
Exploiting Negative Samples: A Catalyst for Cohort Discovery in Healthcare Analytics

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Healthcare analytics, particularly binary diagnosis or prognosis problems, present
2 unique challenges due to the inherent asymmetry between positive and negative
3 samples. While positive samples, representing patients who develop a disease,
4 are defined through rigorous medical criteria, negative samples are defined in an
5 open-ended manner, resulting in a vast potential set. Despite this fundamental
6 asymmetry, previous research has underexplored the role of negative samples,
7 possibly due to the enormous challenge of investigating an infinitely large negative
8 sample space. To bridge this gap, we propose an approach to facilitate cohort
9 discovery within negative samples, which could yield valuable insights into the
10 studied disease, as well as its comorbidity and complications. We measure each
11 sample’s contribution using data Shapley values and construct the Negative Sample
12 Shapley Field to model the distribution of all negative samples. Then we transform
13 this field via manifold learning, preserving the data structure information while
14 imposing an isotropy constraint in data Shapley values. Within this transformed
15 space, we identify cohorts of medical interest through density-based clustering. We
16 empirically evaluate the effectiveness of our approach on our hospital’s electronic
17 medical records. The medical insights revealed in the discovered cohorts are
18 validated by clinicians, which affirms the medical value of our proposal in unveiling
19 meaningful insights consistent with existing domain knowledge, thereby bolstering
20 medical research and well-informed clinical decision-making.

21 1 Introduction

22 Healthcare analytics leverages diverse healthcare data sources to perform many analytic tasks in-
23 cluding diagnosis [28] and prognosis [35]. Electronic Medical Records (EMR) are perhaps the
24 most important of these data sources, since they play a crucial role in recording patients’ essential
25 information and providing a comprehensive view of their health conditions. The recently increasing
26 availability of EMR data has spawned the development of healthcare analytics models for effective
27 patient management and medical resource allocation.

28 Without loss of generality, let us delve into a diagnosis or prognosis problem of predicting whether a
29 patient has developed/will develop a certain disease based on the EMR data. This problem is a binary
30 classification, where patients who develop the disease are “positive samples”, while those who do
31 not are “negative samples”. Notably, we identify the unique nature of such binary classifications in
32 healthcare analytics, as compared to traditional classification tasks. For instance, when classifying
33 cats vs. dogs, both positive and negative samples are based on objective facts. However, in healthcare
34 analytics, positive samples are defined according to rigorous medical criteria, based on medical
35 theories and experience. Contrarily, negative samples are defined in an unrestricted manner, as the
36 complementary set of the positive samples. Consequently, the set of negative samples may encompass
37 a vast number of diverse individuals who are outside the scope of the studied disease or who are
38 healthy. This leads to an inherent asymmetry between positive and negative samples, as positive
39 samples are well-defined and bounded, while negative samples are diverse and open-ended.

40 Despite such fundamental asymmetry in healthcare analytics, previous research has not adequately
41 addressed the role of negative samples. One potential reason for this research gap is the enormous
42 challenge posed by investigating an infinitely large negative sample space, which cannot be easily
43 addressed using existing approaches, e.g., it could be difficult to understand why general healthy
44 individuals do not develop a disease. Nonetheless, it is crucial to probe into negative samples for a
45 more comprehensive investigation of the studied disease. Although it may not have developed in
46 these samples, some may exhibit similar symptoms or even develop related conditions such as its
47 comorbidity or complications. Hence, these negative samples are in urgent need of close medical
48 attention, as they provide an opportunity for clinicians to gain a deeper understanding of the studied
49 disease, leading to more accurate and comprehensive diagnoses, prognoses, and treatment plans.

50 In this paper, we aim to address the gap by exploring negative samples in healthcare analytics.
51 Given the diversity of negative samples, it may not be meaningful to consider them all as one
52 “group”. Instead, we examine the underlying distribution of negative samples to automatically identify
53 medically insightful groups of patients with shared characteristics, referred to as “cohorts” [32, 49].
54 Such cohort discovery among negative samples can provide fresh insights to clinicians on the
55 studied disease, e.g., comprehending the factors contributing to the absence of the disease and the
56 development of related conditions.

57 As front-line clinicians and medical researchers, we bring a unique perspective to guide our method-
58 ology design in effectively discovering cohorts among negative samples. In Sec. 3, we elaborate
59 on our approach with three components. Firstly, we propose to quantify each sample’s contribution
60 to the prediction task using data Shapley values [38, 12]. We then construct the Negative Sample
61 Shapley Field, an inherently existing scalar field describing the distribution and characteristics of all
62 negative samples (Sec. 3.1). Secondly, to effectively discover cohorts, we transform the original field
63 by manifold learning [3] while preserving the original data structure information and ensuring that
64 changes in data Shapley values are isotropic in all orientations (Sec. 3.2). Thirdly, in the transformed
65 manifold space, we identify densely-connected clusters among the negative samples with high data
66 Shapley values through DBSCAN (Sec. 3.3). These clusters help us locate “hot zones”, which are our
67 desired cohorts to discover, exhibiting similar medical characteristics with high data Shapley values.

68 **Our contributions are summarized below:** (i) We bridge the research gap caused by the asymmetry
69 between positive and negative samples in healthcare analytics by exploring negative samples for
70 cohort discovery. (ii) We propose an innovative approach for effective cohort discovery: constructing
71 the Negative Sample Shapley Field, transforming the field by manifold learning with structure
72 preservation and isotropy constraint, and discovering cohorts in the manifold space via DBSCAN.
73 (iii) We empirically evaluate the effectiveness of our approach using our hospital’s EMR (Sec. 4).
74 The experimental results validate the efficacy of each component and demonstrate the capability of
75 our approach for cohort discovery, unveiling meaningful insights that align with existing domain
76 knowledge and have been verified by clinicians. These findings have the potential to benefit medical
77 practitioners by facilitating medical research and clinical decision-making in healthcare delivery.

78 2 Problem and Our Solution

79 **Distinctiveness of negative samples and the unbounded negative sample space.** Let us take
80 hospital-acquired acute kidney injury (AKI), a disease we strive to handle in practice, as an example.
81 AKI is defined according to KDIGO criteria [19] based on a lab test, serum creatinine (sCr). The
82 disease definition has two criteria: absolute AKI and relative AKI. Absolute AKI criterion is met
83 when sCr exhibits a rise exceeding 26.5 $\mu\text{mol/L}$ within the last two days, whereas relative AKI is
84 defined by a rise of sCr 1.5 times or higher over the lowest sCr value within 7 days. In this AKI
85 prediction task, we aim to predict if a patient will develop AKI in the near future. A positive sample
86 is a patient who meets the stringent criteria above, and hence, has a closed definition, whereas a
87 negative sample has an open definition without restrictions. Hence, negative samples in nature form
88 an unbounded space, demonstrating an asymmetry compared to positive samples.

89 **Construction of the Negative Sample Shapley Field for cohort discovery.** To facilitate the analysis
90 of negative samples, we need to investigate their distribution and identify those that are most relevant
91 to the prediction task (e.g., AKI prediction task above) and hence worth exploring. In this regard,
92 we propose to measure the valuation of each negative sample to the task by its data Shapley value.
93 Based on such valuations, we construct a scalar field, the Negative Sample Shapley Field, in which
94 each point is a negative sample, and the point’s value is its data Shapley value. This field depicts the

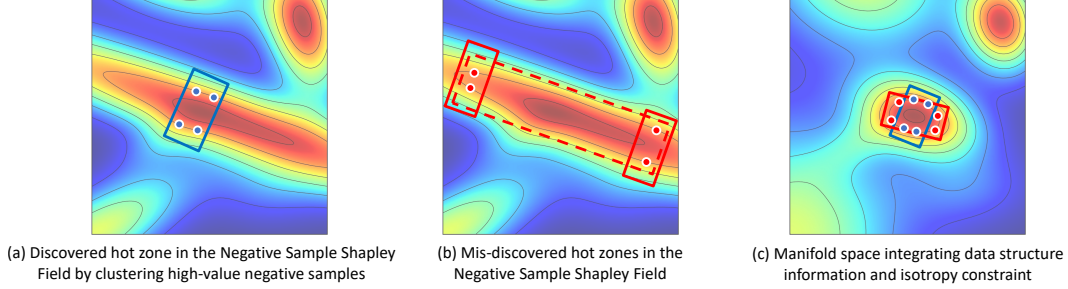


Figure 1: Discovery of hot zones in the Negative Sample Shapley Field.

95 distribution and characteristics of negative samples (see Figure 1(a) for an example). We define “**hot**
 96 **zones**” in this field, identified by points with high data Shapley values, as “**cohorts**”. Our objective
 97 is to automatically detect these cohorts, revealing medically meaningful patterns.

98 **Cohort discovery via manifold learning and density-based clustering.** We note that the vast
 99 number of negative samples renders an exhaustive search infeasible. Although the Negative Sample
 100 Shapley Field is continuously differentiable, the high computational overhead makes it intractable to
 101 find local optima via gradient descent. To overcome this obstacle, we make the assumption that a
 102 subset of negative samples collected in clinical practice carries significant medical value, e.g., patients
 103 who visit hospitals for examinations but do not develop the disease. We posit that these real-world
 104 negative samples should be proximate to our desired hot zones in the space and can effectively sample
 105 our hot zone boundaries, which are hence of medical interest.

106 In Figure 1, we exemplify how to discover hot zones in the Negative Sample Shapley Field. Figure 1(a)
 107 and (b) demonstrate four points situated on the same contour line, indicating their inclusion in the
 108 same hot zone. However, only the former case yields the expected discovered cohort, while the latter
 109 leads to mis-discovery. This highlights that the originally constructed Negative Sample Shapley
 110 Field is suboptimal for cohort discovery among negative samples, due to its anisotropy in data
 111 Shapley values. To overcome this issue, we propose a manifold learning approach. Specifically,
 112 we leverage manifold learning to reduce the dimensionality of the raw sparse EMR data to derive
 113 compact representations that not only preserve the underlying data structure information but also
 114 benefit subsequent spatial clustering analysis. Further, we introduce an isotropy constraint to ensure
 115 uniform changes in data Shapley values across all orientations, which prevents the mis-discovery
 116 as in Figure 1(b). This transformed space, integrating data structure information and the isotropy
 117 constraint, is more suitable for subsequent cohort discovery as illustrated in Figure 1(c).

118 Our objective is then to identify medically meaningful cohorts, specifically dense regions formed by
 119 negative samples with high data Shapley values in the manifold space. We set a data Shapley value
 120 threshold to extract negative samples with high values and employ the DBSCAN algorithm to detect
 121 the hot zones among them. The derived cohorts could shed light on the studied disease, its related
 122 comorbidity, and complications, thereby empowering clinicians in practical healthcare delivery.

123 3 Methodology

124 3.1 Negative Sample Shapley Field Construction

125 Given EMR data $\mathcal{D} = \{d_i\}$, where d_i is a sample with $i \in \{0, \dots, N - 1\}$ and N denotes the
 126 total sample number. We focus on binary classification, and each d_i consists of input features and a
 127 binary label. To investigate negative samples for cohort discovery, we divide \mathcal{D} into \mathcal{D}^+ and \mathcal{D}^- ,
 128 representing positive and negative samples. We denote $\mathcal{D}^- = \{d_i^-\}$, where d_i^- is a negative sample
 129 with $i \in \{0, \dots, N^- - 1\}$ and N^- is the negative sample number.

130 Each negative sample $d_i^- = (\mathbf{x}_i, y_i)$ comprises the input features \mathbf{x}_i and its corresponding binary
 131 label y_i . Our objective is to measure the value of each negative sample by quantifying its contribution
 132 to the prediction performance, which we refer to as data valuation. Data Shapley value [12], stemming
 133 from Shapley value in cooperative game theory, has made significant advances in data valuation [38],
 134 which inspires our proposal to calculate the data Shapley value of each negative sample as its value.
 135 Specifically, let F denote the prediction model and suppose we are interested in evaluating F 's

136 performance on a subset of negative samples $\mathcal{Q} \subseteq \mathcal{D}^-$, along with all the positive samples \mathcal{D}^+ . We
 137 define M as the performance metric function, and then $M(\mathcal{D}^+ \cup \mathcal{Q}, F)$ is the performance achieved
 138 on the combined set of \mathcal{D}^+ and \mathcal{Q} . We define s_i as the data Shapley value for the negative sample
 139 d_i^- . s_i satisfies three properties of Shapley values: (i) null player, (ii) symmetry, and (iii) linearity,
 140 which are the essential properties of an equitable data valuation [12]. We calculate s_i as follows.

141 **Proposition 1** *The data Shapley value s_i for a negative sample d_i^- is given by:*

$$s_i = H \sum_{\mathcal{Q} \subseteq \mathcal{D}^- - \{d_i^-\}} \frac{M(\mathcal{D}^+ \cup \mathcal{Q} \cup \{d_i^-\}, F) - M(\mathcal{D}^+ \cup \mathcal{Q}, F)}{\binom{N^- - 1}{|\mathcal{Q}|}} \quad (1)$$

142 where H is a constant and the summation is taken over all subsets of negative samples, except d_i^- .

143 As the computation of data Shapley value for negative samples has exponential complexity, we
 144 further employ Monte Carlo permutation sampling for approximation [6]. Let Π represent a uniform
 145 distribution of all the permutations among \mathcal{D}^- , s_i can be approximated as the following expectation:

$$s_i = E_{\pi \sim \Pi} [M(\mathcal{D}^+ \cup A_\pi^{d_i^-} \cup \{d_i^-\}, F) - M(\mathcal{D}^+ \cup A_\pi^{d_i^-}, F)] \quad (2)$$

146 where $A_\pi^{d_i^-}$ denotes all the negative samples before d_i^- in a permutation π . By repeating this
 147 approximation, we can derive the estimated data Shapley value s_i efficiently. After computing the
 148 data Shapley value of each negative sample, we define the Negative Sample Shapley Field below.

149 **Definition 1** (*Negative Sample Shapley Field*) *We define the Negative Sample Shapley Field \mathcal{S} as an*
 150 *inherently existing scalar field representing the distribution of data Shapley values across all negative*
 151 *samples in space. In this field, each point denotes a negative sample d_i^- and is associated with its*
 152 *data Shapley value s_i . Therefore, \mathcal{S} is a mathematical function that maps the input of each negative*
 153 *sample to its corresponding data Shapley value: $\mathbf{x}_i \mapsto s_i$.*

154 With this field \mathcal{S} constructed, our goal of cohort discovery within negative samples can be reframed
 155 as the task of identifying “hot zones” - grouped regions within \mathcal{S} exhibiting high data Shapley values.

156 3.2 Manifold Learning with Structure Preservation and Isotropy Constraint

157 As in Figure 1(a) and (b), although we hope to detect a similarly clustered cohort in the Negative
 158 Sample Shapley Field in both scenarios, the anisotropic nature of the space, i.e., the non-uniform
 159 distribution of negative samples with similar data Shapley values, present significant challenges. To
 160 mitigate these challenges, we propose to employ manifold learning [3] to transform the original space
 161 \mathcal{S} into a new geometric space \mathcal{S}' . As elaborated in Sec. 2, to avoid mis-discovery such as Figure 1(b),
 162 we should simultaneously preserve the underlying structural information in the data while imposing
 163 an isotropy constraint on the data Shapley values in \mathcal{S}' . The resulting \mathcal{S}' will be more amenable to
 164 effective cohort discovery, enabling us to identify medically relevant cohorts more accurately.

165 We employ a stacked denoising autoencoder (SDAE) [44] as the backbone model for manifold
 166 learning and integrate the isotropy constraint while preserving the data structure information in \mathbf{x}_i .
 167 Autoencoders (AE) [23, 22] are well-known for capturing data structures by reconstructing input data.
 168 Denoising autoencoders (DAE) [43] are further developed to enhance the learned representations
 169 with the capability of handling input data corruption. By stacking multiple layers of DAE, SDAE can
 170 abstract higher-level robust representations. The model architecture is illustrated in Figure 2.

171 Consider an SDAE consisting of K DAEs. For the k -th DAE ($k \in \{0, \dots, K-1\}$), the encoder takes
 172 $\mathbf{h}_i^{(k)}$ as input, where $\mathbf{h}_i^{(0)} = \mathbf{x}_i$ corresponds to the original input. We define $\tilde{\mathbf{h}}_i^{(k)}$ as the corrupted
 173 version of $\mathbf{h}_i^{(k)}$ with masking noise generated by a stochastic mapping, $\tilde{\mathbf{h}}_i^{(k)} \sim g_{\mathcal{D}}(\tilde{\mathbf{h}}_i^{(k)} | \mathbf{h}_i^{(k)})$, which
 174 randomly sets a fraction of the elements of $\mathbf{h}_i^{(k)}$ to 0. The encoder transforms the corrupted $\tilde{\mathbf{h}}_i^{(k)}$
 175 into an abstract representation $\hat{\mathbf{h}}_i^{(k+1)}$, which is then used by the decoder to recover the uncorrupted
 176 $\mathbf{h}_i^{(k)}$. This process equips the DAE with the capability of extracting useful information for denoising,
 177 which is crucial for healthcare analytics, due to missing data and noise in real-world EMR data [26].

178 **Encoder of the k -th DAE.** The encoder of the k -th DAE transforms the corrupted representation
 179 using an affine transformation followed by a non-linear activation function:

$$\hat{\mathbf{h}}_i^{(k+1)} = f_\theta^{(k+1)}(\tilde{\mathbf{h}}_i^{(k)}) = \sigma(\mathbf{W}_\theta^{(k+1)} \tilde{\mathbf{h}}_i^{(k)} + \mathbf{b}_\theta^{(k+1)}) \quad (3)$$

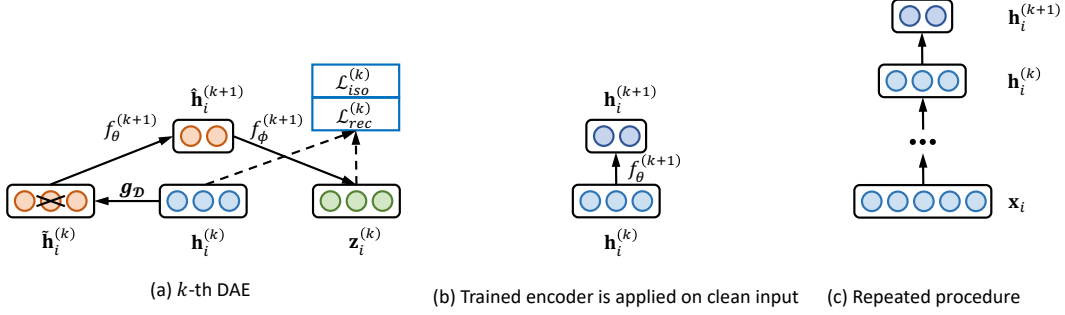


Figure 2: Model architecture of SDAE-based manifold learning.

180 where $f_\theta^{(k+1)}(\cdot)$ is the encoder with $\mathbf{W}_\theta^{(k+1)}$ and $\mathbf{b}_\theta^{(k+1)}$ as the weight matrix and bias vector,
 181 respectively. The rectified linear unit (ReLU) activation function $\sigma(\cdot)$ is used for non-linearity.

182 **Decoder of the k -th DAE.** The derived abstract representation $\hat{\mathbf{h}}_i^{(k+1)}$ is subsequently mapped back
 183 to the original space in the decoder, with the aim of recovering the uncorrupted representation:

$$\mathbf{z}_i^{(k)} = f_\phi^{(k+1)}(\hat{\mathbf{h}}_i^{(k+1)}) = \sigma(\mathbf{W}_\phi^{(k+1)} \hat{\mathbf{h}}_i^{(k+1)} + \mathbf{b}_\phi^{(k+1)}) \quad (4)$$

184 where $f_\phi^{(k+1)}(\cdot)$ is the decoder of the k -th DAE, with $\mathbf{W}_\phi^{(k+1)}$, $\mathbf{b}_\phi^{(k+1)}$ and the ReLU activation.

185 **Structure Preservation.** To attain a stable and robust abstract representation that is resilient to data
 186 corruption, it is crucial to recover the uncorrupted representation as accurately as possible. To achieve
 187 this, we adopt a reconstruction loss that preserves the data structure information. Given a batch of
 188 negative samples \mathcal{B} , the reconstruction loss for this batch is:

$$\mathcal{L}_{rec}^{(k)} = \sum_{i \in \mathcal{B}} \|\mathbf{h}_i^{(k)} - \mathbf{z}_i^{(k)}\|^2 \quad (5)$$

189 **Isotropy Constraint.** In addition to the reconstruction loss, it is essential to enforce an isotropy
 190 constraint to ensure that data Shapley value changes are uniform across orientations. To achieve this,
 191 we introduce a penalty that accounts for the change in data Shapley values relative to the Euclidean
 192 distance between two samples:

$$\mathcal{L}_{iso}^{(k)} = \sum_{i,j \in \mathcal{B}} \left(\frac{s_j - s_i}{\mu_{ij}} \right)^2 \quad (6)$$

193 where i, j are two samples with s_i, s_j as their data Shapley values, μ_{ij} as the distance between $\hat{\mathbf{h}}_i^{(k+1)}$
 194 and $\hat{\mathbf{h}}_j^{(k+1)}$ derived from the encoder. The overall loss is then a weighted sum of the reconstruction loss
 195 and the isotropy penalty, jointly integrating the structural information and the isotropy information:

$$\mathcal{L}^{(k)} = -\frac{1}{|\mathcal{B}|} (\omega_{rec} \mathcal{L}_{rec}^{(k)} + \omega_{iso} \mathcal{L}_{iso}^{(k)}) \quad (7)$$

196 The weights ω_{rec} and ω_{iso} are introduced to address the issue of the two loss terms being on different
 197 scales. This ensures that both losses are decreased at similar rates, leading to a better balance
 198 between the optimization objectives [14, 29]. Specifically, the weights are set to the ratio between the
 199 respective loss in the current iteration (t) and the loss in the previous iteration ($t-1$):

$$\omega_{rec} = \mathcal{L}_{rec}^{(k)}(t) / \mathcal{L}_{rec}^{(k)}(t-1), \quad \omega_{iso} = \mathcal{L}_{iso}^{(k)}(t) / \mathcal{L}_{iso}^{(k)}(t-1) \quad (8)$$

200 We have introduced how to learn the k -th DAE using the loss function in Equation 7, as shown
 201 in Figure 2(a). The corrupted input is only used during the initial training to learn robust feature
 202 extractors. After the encoder $f_\theta^{(k+1)}(\cdot)$ is trained, it will be applied to the clean input as in Figure 2(b):

$$\mathbf{h}_i^{(k+1)} = f_\theta^{(k+1)}(\mathbf{h}_i^{(k)}) = \sigma(\mathbf{W}_\theta^{(k+1)} \mathbf{h}_i^{(k)} + \mathbf{b}_\theta^{(k+1)}) \quad (9)$$

203 $\mathbf{h}_i^{(k+1)}$ will be used as input for the $(k+1)$ -th DAE, as in Figure 2(c), to continue the repeated training
 204 process. When the last DAE, i.e., $(K-1)$ -th DAE, is trained, we obtain the encoded representation
 205 $\mathbf{h}_i^{(K)}$ in the manifold space \mathcal{S}' , which preserves the data structure information in \mathbf{x}_i and integrates
 206 the desired isotropy constraint. $\mathbf{h}_i^{(K)}$ will serve as input for subsequent medical cohort discovery.

207 **3.3 Cohort Discovery Among High Data Shapley Value Negative Samples**

208 We proceed to perform cohort discovery in the encoded manifold space \mathcal{S}' , where each negative
209 sample’s input \mathbf{x}_i is transformed into $\mathbf{h}_i^{(K)}$. We begin by setting a threshold value τ to filter out
210 negative samples with data Shapley values below τ , which focuses our analysis on negative samples
211 with high data Shapley values, i.e., high contributions to the prediction task. Among the remaining
212 negative samples with high data Shapley values, we target to detect the hot zones in \mathcal{S}' , which may
213 represent medically meaningful cohorts of arbitrary shape.

214 To achieve this, we employ DBSCAN, short for density-based spatial clustering of applications with
215 noise [9, 10, 39] on such samples. The core idea of DBSCAN is to group samples that are close
216 to each other in the manifold space \mathcal{S}' into clusters, which could locate potential cohorts, whereas
217 treating the remaining samples as noise or outliers. DBSCAN has three main steps: (i) identify the
218 points within each point’s ε -neighborhood and determine the “core points” with over P_{min} neighbors;
219 (ii) detect the connected components of the core points in the neighbor graph, disregarding any non-
220 core points; (iii) assign each non-core point to the clusters which are the ε -neighborhood of the point;
221 otherwise, label the point as noise. This process results in a set of clusters $\{C_1, C_2, \dots, C_R\}$ and a
222 set of noisy samples Ψ . Given the clusters, we define cohorts as follows.

223 **Definition 2 (Cohorts)** For a dense cluster C_r identified by the DBSCAN algorithm, we consider
224 each of its core points and define a spherical space with the core point as its center and ε as its
225 radius. The joint space of all such spherical spaces is the cohort we aim to discover from this cluster.

226 These discovered cohorts provide a promising avenue for further exploration of medically meaningful
227 patterns in EMR data analytics, potentially revealing important insights.

228 **4 Experimental Evaluation**

229 We evaluate our proposal using our hospital’s EMR data, on which we utilize 709 lab tests to predict
230 whether a patient will develop AKI in each admission in two days (as defined in Section 2). In total,
231 we receive 20,732 admissions, of which 911 develop AKI. We partition the dataset into training
232 data (90%) and testing data (10%). We employ the logistic regression model to compute the data
233 Shapley value for each negative sample as detailed in Section 3.1, using the area under the ROC
234 curve (AUC) as the evaluation metric, and perform Monte Carlo permutation sampling 100,000 times
235 with early stopping. For the manifold learning step, we utilize an SDAE comprising 3 DAEs. The
236 709-dimension inputs are transformed using encoders with dimensions 256, 128, and 64, respectively.

237 **4.1 Cohort Discovery in Clinical Validation**

238 We present the cohort discovery results on our dataset in Figure 3, where we first display the data
239 Shapley value histogram among all the negative samples in Figure 3(a). It is noteworthy that this
240 histogram can be well fitted by a Gaussian mixture model, consisting of three distinct and interesting
241 components. We next examine each component in detail. The first component on the left represents
242 the negative samples with negative data Shapley values. These samples have a negative impact on
243 the prediction task, meaning that they are detrimental to predicting the AKI occurrence. In prior
244 studies, one generally plausible explanation for the presence of such samples is the existence of
245 mislabeled data [12]. However, for a representative acute disease like AKI, these negative samples
246 are highly likely to be positive samples in the future but have not yet exhibited symptoms of AKI
247 within the monitored time duration. Moving on to the second component in the middle, we observe
248 that its data Shapley values are centered around a mean value close to zero. This implies that these
249 negative samples are generally healthy without any apparent AKI-related symptoms. Notably, these
250 healthy samples constitute a relatively significant portion of the data, which is commonly observed in
251 clinical practice and aligns with our initial expectations. The third component on the right represents
252 negative samples that are particularly valuable for the prediction task and merit special attention in
253 our study. To further investigate these samples, we introduce a separation line between the second
254 and third components, i.e., a threshold 60% to exclude the lower 60% negative samples based on
255 their data Shapley values while retaining the remaining 40% for further analysis. Our focus is on
256 these remaining 40% samples for identifying the hot zones, as illustrated in Figure 1.

257 The distribution of all negative samples, in terms of their data Shapley values in the manifold space,
258 is presented in Figure 3(b). Upon performing DBSCAN on the extracted 40% samples with high
259 data Shapley values (points brighter than dark blue), we identify seven distinct cohorts of interest,

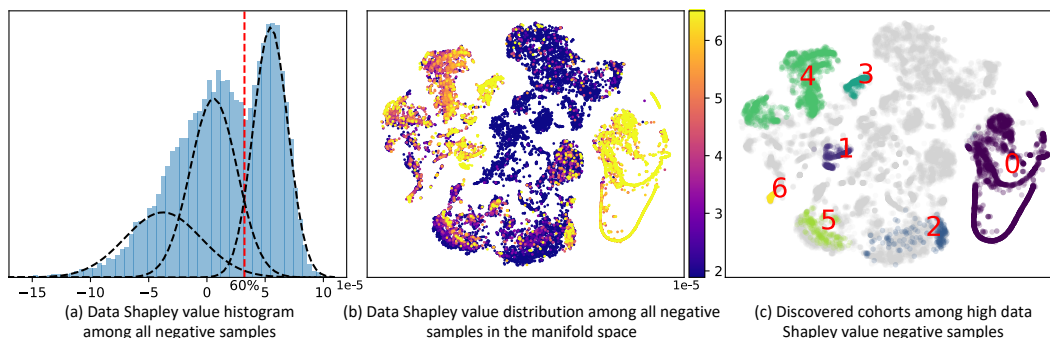


Figure 3: Cohort discovery on our dataset.

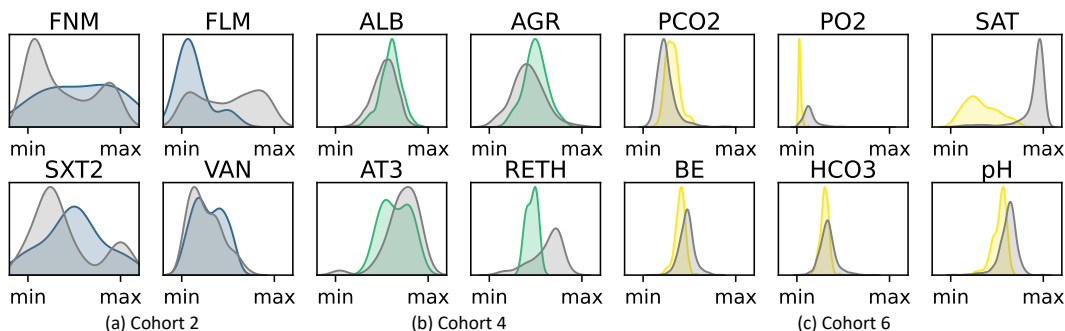


Figure 4: Lab test patterns of discovered Cohorts 2, 4, and 6. In each cohort, the colored region (blue, green, and yellow) represents the lab test value probability density of the samples in the cohort, while the grey region denotes that of all the other samples outside the cohort.

260 which are visually displayed using t-SNE plots in Figure 3(c), in which grey points are either with
 261 low data Shapley values or labeled as noise by DBSCAN. We observe that these discovered cohorts
 262 are distinguishable from one another, potentially corresponding to medically meaningful patterns.

263 4.2 In-depth Analysis of Discovered Cohorts

264 **Cohort 2: inflammatory cohort.** Figure 4(a) indicates a pronounced neutrophil-to-lymphocyte
 265 ratio (NLR) [48] in this patient group, marked by an increase in neutrophils (FNM) and a decrease
 266 in lymphocytes (FLM). This pattern, often tied to infectious, inflammatory, and stress conditions,
 267 suggests an overactive immune response leading to reduced lymphocyte counts [36, 8]. An elevated
 268 NLR, a reliable inflammatory marker, indicates a propensity for invasive infections [16]. Meanwhile,
 269 the levels of Cotrimoxazole (SXT2) and Vancomycin (VAN), both administered to treat infections
 270 including those associated with methicillin-resistant staphylococcus [15], are found to be elevated in
 271 the bodies of these patients. The findings suggest that this patient cohort comprises individuals experi-
 272 encing infections and acute inflammation, and receiving antibiotic treatment. Severe infections can
 273 cause systemic inflammatory response syndrome and kidney injury. Antibiotics like vancomycin can
 274 worsen kidney stress and have nephrotoxic properties [47], potentially leading to kidney dysfunction
 275 during treatment. However, modern medical practice can effectively manage these cases. Infections
 276 are promptly treated with broad-spectrum antibiotics and at appropriate doses within safety limits;
 277 hence, the patients do not develop significant AKI [13].

278 **Cohort 4: hepatic and hematological disorders cohort.** As delineated in Figure 3(c), Cohort 4
 279 exhibits an augmented region and an increased quantity of sampling points, indicative of a more
 280 expansive patient population. A comprehensive analysis of the lab test indicator distribution for
 281 this cohort, portrayed in Figure 4(b), reveals differences in levels of serum proteins. Specifically,
 282 derangements in levels of albumin (ALB) and the albumin-globulin ratio (AGR) signify aberrant
 283 protein synthesis in patients. These may be associated with hepatic dysfunction or hematological
 284 diseases such as myeloma [41, 27]. Hepatic diseases can lead to impaired production of other
 285 proteins such as antithrombin III (AT3) [21]; AT3 may also be lost excessively in nephrotic syndrome

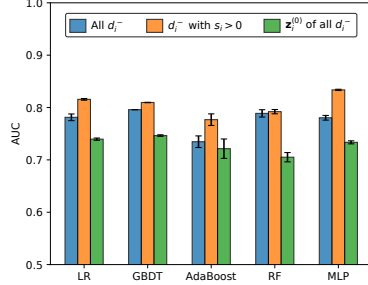


Figure 5: AKI prediction performance of widely adopted classifiers in three different settings.

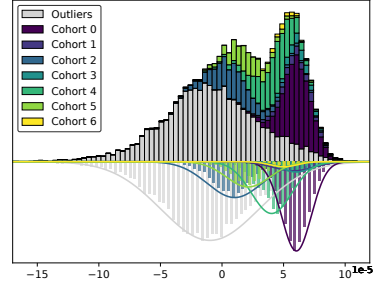


Figure 6: Data Shapley value histogram of the samples within our discovered cohorts.

286 which is a kidney disorder [18], or undergo accelerated consumption in disseminated intravascular
 287 coagulation [34]. Diminished reticulocyte hemoglobin (RETH) is associated with iron deficiency
 288 anemia [2], and could either be linked to hematological disorders or nutritional deficiency. In
 289 addition, imbalances in albumin and globulin may also be associated with dehydration. Therefore,
 290 our observation derived from Cohort 4 may support the pathophysiological relationship that exists
 291 between disorders of the hematological and hepatic systems, which increases the propensity for
 292 kidney disease. Clinicians should exercise vigilance in care when managing these cases.

293 **Cohort 6: respiratory failure and metabolic acidosis cohort.** Figure 4(c) reveals significant
 294 metabolic imbalances in patients, leading to an acid-base imbalance. Specifically, increased carbon
 295 dioxide pressure (PCO₂), reduced oxygen pressure (PO₂), and insufficient blood oxygen saturation
 296 (SAT) suggest respiratory failure [5]. Concurrently, reduced base excess (BE), bicarbonate ion (HCO₃)
 297 levels, and blood pH values hint at metabolic acidosis, indicating possible acute illnesses causing
 298 lactic or ketoacidosis [24]. These results suggest potential severe respiratory complications, such
 299 as advanced pneumonia, heart failure-induced pulmonary edema, or chronic obstructive pulmonary
 300 disease (COPD)[20]. Alternatively, acute conditions like hypoxia, shock, or severe infection could
 301 disrupt aerobic metabolism, leading to anaerobic glucose conversion to lactate, which accumulates in
 302 the bloodstream and causes acidosis. This puts significant strain on the kidneys, potentially resulting
 303 in renal disease symptoms[25]. This cohort of patients under examination does not advance to AKI,
 304 leading to the inference that renal dysfunction may not constitute an end-organ complication. Rather,
 305 this patient cohort appears to exhibit a heightened disposition to respiratory failure.

306 4.3 Validation of Effectiveness for Each Component

307 We validate the effectiveness of each component in our approach for AKI prediction. Specifically,
 308 we evaluate three settings of the negative sample usage in the training data (with positive samples
 309 the same): (i) all d_i^- : use all negative samples; (ii) d_i^- with $s_i > 0$: only use the negative samples
 310 with positive data Shapley values; (iii) $z_i^{(0)}$ of all d_i^- : use the decoded representations from the
 311 SDAE-based manifold learning. $z_i^{(0)}$ is in the same dimension as the raw input but is in the decoding
 312 space after transformation by SDAE. To ensure the credibility of our conclusions across different
 313 settings, we evaluate several widely adopted classifiers: logistic regression (LR), gradient-boosting
 314 decision tree (GBDT), adaptive boosting (AdaBoost), random forest (RF), and multilayer perceptron
 315 (MLP). The experimental results in AUC (mean \pm std) from five repeats are illustrated in Figure 5.

316 **Effectiveness of the Negative Sample Shapley Field.** By comparing settings (i) and (ii), we explore
 317 the effectiveness of our constructed Negative Sample Shapley Field. The results clearly demonstrate
 318 that by removing negative samples with data Shapley values smaller than 0, all the classifiers
 319 exhibit an improvement in AUC. This finding supports the rationale behind our approach of linking
 320 samples of great medical concern with their data Shapley values. Additionally, the effectiveness of
 321 approximating data Shapley values through Monte Carlo permutation sampling is further validated.
 322 Thus, this confirms the efficacy of our constructed Negative Sample Shapley Field.

323 **Effectiveness of Manifold Learning.** By changing the input data from the raw space to the decoder's
 324 output space after our proposed SDAE-based manifold learning (settings (i) vs. (iii)), we observe
 325 a moderate decrease in AUC, approximately 5% in most classifiers. This decrease aligns with our
 326 expectations, as the transformation in SDAE introduces a certain level of information loss. However,

327 the performance degradation remains within an acceptable range. These findings demonstrate that
328 our proposed manifold learning manages to preserve the original data structure information and
329 effectively model the original raw data space, despite a significant reduction in data dimension from
330 709 to 64. Thus, this corroborates our design rationale of employing SDAE for manifold learning
331 with structure preservation and isotropy constraint.

332 **Effectiveness of Cohort Discovery.** We further validate our method’s ability to decompose high
333 data Shapley value samples into distinct, medically relevant cohorts. Figure 6 presents the data
334 Shapley value histogram of our identified cohorts, with the upper part aligned with Figure 3(a)
335 but color-coded by cohort proportion. The lower part shows each cohort’s data Shapley value
336 distribution. We note seven cohorts effectively partition Figure 3(a)’s third component into Gaussian
337 distributions, implying consistent data Shapley values within each cohort. Cohort 2, identified as
338 the inflammatory group, exhibits relatively lower data Shapley values, as immune abnormalities
339 cannot serve as specific features for kidney injury. Conversely, Cohorts 4 and 6, involving critical
340 metabolic systems, display higher data Shapley values, which indicates their significant medical
341 relevance to AKI prediction. These observations confirm the homogeneity within each cohort due
342 to DBSCAN’s detection capability, and similarity in data Shapley values, further substantiating our
343 proposed isotropy constraint in manifold learning. In essence, our approach effectively identifies
344 proximate cohorts with similar data Shapley values, providing valuable medical insights for the
345 prediction task.

346 5 Related Work

347 Shapley value, originally introduced in cooperative game theory [40], offers a solution for the equi-
348 table distribution of a team’s collective value among its individual members [7]. Notable applications
349 of the Shapley value in machine learning encompass data valuation, feature selection, explainable ma-
350 chine learning, etc [38, 12, 46, 31, 31, 30]. Among these applications, data valuation holds particular
351 significance in quantifying the contributions of individual data samples toward predictive models. In
352 this research line, the data Shapley value [12] presents an equitable valuation framework for data value
353 quantification with subsequent research focusing on enhancing computational efficiency [17, 11].

354 Representation learning is a crucial research area contributing to the success of many machine
355 learning algorithms [3]. Among the representation learning methods, manifold learning stands out
356 due to its capability of reducing the dimensionality and visualizing the underlying structure of the
357 data. Traditional manifold learning methods include Isomap [42], locally linear embedding [37], and
358 multi-dimensional scaling [4]. In recent years, AEs have gained significant attention in representation
359 learning, offering efficient and effective representations of unlabeled data. Researchers develop
360 various AE variants for specific application scenarios, e.g., regularized AEs [1], sparse AEs [33],
361 DAEs (denoising AEs) [43]. Specifically, DAEs and their advanced stacked variant SDAEs [44] are
362 highly suitable to tackle EMR data, in which missing and noisy data remains a notorious issue [26].

363 DBSCAN, short for density-based spatial clustering of applications with noise, is introduced to
364 alleviate the burden of parameter selection for users, facilitate the discovery of arbitrarily-shaped
365 clusters, and demonstrate satisfactory efficiency when dealing with large datasets [9, 10, 39].

366 6 Conclusion

367 This paper proposes to examine negative samples for cohort discovery in healthcare analytics, which
368 has not been explored in prior research. In particular, we propose to measure each negative sample’s
369 contribution to the prediction task via its data Shapley value and construct the Negative Sample
370 Shapley Field to model the distribution of all negative samples. To enhance the cohort discovery
371 quality, we transform this original field into an embedded space using manifold learning, incorporating
372 the original data structure information and isotropy constraint. In the transformed space, we manage
373 to identify medically meaningful cohorts within negative samples by DBSCAN. The experiments on
374 our hospital’s EMR data empirically demonstrate the effectiveness of our proposal, and the medical
375 insights derived from our discovered cohorts are validated by clinicians, highlighting the medical
376 value of our approach. Future work includes conducting a long-term validation to further verify the
377 conclusions drawn from cohort discovery. Additionally, more detailed analyses and fine-grained
378 clinical validation are required to explore the detected cohorts that exhibit a hierarchical structure.

379 **References**

- 380 [1] Guillaume Alain and Yoshua Bengio. What regularized auto-encoders learn from the data-
381 generating distribution. *J. Mach. Learn. Res.*, 15(1):3563–3593, 2014.
- 382 [2] Michael Auerbach, Steven J Staffa, and Carlo Brugnara. Using reticulocyte hemoglobin
383 equivalent as a marker for iron deficiency and responsiveness to iron therapy. In *Mayo Clinic*
384 *Proceedings*, volume 96, pages 1510–1519. Elsevier, 2021.
- 385 [3] Yoshua Bengio, Aaron C. Courville, and Pascal Vincent. Representation learning: A review
386 and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, 2013.
- 387 [4] Ingwer Borg and Patrick JF Groenen. *Modern multidimensional scaling: Theory and applica-*
388 *tions*. Springer Science & Business Media, 2005.
- 389 [5] Peter H Breen. Arterial blood gas and ph analysis: clinical approach and interpretation.
390 *Anesthesiology Clinics of North America*, 19(4):885–906, 2001.
- 391 [6] Javier Castro, Daniel Gómez, and Juan Tejada. Polynomial calculation of the shapley value
392 based on sampling. *Comput. Oper. Res.*, 36(5):1726–1730, 2009.
- 393 [7] Georgios Chalkiadakis, Edith Elkind, and Michael J. Wooldridge. *Computational Aspects of*
394 *Cooperative Game Theory*. Synthesis Lectures on Artificial Intelligence and Machine Learning.
395 Morgan & Claypool Publishers, 2011.
- 396 [8] Firdaus S Dhabhar. Enhancing versus suppressive effects of stress on immune function: implica-
397 tions for immunoprotection and immunopathology. *Neuroimmunomodulation*, 16(5):300–317,
398 2009.
- 399 [9] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for
400 discovering clusters in large spatial databases with noise. In *KDD*, pages 226–231. AAAI Press,
401 1996.
- 402 [10] Junhao Gan and Yufei Tao. DBSCAN revisited: Mis-claim, un-fixability, and approximation.
403 In *SIGMOD Conference*, pages 519–530. ACM, 2015.
- 404 [11] Amirata Ghorbani, Michael P. Kim, and James Zou. A distributional framework for data
405 valuation. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 3535–
406 3544. PMLR, 2020.
- 407 [12] Amirata Ghorbani and James Y. Zou. Data shapley: Equitable valuation of data for machine
408 learning. In *ICML*, volume 97 of *Proceedings of Machine Learning Research*, pages 2242–2251.
409 PMLR, 2019.
- 410 [13] Stuart L Goldstein, Theresa Mottes, Kendria Simpson, Cynthia Barclay, Stephen Muething,
411 David B Haslam, and Eric S Kirkendall. A sustained quality improvement program reduces
412 nephrotoxic medication-associated acute kidney injury. *Kidney international*, 90(1):212–221,
413 2016.
- 414 [14] Rick Groenendijk, Sezer Karaoglu, Theo Gevers, and Thomas Mensink. Multi-loss weighting
415 with coefficient of variations. In *WACV*, pages 1468–1477. IEEE, 2021.
- 416 [15] Natasha E Holmes and Benjamin P Howden. What’s new in the treatment of serious mrsa
417 infection? *Current opinion in infectious diseases*, 27(6):471–478, 2014.
- 418 [16] Zhiwei Huang, Zhaoyin Fu, Wujun Huang, and Kegang Huang. Prognostic value of neutrophil-
419 to-lymphocyte ratio in sepsis: A meta-analysis. *The American journal of emergency medicine*,
420 38(3):641–647, 2020.
- 421 [17] Ruoxi Jia, David Dao, Boxin Wang, Frances Ann Hubis, Nick Hynes, Nezihe Merve Gürel,
422 Bo Li, Ce Zhang, Dawn Song, and Costas J. Spanos. Towards efficient data valuation based
423 on the shapley value. In *AISTATS*, volume 89 of *Proceedings of Machine Learning Research*,
424 pages 1167–1176. PMLR, 2019.

- 425 [18] Robert H Kauffmann, Jan J Veltkamp, Nico H Van Tilburg, and Leendert A Van Es. Acquired
426 antithrombin iii deficiency and thrombosis in the nephrotic syndrome. *The American journal of*
427 *medicine*, 65(4):607–613, 1978.
- 428 [19] John A Kellum, Norbert Lameire, Peter Aspelin, Rashad S Barsoum, Emmanuel A Burdmann,
429 Stuart L Goldstein, Charles A Herzog, Michael Joannidis, Andreas Kribben, Andrew S Levey,
430 et al. Kidney disease: improving global outcomes (kdigo) acute kidney injury work group.
431 kdigo clinical practice guideline for acute kidney injury. *Kidney international supplements*,
432 2(1):1–138, 2012.
- 433 [20] Jordan A Kempker, Maria K Abril, Yunyun Chen, Michael R Kramer, Lance A Waller, and
434 Greg S Martin. The epidemiology of respiratory failure in the united states 2002–2017: A serial
435 cross-sectional study. *Critical Care Explorations*, 2(6), 2020.
- 436 [21] E Knot, JW Ten Cate, HR Drijfhout, LH Kahlé, and GN Tytgat. Antithrombin iii metabolism in
437 patients with liver disease. *Journal of clinical pathology*, 37(5):523–530, 1984.
- 438 [22] Mark A Kramer. Nonlinear principal component analysis using autoassociative neural networks.
439 *AIChE journal*, 37(2):233–243, 1991.
- 440 [23] Mark A Kramer. Autoassociative neural networks. *Computers & chemical engineering*,
441 16(4):313–328, 1992.
- 442 [24] Jeffrey A Kraut and Nicolaos E Madias. Metabolic acidosis: pathophysiology, diagnosis and
443 management. *Nature Reviews Nephrology*, 6(5):274–285, 2010.
- 444 [25] Jeffrey A Kraut and Nicolaos E Madias. Lactic acidosis. *New England Journal of Medicine*,
445 371(24):2309–2319, 2014.
- 446 [26] Thomas A Lasko, Joshua C Denny, and Mia A Levy. Computational phenotype discovery
447 using unsupervised feature learning over noisy, sparse, and irregular clinical data. *PloS one*,
448 8(6):e66341, 2013.
- 449 [27] Garrick Edouard Laudin, Peter F Levay, and Buks Coetzer. Globulin fraction and albumin: glob-
450 ulin ratio as a predictor of mortality in a south african multiple myeloma cohort. *International*
451 *Journal of Hematologic Oncology*, 9(3):IJH27, 2020.
- 452 [28] Zachary Chase Lipton, David C. Kale, Charles Elkan, and Randall C. Wetzel. Learning to
453 diagnose with LSTM recurrent neural networks. In *ICLR (Poster)*, 2016.
- 454 [29] Shikun Liu, Edward Johns, and Andrew J. Davison. End-to-end multi-task learning with
455 attention. In *CVPR*, pages 1871–1880. Computer Vision Foundation / IEEE, 2019.
- 456 [30] Zelei Liu, Yuanyuan Chen, Han Yu, Yang Liu, and Lizhen Cui. Gtg-shapley: Efficient and
457 accurate participant contribution evaluation in federated learning. *ACM Trans. Intell. Syst.*
458 *Technol.*, 13(4):60:1–60:21, 2022.
- 459 [31] Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In
460 *NIPS*, pages 4765–4774, 2017.
- 461 [32] Syed S Mahmood, Daniel Levy, Ramachandran S Vasan, and Thomas J Wang. The framingham
462 heart study and the epidemiology of cardiovascular disease: a historical perspective. *The lancet*,
463 383(9921):999–1008, 2014.
- 464 [33] Alireza Makhzani and Brendan J. Frey. k-sparse autoencoders. In *ICLR (Poster)*, 2014.
- 465 [34] Eberhard F Mammen. Antithrombin: its physiological importance and role in dic. In *Seminars*
466 *in thrombosis and hemostasis*, volume 24, pages 19–25. Copyright© 1998 by Thieme Medical
467 Publishers, Inc., 1998.
- 468 [35] DR Mould. Models for disease progression: new approaches and uses. *Clinical Pharmacology*
469 *& Therapeutics*, 92(1):125–131, 2012.
- 470 [36] Carl Nathan. Neutrophils and immunity: challenges and opportunities. *Nature reviews immunol-*
471 *ogy*, 6(3):173–182, 2006.

- 472 [37] Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear
473 embedding. *science*, 290(5500):2323–2326, 2000.
- 474 [38] Benedek Rozemberczki, Lauren Watson, Péter Bayer, Hao-Tsung Yang, Oliver Kiss, Sebastian
475 Nilsson, and Rik Sarkar. The shapley value in machine learning. In *IJCAI*, pages 5572–5579.
476 ijcai.org, 2022.
- 477 [39] Erich Schubert, Jörg Sander, Martin Ester, Hans-Peter Kriegel, and Xiaowei Xu. DBSCAN
478 revisited, revisited: Why and how you should (still) use DBSCAN. *ACM Trans. Database Syst.*,
479 42(3):19:1–19:21, 2017.
- 480 [40] Lloyd S Shapley et al. A value for n-person games. 1953.
- 481 [41] Rosaria Spinella, Rohit Sawhney, and Rajiv Jalan. Albumin in chronic liver disease: structure,
482 functions and therapeutic implications. *Hepatology international*, 10:124–132, 2016.
- 483 [42] Joshua B Tenenbaum, Vin de Silva, and John C Langford. A global geometric framework for
484 nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- 485 [43] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting
486 and composing robust features with denoising autoencoders. In *ICML*, volume 307 of *ACM*
487 *International Conference Proceeding Series*, pages 1096–1103. ACM, 2008.
- 488 [44] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre-Antoine Manzagol.
489 Stacked denoising autoencoders: Learning useful representations in a deep network with a local
490 denoising criterion. *J. Mach. Learn. Res.*, 11:3371–3408, 2010.
- 491 [45] John Burnard West. *Respiratory physiology: the essentials*. Lippincott Williams & Wilkins,
492 2012.
- 493 [46] Brian D. Williamson and Jean Feng. Efficient nonparametric statistical inference on population
494 feature importance using shapley values. In *ICML*, volume 119 of *Proceedings of Machine*
495 *Learning Research*, pages 10282–10291. PMLR, 2020.
- 496 [47] Huizi Wu and Jiaguo Huang. Drug-induced nephrotoxicity: pathogenic mechanisms, biomarkers
497 and prevention strategies. *Current drug metabolism*, 19(7):559–567, 2018.
- 498 [48] R Zahorec et al. Ratio of neutrophil to lymphocyte counts-rapid and simple parameter of
499 systemic inflammation and stress in critically ill. *Bratislavske lekarske listy*, 102(1):5–14, 2001.
- 500 [49] Fei Zhou, Ting Yu, Ronghui Du, Guohui Fan, Ying Liu, Zhibo Liu, Jie Xiang, Yeming Wang,
501 Bin Song, Xiaoying Gu, et al. Clinical course and risk factors for mortality of adult inpatients
502 with covid-19 in wuhan, china: a retrospective cohort study. *The lancet*, 395(10229):1054–1062,
503 2020.