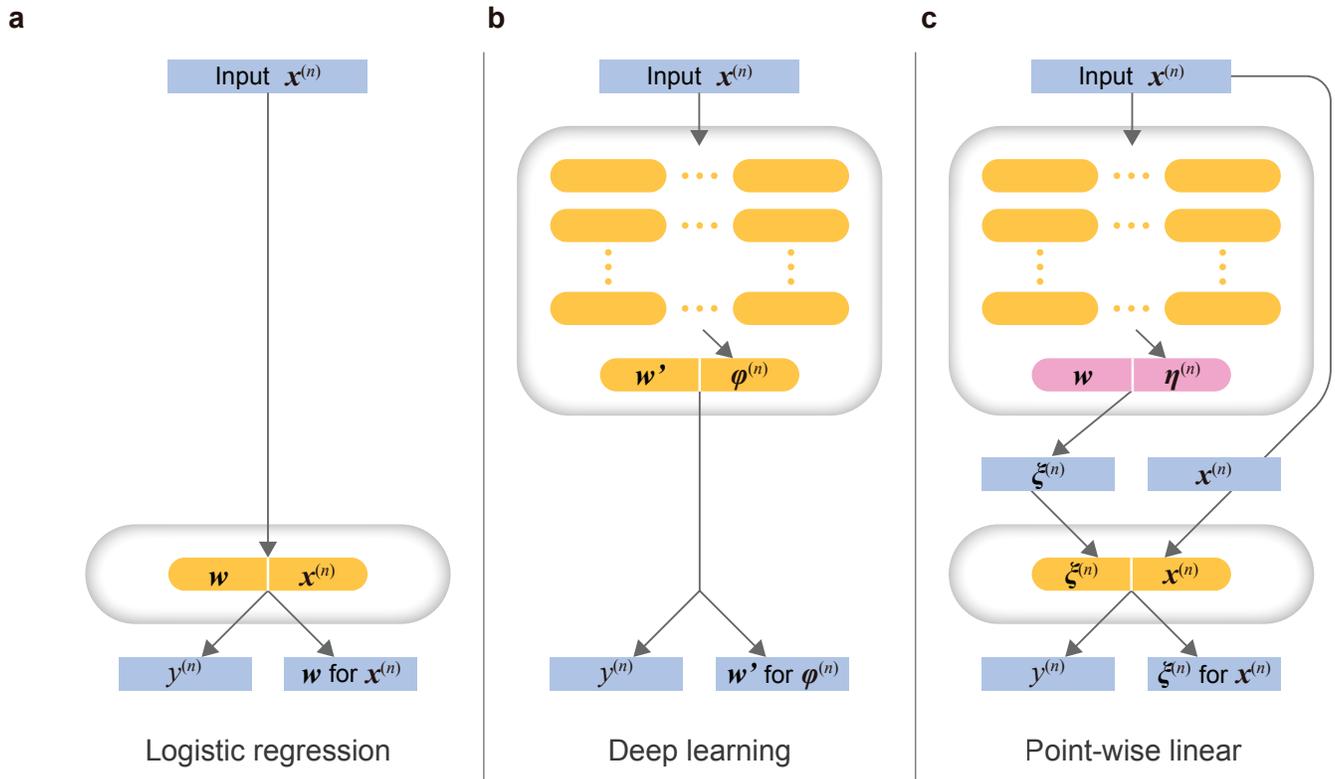


Figure 1. Comparison of network architectures. (a) shows a logistic regression model. $x^{(n)}$ and $y^{(n)}$ are a feature vector and a target value ((n) is sample index), respectively. w is a vector of learning parameters for $x^{(n)}$. (b) shows a fully connected neural network. $\varphi^{(n)}$ and w' are an inner vector and learning parameters, respectively. (c) shows a point-wise linear model. The upper block in (c) is a meta-machine generating a learning parameter $\xi^{(n)}$. The lower block in (c) is a logistic regression model for each feature vector $x^{(n)}$.

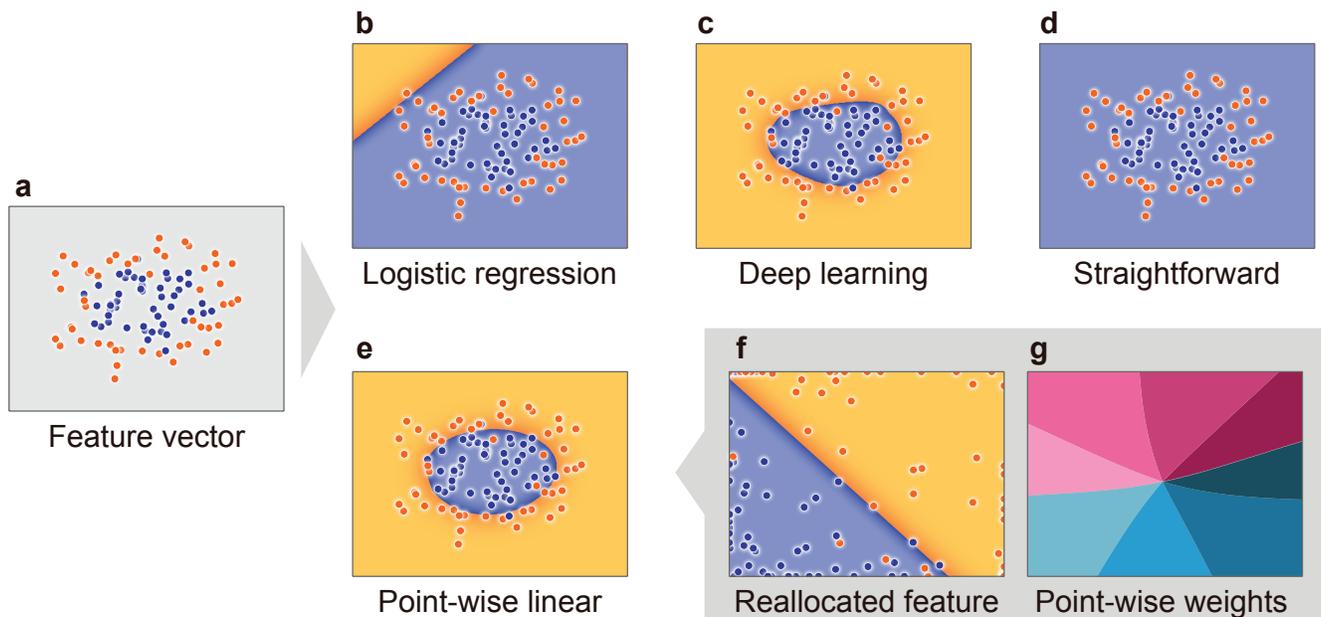


Figure 2. Comparison of learning ability and explainability. (a) is a large circle (orange dots) that contains a smaller circle (blue dots) obtained by *sklearn.datasets.make_circle*. (b) and (c) are the boundaries classified by the logistic regression model and self-normalizing networks (SNNs) model, respectively. (d) and (e) are the boundaries classified by the point-wise linear model of the straightforward manner Eq. (3) and the reallocation function Eq. (4), respectively. (f) is the boundary classified by the reallocated feature vectors ρ . (g) is the arctangent of the angle between the horizontal and vertical elements of the weight vector $\xi^{(n)}$. The weight vector smoothly changes for each data sample.

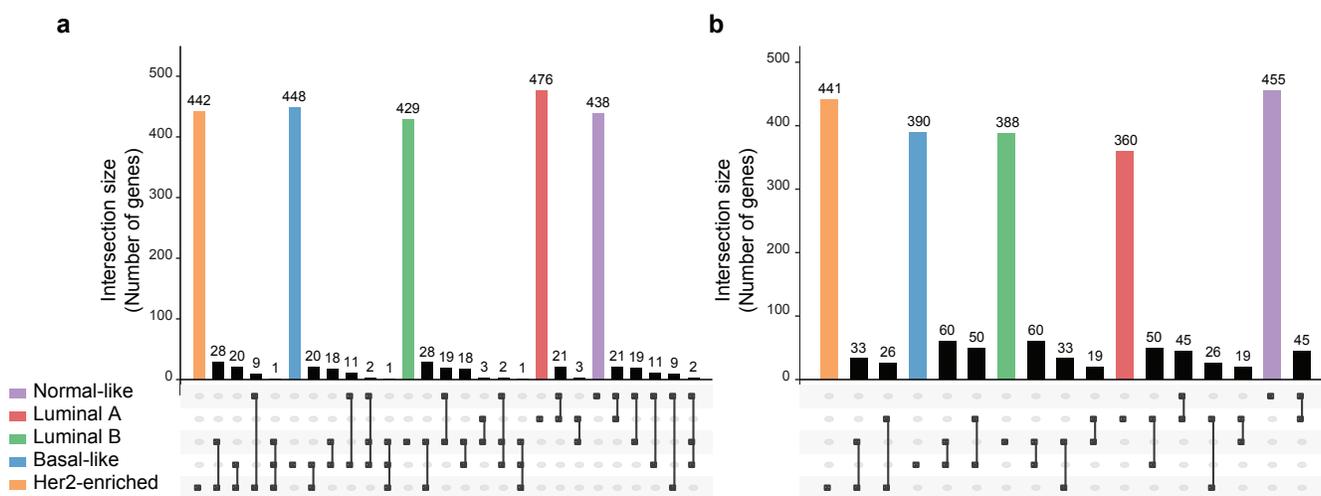


Figure 3. The intersections of the top 500 gene sets calculated by the importance analysis among TCGA subtypes. Left (a) and right (b) panels are for the point-wise linear and logistic regression models, respectively.

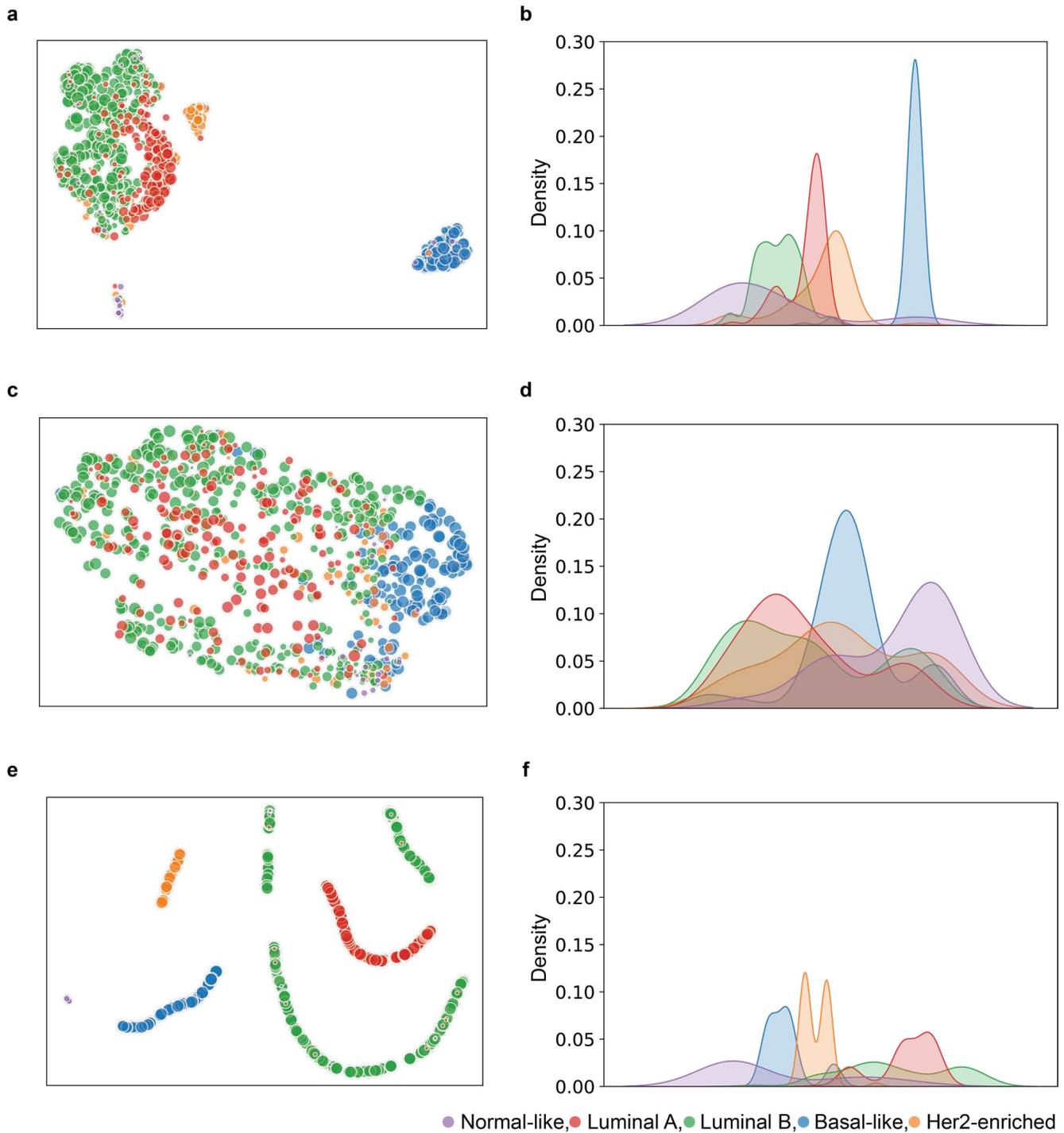


Figure 4. The embeddings of deep learning models projected by UMAP. Top (a), (b), middle (c), (d) and bottom (e), (f) panels are 1D and 2D embeddings of RV η , RF ρ and SNNs, respectively.

Table 1. Number of samples for each breast cancer subtype (Total = 810).

	Normal-like	Luminal A	Luminal B	Basal-like	Her2-enriched
Number of samples	22	406	185	131	66

Table 2. Mean AUCs for 10-fold CV Results.

	Training	Test					
	All	All	Normal-like	Luminal A	Luminal B	Basal-like	Her2-enriched
Logistic regression	0.851	0.853	0.883	0.805	0.746	0.985	0.845
SNNs	0.772	0.745	0.644	0.752	0.615	0.972	0.742
Point-wise linear	0.859	0.862	0.879	0.817	0.779	0.986	0.848

Table 3. Commonalities of the specific feature genes.

		Point-wise linear				
		Normal-like	Luminal A	Luminal B	Basal-like	Her2-enriched
Logistic	Normal-like	37	21	0	0	0
	Luminal A	10	43	39	3	0
	Luminal B	4	0	117	7	4
	Basal-like	2	60	0	82	3
	Her2-enriched	6	3	3	3	206

Logistic: logistic regression.

Table 4. Enriched pathways that the top 500 genes' mRNA expression level were associated with in the 1D embeddings of RV η , RF ρ and the inner vector of SNNs in UMAP.

Canonical Pathways in IPA	$-\log_{10}(p\text{-value})$		
	RV η	RF ρ	SNNs
Cell Cycle: G2/M DNA Damage Checkpoint Regulation	10.57	0.00	0.00
Estrogen-mediated S-phase Entry	10.07	0.36	2.64
Mitotic Roles of Polo-like Kinase	9.91	0.00	0.37
Cell Cycle Control of Chromosomal Replication	9.77	0.00	0.00
Cyclins and Cell Cycle Regulation	8.67	0.28	1.00
Role of CHK Proteins in Cell Cycle Checkpoint Control	8.54	0.89	0.45
Role of BRCA1 in DNA Damage Response	6.81	1.01	0.00
Cell Cycle: G1/S Checkpoint Regulation	6.70	0.00	1.81
tRNA Charging	1.28	4.73	0.24
ErbB Signaling	0.00	1.26	3.67

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryTable8hyperparametersSNNs.xlsx](#)
- [SupplementaryInformations.docx](#)
- [SupplementaryTable10loci.docx](#)
- [SupplementaryAppendix1.pdf](#)
- [SupplementaryTable11IPA.xlsx](#)
- [SupplementaryTable9hyperparameterspointwiselinear.xlsx](#)