
Discerning Decision-Making Process of Deep Neural Networks with Hierarchical Voting Transformation

Anonymous Author(s)

Abstract

1 Neural network based deep learning techniques have shown great success for
2 numerous applications. While it is expected to understand their intrinsic decision-
3 making processes, these deep neural networks often work in a black-box way. To
4 this end, in this paper, we aim to discern the decision-making processes of neural
5 networks through a hierarchical voting strategy by developing an explainable deep
6 learning model, namely Voting Transformation-based Explainable Neural Net-
7 work (VOTEN). Specifically, instead of relying on massive feature combinations,
8 VOTEN creatively models expressive single-valued voting functions between ex-
9 plicitly modeled latent concepts to achieve high fitting ability. Along this line, we
10 first theoretically analyze the major components of VOTEN and prove the relation-
11 ship and advantages of VOTEN compared with Multi-Layer Perceptron (MLP), the
12 basic structure of deep neural networks. Moreover, we design efficient algorithms
13 to improve the model usability by explicitly showing the decision processes of
14 VOTEN. Finally, extensive experiments on multiple real-world datasets clearly
15 validate the performances and explainability of VOTEN.

16 1 Introduction

17 Neural network based deep learning techniques have attracted great attention from both academia and
18 industry in the past decade. Compared with classic machine learning models, deep neural networks
19 have much higher expressiveness and adaptability for complicated data input, and thus have made
20 tremendous success in various application domains, such as Computational Vision [26, 46], Natural
21 Language Processing [44, 18], and Recommender Systems [29, 21]. Nevertheless, since deep neural
22 networks usually have complicated connections of hidden units, a long-standing challenge is how
23 to decipher what’s inside the black box of models for understanding their intrinsic decision-making
24 processes. Indeed, in many real-world scenarios, such as urban management, finance, and medical
25 treatment, the lack of model explainability makes people less likely to be convinced when the
26 decision-making process of the model is not understandable. This prevents a broader application of
27 deep neural networks.

28 While many research efforts have been made in developing explainable deep learning models [40, 9],
29 most of existing works focus on post-hoc explanation [4, 16], i.e., designing metrics to measure
30 the feature relevance/contribution to the outputs of a trained model. Although these methods have
31 made progress on finding important features, the decision-making process of deep learning models
32 is still not available for users. For example, users often need to guess the reason why a prediction
33 is made from the relevant features instead of understanding the decision-making process. Indeed,
34 understanding the intrinsic decision-making process of deep neural networks is a non-trivial task.
35 A major reason is that deep neural networks usually involve massive feature combinations to gain
36 expressiveness on fitting complicated functions. During this process, the effect of features and
37 hidden units may be largely coupled with each other. This indicates that the decision logic of models
38 is inherently buried in the massive feature combinations. Meanwhile, it is difficult for human to
39 understand the intrinsic modeling process of deep learning in a natural manner. Therefore, a key
40 point on improving the explainability of deep neural networks is to decouple the feature combinations
41 and make the modeling process consistent with human decision process. In this way, the model will
42 have an explicit decision-making process and become human-understandable.

43 To this end, in this paper, we propose an explainable deep learning model, namely Voting
44 Transformation-based Explainable Neural Network (VOTEN). Specifically, VOTEN assumes the

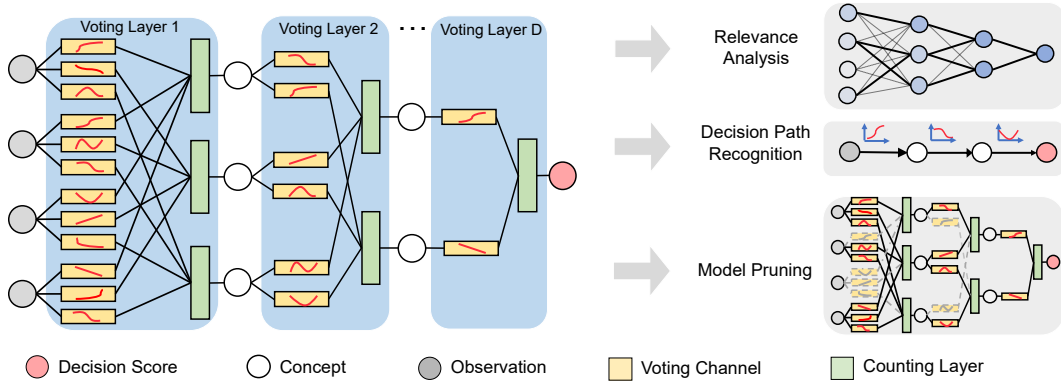


Figure 1: Structure overview of VOTEN.

45 transformation from the input to the output to be a hierarchical voting process. During this process,
 46 lower-level concepts vote for higher-level concepts layer-by-layer in an expressive but explainable
 47 way. Instead of relying on massive feature combinations, VOTEN creatively models expressive
 48 single-valued voting functions between explicitly modeled hidden concepts to gain expressiveness on
 49 fitting complicated functions. This process is explainable for its consistency with human decision-
 50 making process. We first theoretically analyze the major components of VOTEN and prove the
 51 relationship and advantages of VOTEN compared with Multi-Layer Perceptron (MLP), which is the
 52 basic structure of deep neural networks. The results show that MLP can be derived from VOTEN by
 53 using the inexpressive voting functions. Accordingly, we further analyze the effect of inherent votings
 54 and design efficient algorithms for pruning and explaining VOTEN. Finally, we evaluated VOTEN on
 55 multiple real-world datasets with comprehensive experiments. The experimental results demonstrate
 56 that VOTEN generally promotes powerful and explainability deep learning. Specifically, VOTEN (1)
 57 significantly raises prediction performance; (2) exponentially decreases feature combinations; and (3)
 58 supports efficient pruning and effective feature analysis. Meanwhile, we also visualize the intrinsic
 59 decision-making process of VOTEN through case studies, which show VOTEN is explainable and
 60 can discover meaningful latent concepts.

61 2 VOTEN

62 In this section, we introduce the technical details of VOTEN.

63 2.1 Structure

64 The structure overview of VOTEN is shown in Figure 1. Fundamentally, deep neural networks are
 65 difficult to be explained because they work in a totally different way from humans thinking. That
 66 is, while deep neural networks make the decision according to a unified, inseparable complicated
 67 function, humans usually aggregate information to infer intermediate concepts step-by-step. For
 68 example, a student spending a long time in the library is likely to be hard-working. A hard-working
 69 student may master knowledge well. Then, a student mastering knowledge well may get a high GPA.
 70 In this process, we use the concept “long time in the library” to infer the concept “hard-working”,
 71 then use “hard-working” to infer the concept “master knowledge well”. Finally, we infer “high GPA”
 72 with “master knowledge well”. Along this line, human has specific way on raising our decision
 73 model’s complexity and making more comprehensive decisions. Usually, we learn to infer from more
 74 views with more observations and hidden concepts. During this process, a concept may be inferred
 75 from multiple observations, and an observation may be used to infer multiple concepts. For example,
 76 “little missing of classes” may also imply “hard-working”.

77 Existing works have made significant progress in designing tools to transform model logic to human-
 78 understandable logic. However, we believe a more natural way to raise model explainability is
 79 to make its intrinsic decision-making process consistent with humans. Therefore, motivated by
 80 human decision-making process, VOTEN explicitly models a small number of concepts and focuses
 81 on seizing the quantitative relationship between different levels of individual concepts. For better
 82 understanding, observations and concepts in VOTEN can be regarded as *voters* organized in a
 83 hierarchical structure. Each voter independently votes for some higher-level concepts based on their
 84 value. The higher-level concepts aggregate the votes to get their value estimation and further vote for
 85 the next level. During training, the model builds intermediate concepts and learns the quantitative
 86 relationship between concepts in adjacent levels for making proper votes.

87 Formally, for a VOTEN model with D levels of concepts, we use \mathcal{C}_i^d to denote the i -th concept in the
 88 d -th layer, where \mathcal{C}_i^0 denotes an input feature. We refer to transformations between adjacent levels
 89 of concepts as a *voting layer*. In the d -th voting layer, each concept \mathcal{C}_i^d votes for each higher-level
 90 concept \mathcal{C}_j^{d+1} with an independent *voting channel* $\mathcal{V}_{i,j}^d$. $\mathcal{V}_{i,j}^d$ takes the value of \mathcal{C}_i^d as the input and
 91 votes with a single-valued nonlinear function $f_{i,j}^d : \mathbb{R} \rightarrow \mathbb{R}$. Then, a counting layer gets weighted
 92 average of the votes and estimates the value of \mathcal{C}_j^{d+1} as

$$x_j^{d+1} = \sum_{k=1}^{n_d} a_{k,j}^d f_{k,j}^d(x_k^d) \quad \text{s.t. } \forall d, j, \sum_{k=1}^{n_d} a_{k,j}^d = 1, \quad (1)$$

93 where n_d denotes the number of concepts in the d -th layer, x_k^d denotes the value of \mathcal{C}_k^d , $a_{k,j}^d$ denotes
 94 the weight of $\mathcal{V}_{i,j}^d$. The concepts of the last layer is regarded as the decision score, which generates
 95 the model output with $\mathbf{o} = F^{VOTEN}(\mathbf{x}^D)$. In particular, to assure the ability of the voting functions
 96 on effective concept transformation, VOTEN models each voting function with a voting network.
 97 The voting network can be designed in complicated ways without influencing model explainability,
 98 as long as the function is still single-valued. For example, we can adopt weight-sharing structures to
 99 reduce model complexity.

100 2.2 Why is VOTEN more explainable than MLP?

101 In this part, we theoretically discuss VOTEN’s advantages over MLP. Specifically, we claim that
 102 voting expressiveness is the core proposition of VOTEN that raises explainability.

103 **Theorem 1** *MLP can be derived from a subset of degenerated VOTEN models whose voting functions*
 104 *in the form of $f_{i,j}^d(x_i^d) = w_{i,j}^d \sigma(x_i^d) + b_{i,j}^d$, where σ is a predefined activation function, the scalars*
 105 *$w_{i,j}^d, b_{i,j}^d \in \mathbb{R}$ are trainable parameters.*

106 *Proof.* Please refer to Appendix A.

107 From this point of view, MLP also can be regarded as consistent with human decision-making process.
 108 However, MLP is still recognized as difficult to explain. Empirically, concept-transformation-based
 109 decision-making is easier to understand with (1) fewer concepts, (2) shorter concept transformation
 110 paths, and (3) fewer decision-making patterns. Under VOTEN schema, we show how inexpensive
 111 voting channels make MLP disobey these three principles.

112 **Corollary 1** *In MLP, votings from the same concept are linearly correlated.*

113 This means individual voters in MLP are weak in distinguishing different concepts in the next level.
 114 Therefore, MLP relies on highly complicated feature combinations of a large number of deeply
 115 tiled voters to achieve high fitting ability. During this process, necessary intermediate information
 116 is inherently modeled through combinatorial effects of hidden concepts with inexplicit meanings.
 117 Indeed, previous works have proved that human-understandable concepts are inherently mounted in
 118 the hidden units of neural network models [25].

119 **Corollary 2** *In MLP, the hypothesis space for voting distribution is limited to scaling and shifting an*
 120 *input distribution.*

121 This means the voters in MLP may not always achieve significant transformation between different
 122 levels of concepts. Then, the input needs to go through a long path of transformations before it
 123 contributes to the output. Moreover, the deeply tiled massive concepts make each input feature
 124 has massive paths to the output, which brings a large number of possible decision-making patterns
 125 of the model. In addition, the effect of votings can easily couple and cancel each other out in the
 126 downstream calculations. As a result, separately analyzing the role of individual concepts or decision
 127 paths becomes meaningless.

128 Different from MLP, VOTEN has far more expressive voting functions. It directly models nonlinear
 129 transformations between essential intermediate concepts without relying on massive feature combina-
 130 tions and naive transformations. As a result, VOTEN’s decision-making process is explicit with only
 131 a small number of meaningful concepts, thus is explainable.

132 2.3 Explaining the effect of votings

133 In VOTEN, individual voting channels play explicit roles in influencing the model decision. In this
 134 part, we analyze the effect of votings. First, we use a concept function $g_i^d : \mathbb{R}^{n_0} \rightarrow \mathbb{R}$ to represent the
 135 transformation from model inputs $\mathbf{x}^0 \in \mathbb{R}^{n_0}$ to \mathcal{C}_i^d , which decides the meaning of the concept.

136 **Definition 1** In VOTEN, we refer to two concept functions g and g' as equivalent iff. there exists an
 137 invertible function Φ , so that $\forall \mathbf{x}^0 \in \mathbb{R}^{n_0}$, $\Phi(g(\mathbf{x}^0)) = g'(\mathbf{x}^0)$.

138 Since Φ is invertible, the outputs of two equivalent concept functions have one-to-one correspondence
 139 over all the possible inputs. Then, they can effect equally in the decision-making process.

140 **Theorem 2** In VOTEN, if replacing a concept function with an equivalent form, there exists a way to
 141 replace its voting functions so that all the downstream concept functions stay unchanged.

142 *Proof.* Please refer to Appendix A.

143 To analyze how voting channels effect on concept functions, we can reformulate Equation 1 as

$$x_j^{d+1} = g_j^{d+1}(\mathbf{x}^0) = \sum_{i=1}^{n_d} a_{i,j}^d (f_{i,j}^d(g_i^d(\mathbf{x}^0)) - \mathbb{E}_{\mathbf{x}}[f_{i,j}^d(g_i^d(\mathbf{x}))]) + b_j^{d+1}, \quad (2)$$

144 where $\mathbb{E}_{\mathbf{x}}[f_{i,j}^d(g_i^d(\mathbf{x}))]$ denotes the expectation of $f_{i,j}^d$ over all the instances, b_j^{d+1} is a sample-
 145 independent bias. In particular,

$$b_j^{d+1} = \sum_{i=1}^{n_d} a_{i,j}^d \mathbb{E}_{\mathbf{x}}[f_{i,j}^d(g_i^d(\mathbf{x}))] = \mathbb{E}_{\mathbf{x}}[g_j^{d+1}(\mathbf{x})]. \quad (3)$$

146 For simplicity, we use $\overline{h_{i,j}^d}$ to denote $\mathbb{E}_{\mathbf{x}}[f_{i,j}^d(g_i^d(\mathbf{x}))]$. Notably, by adding an arbitrary bias to g_j^d , we
 147 obtain an equivalence of the original concept. According to Theorem 2, we can construct a model
 148 with exactly the same expression (i.e., equivalent concepts and the same predictions) as the previous
 149 one. This implies that VOTEN may converge to models with differed internal bias but exactly the
 150 same decision-making process, indicating VOTEN explanation should be independent of concept
 151 bias. According to Equation 2 and 3, voting channels' average only influence concept bias while the
 152 voting deviation $f_{i,j}^d(g_i^d(\mathbf{x}^0)) - \overline{h_{i,j}^d}$ reveals the effect of $\mathcal{V}_{i,j}^d$ for the decision. This can be intuitively
 153 explained as each voter can vote with different basic scores and only the deviation from the basic
 154 score reflects their judgement for a specific instance. Similarly, from the global view, we reformulate

155 the concept function as $x_j^{d+1} = \sum_i a_{i,j}^d \sigma_{i,j}^d K_{i,j}^d(\mathbf{x}^0) + b_j^{d+1}$, where $K_{i,j}^d(\mathbf{x}^0) = \frac{f_{i,j}^d(g_i^d(\mathbf{x}^0)) - \overline{h_{i,j}^d}}{\sigma_{i,j}^d}$.

156 Notably, $K_{i,j}^d(\cdot)$ generates a distribution with mean 0 and variance 1 over all the instances. Therefore,
 157 $a_{i,j}^d$ and $\sigma_{i,j}^d$ jointly decide the overall effect of votings. Specifically, the counting layer explicitly
 158 adjusts $a_{i,j}^d$ so that reliable voters have stronger influences. Meanwhile, the voter implicitly adjusts
 159 $\sigma_{i,j}^d$. When the concept is less related to the target concept and cannot support proper votes, they
 160 decrease $\sigma_{i,j}^d$ and tend to always vote the basic score to avoid disturbing the model. Otherwise, they
 161 increase $\sigma_{i,j}^d$ and vote confidently to lead the model to correct estimation.

162 Based on the above analysis, we can easily design algorithms to ease both VOTEN local and global
 163 explanation, such as recognizing decision paths and quantifying the concept/feature relevance to the
 164 prediction, which can be found in Appendix B and Appendix C.

165 2.4 VOTEN supports effective model pruning

166 In VOTEN, concepts are estimated by averaging the votes. This means we can delete a voting channel
 167 while keeping the physical meaning of the target concept unchanged. This supports effective link
 168 pruning. Specifically, since only the voting deviation from the average decides the effect, we can
 169 assume the absent channel votes the basic score regardless of the input. Formally, the value of C_j^{d+1}

170 is estimated as $\hat{x}_j^{d+1} = \sum_{i=1}^{n_d} I_{i,j}^d a_{i,j}^d (f_{i,j}^d(x_i^d) - \overline{h_{i,j}^d}) + b_j^{d+1}$, where $I_{i,j}^d \in \{0, 1\}$ indicates if $\mathcal{V}_{i,j}^d$
 171 is not absent. Furthermore, since only involving a small number of voting channels, we can achieve
 172 network pruning on VOTEN by exhaustively exploring how the performance will get influenced if
 173 some voting channels are absent, which is infeasible for MLP. In MLP, we usually prune unimportant
 174 hidden units to reduce model complexity. However, usually not all the voting channels from an
 175 important concept are necessary. In VOTEN, these unnecessary concept transformations can be
 176 further eliminated to not only reduce model complexity but also raises the explainability of the model.
 177 In Appendix D, we give an efficient VOTEN pruning algorithm with a lazy updating strategy.

178 3 Experiment

179 To evaluate the effectiveness and explainability of VOTEN for seizing comprehensive decision-
 180 making patterns. We conducted experiments with 4 large public datasets, including Context-aware

Table 1: Model Performance. We conducted 10 independent runs on each dataset and show the average \pm standard deviation of AUC and AP. In particular, for multi-classifications, we estimated the macro average of each metric. We also did significance test, where * and ** denote significantly (i.e., p -value ≤ 0.05) and very significantly (i.e., p -value ≤ 0.01) worse than VOTEN.

Model	MR		RP		CT		CI	
	AP	AUC	AP	AUC	AP	AUC	AP	AUC
DecisionTree	0.2557**	0.6762**	0.0177**	0.5134**	0.8119**	0.9396**	0.2514**	0.7250**
\pm	0.0007	0.0096	0.0001	0.0014	0.0010	0.0005	0.0022	0.0018
RandomForest	0.4624**	0.9007**	0.0349**	0.6566**	0.9763**	0.9979**	0.6348**	0.9388**
\pm	0.0010	0.0006	0.0011	0.0052	0.0003	0.0001	0.0014	0.0005
LightGBM	0.4817**	0.9209**	0.0525*	0.7305*	0.9753**	0.9965**	0.6972	0.9566
\pm	0.0007	0.0001	0.0017	0.0016	0.0003	0.0001	0.0001	0.0001
MLP	0.4870**	0.9199**	0.0516**	0.7250**	0.9662**	0.9965**	0.6215**	0.9454**
\pm	0.0040	0.0008	0.0014	0.0017	0.0018	0.0002	0.0030	0.0004
VOTEN	0.5007	0.9237	0.0550	0.7322	0.9783	0.9983	0.6522	0.9508
\pm	0.0041	0.0005	0.0001	0.0003	0.0013	0.0001	0.0019	0.0002

181 Multi-Modal Transportation Recommendation (MR) [1, 48], IJCAI-18 Search Conversion Rate
 182 Prediction (RP) [2], Forest Cover Type Prediction (CT) [10], and Census-Income Prediction (CI) [34].
 183 The detailed descriptions of experimental setup can be found in Appendix E.

184 3.1 Performance Evaluation: Can VOTEN achieve higher performance than MLP?

185 We used two widely adopted metrics for imbalanced classification performance evaluation, including
 186 Area Under ROC Curve (AUC) [12] and average precision (AP) [17], whose higher value means
 187 higher performance. In Table 1, we compare the performance of VOTEN with several baselines
 188 (see Appendix E.2). For each dataset, we have carefully tuned the parameters of the baselines
 189 to achieve their best performance. Especially, detailed analysis on MLP parameters can be found in
 190 Appendix E.3. It can be observed that, while deep neural networks are powerful when incorporated
 191 with purposely designed modules or prior knowledge for specific tasks, its standard form (i.e., MLP)
 192 without task-specific structures may perform worse than LightGBM, which have also been shown by
 193 many previous studies [41, 7, 3]. Indeed, LightGBM is believed to be powerful in handling structured
 194 data and often appears as the major model of the top solutions in data-mining competitions [48, 31].
 195 In contrast, VOTEN, also in its standard form, significantly outperforms MLP for all these tasks, even
 196 better than LightGBM (except for CI, which contains mostly discrete features and is naturally suitable
 197 for tree-based models). This shows the effectiveness of VOTEN in terms of handling real-world
 198 problems. Indeed, VOTEN is more suitable for handling data-mining tasks since it is consistent
 199 with the practical decision-making process. It should be noticed that, similar to MLP, VOTEN is
 200 a standard and generic model that can be easily expanded to task-specific networks to raise model
 201 performance (for example, we can replace MLP with VOTEN in DeepFM [21]). While this paper
 202 focuses on the generic performance of standard VOTEN, it shows the possibility of building more
 203 powerful deep learning solutions for a wide range of applications.

204 3.2 Explanation Complexity: Is the decision-making process of VOTEN recognizable?

205 In Table 2, we compare the explanation complexity of MLP and VOTEN. As we have dis-
 206 cussed in Section 2.2, we use the number of
 207 feature combinations, the length of decision
 208 paths, and the number of possible decision-
 209 making patterns to show the explanation com-
 210 plexity of a model. In particular, we trained
 211 two VOTEN models with different settings for
 212 each dataset. Specifically, “VOTEN⁻” is com-
 213 parable to the best performance of MLP, with
 214 the least feature combinations. “VOTEN” is
 215 the one with the best performance. It can be
 216 observed that VOTEN greatly reduces feature
 217 combinations and decision paths. For example,
 218 on the MR dataset, MLP needs 7 voting layers
 219 that each contain 128 concepts. This brings an exponentially large number of long decision paths. In
 220 contrast, VOTEN achieves comparable performance with 16 hidden concepts in total. Then, each
 221 feature only has 16 possible paths with a length of 2 to reach the output. This significantly eases
 222 the understanding of the decision-making process. For some datasets, VOTEN with no intermediate
 223

Table 2: Explanation complexity. “#C/L” counts hidden concepts in each layer. “#P/F” counts possible decision paths from each feature.

Data	Model	Performance		#C/L	Depth	#P/F
		AP	AUC			
MR	MLP	0.487	0.920	128	7	2 ⁴⁹
	VOTEN ⁻	0.497	0.922	16	1	16
	VOTEN	0.501	0.924	16	2	256
RP	MLP	0.052	0.725	12	3	1,728
	VOTEN ⁻	0.055	0.732	0	0	1
	VOTEN	0.055	0.732	8	2	64
CT	MLP	0.966	0.996	128	7	2 ⁴⁹
	VOTEN ⁻	0.969	0.997	32	2	1,024
	VOTEN	0.978	0.998	64	2	4,096
CI	MLP	0.622	0.945	32	2	1,024
	VOTEN ⁻	0.652	0.950	0	0	1
	VOTEN	0.652	0.951	4	2	16

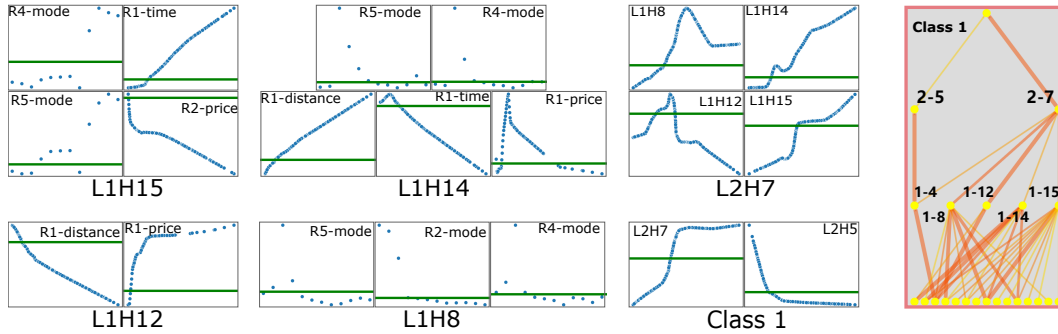


Figure 2: Case study of VOTEN global decision-making process in MR dataset. We show the important paths from inputs to one output, where wide lines indicate important voting channels. The concept $LxHy$ indicates the y -th concept in the x -th layer. We visualize the important voting functions for each concept along the paths relevant to L2H7, where blue lines show the function while green lines show the average vote.

224 concepts (i.e., features directly vote for the prediction) can achieve comparable performance to MLP.
 225 Moreover, even when reaching its best performance, VOTEN still has a much smaller explanation
 226 complexity than MLP. It should be noticed that we have listed the maximum possible number of
 227 decision paths for convincing illustration. In practice, we can easily distinguish a fewer number of
 228 important decision paths, which will be shown in the following experiments.

229 3.3 Case Study: How to explain a VOTEN model?

230 We show how to explain a VOTEN model with an example in the MR dataset. Specifically, we
 231 first analyze the global decision-making process and discover the meaning of concepts by observing
 232 voting functions on important decision paths. Then, we locally explain how the inputs of an instance
 233 hierarchically vote the final prediction.

234 **Task Description.** The task is to recommend the transport mode for online map app users, given a
 235 user and an Origin-Destination (OD) pair. The features mainly contain user portraits and an ordered
 236 list of recommended plans of the map app. Each plan consists of transport mode, time, distance, and
 237 price. We deleted the first recommended plan’s mode information since it is too strongly correlated
 238 with the label (many users choose the first recommendation as default). In this way, we can better
 239 observe how the model incorporates complicated information for meaningful decisions.

240 **Global Decision-Making Process.** We focus on the decision-making process for the class “Subway”.
 241 First, we discover the important transformations from observations to the prediction with our decision
 242 path recognition algorithm. The results are shown in Figure 2. Next, we analyze the meaning of
 243 concepts from the bottom to the top. L1H8 gets larger when modes of more recommended plans
 244 are “Subway”, which we regard as “*OD pair with flexible subway-based plans*”. L1H12 decreases
 245 with higher prices and lower distances, which we regard as “*OD pair’s distance-cost performance*”.
 246 L1H14 observes if the OD pair is distant but still available with inexpensive and fast transportation.
 247 Besides, it also observes if modes of public transportation ever appear in the plan list. Therefore, we
 248 regard L1H14 as “*distant OD pair with convenient and economical public transportation*”. L1H15
 249 is sensitive to a time-consuming top-1 plan. It also observes the other plans’ modes and prices to
 250 estimate if trading money for efficiency is infeasible. Therefore, we regard L1H15 as “*no choice*
 251 *but a time-consuming transportation*”. In the second layer, L2H7 is estimated based on L1H8,
 252 L1H12, L1H14, and L1H15. It gets higher if many subway-based plans available (higher L1H8),
 253 transportation with balanced distance-cost is recommended (has a peak for L1H12), distant OD pair
 254 but still has convenient public transportation (higher L1H14), or costly time-saving transportation is
 255 infeasible (higher L1H15). Comprehensively considering these reasons, it implies “*OD pair suitable*
 256 *for subway transportation*, which votes for the score of “Subway” with a monotonic transformation.
 257 Along the other path, L1H4 is a concept “*OD pair with inexpensive transportation*”. Then, L2H5 is
 258 also about the price since it is mainly based on L1H4. It should be noticed that its estimation may
 259 still be enhanced by other lower-level concepts when dealing with specific instances, although they
 260 are less important from the global view. Finally, L2H5 and L2H7 vote nonlinearly to the score so that
 261 the model can make accurate quantitative predictions. With the above analysis, we can qualitatively
 262 understand the logic of VOTEN on recommending “Subway”. Actually, the decision-making process
 263 is quantitatively more complicated and can handle more special cases. In practice, domain experts

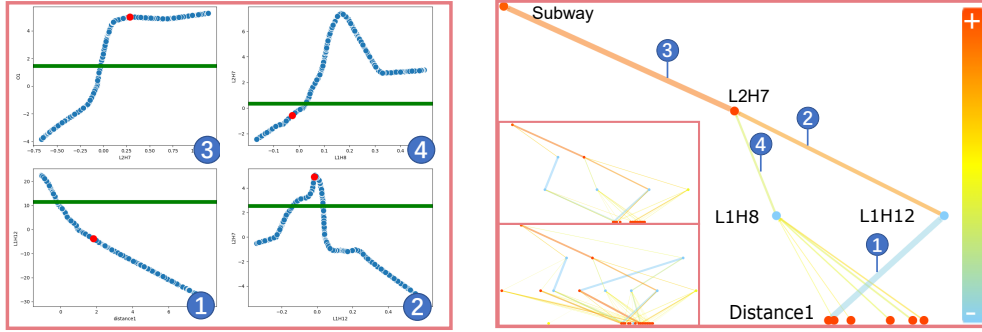


Figure 3: Case study of VOTEN local decision-making process in MR dataset, where wider and darker lines indicate stronger influence (negatively in blue and positively in red). The bottom left shows filtered paths when we gradually bring up the thresholds. We also show the important voting functions, where the green lines show the average and red points show vote for the current instance.

264 can thoroughly analyze the shape and gradients of the voting functions to get more insights into the
 265 concepts. This may help researchers to find new concepts and develop new theories, especially in
 266 fields such as psychology and management, where scientists work on finding mechanisms linking
 267 observations to outcomes. Extra visualizations on global decision path and voting functions can be
 268 found in Appendix F and Appendix G.

269 **Local Decision-Making Process.** Then, we analyze the decision paths for specific instances. In
 270 Figure 3, we visualize the most important paths for a sample with a high score for “Subway”,
 271 which are filtered with a small threshold in our local decision path recognition algorithm. Based on
 272 the global analysis, we can easily tell how the observations gradually transformed to higher-level
 273 concepts. Specifically, for L1H12, the vote from price is filtered out, showing the OD pair’s price is
 274 normal over the dataset. However, it finds the OD pair to be distant, which contributes to a relatively
 275 larger *distance-cost performance* than average, given the normal price. This indicates the plan to
 276 be relatively more economical. L2H7 observes the value of L1H12, and finds that the OD pair can
 277 choose transportation whose *distance-cost performance* more balanced than average cases (reaches
 278 the peak of voting). In this case, L2H7 votes high for subway, which is an economical and balanced
 279 transportation mode. On the other decision path, L1H8 thinks “Subway” may not be suitable since
 280 subway-based plans seldom appear in the list, thus votes negatively for L2H7. But L1H12 makes a
 281 very confident judgment based on the balance of distance-price performance, which dominates L2H7
 282 and makes the model predict correctly. Interestingly, we find one of the most important strategies
 283 in this task is to guess the transportation mode most recommended by the app (the information that
 284 we hide in prior), which is reasonable. On the one hand, the app’s recommender system trained
 285 with abundant information can naturally achieve high performance. On the other hand, many users
 286 will click the first recommendation as default. In addition to this strategy, the model will use more
 287 important decision paths to achieve more accurate predictions. As we gradually increase the threshold,
 288 more decision-making patterns appeared. Extra visualizations on local explanation of other datasets
 289 can be found in Appendix H.

290 3.4 Relevance Analysis: Can VOTEN help quantify feature relevance?

291 **Propagation-Based Relevance.** Motivated by relevance propagation [9], we propose an algorithm
 292 (see Appendix D) to quantify the relevance of features and concepts, based on the important decision
 293 paths. The short decision paths of VOTEN decreases error accumulation during the propagation and
 294 enables more accurate relevance estimation. In Figure 4, we visualize the propagation-based feature
 295 relevance with heatmaps. Obviously, features from #69 to #86 are generally important in the MR
 296 dataset, among which the first several features (information about the top-1 plan in the list) are the
 297 most relevant. Furthermore, the relevant features vary for different samples in terms of different
 298 classes, which indicates VOTEN to predict in multiple patterns. For example, Figure 4(b) shows that
 299 VOTEN assigns high score to class 2 for sample 4 and sample 5 with different reasons. Specifically,
 300 feature #3 and #53 contribute positively for sample 5 while negatively for sample 4. Instead, #50
 301 is more positively relevant for sample 4. Extra visualizations can be found in Appendix I.

302 **Single-Sighted Prediction Strength.** We can also estimate the feature relevance from the view of
 303 model performance when the prediction is supported by a single voter in some layer. Specifically,
 304 we disable the other voters in a similar way as we do in model pruning. Table 3 shows 5 MR
 305 features achieving the highest AUC when conducting single-sighted prediction for class 0 and

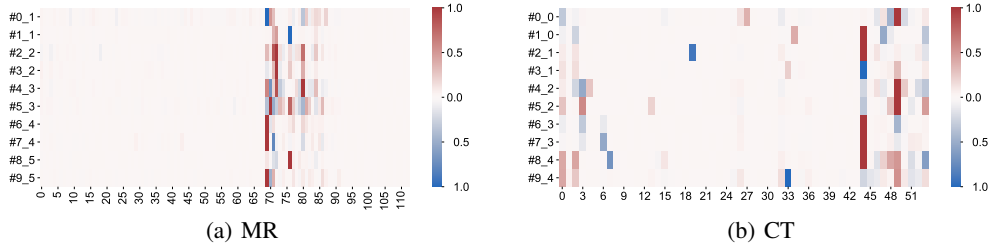


Figure 4: Heatmap for the propagation-based relevance in VOTEN. The x-axis represents features and the y-axis represents the instance-output pairs, where $\#I_O$ indicates the feature relevance of the I -th instance to class O . Red means positive relevance while blue means negative relevance.

Table 3: Features with top-5 single-sighted prediction strength for class 0 and class 1 of MR dataset in VOTEN. ‘‘Input’’ indicates the AUC of ranking the samples according to the feature’s value.

Model	Class 0					Class 1				
	Top-1	Top-2	Top-3	Top-4	Top-5	Top-1	Top-2	Top-3	Top-4	Top-5
VOTEN	0.675	0.653	0.647	0.606	0.592	0.747	0.632	0.625	0.587	0.566
Input	0.596	0.533	0.572	0.600	0.586	0.535	0.527	0.534	0.555	0.525

306 class 1 in VOTEN. As a comparison, we also show the AUC of directly ranking with the feature,
 307 which reveals the feature’s linear correlation with the prediction. In particular, since there can be
 308 negative correlations, we evaluate AUC for the rank in increasing and decreasing order, and use
 309 the larger one as the performance. Moreover, Figure 5 shows the results on all the features. It
 310 can be observed that, even if approximated into single-sighted, VOTEN significantly raises AUC,
 311 showing VOTEN to recognize important features and strengthen their effectiveness with nonlinear
 312 transformations. Especially, while feature #85 originally seems not correlated with the output,
 313 VOTEN finds it in practice nonlinearly very relevant to the output. This proves the effectiveness
 314 of VOTEN on quantifying the inherent nonlinear relationships between the observations and the
 315 prediction. Interestingly, VOTEN weakens the effect of some features (e.g., feature #3 for RP) to
 316 prevent them from disturbing the prediction. Extra visualizations can be found in Appendix J.

317 3.5 Pruning Experiment

318 We conducted pruning experiments on
 319 VOTEN for MR and CT datasets with
 320 our pruning algorithm. As we gradually
 321 deleted voting channels, we monitored the
 322 change of AUC during this process, which
 323 are shown in Figure 6. For MR, AUC is
 324 still near 0.924 after pruning nearly half
 325 of the voting channels. For CT, AUC is
 326 still near 0.998 after pruning a quarter
 327 of the channels. Interestingly, proper prun-
 328 ing may slightly raise model performance,
 329 which is reasonable as a simple model has
 330 less chance of over-fitting. In practice, operations like fine-tuning can be adopted to further raise
 331 the performance of the pruned VOTEN model. Then, the model can be further compressed without
 332 affecting the prediction much. These results prove VOTEN to support effective pruning, which is
 333 helpful. We can use complicated information for training and easily prune the model to decrease the
 334 complexity both for storage, calculation, and explanation.

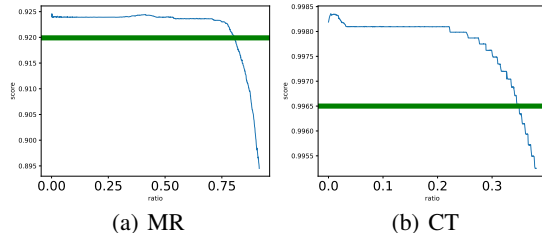


Figure 6: VOTEN pruning experiments. The x-axis shows the ratio of pruned channels while the y-axis shows AUC. The green line shows MLP’s performance.

335 4 Related Work

336 **Post-Hoc Deep Learning Explanation.** Post-Hoc explaining algorithms analyze the relevance of
 337 features in a model-free way, mainly including propagation-based methods [27, 5] and perturbation-
 338 based methods [49, 46]. Propagation-based methods propagate the relevance score backward to the
 339 inputs. For example, Simonyan *et al.* [40] generates saliency maps with the gradients of the output
 340 category with respect to the inputs. Bach *et al.* [9] proposed Layer-wise Relevance Propagation
 341 (LRP), which designed effective rules for the propagation. Perturbation-based methods explain
 342 model behavior by observing how the output reacts to purposely perturbed or constructed inputs. For
 343 example, Local Interpretable Model-Agnostic Explanations (LIME) [38] trains a local explainable

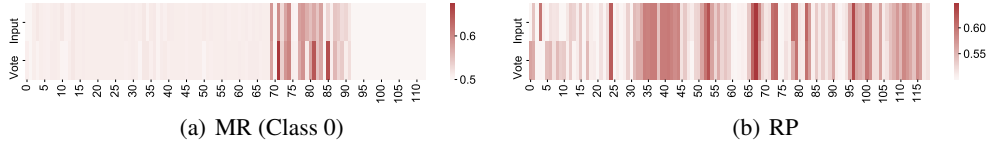


Figure 5: Single-sighted prediction strength in VOTEN. We also show the AUC of linearly ranking the samples as comparisons. Darker color means higher AUC.

344 approximation model around the prediction with randomly perturbed features and the corresponding
 345 outputs. SHapley Additive exPlanations (SHAP) [30] estimates the Shapley value of features to
 346 measure their contribution to model performance. In addition to features, some works estimate concept
 347 importance for a model [20, 43, 25]. For example, Kim *et al.* [25] learn the representation of human-
 348 understandable concepts with labeled concept-relevant examples and estimate concept sensitivity
 349 according to the directional derivative towards the concepts. Along this line, abundant works have
 350 been proposed to further raise the effectiveness of post-hoc interpretation algorithms [28, 45, 13].
 351 However, these algorithms regard models as blackboxes and heuristically explain with their own
 352 metrics, which cannot give explicit understandings of the actual decision-making process. Different
 353 from existing works, we proposed a naturally understandable neural network model.

354 **Intrinsically Explainable Machine Learning Techniques.** Intrinsically explainable models can be
 355 explained without relying on post-hoc algorithms [8, 4], mainly including classic models such as
 356 logistic regression [33], linear support vector machine [23], decision trees [39], generalized additive
 357 models [22] and Bayesian models [37]. Recently, Agarwal *et al.* [6] proposed Neural Additive Model
 358 that predicts with a linear combination of neural networks. However, all these models usually have
 359 tight restrictions on the hypothesis space, which limits their fitting ability on complicated real-world
 360 problems. Based on these methods, complicated models are developed for higher performance.
 361 However, even with intrinsically explainable base models, these complicated models still need post-
 362 hoc algorithms for explanation [8]. For example, ensemble tree models [14, 24] predict with a large
 363 number of weak learners. However, the joint decision-making process of massive decision trees is
 364 difficult to understand. Kernel functions [35] are incorporated in support vector machines to seize
 365 high dimensional feature interactions. However, the dimension transformation is implicit and not
 366 understandable. In recent years, researchers also try to design explainable neural network models by
 367 incorporating purposely designed task-specific constraints or structures [15, 47, 36]. For example,
 368 Qin *et al.* [36] proposed a person-job fit network, which used attention layers to find important talent
 369 working experience for different job requirements. However, these models cannot be adopted by the
 370 general tasks. Besides, they only provide heuristic and domain-specific intermediate information
 371 instead of telling the complete decision-making process. Different from these works, we aim at a
 372 general explainable neural network model, which has an intrinsically explainable decision-making
 373 process while retaining the high fitting ability.

374 5 Concluding Remarks

375 In this paper, we have proposed an explainable deep learning model, VOTEN. Specifically, we
 376 theoretically analyzed the major components of VOTEN and discussed its priority over MLP, and
 377 accordingly proposed some efficient algorithms to raise the model usability. Experimental results
 378 on multiple real-world datasets clearly demonstrated that VOTEN can significantly improve the
 379 explainability and performance of deep learning.

380 **Limitations.** In this paper, we focused on comparing VOTEN with MLP, which is the generic
 381 and basic structure of deep learning models. Many powerful problem-specific structures can be
 382 derived from MLP by adding operations such as weight sharing (e.g., CNN). Similar to MLP,
 383 VOTEN is a basic and generic structure. It can be adopted to problem-specific models (e.g., we
 384 can simply use VOTEN to replace MLPs in deepFM [21] or MMoE [32]). Indeed, recent studies
 385 show that if properly designed, simple MLP-based structure achieves comparable performance to
 386 complicated SOTA models [19, 42]. VOTEN’s advantages over MLP provides great possibility on
 387 further improving a wide range of deep learning applications. In the future, we will also explore
 388 building VOTEN-based task-specific structures. In addition, since VOTEN automatically extracts
 389 concepts during training, human effort is needed to observe the voting functions for understanding
 390 the concepts, which is a common issue in unsupervised concept modeling, such as Latent Dirichlet
 391 Allocation [11]. In the future, we will work on easing the concept understanding of VOTEN, such as
 392 recognizing concept-related samples or aligning VOTEN with human-understandable concepts.

393 References

- 394 [1] Context-aware multi-modal transportation recommendation. <https://dianshi.bce.baidu.com/competition/29/question>.
395
- 396 [2] Ijcai-18 alimama sponsored search conversion rate(cvr) prediction contest. <https://tianchi.aliyun.com/competition/entrance/231647/information>.
397
- 398 [3] Lightgbm vs neural network. <https://mljar.com/machine-learning/lightgbm-vs-neural-network/>.
399
- 400 [4] Amina Adadi and Mohammed Berrada. Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access*, 6:52138–52160, 2018.
401
- 402 [5] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. *arXiv preprint arXiv:1810.03292*, 2018.
403
- 404 [6] Rishabh Agarwal, Nicholas Frosst, Xuezhou Zhang, Rich Caruana, and Geoffrey E Hinton. Neural additive models: Interpretable machine learning with neural nets. *arXiv preprint arXiv:2004.13912*, 2020.
405
406
- 407 [7] Anesh Alvanpour, Sumit Kumar Das, Christopher Kevin Robinson, Olfa Nasraoui, and Dan Popa. Robot failure mode prediction with explainable machine learning. In *2020 IEEE 16th International Conference on Automation Science and Engineering (CASE)*, pages 61–66. IEEE, 2020.
408
409
410
- 411 [8] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58:82–115, 2020.
412
413
414
- 415 [9] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one*, 10(7):e0130140, 2015.
416
417
- 418 [10] Jock A Blackard and Denis J Dean. Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables. *Computers and electronics in agriculture*, 24(3):131–151, 1999.
419
420
- 421 [11] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
422
- 423 [12] Andrew P Bradley. The use of the area under the roc curve in the evaluation of machine learning algorithms. *Pattern recognition*, 30(7):1145–1159, 1997.
424
- 425 [13] Steven Bramhall, Hayley Horn, Michael Tieu, and Nibhrat Lohia. Qlime-a quadratic local interpretable model-agnostic explanation approach. *SMU Data Science Review*, 3(1):4, 2020.
426
- 427 [14] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
428
429
- 430 [15] Xu Chen, Hanxiong Chen, Hongteng Xu, Yongfeng Zhang, Yixin Cao, Zheng Qin, and Hongyuan Zha. Personalized fashion recommendation with visual explanations based on multimodal attention network: Towards visually explainable recommendation. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 765–774, 2019.
431
432
433
434
- 435 [16] Arun Das and Paul Rad. Opportunities and challenges in explainable artificial intelligence (xai): A survey. *arXiv preprint arXiv:2006.11371*, 2020.
436
- 437 [17] Jesse Davis and Mark Goadrich. The relationship between precision-recall and roc curves. In *Proceedings of the 23rd international conference on Machine learning*, pages 233–240, 2006.
438

- 439 [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training
440 of deep bidirectional transformers for language understanding. In *Proceedings of the 2019*
441 *Conference of the North American Chapter of the Association for Computational Linguistics:*
442 *Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019.
- 443 [19] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Repmlp: Re-parameterizing
444 convolutions into fully-connected layers for image recognition. *arXiv e-prints*, pages arXiv-
445 2105, 2021.
- 446 [20] Amirata Ghorbani, James Wexler, James Zou, and Been Kim. Towards automatic concept-based
447 explanations. *arXiv preprint arXiv:1902.03129*, 2019.
- 448 [21] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. Deepfm: a
449 factorization-machine based neural network for ctr prediction. *arXiv preprint arXiv:1703.04247*,
450 2017.
- 451 [22] Trevor J Hastie and Robert J Tibshirani. *Generalized additive models*, volume 43. CRC press,
452 1990.
- 453 [23] Shujun Huang, Nianguang Cai, Pedro Penzuti Pacheco, Shavira Narrandes, Yang Wang, and
454 Wayne Xu. Applications of support vector machine (svm) learning in cancer genomics. *Cancer*
455 *Genomics-Proteomics*, 15(1):41–51, 2018.
- 456 [24] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and
457 Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural*
458 *information processing systems*, 30:3146–3154, 2017.
- 459 [25] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al.
460 Interpretability beyond feature attribution: Quantitative testing with concept activation vectors
461 (tcav). In *International conference on machine learning*, pages 2668–2677. PMLR, 2018.
- 462 [26] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
463 convolutional neural networks. *Advances in neural information processing systems*, 25:1097–
464 1105, 2012.
- 465 [27] Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Montavon, Wojciech
466 Samek, and Klaus-Robert Müller. Unmasking clever hans predictors and assessing what
467 machines really learn. *Nature communications*, 10(1):1–8, 2019.
- 468 [28] Heyi Li, Yunke Tian, Klaus Mueller, and Xin Chen. Beyond saliency: understanding convolu-
469 tional neural networks from saliency prediction on layer-wise relevance propagation. *Image*
470 *and Vision Computing*, 83:70–86, 2019.
- 471 [29] Jianxun Lian, Xiaohuan Zhou, Fuzheng Zhang, Zhongxia Chen, Xing Xie, and Guangzhong
472 Sun. xdeepfm: Combining explicit and implicit feature interactions for recommender systems.
473 In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery*
474 *& Data Mining*, pages 1754–1763, 2018.
- 475 [30] Scott Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *arXiv*
476 *preprint arXiv:1705.07874*, 2017.
- 477 [31] Zhipeng Luo, Jianqiang Huang, Ke Hu, Xue Li, and Peng Zhang. Accuair: Winning solution to
478 air quality prediction for kdd cup 2018. In *Proceedings of the 25th ACM SIGKDD International*
479 *Conference on Knowledge Discovery & Data Mining*, pages 1842–1850, 2019.
- 480 [32] Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H. Chi. Modeling task
481 relationships in multi-task learning with multi-gate mixture-of-experts. In *Proceedings of the*
482 *24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD
483 '18, page 1930–1939, 2018.
- 484 [33] Edward C Norton, Bryan E Dowd, and Matthew L Maciejewski. Marginal effects—quantifying
485 the effect of changes in risk factors in logistic regression models. *Jama*, 321(13):1304–1305,
486 2019.

- 487 [34] Nikunj C Oza and Stuart Russell. Experimental comparisons of online and batch versions of
488 bagging and boosting. In *Proceedings of the seventh ACM SIGKDD international conference*
489 *on Knowledge discovery and data mining*, pages 359–364, 2001.
- 490 [35] Luis Carlos Padierna, Martin Carpio, Alfonso Rojas-Dominguez, Hector Puga, and Hector
491 Fraire. A novel formulation of orthogonal polynomial kernel functions for svm classifiers: the
492 gegenbauer family. *Pattern Recognition*, 84:211–225, 2018.
- 493 [36] Chuan Qin, Hengshu Zhu, Tong Xu, Chen Zhu, Liang Jiang, Enhong Chen, and Hui Xiong.
494 Enhancing person-job fit for talent recruitment: An ability-aware neural network approach.
495 In *The 41st international ACM SIGIR conference on research & development in information*
496 *retrieval*, pages 25–34, 2018.
- 497 [37] Adrian E Raftery. Bayesian model selection in social research. *Sociological methodology*, pages
498 111–163, 1995.
- 499 [38] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining
500 the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international*
501 *conference on knowledge discovery and data mining*, pages 1135–1144, 2016.
- 502 [39] S Rasoul Safavian and David Landgrebe. A survey of decision tree classifier methodology.
503 *IEEE transactions on systems, man, and cybernetics*, 21(3):660–674, 1991.
- 504 [40] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks:
505 Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*,
506 2013.
- 507 [41] Ying Sun, Fuzhen Zhuang, Hengshu Zhu, Qi Zhang, Qing He, and Hui Xiong. Market-oriented
508 job skill valuation with cooperative composition neural network. *Nature Communications*,
509 12(1):1–12, 2021.
- 510 [42] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas
511 Unterthiner, Jessica Yung, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, et al. Mlp-mixer: An
512 all-mlp architecture for vision. *arXiv preprint arXiv:2105.01601*, 2021.
- 513 [43] Weibin Wu, Yuxin Su, Xixian Chen, Shenglin Zhao, Irwin King, Michael R Lyu, and Yu-Wing
514 Tai. Towards global explanations of convolutional neural networks with concept attribution. In
515 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
516 8652–8661, 2020.
- 517 [44] Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. Hierarchical
518 attention networks for document classification. In *Proceedings of the 2016 conference of the*
519 *North American chapter of the association for computational linguistics: human language*
520 *technologies*, pages 1480–1489, 2016.
- 521 [45] Muhammad Rehman Zafar and Naimul Mefraz Khan. Dlime: a deterministic local interpretable
522 model-agnostic explanations approach for computer-aided diagnosis systems. *arXiv preprint*
523 *arXiv:1906.10263*, 2019.
- 524 [46] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In
525 *European conference on computer vision*, pages 818–833. Springer, 2014.
- 526 [47] Yingying Zhang, Shengsheng Qian, Quan Fang, and Changsheng Xu. Multi-modal knowledge-
527 aware hierarchical attention network for explainable medical question answering. In *Proceedings*
528 *of the 27th ACM International Conference on Multimedia*, pages 1089–1097, 2019.
- 529 [48] Wenjun Zhou, Taposh Dutta Roy, and Iryna Skrypnik. The kdd cup 2019 report. *ACM SIGKDD*
530 *Explorations Newsletter*, 22(1):8–17, 2020.
- 531 [49] Luisa M Zintgraf, Taco S Cohen, Tameem Adel, and Max Welling. Visualizing deep neural
532 network decisions: Prediction difference analysis. *arXiv preprint arXiv:1702.04595*, 2017.

533 **Checklist**

- 534 1. For all authors...
- 535 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
536 contributions and scope? [Yes]
- 537 (b) Did you describe the limitations of your work? [Yes] We have described the limitations
538 and future works in Section 5.
- 539 (c) Did you discuss any potential negative societal impacts of your work? [N/A] Our work
540 aims to raise the explainability of deep learning. There is no potential negative societal
541 impact.
- 542 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
543 them? [Yes]
- 544 2. If you are including theoretical results...
- 545 (a) Did you state the full set of assumptions of all theoretical results? [N/A] No assumptions
546 for theoretical analysis in this paper.
- 547 (b) Did you include complete proofs of all theoretical results? [Yes] We have provided
548 complete proofs for our theoretical results in Appendix A.
- 549 3. If you ran experiments...
- 550 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
551 mental results (either in the supplemental material or as a URL)? [Yes] The supplemental
552 material contains code and instructions to reproduce the main experimental results
553 of this paper. The original public dataset can be downloaded by readers on their own,
554 which we have provided the URLs in the code package.
- 555 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
556 were chosen)? [Yes] We have specified these training details, please refer to Appendix
557 E.
- 558 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
559 ments multiple times)? [Yes] We have reported error bars after running experiment
560 multiple times. In addition, we have conducted significant test. Related information
561 can be found in Table 1.
- 562 (d) Did you include the total amount of compute and the type of resources used (e.g., type
563 of GPUs, internal cluster, or cloud provider)? [Yes] We have reported these information,
564 please refer to Appendix E.
- 565 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 566 (a) If your work uses existing assets, did you cite the creators? [Yes] Our work use public
567 dataset. We have cited the creators.
- 568 (b) Did you mention the license of the assets? [No] We are sorry that we cannot find
569 the license of the datasets. Specifically, MR and PR can be downloaded from the
570 competition and CT and CI are publicly available in UCI repository. The licenses are
571 not specified.
- 572 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
573 We have submitted our code as supplemental material.
- 574 (d) Did you discuss whether and how consent was obtained from people whose data you're
575 using/curating? [Yes] We have mentioned that we use public datasets, please refer to
576 Section 3 and Appendix E.
- 577 (e) Did you discuss whether the data you are using/curating contains personally identifiable
578 information or offensive content? [Yes] Our data contain no such information. We have
579 emphasized that all the data are anomalous, please refer to Appendix E.
- 580 5. If you used crowdsourcing or conducted research with human subjects...
- 581 (a) Did you include the full text of instructions given to participants and screenshots, if
582 applicable? [N/A] No human subject involved.
- 583 (b) Did you describe any potential participant risks, with links to Institutional Review
584 Board (IRB) approvals, if applicable? [N/A] No human subject involved.
- 585 (c) Did you include the estimated hourly wage paid to participants and the total amount
586 spent on participant compensation? [N/A] No human subject involved.