

Efficient Shapley Explanation For Features Importance Estimation Under Uncertainty

Xiaoxiao Li^{*}, Yuan Zhou[†], Nicha C. Dvornek^{†*}, Yufeng Gu[◇], Pamela Ventola[‡]
and James S. Duncan^{**†}

^{*} Biomedical Engineering, Yale University, New Haven, CT, USA

^{*} Electrical Engineering, Yale University, New Haven, CT, USA

[†] Radiology & Biomedical Imaging, Yale School of Medicine, New Haven, CT, USA

[‡] Child Study Center, Yale School of Medicine, New Haven, CT, USA

[◇] College of Information Science & Electronic Engineering, Zhejiang University,
Hangzhou, China

Abstract. Complex deep learning models have shown their impressive power in analyzing high-dimensional medical image data. To increase the trust of applying deep learning models in medical field, it is essential to understand why a particular prediction was reached. Data feature importance estimation is an important approach to understand both the model and the underlying properties of data. Shapley value explanation (SHAP) is a technique to fairly evaluate input feature importance of a given model. However, the existing SHAP-based explanation works have limitations such as 1) computational complexity, which hinders their applications on high-dimensional medical image data; 2) being sensitive to noise, which can lead to serious errors. Therefore, we propose an uncertainty estimation method for the feature importance results calculated by SHAP. Then we theoretically justify the methods under a Shapley value framework. Finally we evaluate our methods on MNIST and a public neuroimaging dataset. We show the potential of our method to discover disease related biomarkers from neuroimaging data.

1 Introduction

Deep learning models, trained on extremely large data sets, have even exceeded human performance in many tasks. Even in the medical field, there are impressive results. Due to the difficulty of interpreting complex deep learning models, some simple ones, such as linear regression and random forest, are often preferred in clinical usage. One common property of linear regression and random forest is their interpretability on feature importance. Without sacrificing the benefits of using deep learning models to improve task performance, many efforts have been made to estimate feature importance scores for deep learning models. There are three main approaches for feature importance estimation: 1) gradient-based methods, such as Simple Gradient (SG) [1], Integrated Gradient (IG) [2], LRP

¹ This work was supported by NIH Grant [R01NS035193, R01MH100028].

² Our code is publicly available at: <https://github.com/xxlya/DistDeepSHAP/>

[3], and DeepLIFT [4]; 2) sensitivity based methods, such as LIME [5] and SHAP [6]; 3) methods that mimic deep learning models using tree- or rule-based models [7]. The Shapley value is a means of fairly portioning the collective profit attained by a coalition of players, based on the relative contributions of the players in a game [8]. In this work, we focus on SHAP-based explanation [6].

Although reliability is necessary for model explanations to be trustworthy, relatively few studies have focused on quantifying the uncertainty and robustness of explanation methods. For example, it has been shown that multiple importance estimation methods incorrectly attribute the scores when a constant vector shift is applied to the input [9]. The attributions provided by interpretation methods may themselves contain significant uncertainty [10] and imperceptibly small perturbations of the input can significantly alter the explanations [11].

In order to apply sampling based uncertainty estimation, we modified the original formulation of DeepSHAP [6]. Different from DeepSHAP that back-propagates the prediction difference between input and the a point estimate of references, we consider the reference values as a distribution and show that the Shapley value can be estimated by bootstrap sampling. Therefore, uncertainty of feature importance scores could be measured. Our experiments quantify the performance of different uncertainty estimation methods and their impact on uncertainty-related error reduction. Our key contributions are summarized as follows: 1) Propose a Shapley value estimation framework with uncertainty estimation; 2) Evaluate uncertainty estimation using both human interpretation and quantitative methods; and 3) Apply our method to a subgroup biomarker detection problem in neuroimaging.

2 Preliminaries

2.1 Shapley Value For Feature Importance Estimation

Consider a cooperative game with N players aiming at maximizing a payoff, and let $S \subset \mathcal{N} = \{1, \dots, N\}$ be a subset consisting of $|S|$ players. Suppose the prediction function $f : \mathbb{R}^N \rightarrow \mathbb{R}$ is the model to be explained, which outputs an importance score for a specific class. Denote $v_x : 2^{\mathcal{N}} \rightarrow \mathbb{R}$ (which maps a subset of \mathcal{N} to a real number) as the importance score evaluation function given a subset of features of input $x = [x_1, \dots, x_N]^T \in \mathbb{R}^N$. We have $f(x) = v_x(\mathcal{N})$. The prediction power for the i th feature is the weighted sum of all possible marginal contributions:

$$\phi_i^x = \frac{1}{N} \sum_{S \subseteq \mathcal{N} \setminus \{i\}} \binom{N-1}{|S|}^{-1} (v_x(S \cup \{i\}) - v_x(S)). \quad (1)$$

While Shapley values give a more accurate interpretation of the importance of each player in a coalition, their calculation is expensive. When the number of features (i.e., players in the game) is a large N , the computational complexity will be 2^N , which is especially expensive. Therefore, computing accurate Shapley values is a challenging task.

In classification scenario, f is the output score for a specific class and $v_x(S)$ is computed by feeding x into f with $\{x_i : i \notin S\}$ replaced by some reference values, which come from samples of some other reference classes.

2.2 Propagating Shapley values

DeepSHAP [6], which is built on DeepLIFT [4], is a fast approximation to Shapley value by recursively passing DeepLIFT’s multipliers. With uninformative reference $r \in \mathbb{R}^N$ (i.e. background of images) and prediction model f and $y = f(x)$, we define $\Delta y = f(x) - f(r)$, $\Delta x_i = x_i - r_i$. DeepLIFT has $\sum_{i=1}^N m_{x_i f} \Delta x_i = \Delta y$, where $m_{x_i f}$ is the multiplier and the importance score of x_i is $m_{x_i f} \Delta x_i$. Suppose $f = g \circ h$, and $h : \mathbb{R}^N \rightarrow \mathbb{R}^M : x \mapsto [h_1(x), \dots, h_M(x)]^T$ is a hidden layer with neurons h_1, \dots, h_M , where $h_j : \mathbb{R}^N \rightarrow \mathbb{R}$ (corresponding outputs are z_1, \dots, z_M and $z_j = h_j(x)$), by the chain rule, the multiplier $m_{x_i f}$ is calculated as :

$$m_{x_i f} = \sum_j m_{x_i h_j} m_{z_j g}. \quad (2)$$

Therefore, the output prediction can be decomposed and backpropagated to each feature dimension. DeepSHAP approximation replaces reference input r in DeepLIFT with $E[x]$ and equates ϕ_i to $m_{x_i f} \Delta x_i$. By chain rule and linear approximation, DeepSHAP calculates Shapley value as:

$$\begin{aligned} m_{x_i f} &= \sum_{j_1} m_{x_i h_{j_1}^{(1)}} \sum_{j_2} m_{z_{j_1}^{(1)} h_{j_2}^{(2)}} \cdots \sum_{j_L} m_{z_{j_{L-1}}^{(L-1)} h_{j_L}^{(L)}} \quad (\text{Chain rule}), \\ \phi_i^x &= m_{x_i f} (x_i - E[x_i]), \end{aligned} \quad (3)$$

where $h_{j_l}^{(l)}$ and $z_{j_l}^{(l)}$ are the j_l th hidden neuron function and output at the l th layer, for $l \in \{1, \dots, L\}$ and $\sum_{j_L} h_{j_L}^{(L)}(z_{j_{L-1}}^{(L-1)}) = y = f(x)$. The *Rescale* and *RevealCancel* mechanisms [4] in DeepLIFT can make sure that the attributes are correctly propagated to the input. Since the Shapley values for the simple network components can be efficiently solved analytically if they are either linear, max pooling, or an activation function with just one input, this composition rule enables a fast approximation to the original Shapley value for the whole model.

3 Proposed Approach

3.1 DeepSHAP with Reference Distribution

In DeepSHAP, the “missing” features are set to the sample mean of the dataset. Hence, the estimated Shapley value may change when the given dataset changes. When the reference distribution (empirical distribution of the dataset) is given, we show that a more accurate approach is to obtain Shapley values for samples from the reference distribution and name it as distribution-based DeepSHAP (DistDeepSHAP). The improvement over the original DeepSHAP can be easily

seen, because we can rewrite the Shapley value as an average over contributions from all the reference samples.

Next, we introduce an alternative formulation of the Shapley value. Let $\pi(\mathcal{N})$ be the set of all ordered permutations of \mathcal{N} . Let $Pre^i(O)$ be the set of players which are predecessors of player i in the order $O \in \pi(\mathcal{N})$, we have

$$\phi_i^x = \frac{1}{N!} \sum_{O \in \pi(\mathcal{N})} (v_x(Pre^i(O) \cup \{i\}) - v_x(Pre^i(O))). \quad (4)$$

Given a single reference \hat{x} sampled from data distribution \mathcal{D} and using $f(x)$ to replace v_x , we define the single reference SHAP as

$$\phi_i^{x|\hat{x}} = \frac{1}{N!} \sum_{O \in \pi(\mathcal{N})} \{f(\tau(x, \hat{x}, Pre^i(O) \cup \{i\})) - f(\tau(x, \hat{x}, Pre^i(O)))\}, \quad (5)$$

where

$$\tau(x, \hat{x}, P) = (v_1, v_2, \dots, v_N), \quad v_j = \begin{cases} x_j, & j \in P \\ \hat{x}_j, & j \notin P \end{cases} \quad (6)$$

Since the reference \hat{x} is a random variable, $\phi_i^{x|\hat{x}}$ is an induced random variable by \hat{x} . We estimate Eq. (4) by:

$$\phi_i^x = \mathbb{E}_{\hat{x} \in \mathcal{D}} [\phi_i^{x|\hat{x}}], \quad (7)$$

It is obvious that Eq. (5) is a typical Shapley value format as Eq. (4). In practice, we can borrow the efficient approximation to Shapley value based on DeepSHAP (Eq. (3), (6)). The, a single reference SHAP value (Eq. (5)) can be efficiently computed as:

$$\phi_i^{x|\hat{x}} = m_{x_i f}(x_i - \hat{x}_i). \quad (8)$$

3.2 Uncertainty of Shapley Values

To estimate the uncertainty of Shapley values, we collect samples of $\phi_i^{x|\hat{x}}$ based on randomly drawn \hat{x} and use the percentiles to measure its uncertainty, u_i , which is associated with each individual Shapley value ϕ_i^x produced by Eq. (8). In particular, we calculate confidence intervals $CI_{i,\gamma} = [c_{i,\frac{\gamma}{2}}, c_{i,1-\frac{\gamma}{2}}]$ with lower bounds $c_{i,\frac{\gamma}{2}}$ and upper bounds $c_{i,1-\frac{\gamma}{2}}$ at confidence level $\gamma = 1 - \alpha$ for each assigned Shapley value. The detailed methods are presented in Algorithm 1. When the training dataset is not available, we can use the testing dataset to approximate the reference distribution.

3.3 Evaluation of Uncertainty

In general settings, it is difficult to evaluate uncertainty estimates for feature importance estimation methods, since we typically do not have per-feature ground-truth to evaluate against. To quantitatively and qualitatively assess the accuracy of the uncertainty estimates provided by DistDeepSHAP, we propose a

Algorithm 1 Estimating Shapley value ϕ_i^x and its uncertainty u_i^x for the i th feature given input x

Input: x , a given instance; R , number of repeats; y , prediction model output with respect to x ; \hat{X} , a set of \hat{x} (i.e. training data); and Φ_i^x , a list to store $\phi_i^{x|\hat{x}}$ calculated by different samples; $1 - \alpha$, confidence level.

```

1:  $\Phi_i^x \leftarrow$  empty list
2: for  $r = 1$  to  $R$  do
3:   choose a random instance  $\hat{x}$  from  $\hat{X}$ 
4:    $\phi_i^{x|\hat{x}} \leftarrow m_{x_i y}(x_i - \hat{x}_i)$   $\triangleright m_{x_i y}$  is calculated by Eq. (3)
5:   Add  $\phi_i^{x|\hat{x}}$  to list  $\Phi_i^x$ 
6: end for
7:  $\phi_i^x \leftarrow \text{MEAN}(\Phi_i^x)$ 
8:  $u_i^x \leftarrow c_{i, 1-\frac{\alpha}{2}} - c_{i, \frac{\alpha}{2}}$ 
9: Output:  $\phi_i^x, u_i^x$ 

```

calibration-based method. The intuition is that the uncertainty estimated on a single input x should reflect the distribution of important scores of all the instances within the class of x . Specifically, we randomly sample m instances within the same class of x and estimate the feature importance scores of those instances. Intuitively, calibration means that given the $1 - \alpha$ confidence interval of a feature, the feature importance scores should occur within this interval with probability $1 - \alpha$.

4 Experiments and Results

4.1 Validation on MNIST Dataset

In order to show the feasibility of our proposed method, we test the explanation results on the MNIST dataset [12], which is intuitive for human judgment. We flattened the MNIST images to vectors. Denote dropout layer with ratio 0.5 as 'Drop', Relu non-linear activation as 'Relu' and fully-connected layer as 'FC'. Then, we trained a Multilayer Perceptron (MLP) classifier (Input(784) \rightarrow Drop(Relu(FC(512))) \rightarrow Drop(Relu(FC(512))) \rightarrow FC(10)) achieving 97.12% accuracy. Given the pre-trained classifier, we compare the feature importance estimation results using DistDeepSHAP with the alternative methods (including GuideBackProp [13], Integrated Gradient [2], DeepLIFT [4] and DeepSHAP [6]) and the corresponding uncertainty estimation results using our proposed DistDeepSHAP in Fig. 1. We set $\alpha = 0.1$ and repeat sampling times $R = 100$ ¹. DeepLIFT uses value 0 as reference. Different from DeepLIFT, DeepSHAP and DistDeepSHAP, GuideBackProp and Integrated Gradient do not need a reference. As we expected, our proposed DistDeepSHAP achieved similar feature estimation results to DeepSHAP. However, DistDeepSHAP can provide additional uncertainty estimations of importance score as shown in Fig. 1 (h) and

¹ The investigation on the repeat sampling times is left in the Appendix

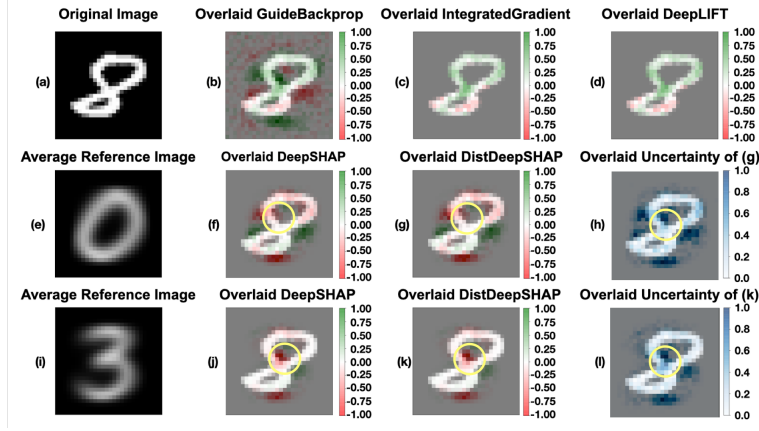


Fig. 1: Explanation results on one digit ‘8’ shown in (a). The color-bar next to each subfigure indicates the scale. (b) (c) (d) show the overlaid pixel importance score to identify this ‘8’ using GuideBackProp [13], Integrated Gradient [2] and DeepLIFT [4] separately. For SHAP-based methods, we can choose another class as the reference. With digits ‘0’s (e) as reference, (f) shows the overlaid Shapley values using DeepSHAP and (g) shows the mean of SHAP scores estimated by our proposed DistDeepSHAP, (h) provides the uncertainty of the Shapley values computed by DistDeepSHAP. With digits ‘3’s (i) as reference, (j) shows the overlaid Shapley values using DeepSHAP and (k) shows the mean of SHAP scores estimated by our proposed DistDeepSHAP, (l) provides the uncertainty of the Shapley values computed by DistDeepSHAP. All values are normalized to range [0,1].

Fig. 1(l). For human interpretation, “central cross” (circled in Fig. 1) is an important pattern to identify digit 8. However, SHAP-based methods marked the central part as negative evidence, which disagrees with human interpretation, and also was different from the results using GuideBackProp, Integrated Gradient and DeepLIFT. We noticed that Fig. 1(h) and Fig. 1(l) assigned the central part high uncertainty values. Therefore, we could use this additional uncertainty information to help judge whether the interpretation result is reliable.

4.2 Application on Autism Resting-state fMRI

Data and Preprocessing The study was carried out using resting-state fMRI (rs-fMRI) data from the Autism Brain Imaging Data Exchange dataset (ABIDE I preprocessed, [14]). We downloaded Regions of Interests (ROIs) fMRI series of the top four largest sites (UM1, NYU, USM, UCLA1) from the preprocessed ABIDE dataset with Configurable Pipeline for the Analysis of Connectomes (CPAC), band-pass filtering (0.01 - 0.1 Hz), no global signal regression, parcellated by Harvard-Oxford (HO) atlas. Skipping subjects with missing files, we downloaded 106, 175, 72, 71 subjects from UM1, NYU, USM, UCLA1 separately. For data augmentation, we used sliding window with length 32 and stride 1 to truncate fMRI series.

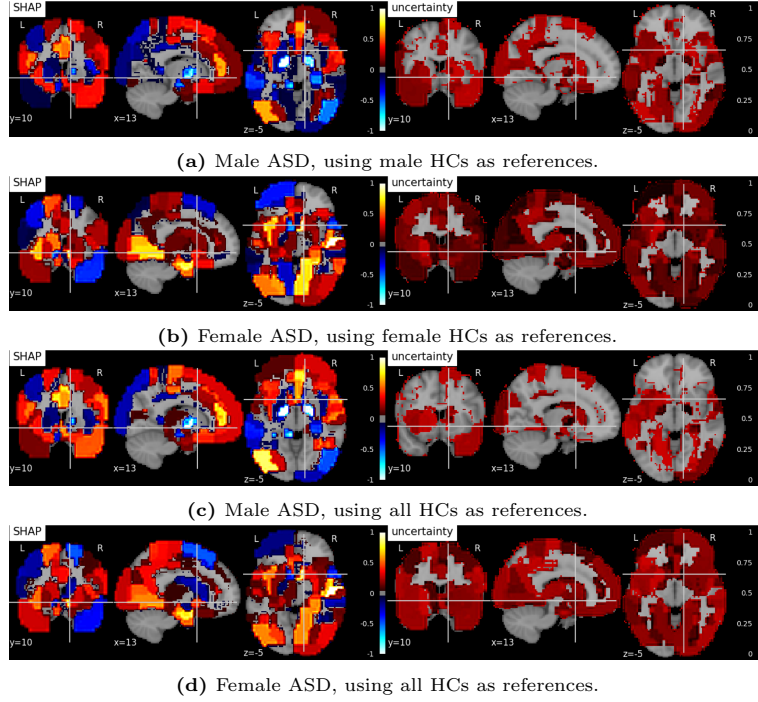


Fig. 2: Shapley values and their corresponding uncertainty (normalized to range $[0,1]$).

Implementation Details The task we performed on the ABIDE datasets was to identify autism spectrum disorders (ASD) or healthy control (HC). We used the mean time sequences of ROIs to compute the correlation matrix as functional connectivity. The functional connectivity provided an index of the level of co-activation of brain regions based on the time series of rs-fMRI brain imaging data. Each element of the correlation matrix was calculated using Pearson correlation coefficient and Fisher transformation. As the correlation matrices were symmetric, we only kept the upper-triangle of the matrices and flattened the triangle values to vectors. The number of resultant features was defined by $N(N-1)/2$, where N was the number of ROIs. Under the HO atlas (110 ROIs), the procedure resulted in 5995 features. We designed an MLP (Input(5995) \rightarrow Drop(ReLu(FC(8))) \rightarrow FC(2)) to classify the inputs. We randomly split 80% of the data as training data and the remaining 20% as testing data based on subjects. Adam optimization was applied with initial learning rate $1e-5$ and reduced by $1/2$ for every 20 epochs and stopped at the 50th epoch. We achieved 69.04% accuracy on the testing set.

Feature Importance Estimation Results Given the pre-trained ASD vs. HC classifier, we showed the ROI importance results on a male ASD and a female ASD instance with different reference distributions ($\alpha = 0.1$ and $R = 100$)

in Fig. 2. We interpreted the ROIs with high Shapley values as biomarkers for a certain group. Specifically, Fig. 2 showed the biomarkers to distinguish ASD in a certain sex group from different reference population. The results indicate that an ROI may have different importance scores when compared to different references. The uncertainty results presented on the right column of Fig. 2 gave us the confidence of trusting the biomarker detection results. For example, both Fig. 2(a) and Fig. 2(c) indicated high Shapley values and low uncertainty on frontal cortex. Existing clinical studies [15,16,17] have shown that frontal cortex is a salient biomarker in ASD reflected by reduced levels of GABA and reduced MeCP2 expression. In neurological biomarker detection studies, investigators could take the uncertainty information and assign higher priority to the trustable important ROIs (with large Shapley score and low uncertainty) in clinical investigation, as conducting clinical biomarkers evaluation without prior knowledge is costly and time consuming. With the flexibility of choosing subgroup references, our proposed DistDeepSHAP could explore the fine-grained neurological biomarkers of subgroups. The almost opposite importance scores for the male and the female instances showed the heterogeneous ASD neurological patterns in different genders [18,19].

4.3 Evaluation of Uncertainty on the Two Experiments

As we described in Section 3.3, a reasonable uncertainty measurement has calibration property. Namely, the uncertainty estimation can imply how the other instances agree on the importance score calculated from x , and how certain the feature importance estimates are on previously unseen sample images. Hence uncertainty can be used as out of distribution data importance score calibration, and an accurate uncertainty estimate of a given testing sample should be correlated with the Shapley value estimates on the **held-out** testing set. For example, the best uncertainty estimates with given significance level $\alpha = 0.1$ will include $1 - \alpha = 90\%$ of the testing samples whose Shapley values fall into this $1 - \alpha$ confidence interval. Given the estimation on a digit ‘8’ input and a male ASD input, we sampled another 100 digit ‘8’s and 10 male ASDs from the testing set of the two experiments separately. We used digit ‘3’s as the references for digit ‘8’s and used male HCs as the references for male ASDs. We applied the uncertainty evaluation method

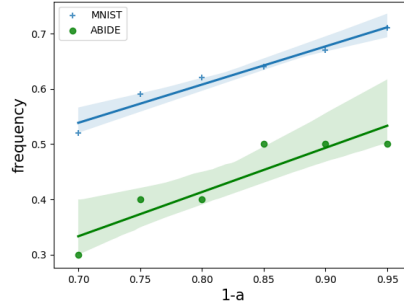


Fig. 3: Relationship between $1 - \alpha$ and percentage of held-out testing samples falling into the confidence interval for uncertainty evaluation.

proposed in Section 3.3 for the given testing samples. The top 10% uncertain features were selected. For a held-out testing sample, only if its Shapley value of the top 10% uncertain features fall in the estimated confidence interval of the previous given testing sample, we regarded the held-out testing sample as success. We varied $1 - \alpha \in \{0.7, 0.75, 0.8, 0.85, 0.9, 0.95\}$. The relationship between $1 - \alpha$ and the percentage of the held-out testing samples that are located in the confidence interval, was shown in Fig. 3. The correlations of $1 - \alpha$ and this percentage were 0.981 for MNIST and 0.916 for ABIDE respectively. The uncertainty estimated for the classification model on the MNIST dataset had higher correlation and closer percentage to $1 - \alpha$. The reason could be that the classification model on the MNIST dataset achieved higher accuracy and we had more held-out testing samples than those of ABIDE. Overall, the calibration patterns could validate the reliability of our proposed uncertainty estimates.

5 Conclusion

In this work, we propose DistDeepSHAP, a post-hoc feature importance estimation method with uncertainty evaluation for deep learning models. DistDeepSHAP is based on the idea of DeepSHAP, but improves DeepSHAP by sampling the references from a distribution and calculating Shapley values for these references. Our proposed DistDeepSHAP has several advantages over DeepSHAP. First, it can obtain uncertainty estimates for the provided feature importance scores. Second, it can better utilize the empirical reference distribution and has the potential for better feature importance score estimation. Last but not least, it can be calculated with arbitrary subgroup references and interpret salient features with respect to a subgroup, which is crucial for neuroscience study.

References

1. K. Simonyan *et al.*, “Deep inside convolutional networks: Visualising image classification models and saliency maps,” *arXiv preprint arXiv:1312.6034*, 2013.
2. M. Sundararajan *et al.*, “Axiomatic attribution for deep networks,” in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 3319–3328, JMLR. org, 2017.
3. G. Montavon *et al.*, “Methods for interpreting and understanding deep neural networks,” *Digital Signal Processing*, vol. 73, pp. 1–15, 2018.
4. A. Shrikumar *et al.*, “Learning important features through propagating activation differences,” in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 3145–3153, JMLR. org, 2017.
5. M. T. Ribeiro *et al.*, “‘‘ why should i trust you?’’ explaining the predictions of any classifier,” in *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144, 2016.
6. S. M. Lundberg and S.-I. Lee, “A unified approach to interpreting model predictions,” in *Advances in neural information processing systems*, pp. 4765–4774, 2017.

7. P. Schwab and H. Hlavacs, "Capturing the essence: Towards the automated generation of transparent behavior models," in *Eleventh Artificial Intelligence and Interactive Digital Entertainment Conference*, 2015.
8. L. S. Shapley, "A value for n-person games," *Contributions to the Theory of Games*, vol. 2, no. 28, pp. 307–317, 1953.
9. P.-J. Kindermans *et al.*, "The (un) reliability of saliency methods," in *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pp. 267–280, Springer, 2019.
10. H. Fen *et al.*, "Why should you trust my interpretation? understanding uncertainty in lime predictions," *arXiv preprint arXiv:1904.12991*, 2019.
11. J. Adebayo *et al.*, "Sanity checks for saliency maps," in *Advances in Neural Information Processing Systems*, pp. 9505–9515, 2018.
12. Y. LeCun *et al.*, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
13. J. T. Springenberg *et al.*, "Striving for simplicity: The all convolutional net," *arXiv preprint arXiv:1412.6806*, 2014.
14. A. Di Martino *et al.*, "The autism brain imaging data exchange: towards a large-scale evaluation of the intrinsic brain architecture in autism," *Molecular psychiatry*, vol. 19, no. 6, p. 659, 2014.
15. A. A. Goldani *et al.*, "Biomarkers in autism," *Frontiers in psychiatry*, vol. 5, p. 100, 2014.
16. R. Nagarajan *et al.*, "Reduced mecp2 expression is frequent in autism frontal cortex and correlates with aberrant mecp2 promoter methylation," *Epigenetics*, vol. 1, no. 4, pp. 172–182, 2006.
17. T. Watanabe *et al.*, "Mitigation of sociocommunicational deficits of autism through oxytocin-induced recovery of medial prefrontal activity: a randomized trial," *JAMA psychiatry*, vol. 71, no. 2, pp. 166–175, 2014.
18. T. T. Rivet and J. L. Matson, "Review of gender differences in core symptomatology in autism spectrum disorders," *Research in Autism Spectrum Disorders*, vol. 5, no. 3, pp. 957–976, 2011.
19. A. K. Halladay *et al.*, "Sex and gender differences in autism spectrum disorder: summarizing evidence gaps and identifying emerging areas of priority," *Molecular autism*, vol. 6, no. 1, pp. 1–5, 2015.