

Reconnoitering the class distinguishing abilities of the features, to know them better

Payel Sadhukhan^a, Sarbani Palit^b, Kausik Sengupta^a

^aInstitute for Advancing Intelligence, TCG CREST, Kolkata, 700091, West Bengal, India

^bComputer Vision and Pattern Recognition Unit, Indian Statistical Institute, Kolkata, 700108, West Bengal, India

Abstract

The relevance of machine learning (ML) in our daily lives is closely intertwined with its explainability. Explainability can allow end-users to have a transparent and humane reckoning of a ML scheme's capability and utility. It will also foster the user's confidence in the automated decisions of a system. Explaining the variables or features to explain a model's decision is a need of the present times. We could not really find any work, which explains the features on the basis of their class-distinguishing abilities (specially when the real world data are mostly of multi-class nature). In any given dataset, a feature is not equally good at making distinctions between the different possible categorizations (or classes) of the data points. In this work, we explain the features on the basis of their class or category-distinguishing capabilities. We particularly estimate the class-distinguishing capabilities (scores) of the variables for pair-wise class combinations. We validate the explainability given by our scheme empirically on several real-world, multi-class datasets. We further utilize the class-distinguishing scores in a latent feature context and propose a novel decision making protocol. Another novelty of this work lies with a *refuse to render decision* option when the latent variable (of the test point) has a high class-distinguishing potential for the likely classes.

Keywords: feature importance, explainability, class distinguishing ability, decision to refuse or render

1. Crux of the work—A toy example

We have images of hand-written digits. Given that we have 10 digits, we define 10 classes or categories. We assign each digit to a unique class – class 0 corresponds to the images with hand-written 0, class 1 corresponds to the images with hand-written 1. We assume that, in an image we will find only one hand-written digit. The task is to classify each image to its respective class. For example, an image with a handwritten 1 should be classified to class 1 and an image with a handwritten digit 2 should be classified to 2. Each image in a sub-figure has four quadrants (top-left, top-right, bottom-left and bottom-right). Each quadrant in a image represents a feature. Now, to better understand the problem being addressed, consider Figure 1. There are six sub-figures in Figure 1, each of which shows 2 images of hand-written digits which we want to distinguish. In subfigures (a) and (d), the full images are visible. In subfigures

(b) and (e), we conceal the upper-right quadrants of images in (a) and (d) respectively. In subfigures (c) and (f), we conceal the bottom-left quadrants of images (a) and (d) respectively. For the sake of explaining the problem being addressed, we deal with two particular cases here— i] **Case 1:** subfigures (a), (b) and (c) – distinguishing 7 and 1 and ii] **Case 2:** subfigures (d), (e) and (f) – distinguishing 5 and 6.

In Case 1, we have to make a choice between 7 and 1. Sub-figure (b) shows that *top right quadrant* plays a decisive role in making that distinction. Sub-figure (b) shows that the absence of top-right quadrant information makes 7 and 1 indistinguishable to a large extent. If, for a given image categorization, the task is to decide between 1 and 7 and the information from the top-right quadrant is unavailable, *the judicious choice would be to refuse to predict a category*. On the contrary, a missing information from the bottom right quadrant (subfigure (c)) is not bothersome and does not pose any hindrance in distinguishing 1 and 7, .

In Case 2 (subfigures (d)-(f)), we have to make a distinction between 5 and 6. An analysis similar to Case 1 shows that the missing information from the bottom-left quadrant (sub-figure (f)) is bothersome here.

Case 1 and Case 2 reveal that the information from the top-right quadrant is instrumental in distinguishing 1 and 7, but is not as effective in distinguishing 5 and 6. On the contrary, the feature from the bottom-left quadrant is instrumental in distinguishing 5 and 6, but not very effective in discriminating between 1 and 7. The crux is – the efficaciousness of a feature in a classification task is dependent on the classes which it has

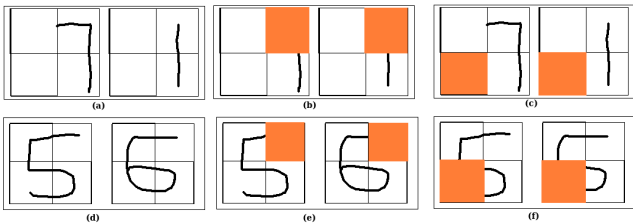


Figure 1: A toy example to present a visual illustration of the research problem which we have addressed in this work.

to distinguish. To address this aspect, we need an explainability and quantification (or an explainable quantification) of the class-distinguishing abilities of the features. *The basic objective of our work is to obtain the explainability of the features in terms of their class distinguishing capabilities.*

2. Introduction

Nowadays, *machine learning* is an inevitability in our daily life – and its relevance is increasing with each passing day [1, 2, 3]. A machine learning model usually learns from prior or available data, generalizes them to learn a model and uses it to make predictions and categorization of an unknown instance [4]. In particular, the research community has been instrumental in tailoring the ML models and real-world datasets for one another to enable the extraction of meaningful information from the data [5, 6, 7]. Sincere efforts led to the digital revolution and ML models have become indispensable in various spheres of our lives from medical applications [8, 9, 10], construction purposes [11], investigation of automobile crashes [12], finance sector [13] and analyzing student attrition [14]. A key area of application of ML is the automated decision making or decision support systems. To increase the trust and faith of the users in automated decision making, we need to work on the explainability of the decisions and interpretations of the models. The impact of the wrong decisions by a model is a stakeholder in this application. As the automated decision making extends to sensitive domains like clinical decision support systems [8, 9, 10] and finance-banking sectors [15] the wrong decisions of ML model will be critical. Consequently, the veracity of a model’s decision plays a key role in its acceptance. In such a scenario, a human reckoning or an interpretability of the decision making systems is highly desirable to maintain and boost the trust on the systems.

This gave rise to a new need — to study the rationale behind the working of the ML models and *explainable machine learning* has come up as the call for the day. In general, ML models are perceived as *black boxes* – where we input some unseen data (that we want to predict), learning goes on inside the *opaque* black box and we get an output which is the prediction or categorization of the input data [16, 17]. To build the credence and confidence of the end users on the ML models, we will need an interpretable framework where we will have reduced opacity and can explain the causation of the decisions rendered by the models.

One way to understand an ML model is through the features on which it is trained – how each feature is influencing the decisions of the model? In recent years, a number of works have focused on explaining the association of the range of features and the outcomes from the models [18, 19, 20, 21, 22]. They have particularly analyzed the influence of the feature values on the predictive outcomes of the test points. In recent years, studies have been carried out to figure out the *interpretability* of the interpretable or explainable models [23, 24]. The focus is on the degree of interpretability achieved from the interpretable models by the data scientists. The studies report the over-trust and misuse of some of the popular interpretable tools [22, 21]

to a considerable extent. This indicates that we require some more intuition into the explainabilities of the features. A more detailed analysis of the features can be more bearing in this regard.

In this work, we present a novel and intuitive protocol of explaining the features through their class-distinguishing abilities. To the best of our knowledge, we could not find any extant work which has followed this line of thought. We are particularly interested to study the *role or capabilities* of each feature in discriminating between the different classes for a dataset. It is because – *the distinguishing capabilities of the features will influence the distinguishing capability of the model.* Hence, if we can explain and estimate the distinguishing capabilities of the features, we will be able to explain the decisions of the models. For each feature, we inspect its distinguishing capability for all possible pair-wise class combinations. For c classes, there can be cC_2 pairwise classes and we explore and quantify cC_2 class combinations. Technically, we employ a latent feature framework to obtain the class-distinguishing scores of the features.

For example, let us have 4 features a, b, c and d and 3 classes 1, 2 and 3. We are looking for the following answers — *How effective is feature a in distinguishing between class 1 and class 3? Is it equally effective in distinguishing between class 1 and class 2 also?* We look for similar answers for all features. Once we have the answers, we have explainability on the class-distinguishing capacities of the features. In this paper, we work towards a obtaining a quantitative estimate of the class distinguishing abilities of the features. In our scheme, it is possible to compute the gross importance of each individual feature through their class-distinguishing scores. We may emphasize here that the focus of this work is beyond the gross importance computation and we want to delve into their class-distinguishing aspects.

After addressing the research aspect stated above, we explore another perspective which we believe is intertwined with the one that we have already talked about. The main objective of the second research problem is – validating the correctness of the explainable class distinguishing scores (which we obtain from the first exploration) through a classification task. We use the same framework (operating under latent or missing feature constraint) to design a novel classifier model. The novelty of this classifier is – *it will render a decision on a test instance only when it is confident of doing so. It will refuse to render a decision when it is not confident enough.* The framework is explained briefly as follows. Let feature a be absent (latent feature) for some test point \mathbf{q} . We obtain (have the provision in our framework) an intermediate output from the model which tells us that \mathbf{q} either belongs to class 1 or to class 3. The question is – **Would we like to classify it?** If feature a is instrumental in distinguishing between class 1 and class 3, we would probably like to deter ourselves from rendering a prediction. On the contrary, if feature a not that instrumental in distinguishing between class 1 and class 3, we would like to proceed with the prediction. We incorporate this line of thought into our scheme. We have carried out the experimental study on 5 real-world datasets. The empirical findings indicate that our scheme’s decision on withholding or rendering a decision is mostly right

across all datasets. The difference in accuracies of the two scenarios *rendering a prediction* and *withholding the predictions* is a significant figure for all five datasets – the accuracy has been much more in the case where the model has rendered a decision than that of the cases where the model withheld the decision. This in turn validates the explainability of the class-distinguishing capabilities of the features. We also present case studies on two datasets – *maternal health* and *contraceptive*. We semantically analyze the features of these datasets and correlate them with our quantitative finding.

The rest of this article is arranged in the following manner. In Section 3, we describe the extant works in this domain. We elaborate on the key aspect of this work, *class-discernible scores* in the Section 4 followed by a detailed presentation of the Proposed Approach in Section 5. We describe the Empirical Setup in Section 6 and discuss the Results and Analysis in Section 7. We wrap up this article with Conclusion and scope of Future Work in Section 8.

3. Related Works

The issue which is being addressed in this work is the interpretation of the features in terms of their class distinguishing capabilities. As mentioned earlier, we could not really find any work which addresses this matter but a number of researches have been carried out in feature explainability and interpretability in general. We present an account of the prior works in explainability in machine learning followed by a focus on feature explainability in the later part.

The inevitability of ML has lead to the requirement of its explanations which can be understood by the humans. Consequently, explaining the decisions of algorithms which can influence human lives has come up as the call for the day [25, 26]. Recent researches in diverse domain like property-value prediction, breast cancer profiling [8], fake news detection [27], glass transition temperature prediction [28], point cloud classification [29] has particularly focused on explainability of an algorithm’s decision besides solving the problem. In [30], an interactive machine learning framework is adopted to cater to the need for explainability. From a ML perspective, feature explainability is usually of one of the two types – i] *local*, or ii] *global*. A local explanation deals with the rationale of a model’s decision on some particular input. On the other hand, a global explanation usually renders the rationale behind the decision through a summarization, which is usually independent of any individual input [31]. A key work in the domain of local explanation is LIME [20] and in the domain of global explanation is SHAP [21] – these works deal with the importance or contributions of the variables while rendering a decision. Other notable works are [32]. These methods are model-agnostic and is not specific to any classifier model. Apart from these, model-specific research has also been catered to – a deep taylor based decomposition for neural networks and tree interpreter for random forests [33], understanding deep neural network through prototype and typicality [34]. In [35, 36], explainable classifiers are proposed which is the need of the present times. The nature of explainability of ML models and the features across

domains. In recent years, a number of works on feature explainability have focused on specific domains like loan underwriting [15] and business model design and sustainability [37]. Explainability of features is particularly essential in critical areas like decision support systems in clinical domain. A few notable works in this area are – diabetes mellitus prediction [9], non-communicable disease prediction [38], clinical decision support system in medical imaging [10], renal mass calcification [39]. [40] has worked towards reducing the opacity of the features and obtained interpretable features with the assistance from the prediction of the black-box models.

In this work, we learn the class-distinguishing abilities of the features through a framework of imputed features. The paradigm of imputed features has existed in the machine learning domain for long and it has been perceived as a constraint in the usual learning [41, 42]. This learning paradigm has also been used for in some works which deal with explainability. Some notable ones are on explainable anatomical shape analysis [43], explainable prediction of electrical energy demand [44] and explainable recommendations [45]. In this work, we use the paradigm of latent variable or missing feature in a novel way – to learn the class-distinguishing capabilities of the features.

4. Class-distinguishing scores

We are talking about classifiers in this work. We are going to analyse the class-distinguishing capabilities of the classifier models. The motivation of this work is – a classifier model may not be equally good at distinguishing the pairs of possible classes in a dataset. On a micro-scale, we will analyze the class-distinguishing capabilities of the features individually. For example, let there be three possible classes of *cat*, *fish* and *boat* in a dataset. We may find that the classifier is effective towards distinguishing between *boat* and *cat* but not as good while distinguishing *cat* and *fish*. Instead of providing the explainability of a model’s decision as a function of the feature aggregates, we will try to explain and quantitatively capture the class-distinguishing capabilities of individual features. It is because – features are the building blocks of a classification model.

4.1. Classifiers used in this work

We assume the number of features in our data to be γ . We will require $(\gamma + 1)$ classifiers in our scheme. We will have a dedicated classifier for each feature, γ classifiers accounting from that. The remaining one classifier will be a generic one, learnt from all the features. Let us denote the classifiers as $M_0, M_{-1}, M_{-2}, \dots, M_{-\gamma}$. We use all γ features to model classifier M_0 . Classifiers M_{-a} , ($0 < a \leq \gamma$) is modelled by imputing or removing feature a . That is, in classifier M_{-a} , a is the latent feature. From the given dataset, we remove or impute feature a , and use the remaining $(\gamma - 1)$ features to train model M_{-a} , ($0 < a \leq \gamma$).

4.2. Obtaining the class-distinguishing scores

Let us assume that there are c classes. We will obtain a class-distinguishing score for each pair of classes α and β , $\alpha \neq \beta$, $1 \leq \alpha, \beta \leq c$. Let $CS_{\alpha\beta}^{-a}$ be the class-distinguishing score of feature a with respect to classes α and β . We would like to compute the scores from the classification output of n instances. We denote class_i to be the true class of instance i and pred_i^{-a} denotes the prediction of instance i obtained from model M_{-a} .

For each classifier M_{-a} , $0 \leq a \leq \gamma$ and a pair of classes α and β , we will first obtain T_{α}^{-a} , it denotes that number of instances correctly classified to class α by classifier M_{-a} . Similarly, T_{β}^{-a} or T_{β}^{-a} denotes the number of instances correctly classified to class β by classifier M_{-a} respectively. $F_{\alpha\beta}^{-a}$ denotes the number of instances which belong to class β but has wrongly been classified to class α by classifier M_{-a} . Similarly, $F_{\alpha\beta}^{-a}$ denotes the number of instances which belong to class α but has wrongly been classified to class β by classifier M_{-a} .

$$\begin{aligned} T_{\alpha}^{-a} &= \sum_{i=1}^n \{\text{pred}_i^{-a} = \alpha \cap \text{class}_i = \alpha\} \\ T_{\beta}^{-a} &= \sum_{i=1}^n \{\text{pred}_i^{-a} = \beta \cap \text{class}_i = \beta\} \\ F_{\alpha\beta}^{-a} &= \sum_{i=1}^n \{\text{pred}_i^{-a} = \alpha \cap \text{class}_i = \beta\} \\ F_{\alpha\beta}^{-a} &= \sum_{i=1}^n \{\text{pred}_i^{-a} = \beta \cap \text{class}_i = \alpha\} \end{aligned} \quad (1)$$

Based on these four sets of values, we will compute the class-distinguishing score, $CS_{\alpha\beta}^{-a}$ for each pair of classes (α, β) , $\alpha \neq \beta$, $1 \leq \alpha, \beta \leq c$ specific to each feature a , $1 \leq a \leq \gamma$. Two intermediate scores, $\text{pre}_{\alpha\beta}^{-a}$ and $\text{rec}_{\alpha\beta}^{-a}$ are obtained from the parameters obtained in Equation (1). We calculate the class-discernibility score, $CS_{\alpha\beta}^{-a}$ from these two intermediate values. The calculation is substantially similar to that of F_1 . The difference here is, we account for the *true positives* of the two classes together.

$$\begin{aligned} \text{pre}_{\alpha\beta}^{-a} &= \frac{T_{\alpha}^{-a} + T_{\beta}^{-a}}{T_{\alpha}^{-a} + T_{\beta}^{-a} + F_{\alpha\beta}^{-a}} \\ \text{rec}_{\alpha\beta}^{-a} &= \frac{T_{\alpha}^{-a} + T_{\beta}^{-a}}{T_{\alpha}^{-a} + T_{\beta}^{-a} + F_{\beta\alpha}^{-a}} \\ CS_{\alpha\beta}^{-a} &= \frac{2 \times \text{pre}_{\alpha\beta}^{-a} \times \text{rec}_{\alpha\beta}^{-a}}{\text{pre}_{\alpha\beta}^{-a} + \text{rec}_{\alpha\beta}^{-a}} \end{aligned} \quad (2)$$

When $a = 0$, training happens on the full and complete dataset. We denote the class-distinguishing score for (α, β) of the classifier model trained on the full complete data as $CS_{\alpha\beta}^0$. For $1 \leq a \leq \gamma$, we remove feature a and train the model. Subsequently, while using model M_{-a} for predicting a test point, feature a has to be imputed. As we have said earlier, in a similar fashion, we train γ classifiers $M_{-1}, M_{-2}, \dots, M_{-\gamma}$, each of

which is trained by imputing one feature at a time. We calculate their class distinguishing scores as well. $CS_{\alpha\beta}^{-a}$ denotes the class-distinguishing score of classifier model M_{-a} w.r.t. classes (α, β) .

For c classes, the class distinguishing score for each feature will be a $c \times c$ matrix. The symmetricity of the matrix will depend on whether we select a symmetric score function or not. Let us assume a symmetric score function as of now.

In the next section, we present an detailed explanation of the proposed methodology.

5. Proposed Approach

In this work, we are motivated to estimate the goodness and utility of the features of a dataset. It boils down to — what role each feature is playing in the predictions. The class-distinguishing powers of the features are a good indicator of their predictive powers. To answer this question, we have to engage in a micro-analysis of the features with respect to their class-distinguishing capabilities. In the next subsection, we will introduce and explain a scheme for estimating the predictive capability or *goodness* of a feature.

5.1. Goodness of features

As in the previous subsection, we will assume that our dataset has γ features and c classes. Let the features be denoted as $f_1, f_2, \dots, f_{\gamma}$. Let \mathcal{D} be the full and complete dataset and \mathcal{D}_{-a} denote the dataset without feature a . \mathcal{D}_{-a} is the entire dataset consisting of all instances without feature a . We will train a model M_0 with the full and complete dataset \mathcal{D} , evaluate the class-distinguishing scores and use it as yardstick for measuring the goodness and predictive power of each feature. We denote the class-distinguishing score of \mathcal{D} as $CS_{\alpha\beta}^0$.

To evaluate the class-distinguishing capability of feature a , we will train a classifier model M_{-a} with \mathcal{D}_{-a} . Essentially, we are training classifier model M_{-a} without feature a and we will obtain the class-distinguishing scores from the model. The disparity or difference in the class-distinguishing capability of M_0 and M_{-a} will be indicative of the predictive power of feature a . We will explore this capability for all features $f_1, f_2, \dots, f_{\gamma}$. We denote the class-distinguishing score with respect to \mathcal{D}_{-a} for feature a with $CS_{\alpha\beta}^{-a}$.

The normalized class-distinguishing score for feature a with respect to classes α and β is denoted with $NCS_{\alpha\beta}^a$. Note that $1 \leq a \leq \gamma$ where γ denotes the number of features. We compute it as follows.

$$NCS_{\alpha\beta}^a = CS_{\alpha\beta}^0 - CS_{\alpha\beta}^{-a} \quad (3)$$

We may note that $NCS_{\alpha\beta}^a < 0$ is possible. It signifies that removal of feature a causes betterment of the classification performance of the model with respect to classes α and β . In this particular scenario, we should not bother about the absence of feature a for distinguishing between classes α and β . To be precise, these cases

- If $NCS_{\alpha\beta}^a \approx 0$, $\forall \alpha\beta, \alpha \neq \beta$, feature a has no role in differentiating between the classes of a dataset. It is so because the distinguishing capability of the models do not deteriorate without the removal of feature a across all class combinations, which signifies that feature a is rather redundant or unimportant in learning from the dataset.
- Alternatively, let us have $NCS_{\alpha\beta}^a \geq \varepsilon$, $\forall \alpha\beta, \alpha \neq \beta$ where $0 < \varepsilon \leq 1$ is a significant amount with respect to the problem. In this particular case, feature a plays a decisive role in differentiating between all pairs of classes.

We explain the basic intuition of our scheme above. The basic idea is – if the class distinguishing reduce change after removal of some feature a , feature a is important. In the previous two items, we have assumed for all pairs of classes α and β , where $\alpha \neq \beta$. But from Figure 1, we have seen that class distinguishing scores of the are specific to the pair of classes. The feature from the top-right quadrant is effective for distinguishing 1 and 7 but not as effective for distinguishing 5 and 6. The reverse holds for the feature from the bottom-left quadrant. Hence, we will need a class-specific (pair wise class) analysis to have a

Remarks:

- We can calculate the overall importance of a feature from its normalized class-distinguishing scores for all possible classes. The overall importance of a feature is the cumulative sum of all its normalized class-distinguishing scores. Let $I_{overall}^a$ be the overall importance of feature a .

$$I_{overall}^a = \sum_{\alpha=1}^{\gamma} \sum_{\beta=1, \alpha \neq \beta}^{\gamma} NCS_{\alpha\beta}^{-a} \quad (4)$$

- The real-world datasets are usually multi-class (more than two classes). In such cases class-distinguishing scores for a feature is usually more than 2 (3 for 3 classes, 6 for 4 classes and so on). We have to look at the normalized scores, $NCS_{\alpha\beta}^{-a}$ ($\forall a$ features, $\forall \alpha, \beta, \alpha \neq \beta$) features to have a detailed picture. We do that in the next subsection.

5.2. Class-distinguishing capabilities

Now that we know how to evaluate the importance of each feature, we want to use it to evaluate the class-discerning capabilities of the features for pair-wise class combinations. $(CS_{\alpha\beta}^0 - CS_{\alpha\beta}^{-a})$ denotes that capability of feature a to discriminate classes α and β . We have to accumulate this information for all pairs of classes α and β , where $\alpha \neq \beta$ to get a micro-level information. We will remove each feature (making it a latent variable), and for each feature, we will maintain a matrix to organise this information. We may note that the diagonal elements signify the class-distinguishing scores of one class to itself and it do not bear any relevance (hence we mark them as don't care).

Note that: If we consider a symmetric function to calculate the class-distinguishing scores $(CS_{\alpha\beta}^{-a})$ will be equal to $(CS_{\beta\alpha}^{-a})$, and the class-distinguishing matrix for each feature will be a symmetric one.

	class 1	class 2	class 3
class1	Don't care	CS_{12}^0	CS_{13}^0
class 2	CS_{21}^0	Don't care	CS_{23}^0
class 3	CS_{31}^0	CS_{32}^0	Don't care

Table 1: Class-distinguishing scores for the full, complete dataset. CS_{ii}^0 or the diagonal elements do not bear any significance hence it is written as don't care. CS_{12}^0 denotes the capability of M_{-0} to distinguish between classes 1 and 2.

	class 1	class 2	class 3
class1	Don't care	CS_{12}^{-a}	CS_{13}^{-a}
class 2	CS_{21}^{-a}	Don't care	CS_{23}^{-a}
class 3	CS_{31}^{-a}	CS_{32}^{-a}	Don't care

Table 2: Class-distinguishing scores for some a , CS_{ii}^{-a} or the diagonal elements do not bear any significance hence it is written as don't care. CS_{12}^{-a} denotes the capability of M_{-a} to distinguish between classes 1 and 2 (the capability of a model to distinguish between class 1 and class 2 without feature a). If $NCS_{12}^{-a} \geq \varepsilon$, we may say that feature a is instrumental in distinguishing between classes 1 and 2. In such a scenario, if we have a point whose feature a information is missing and we know that the two likely classes are 1 and 2, it is judicious to refuse a decision for that instance.

Further note: The procedure explained above is carried out on the training data. From $(CS_{\alpha\beta}^0 - CS_{\alpha\beta}^{-a})$, we will evaluate the importance of feature a in distinguishing classes α and β . A higher value indicates that feature a plays a decisive role in distinguishing classes α and β . Now, we would like to use this information in a setting where we have to predict test points with some missing features. How do we do that?

5.3. Handling the feature importance to make predictions on test points which may have some missing features

Let us ponder over our last statement in the previous subsection.

- We assume that $(CS_{\alpha\beta}^0 - CS_{\alpha\beta}^{-a})$ is a significant value, which manifests the decisive role of feature a in distinguishing classes α and β .
- The other assumption is — we have to classify a test point \mathbf{q} for which feature a is latent variable.
- We invoke classifier model M_{-a} to classify \mathbf{q} . After obtaining the prediction from M_{-a} , we find that the two most probable classes for \mathbf{q} are α and β in some order (either α more probable than β).

Now that we already know that feature a is decisive in distinguishing between α and β , will it be a judicious decision to distinguish between the two in absence of feature a information? It is certainly not.

We motivate our next line of action on this thought. While classifying or predicting the test points under latent variable constraint (single latent variable), we will obtain the two most probable classes for each point and also the difference in their likelihood according to the classifier. Let the latent variable be

denoted by a and we invoke classifier M_{-a} to classify \mathbf{q} . Let the two probable classes for \mathbf{q} be α and β and their respective likelihood for instance \mathbf{q} be $\mathcal{L}_{\alpha}^{-a}(\mathbf{q})$ and $\mathcal{L}_{\beta}^{-a}(\mathbf{q})$ respectively. Let us assume $\mathcal{L}_{\alpha}^{-a}(\mathbf{q}) > \mathcal{L}_{\beta}^{-a}(\mathbf{q})$.

We use $(CS_{\alpha\beta}^0 - CS_{\alpha\beta}^{-a})$ as the class-distinguishing score of feature a for distinguishing class α and class β . We desegregate this information from the training phase along with $\mathcal{L}_{\alpha}^{-a}(\mathbf{q})$ and $\mathcal{L}_{\beta}^{-a}(\mathbf{q})$ obtained about \mathbf{q} from the test phase to decide whether we make a judgement about the class of \mathbf{q} or not. Let **Prediction**(\mathbf{q}) be the prediction for \mathbf{q} .

Note that, in the first part of Equation (1), we compare two types of terms i] $(\mathcal{L}_{\alpha}^{-a}(\mathbf{q}), \mathcal{L}_{\beta}^{-a}(\mathbf{q}))$ – which is obtained from the test phase and is dependent on test instance \mathbf{q} , and, ii] $(CS_{\alpha\beta}^0 - CS_{\alpha\beta}^{-a})$ – which is obtained solely from the training phase.

The equation signifies that — i] when the difference in likelihood of the two most probable classes is more than the class-distinguishing scores (with respect to those two classes) for the missing feature, *the scheme makes the prediction. The class with higher likelihood is predicted*, ii] When the difference in likelihood is less than the class-distinguishing scores (which is obtained in the training phase), *the scheme refuses to make a prediction*. Hence, we will have two principal classes of prediction — i] where we make a prediction or deliver a decision ii] where we refuse to make a decision. We may note that, when we have a test point with all the features, we do not consider the *refusal to decision* option and we always deliver a decision.

5.4. Were we right to refuse a decision? — and the goodness of the class- distinguishing scores estimated in the training phase

We have the test set denoted by \mathcal{D}_{te} . It has points with full and complete features (no missing features). We consider test points $\mathbf{q}, \mathbf{q} \in \mathcal{D}_{te}$ and impute its feature a , to form a dataset \mathcal{D}_{te}^{-a} , $a \in \{f_1, f_2, \dots, f_{\gamma}\}$. Each point in \mathcal{D}_{te}^{-a} is a point in \mathcal{D}_{te} with feature a removed. We invoke classifier model M_{-a} to obtain the predictions of points in \mathcal{D}_{te}^{-a} .

We consider two classes of cases which we described in the previous subsection and compute their respective accuracies.

- Case 1- where we deliver a decision in spite of having missing features – decisions are made and we compute the accuracies from the decisions (predictions) and true class (category) information of the points. We denote it with $acc_{decision}$.
- Case 2 - where we refuse to deliver a render in presence of a missing feature (because we were not confident of deciding in absence of the feature), we consider the most probable choice of the classifier model as the *tentative prediction*. We compute the accuracy on the basis of the tentative predictions and the true class (category) information of the points. We denote the accuracy with acc_{reject} .

Intuition: We have accuracies from two distinct cases where i] the classifier makes a decision, ii] the classifier is not confident of passing a decision and it refuses a decision.

The decisions to predict or reject are made on the basis of the explainability and class-distinguishing capacities of the features. The explainability or class-discernibility quantification is also a part of the model. *Given this scenario, it is logical to deduce that if the model has a good feature explainability, it should be correct about its decision to predict or reject. Or, in other words, the model's classification accuracy should be more in situations where it is predicting than situations where it is refusing to predict (reject).*

6. Experimental Setup

The empirical study is designed to evaluate two primary goals — i] Compute the class distinguishing capabilities (scores) of the features. ii] Validate and check the goodness of the computed scores by employing a missing feature scenario. We have employed 10 datasets of diverse specifications — i] the number of points in the datasets vary from 336 to 58000, ii] the number of features vary from 7 to 36 and iii] the number of classes vary from 3 to 9.

In addition to that, we perform micro-analysis of the performance of the proposed scheme on credit-r dataset.

Table 3: Description of datasets. n denotes the number of points in a dataset. f and c denotes the number of features and number of classes of a dataset respectively. d denotes the number of pair-wise combinations of classes to be distinguished by each feature (it is essentially cC_2).

Dataset	n	f	c	d
maternal health risk	1014	7	3	3
contraceptive	1473	9	3	3
segment	2310	19	7	21
balance	625	4	3	3
marketing	6876	13	9	36

6.1. Data partitioning

Data partitioning is an important component any learning procedure. The scheme presented in this work requires us to delve beyond the usual *training-test* partitioning where we have a *three-way* partitioning of data. Before describing the data partitioning, let us recapitulate the assumptions that we have made throughout this work. We assume that a dataset be denoted by \mathcal{D} and there are c classes, γ features. We train $(\gamma + 1)$ classifiers. M_0 is trained on the full and complete dataset, while M_{-i} , $(1 \leq i \leq \gamma)$ is trained without (permuted) i^{th} feature. The training phase of our work has two principal elements — i] we train the $(\gamma + 1)$ classifiers, data used in this phase is denoted with \mathcal{D}_{tr} and ii] we evaluate the class-distinguishing capabilities of the features by comparing the performances of the $(\gamma + 1)$ classifiers using a *test dataset which is also a part of the training phase*. We denote this set with \mathcal{D}_{trte} (it is essentially the test set of the training phase). The remaining component pertains to the usual test phase and we require a test set on which we evaluate the performance of the proposed learning scheme. We denote this set as \mathcal{D}_{te} .

We may note that \mathcal{D}_{tr} , \mathcal{D}_{trte} and \mathcal{D}_{te} are mutually exclusive and

there is no common data points between these sets. Additionally, there are exhaustive with respect to \mathcal{D} as well.

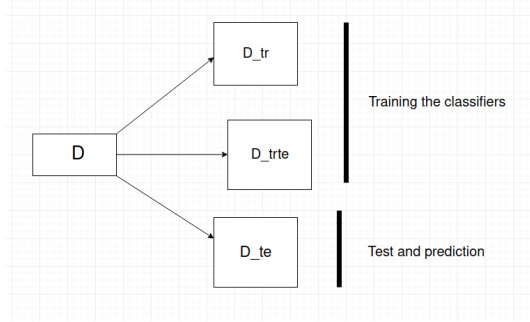


Figure 2: Data partitioning

$$\begin{aligned}
 \mathcal{D} &= \mathcal{D}_{tr} \cup \mathcal{D}_{trte} \cup \mathcal{D}_{te}, \\
 \mathcal{D}_{tr} \cap \mathcal{D}_{trte} &= \phi \\
 \mathcal{D}_{trte} \cap \mathcal{D}_{te} &= \phi \\
 \mathcal{D}_{tr} \cap \mathcal{D}_{te} &= \phi.
 \end{aligned} \tag{5}$$

For all the datasets, we have set the ratio of the number of points in \mathcal{D}_{tr} , \mathcal{D}_{trte} and \mathcal{D}_{te} to 0.5, 0.3 and 0.2 respectively.

6.2. Training and Testing

- **Training phase:** The model is trained on \mathcal{D}_{tr} and the performance scores on \mathcal{D}_{trte} are used to compute the class-distinguishing scores of the features for each pair of classes.
- **Test phase:** The model trained on \mathcal{D}_{tr} is used to get the initial predictions of a test point from \mathcal{D}_{te} . While making the final decision on rendering or refusing the prediction for a test instance, the class-distinguishing scores obtained via \mathcal{D}_{trte} is also taken into account.

6.3. Metrics used for evaluation

We have evaluated the classification performances on *accuracy*. Accuracy evaluates the number of correct predictions made by a classifier overall.

6.4. Metrics used for evaluation

Random forest classifier is used in this work for all classification purposes.

7. Results and analysis

The results and analysis of the work is presented in the following manner. Our work has two major analyses. At first, we evaluate the class-distinguishing scores for two real-world datasets — i] *Maternal health* and ii] *Contraceptive*. The scores are obtained by our model. We present a semantic analysis and explanation of the features of each dataset in terms of their class-distinguishing abilities (captured through their scores).

The class-distinguishing scores provides an way of computing the overall importance of a feature. SHAP [21] is a notable work on feature importance. We compute the feature importance for these datasets using SHAP and show the outcomes graphically. We also show the overall class-distinguishing scores obtained by our method and show the correspondance between the two.

Secondly, we report the utility of the class-discernible scores of our scheme in a classification task.

7.1. Analysis using *Maternal health* dataset

Dataset description: The dataset provides an insight into the risk factors in maternal mortality, the data is collected from rural hospitals in Bangladesh. It reports age and five health parameters of 1014 pregnant women and calculates the risk factor of *pregnancy* on the basis of these parameters [46]. The health parameters are collected through wearable IoT enabled devices. The parameters reported in this study are — *systolic BP*, *diastolic BP*, *blood sugar level*, *body temperature* and *heart rate*. The risk factor is categorized as — *low risk* (1), *medium risk* (2) and *high risk* (3). So, there are *three* classes and we have ${}^3C_2 = 3$ cases of pair-wise class distinctions. It is worth noting that in quite a few cases, the maternal age provided in the dataset was as low as 10 years to as high as 70 years. To operate in a feasible scenario, we have removed all the points where the age was found to be less than 12 years or more than 55 years. The study is carried out on this reduced dataset.

In this section, we will analyse the class-distinguishing scores for each feature. The class-distinguishing scores for each of the features are reported in Figure 6. The bars in the figure indicates the class distinguishing capability of each figure — more the height of a bar, more is the distinguishing capability of a feature with respect to that pair of classes. The class distinguishing scores for a feature can be interpreted as a measure of its importance (specific to pair-wise classes). For example, the blue bar corresponding to the left-most item (age) indicates the — class-distinguishing ability of *maternal age* in distinguishing between class 1 (*low risk*) and class 2 (*medium risk*). Similarly, the green bar corresponding to the left-most item (age) indicates the — class-distinguishing ability of *maternal age* in distinguishing between class 1 (*low risk*) and class 3 (*high risk*). The height of the green bar is more than that of the blue bar, and it indicates that *maternal age* is more instrumental in distinguishing between *low risk* and *high risk* cases than that of *low risk* and *medium risk* cases. Let us analyze the key observations from the plot of the class-distinguishing scores Let CS_{12}^{age} be the class distinguishing score with respect to age (feature 1 in this dataset).

- For each feature, the highest class-discernible score is obtained in **medium risk-high risk cases**. For feature 2 (Systolic BP), feature 4 (Blood Sugar) and feature 5 (Body temperature), the score for medium risk-high risk case has surpassed the remaining two cases (low risk-medium risk cases and low risk-high risk cases) by a considerable amount. It is of the significance that — when the classifier predicts high risk (class 3) and medium risk (class 2) as

two most likely classes for a point, the absence of any one of the above three features can lead us to a *refusal in prediction* (if the classifier’s confidence on the mostly likely class is not very high with respect to the second most probable class). For the remaining three features (Age, Diastolic BP and Heart Rate), this difference is not much pronounced. In these cases, when the classifier predicts high risk and medium risk (in either order) as the two most likely classes and any one of the above three features is absent – the scheme will be confident of rendering a decision even when the likelihood of the two most probable classes varies by a narrow margin. It is because of the low class-distinguishing scores.

- We have expected **age (Feature 1)** to play a decisive role as a feature for learning the pregnancy. But the outcome that we have obtained is not in congruence with our expectations (with respect to the other features). The class-distinguishing scores of *age* as a feature is highest for *medium risk vs high risk* cases, followed by *low risk vs high risk cases* and the class-distinguishing scores of *low risk vs high risk* scores is the least.
- The bargraphs of **Systolic Blood Pressure (BP) (Feature 2)** and **Diastolic Blood Pressure (BP) (Feature 3)** are of similar nature, with *Systolic BP* being more informative than *Diastolic BP*. The respective bars for all three cases is taller for Systolic BP than that of Diastolic BP. The BP parameters’ capability of distinguishing *low-risk vs medium-risk* cases is more than that of *low-risk vs high-risk* cases. The class-distinguishing score in *medium risk vs high risk* case is more pronounced for Systolic BP than that of Diastolic BP. The class-discernible
- **Blood Sugar Level (Feature 4)**: Blood sugar levels indicates our body’s metabolism rate, which, in turn gives a crucial insight into our overall health. The findings of our study also indicates the same and barplot of Figure 3 indicates that **Blood Sugar** is key feature for distinguishing the high risk cases from low risk cases and the medium risk cases. The class discernible scores of this feature is highest for medium risk-high risk cases followed by low risk-high risk cases. It indicates that it will be difficult for the scheme to choose between *high risk vs low risk/ medium risk* cases in absence of the blood sugar information. The class discernible score for the *low risk vs medium risk* is not much higher. This indicates that, in absence of Blood Sugar information, the scheme might be able to render a decision for a point when the classifier predicts low risk and medium risk as the two most likely classes.
- **Body Temperature (Feature 5)**: Basal Body temperature’s class-distinguishing capacity is low in *low risk vs medium risk* and *low risk vs high risk* cases. But the scores are high in *medium risk vs high risk* cases which signifies that it plays a crucial role for distinguishing between *medium risk and high risk* pregnant women.

- **Heart Rate (Feature 6)**: This class-distinguishing capacity of this feature is reported as nominal from Figure 6. The bar of class-distinguishing score for *medium risk vs high risk* case is slightly better as compared to that of *low risk vs medium risk* and *low risk vs high risk* cases. But the absolute values are much less compared to other features.
- **Overall feature importance**: We show the overall importance of the features obtained via SHAP and the proposed method in Figure 4. In our method, the overall importance of a feature is obtained by summing up its class-distinguishing scores for all possible classes. The outcomes given our scheme is mostly congruent with the outcomes given by SHAP. The top four features are same in both cases. Blood sugar is the most important feature followed by Systolic blood pressure, age and Diastolic blood pressure. In case of SHAP, body temperature is least important feature. But, the proposed method indicates heart rate as the least important feature. However, we can account this difference (very minute) to the randomness associated with a Random Forest Classifier.

7.2. Analysis using *Contraceptive*

Dataset description: This dataset is a subset of the 1987 National Indonesia Contraceptive Prevalence Survey. The goal of this data collection was to understand and correlate the socio-economic and demographic factors with the choice of contraceptive method. The choice of contraceptive were — i] no choice, ii] long-term method and iii] short term method. The data was collected from 1473 non-pregnant or not known to be pregnant women in Indonesia. The data was collected about the socio-economic and demographic situation these served as the features) of each woman and her choice of the contraceptive method (it served as the class).

There are two numerical features (age of wife and number of children ever born) and seven categorical features (wife’s education, husband’s education, wife’s religion, wife is now working or not, husband’s occupation, standard of living index and media exposure). The dataset is multi-class and we have to classify each point (supposedly woman) to three classes – *no choice (class 1)*, *long-term method (class 2)* and *short-term method (class 3)*. Figure 5 shows the class-distinguishing scores of the features for the pair-wise classes. We discuss the key findings below.

- **Wife’s age (Feature 1)** shows strong class-discernible capabilities for two cases — i] *no choice vs short term method* (classes 1 and 2) and ii] *long term method vs short term method* (classes 2 and 3). In particular, the class discernible scores of the latter case is the highest for all cases (all features and all pair-wise classes). The likely explanation for this is the prevalence of long-term contraceptive method among older female (or male), and the prevalence of short-term or reversible methods among younger population (specially if they are undecided about having more

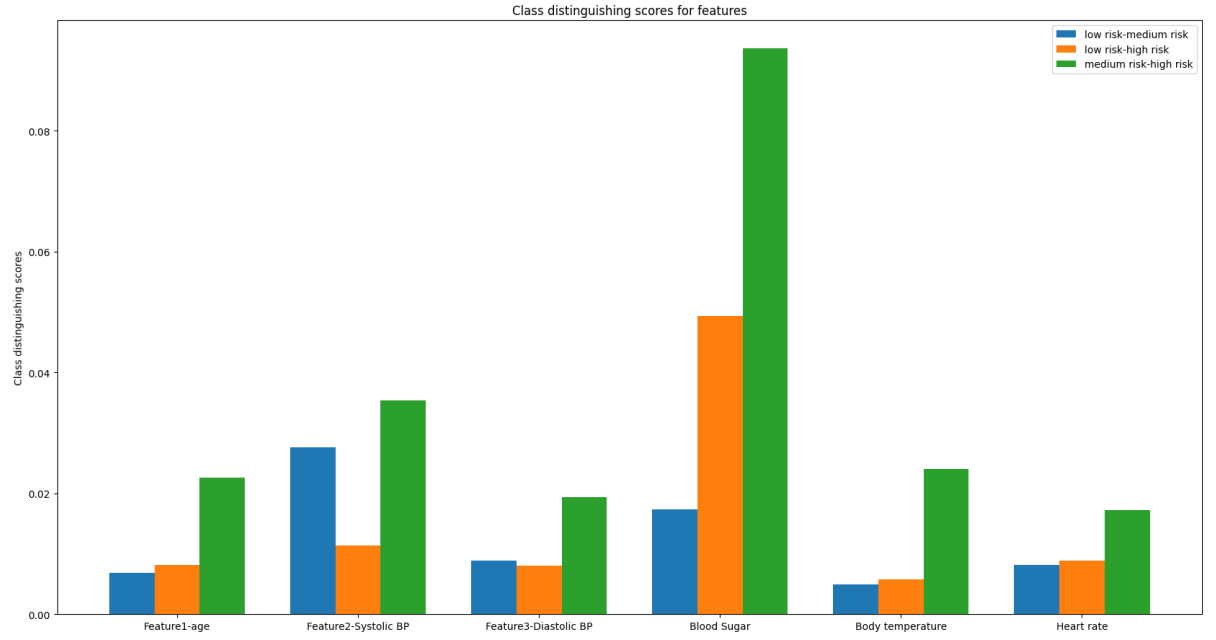


Figure 3: This figure shows the class-distinguishing scores of features of the *maternal health* dataset

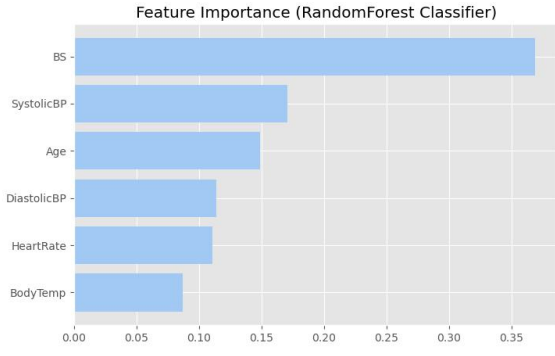
children in future or losing fertility). Age is not as effective for distinguishing *no use vs long term* cases. A possible cause may be the prevalence of higher age group in the dataset, for whom the no-use and long-term method are equally likely.

- **Number of children ever born (Feature 4)** has emerged as a key feature from our empirical study. It has good class-discernible scores for all three cases of pair-wise classes. This feature shows a particularly strong class-discernible score for distinguishing *no contraceptive vs long term contraceptive* cases. Note that, this feature gives the highest class distinguishing score for distinguishing this pair of classes. It is expected that couples with fewer children (this was back in 1987) would be willing to have more child (hence no contraceptive) and couples with more number of children (or a figure according to their own optimality) will resort to long-term contraceptive methods. Class discernibility score for *long term vs short term methods* is also a significant figure with respect to this feature. An explanation for this can be – couples with lesser number of children will prefer short term methods while couples with more number of children will prefer long term methods. The bar-graph corresponding to *no contraceptive vs short term contraceptive* case is also high. The same explanation for the no use-long term case is also applicable here.
- **Wife’s education (Feature 2)** shows weighty class-distinguishing figures for two cases i] *no contraceptive and long term contraceptive* cases, and ii] *short term contraceptive and long term contraceptive* cases. In context of

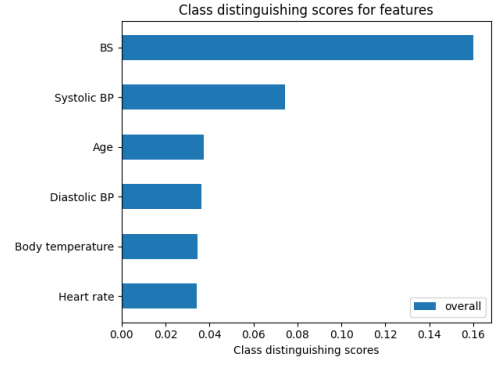
1987, this finding seems pretty reasonable – wife’s education leads to an awareness, which promotes the use of contraceptive. The class-discernibility score of *no contraceptive vs short term contraceptive* cases is not as high as the other cases. A likely explanation for this may be the prevalence of short term contraceptive method in the lower age group, for whom no contraceptive is also a viable choice (wanting children). And the wives from the lower age group is likely to be more educated — leading to two similar choices.

A technical implication of the above: Suppose, we want to classify two women, \mathbf{p} and \mathbf{q} on the basis of all features but one - *wife’s education*. We find that two most probable classes for each woman from the classifier. Let us assume that for \mathbf{p} , the two most probable classes are *no contraceptive method (class 1)* and *short term method (class 3)*. Since *wife’s education* has good class-discernible scores for classifying these two cases (and we do not have access to this information), there is a good chance that the scheme will *refuse* to render a decision in this case. Let us assume that for \mathbf{q} , the two most probable classes are *no contraceptive (class 2)* and *short term method (class 3)*. Since *wife’s education* has less class-discernible scores for classifying these two cases, there is a good chance that the classifier will not bother about this missing information and the scheme will render a decision.

- **Media exposure (Feature 9):** The class-distinguishing scores of this feature is more for *no contraceptive vs long term method* and *long term method vs short term method* cases. Its competence in the remaining case is much



(a) SHAP



(b) Proposed method

Figure 4: Importance of the features of *Maternal health* dataset.

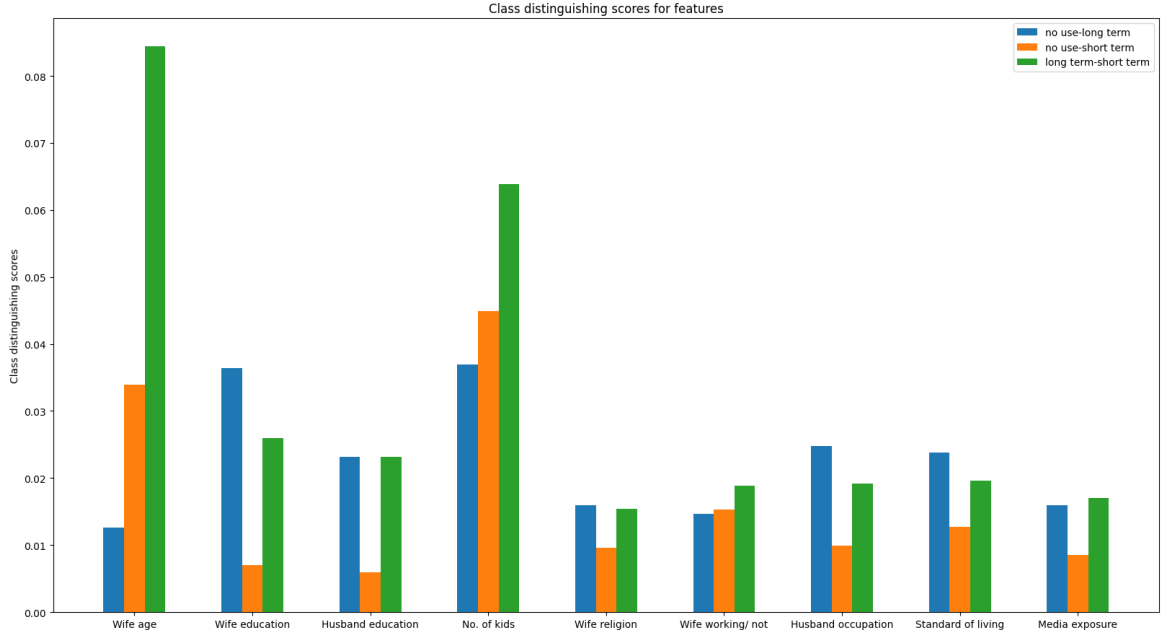


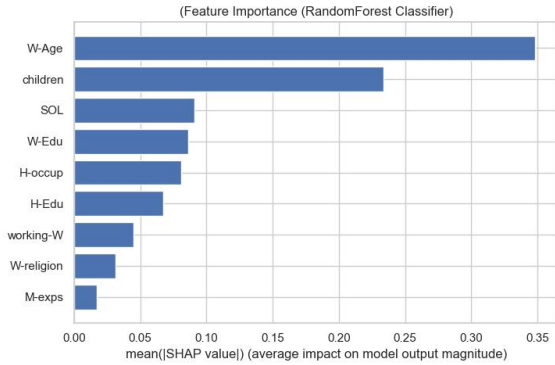
Figure 5: This figure shows the class-distinguishing scores of features of the *contraceptive* dataset

less though. But, we may note that this feature binary-categorical feature is severely imbalanced. That might affect the learning and we might not have a true picture of the socio-economic correlation with the classes.

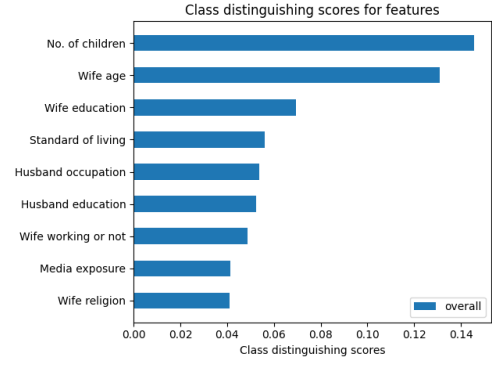
- **Wife’s religion, Wife’s education, Husband’s occupation and Standard of living, Husband’s education** – the trends of these five features are mostly similar. For each of the feature, the class-discernability scores are in following order – *no contraceptive vs long term* cases > *long term vs short term* cases > *no use vs short term* cases. In all five of these cases, the class-distinguishing ability of these features is much less in *no use vs short term* cases with respect to the remaining two cases. The outcome can be accounted to one aspect — these features may have an underlying correlation. We may further note that, of all the fea-

tures, wife’s religion shows the lowest class-distinguishing scores.

- **Overall feature importance:** We report the overall importance of the features of *contraceptive* dataset in Figure 6. In Figure (6a) and Figure (6B), we show the importance obtained via SHAP and the proposed method respectively. The findings from the two are partially in congruence with each other, in essence they are largely similar. *Wife’s education* and *No. of children* are top two features in either case but SHAP outputs *Wife’s education* as the most important feature whereas the proposed method’s plot shows *No. of children* as the most important feature. The same holds for the 3rd most and 4th most important features with respect to the two schemes. According to SHAP, the 3rd and 4th most important features are Standard-of-living



(a) SHAP



(b) Proposed method

Figure 6: Importance of the features of *Contraceptive* dataset.

and Wife’s education respectively. On the contrary, the outcomes from the proposed method indicates just the reverse. The 5th, 6th and 7th features are same for both the methods. We again have a reversal in ranking for the two least important features. We may further note that barring the top two features (obtained in either case), the importances of the remaining features is much lesser.

7.3. Classification and utility of the class-discernible scores

In this subsection, we study the utility of the proposed *class-distinguishing scores* in rendering a decision. We dedicate one figure to each dataset. For each dataset, we report feature-specific performance of the classifiers and plot four bar-graphs for each — i] Plot 1 (blue coloured bar)- all features are present (this serves as the baseline), ii] Plot 2 (orange coloured bar)- a particular feature a is absent, iii] Plot 3 (green coloured bar)- feature a is present and our scheme renders a decision, and iv] Plot 4 (red coloured bar)- feature a is present and our scheme refuses to render a decision. In this particular case, the class with highest probability estimate is chosen as the decision. We plot classification accuracy in all the cases.

In each figure, plot 3 and plot 4 are the outcomes of our scheme. Plot 3 cumulates the cases in which our scheme is confident of rendering a decision. On the contrary, plot 4 cumulates the cases in which our scheme does not want to render a decision. The **decision** about *rendering or refusing a decision* is an integral part of the proposed scheme and indicative of the correctness of the explainability given the scheme. We may note that – if we get a difference in accuracies of these two cases, we can say that decision of our scheme bears some significant. To be precise, if the accuracy in plot 3 case is more than the respective accuracy in plot 4 cases, the scheme is right about its *decision*. This shows that the scheme is refusing a decision which could lead to an incorrect one. We may further note that our scheme has achieved this for all four datasets and across all features. The red bar (plot 4) has much lower height than its corresponding green bar (plot 3) in all the cases. This signifies that our scheme is indeed right about refusing the decision — had it decided, it would have been wrong in the prediction

(as indicated from the accuracy). This also indicates the correctness of the class-discernible abilities (explainability) of the features — on the basis of which the decision about rendering or refusal has been taken.

8. Conclusion and Future Work

In this work, we analyze and explain the features in micro-scale. We study and quantify their class-distinguishing abilities and obtain a set of feature-specific class-distinguishing scores. We use the computed scores to aid the decision making in latent feature scenarios. A key characteristic of this decision making model is — it can refuse or render a decision (categorization) for a test point. The decision of refusal or rendering is dependent on two things — i] the two most probable classes or categories of the test point, and ii] the latent feature’s ability to distinguish these two likely classes.

The empirical results from the entire procedure indicates the correctness of the class-distinguishing scores as well as its suitability in the latent feature framework. Detailed analyses carried on two real-world datasets also manifests the same. The proposed framework has the provides for the overall importances of the features as well as their class-specific importances. This proposed framework for obtaining feature importance and explainability can aid in the automated decision making of real-world systems. In our future work, we will like to explore two more aspects of real-world data, namely class-imbalance and feature imbalance. We will specifically focus on how they affect the class-distinguishing abilities of the features and the overall feature importances.

References

- [1] S. Chancellor, E. P. Baumer, M. De Choudhury, Who is the” human” in human-centered machine learning: The case of predicting mental health from social media, *Proceedings of the ACM on Human-Computer Interaction* 3 (CSCW) (2019) 1–32.
- [2] L. Ma, B. Sun, Machine learning and ai in marketing—connecting computing power to human insights, *International Journal of Research in Marketing* 37 (3) (2020) 481–504.

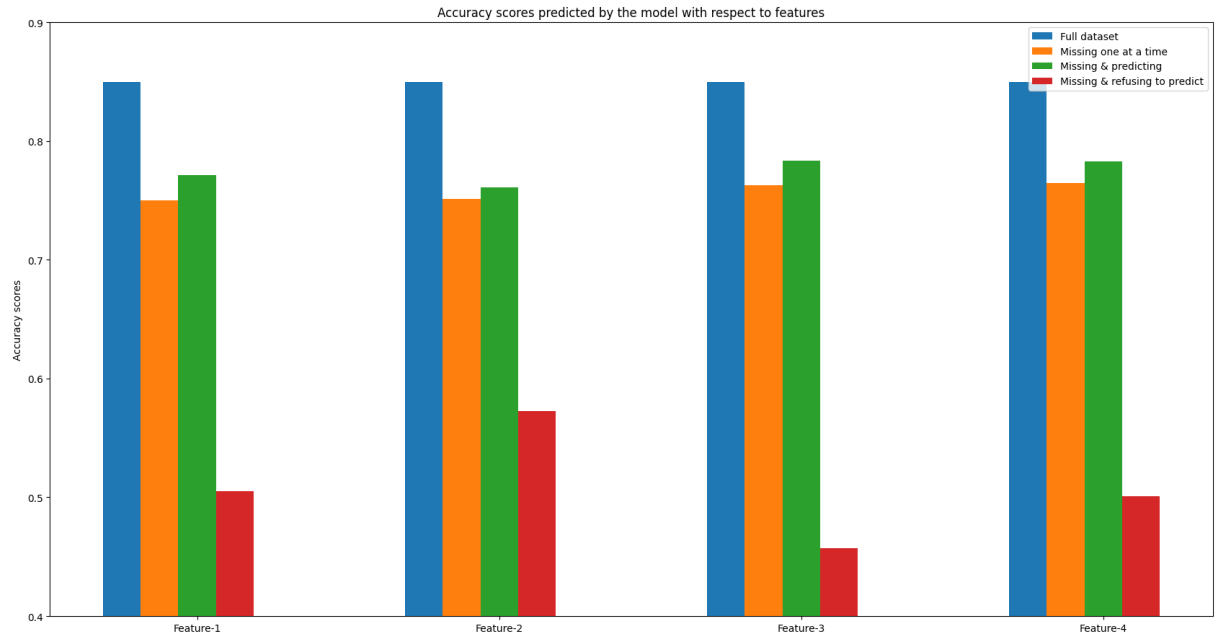


Figure 7: Balance

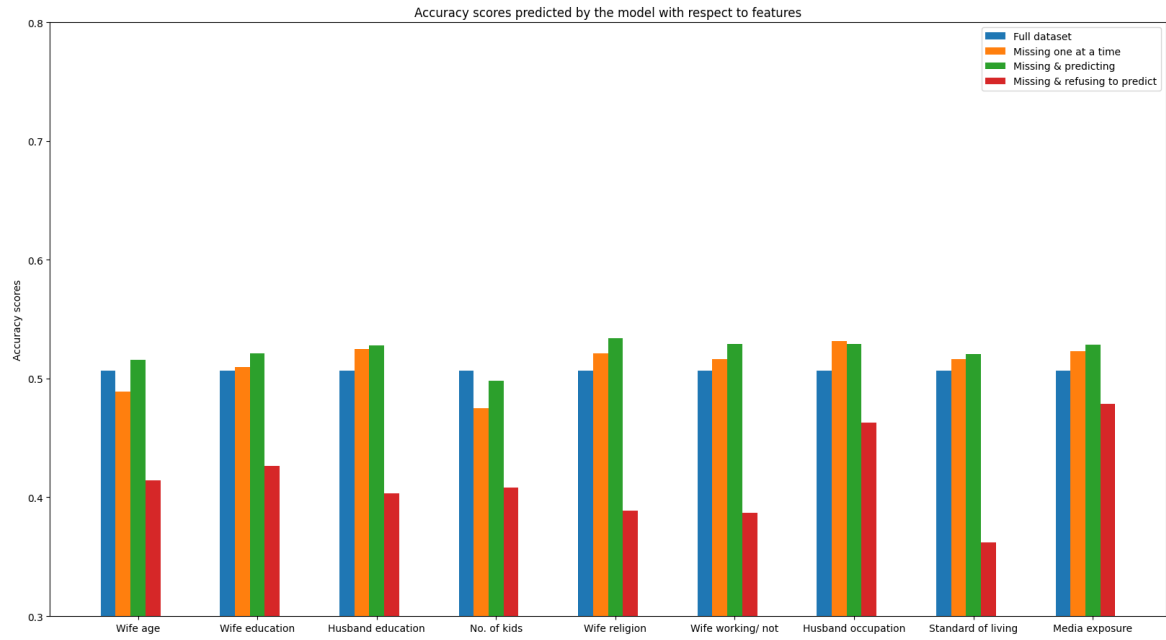


Figure 8: Contraceptive

- [3] I. Poola, How artificial intelligence in impacting real life everyday, International Journal for Advance Research and Development 2 (10) (2017) 96–100.
- [4] M. I. Jordan, T. M. Mitchell, Machine learning: Trends, perspectives, and prospects, Science 349 (6245) (2015) 255–260.
- [5] I. H. Sarker, Machine learning: Algorithms, real-world applications and research directions, SN Computer Science 2 (3) (2021) 1–21.
- [6] L. E. Lwakatare, A. Raj, I. Crnkovic, J. Bosch, H. H. Olsson, Large-scale machine learning systems in real-world industrial settings: A review of challenges and solutions, Information and software technology 127 (2020) 106368.
- [7] D. Gómez, A. Rojas, An empirical overview of the no free lunch theorem and its effect on real-world machine learning classification, Neural computation 28 (1) (2016) 216–228.
- [8] A. Binder, M. Bockmayr, M. Hägele, S. Wienert, D. Heim, K. Hellweg, M. Ishii, A. Stenzinger, A. Hocke, C. Denkert, et al., Morphological and

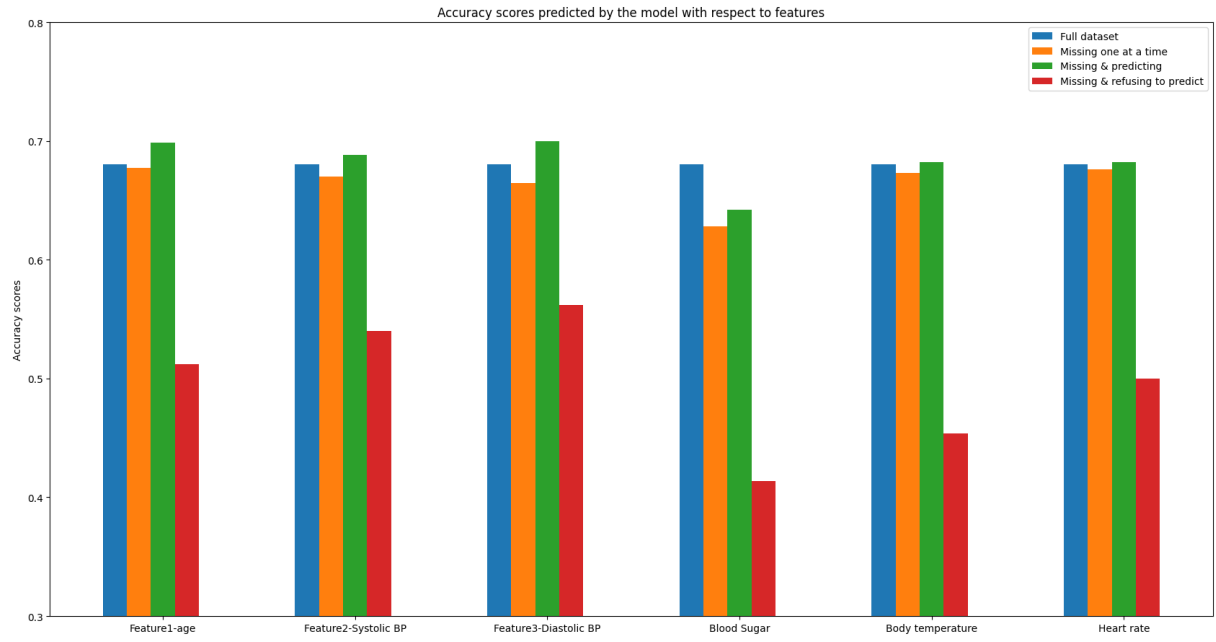


Figure 9: Maternal health

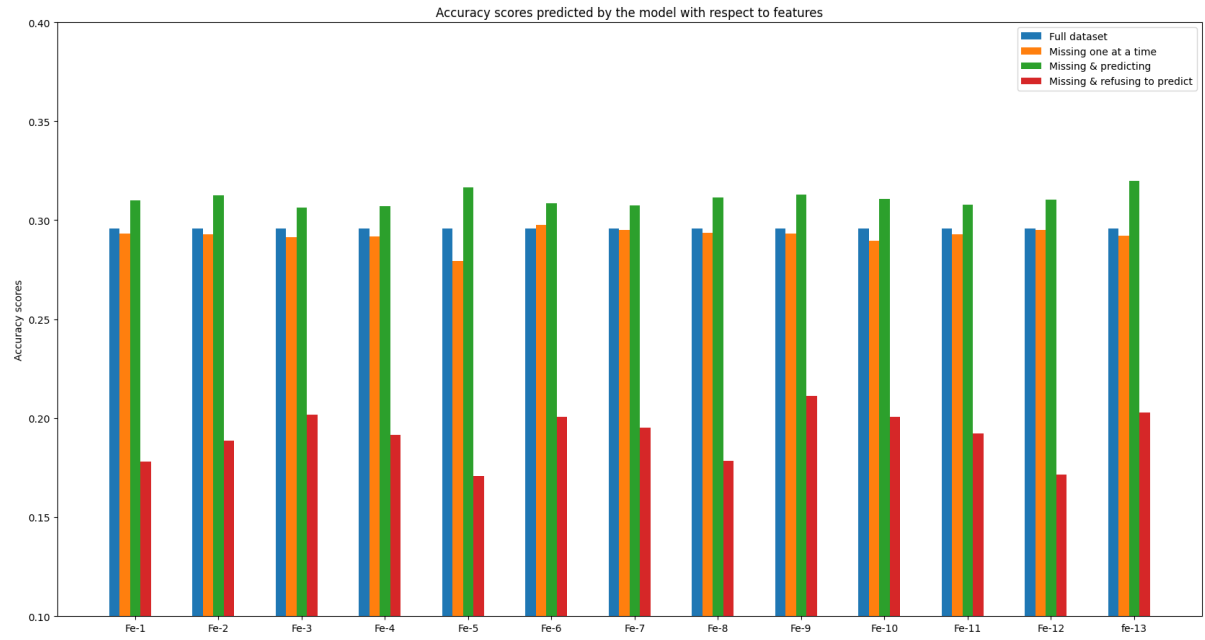


Figure 10: Marketing

- molecular breast cancer profiling through explainable machine learning, *Nature Machine Intelligence* 3 (4) (2021) 355–366.
- [9] Y. Du, A. R. Rafferty, F. M. McAuliffe, L. Wei, C. Mooney, An explainable machine learning-based clinical decision support system for prediction of gestational diabetes mellitus, *Scientific Reports* 12 (1) (2022) 1–14.
- [10] A. Katzmann, O. Taubmann, S. Ahmad, A. Mühlberg, M. Sühling, H.-M. Groß, Explaining clinical decision support systems in medical imaging using cycle-consistent activation maximization, *Neurocomputing* 458 (2021) 141–156.
- [11] M. Abedin, S. Mokhtari, A. B. Mehrabi, Bridge damage detection using machine learning algorithms, in: P. Fromme, Z. Su (Eds.), *Health Monitoring of Structural and Biological Systems XV*, Vol. 11593, International Society for Optics and Photonics, SPIE, 2021.
- [12] K. Topuz, D. Delen, A probabilistic bayesian inference model to investigate injury severity in automobile crashes, *Decision Support Systems* 150

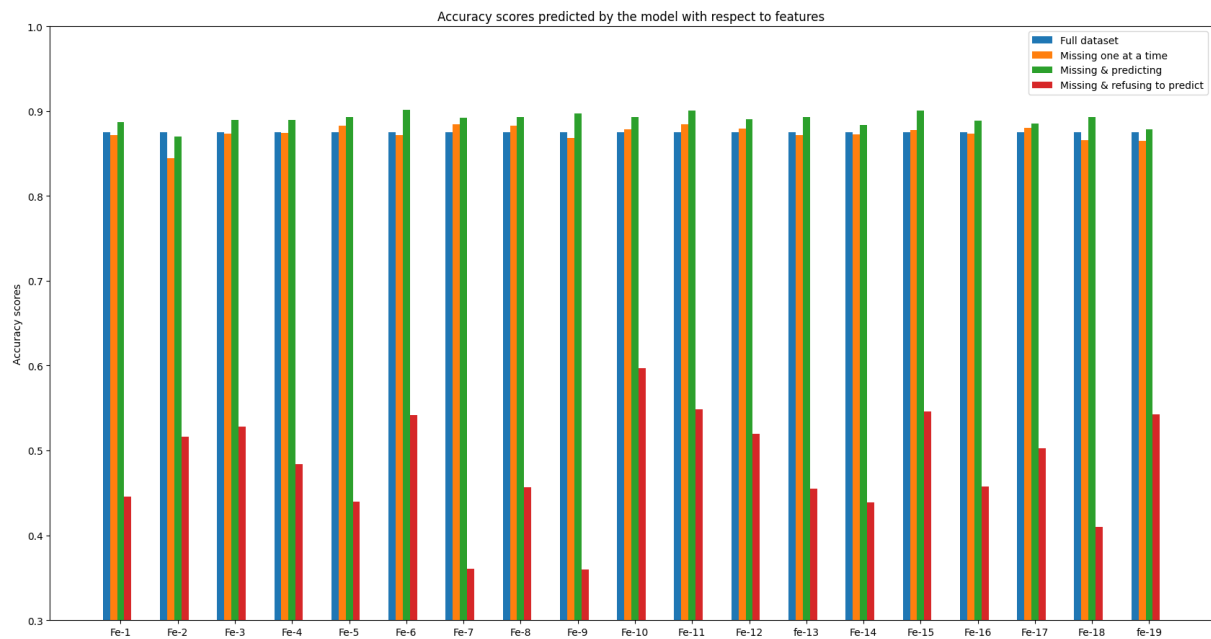


Figure 11: Segment

- (2021) 113557, interpretable Data Science For Decision Making.
- [13] F. Vandervorst, W. Verbeke, T. Verdonck, Data misrepresentation detection for insurance underwriting fraud prevention, *Decision Support Systems* 159 (2022) 113798.
 - [14] D. Delen, K. Topuz, E. Eryarsoy, Development of a bayesian belief network-based dss for predicting and understanding freshmen student attrition, *European Journal of Operational Research* 281 (3) (2020) 575–587, featured Cluster: Business Analytics: Defining the field and identifying a research agenda.
 - [15] S. Sachan, J.-B. Yang, D.-L. Xu, D. E. Benavides, Y. Li, An explainable ai decision-support-system to automate loan underwriting, *Expert Systems with Applications* 144 (2020) 113100.
 - [16] C. Rudin, Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead, *Nature Machine Intelligence* 1 (5) (2019) 206–215.
 - [17] C. B. Azodi, J. Tang, S.-H. Shiu, Opening the black box: interpretable machine learning for geneticists, *Trends in genetics* 36 (6) (2020) 442–455.
 - [18] S. R. Islam, W. Eberle, S. K. Ghafoor, M. Ahmed, Explainable artificial intelligence approaches: A survey, *arXiv preprint arXiv:2101.09429* (2021).
 - [19] S. Luo, H. Ivison, C. Han, J. Poon, Local interpretations for explainable natural language processing: A survey, *arXiv preprint arXiv:2103.11072* (2021).
 - [20] M. T. Ribeiro, S. Singh, C. Guestrin, "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*, 2016, pp. 1135–1144.
 - [21] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, *Advances in neural information processing systems* 30 (2017).
 - [22] R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, N. Elhadad, *Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission*, KDD '15, Association for Computing Machinery, New York, NY, USA, 2015. doi:10.1145/2783258.2788613. URL <https://doi.org/10.1145/2783258.2788613>
 - [23] H. Kaur, H. Nori, S. Jenkins, R. Caruana, H. Wallach, J. Wortman Vaughan, Interpreting interpretability: Understanding data scientists' use of interpretability tools for machine learning, in: *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, CHI '20, Association for Computing Machinery, New York, NY, USA, 2020, p. 1–14. doi:10.1145/3313831.3376219. URL <https://doi.org/10.1145/3313831.3376219>
 - [24] D. Rengasamy, B. C. Rothwell, G. P. Figueredo, Towards a more reliable interpretation of machine learning outputs for safety-critical systems using feature importance fusion, *Applied Sciences* 11 (24) (2021). doi:10.3390/app112411854. URL <https://www.mdpi.com/2076-3417/11/24/11854>
 - [25] E. Hickman, M. Petrin, Trustworthy ai and corporate governance: The eu's ethics guidelines for trustworthy artificial intelligence from a company law perspective, *European Business Organization Law Review* 22 (4) (2021) 593–625. doi:10.1007/s40804-021-00224-0. URL <https://doi.org/10.1007/s40804-021-00224-0>
 - [26] B. Goodman, S. Flaxman, European union regulations on algorithmic decision-making and a "right to explanation", *AI Magazine* 38 (3) (2017) 50–57.
 - [27] J. C. Reis, A. Correia, F. Murai, A. Veloso, F. Benevenuto, Explainable machine learning for fake news detection, in: *Proceedings of the 10th ACM conference on web science*, 2019, pp. 17–26.
 - [28] E. Alcobaca, S. M. Mastelini, T. Botari, B. A. Pimentel, D. R. Cassar, A. C. P. de Leon Ferreira, E. D. Zanotto, et al., Explainable machine learning algorithms for predicting glass transition temperatures, *Acta Materialia* 188 (2020) 92–100.
 - [29] M. Zhang, H. You, P. Kadam, S. Liu, C.-C. J. Kuo, Pointhop: An explainable machine learning method for point cloud classification, *IEEE Transactions on Multimedia* 22 (7) (2020) 1744–1755.
 - [30] L. Jiang, S. Liu, C. Chen, Recent research advances on interactive machine learning, *Journal of Visualization* 22 (2) (2019) 401–417.
 - [31] M. Danilevsky, K. Qian, R. Aharonov, Y. Katsis, B. Kawas, P. Sen, A survey of the state of explainable ai for natural language processing, *arXiv preprint arXiv:2010.00711* (2020).
 - [32] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, S.-I. Lee, From local explanations to global understanding with explainable ai for trees, *Nature machine intelligence* 2 (1) (2020) 56–67.
 - [33] G. Montavon, S. Lapuschkin, A. Binder, W. Samek, K.-R. Müller, Explaining nonlinear classification decisions with deep taylor decomposition, *Pattern Recognition* 65 (2017) 211–222. doi:<https://doi.org/10.1016/j.patcog.2016.11.008>

- [34] P. Angelov, E. Soares, Towards explainable deep neural networks (xdnn), *Neural Networks* 130 (2020) 185–194.
- [35] J. M. Alonso Moral, C. Castiello, L. Magdalena, C. Mencar, Design and validation of an explainable fuzzy beer style classifier, in: *Explainable Fuzzy Systems*, Springer, 2021, pp. 169–217.
- [36] P. Sabol, P. Sinčák, K. Ogawa, P. Hartono, explainable classifier supporting decision-making for breast cancer diagnosis from histopathological images, in: *2019 International Joint Conference on Neural Networks (IJCNN)*, IEEE, 2019, pp. 1–8.
- [37] B. Hamrouni, A. Bourouis, A. Korichi, M. Brahmi, Explainable ontology-based intelligent decision support system for business model design and sustainability, *Sustainability* 13 (17) (2021) 9819.
- [38] K. Davagdorj, J.-W. Bae, V.-H. Pham, N. Theera-Umporn, K. H. Ryu, Explainable artificial intelligence based framework for non-communicable diseases prediction, *IEEE Access* 9 (2021) 123672–123688.
- [39] G. Kunapuli, B. A. Varghese, P. Ganapathy, B. Desai, S. Cen, M. Aron, I. Gill, V. Duddalwar, A decision-support tool for renal mass classification, *Journal of Digital Imaging* 31 (6) (2018) 929–939.
- [40] A. Gosiewska, A. Kozak, P. Biecek, Simpler is better: Lifting interpretability-performance trade-off via automated feature engineering, *Decision Support Systems* 150 (2021) 113556, interpretable Data Science For Decision Making.
- [41] A. N. Baraldi, C. K. Enders, An introduction to modern missing data analyses, *Journal of school psychology* 48 (1) (2010) 5–37.
- [42] T. D. Little, T. D. Jorgensen, K. M. Lang, E. W. G. Moore, On the joys of missing data, *Journal of pediatric psychology* 39 (2) (2014) 151–162.
- [43] C. Biffi, J. J. Cerrolaza, G. Tarroni, W. Bai, A. De Marvao, O. Oktay, C. Ledig, L. Le Folgoc, K. Kamnitsas, G. Doumou, et al., Explainable anatomical shape analysis through deep hierarchical generative models, *IEEE transactions on medical imaging* 39 (6) (2020) 2088–2099.
- [44] J.-Y. Kim, S.-B. Cho, Explainable prediction of electric energy demand using a deep autoencoder with interpretable latent space, *Expert Systems with Applications* 186 (2021) 115842.
- [45] N. Yang, Y. Ma, L. Chen, P. S. Yu, A meta-feature based unified framework for both cold-start and warm-start explainable recommendations, *World Wide Web* 23 (1) (2020) 241–265.
- [46] M. Ahmed, M. A. Kashem, M. Rahman, S. Khatun, Review and analysis of risk factor of maternal health in remote area using the internet of things (iot), in: *InECCE2019*, Springer Singapore, Singapore, 2020, pp. 357–365.