# Explaining a Series of Models by Propagating Shapley Values

Hugh Chen[1], Scott M. Lundberg[2], and Su-In Lee[1,*]

[1]Paul G. Allen School of Computer Science and Engineering, University of Washington

[2]Microsoft Research

[*]Corresponding: suinlee@cs.washington.edu

## Abstract

Local feature attribution methods are increasingly used to explain complex machine learning models. However, current methods are limited because they are extremely expensive to compute or are not capable of explaining a distributed series of models where each model is owned by a separate institution. The latter is particularly important because it often arises in finance where explanations are mandated. Here, we present DeepSHAP, a tractable method to propagate local feature attributions through complex series of models based on a connection to the Shapley value. We evaluate DeepSHAP across biological, health, and financial datasets to show that it provides equally salient explanations an order of magnitude faster than existing model-agnostic attribution techniques and demonstrate its use in an important distributed series of models setting.

## 1 Introduction

With the widespread adoption of machine learning (ML), *series of models* (i.e., where the outputs of predictive models are used as inputs to separate predictive models) are increasingly common. Examples include: (1) *stacked generalization*, a widely used technique [1–5] to improve generalization performance by ensembling the predictions of many models (called base-learners) using another model (called a meta-learner) [6], (2) *neural network feature extraction*, where models are trained on features extracted using neural networks [7, 8], typically for structured data [9–11], and (3) *consumer scores*, where predictive models that describe a specific behavior (e.g., credit scores [12]) are used as inputs to downstream predictive models. For example, a bank may use a model to predict customers' loan eligibility on the basis of their bank statements and their credit score, which itself is often a predictive model [13].

Explaining a series of models is crucial for debugging and building trust, even more so because a series of models is inherently harder to explain compared to a single model. One popular paradigm for explaining models are *local feature attributions*, which explain why a model makes a prediction for a single sample (known as the "explicand" [14]). Existing *model-agnostic* local feature attribution methods (e.g., IME [15], LIME [16], KernelSHAP [17]) work regardless of the specific model being explained. They can explain a series of models, but suffer from two distinct shortcomings: (1) their sampling-based estimates of feature importance are inherently variable, and (2) they have high computational cost which may not be tractable for large pipelines. Alternatively, *model-specific* local feature attribution methods (i.e. attribution methods that work for specific types of models) are often much faster than model-agnostic approaches, but generally cannot be used to explain a series of models. Examples include those for (1) deep models (e.g., DeepLIFT [18], Integrated Gradients [19]) and (2) tree models (e.g., Gain/Gini Importance [20], TreeSHAP [21]).

In this paper, we present DeepSHAP – a local feature attribution method that is *faster than model-agnostic methods* and *can explain complex series of models that pre-existing model-specific methods cannot*. DeepSHAP is based on connections to the Shapley value, a concept from game theory that satisfies many desirable axioms. We make several important contributions:

1. We propose a theoretical framework (Methods Section 6.6) that connects the rules introduced in Shrikumar *et al.* to the Shapley value with an interventional conditional expectation set function[1] (ICE

---

[1]with a flat causal graph

Shapley value) (Methods Section 6.1).

2. We show that the ICE Shapley value decomposes into an average over "single baseline attributions"[2] (Methods Section 6.2), where a single baseline attribution explains the model for a single sample (explicand) by comparing to a single sample (baseline).

3. We propose a *generalized rescale rule* to explain a complex series of models by propagating attributions while enforcing efficiency at each layer (Figure 1b, Methods Section 6.4). This framework extends DeepSHAP to explain any series of models composed of linear, deep, and tree models.

4. We propose a *group rescale rule* to propagate local feature attributions to groups of features (Methods Section 6.7). We show that these group attributions better explain models with many features.

Many feature attribution methods must define the absence of a feature, often by masking features according to a single baseline sample (single baseline attribution) [14, 18, 19]. In contrast, we show that under certain assumptions, the correct approach is to use many baseline samples instead (Appendix Section A.5.2). Qualitatively, we show that using many baselines avoids bias that can be introduced by single baseline attributions (Section 3.1). Additionally, we show that the choice of baseline samples is a useful parameter which changes the question answered by the attributions (Figure 1c, Section 3.2).

We qualitatively and quantitatively evaluate DeepSHAP in real-world datasets including biological, health, image, and financial data sets. In the biological datasets [23–26], we qualitatively assess group feature attributions based on gene sets identified in prior literature (Section 4.1). In the health, image, and financial datasets [27–29], we quantitively show that DeepSHAP provides useful explanations and is drastically faster than model agnostic approaches using an ablation test, where we hide features according to their attribution values (Sections 4.2, 4.3, 4.4).

In addition, DeepSHAP is the only approach we are aware of that enables explanations of a distributed series of models (where each model belongs to a separate institution). Model-agnostic approaches do not work because they need access to every model in the series, but institutions cannot share models because they are proprietary. One extremely prevalent example of distributed models are *consumer scores* which exist for nearly every American consumer [12] (Section 4.4). In this setting, transparency is a critical issue, because opaque scores can hide discrimination or unfair practices.

## 2 Generalizing DeepSHAP local explanations

A closely related method to DeepSHAP was designed to explain deep models ($f : \mathbb{R}^m \to \mathbb{R}$) [17], by performing DeepLIFT [18] using the average as a baseline. However, using a single average baseline is not the correct approach to explain non-linear models based on connections to Shapley values with an interventional conditional expectation set function and a flat causal graph [30]. Instead, we show that the correct way to obtain the interventional Shapley value local feature attributions (denoted as $\phi(f, x^e) \in \mathbb{R}^m$) based on an explicand ($x^e \in \mathbb{R}^m$), or sample being explained, is to average over single baseline feature attributions (denoted as $\phi(f, x^e, x^b) \in \mathbb{R}^m$) where baselines are $x^b \in \mathbb{R}^m$ and $D$ is the set of all baselines (details in Methods Section 6.2):

$$\phi(f, x^e) = \frac{1}{|D|} \sum_{x^b \in D} \phi(f, x^e, x^b) \tag{1}$$

DeepLIFT [18] explains deep models by propagating feature attributions at each layer of the deep model. Here, we extend DeepLIFT by generalizing DeepLIFT's rescale rule to accommodate more than neural network layers while guaranteeing layer-wise efficiency (details in Methods Section 6.4). For a series of models which can be represented as a composition of functions ($f_k(x) = (h_k \circ \cdots \circ h_1)(x)$, where $h_i : \mathbb{R}^{m_i} \to \mathbb{R}^{o_i}$, $h_i = o_{i-1} \forall i \in 2, \cdots k$, $h_1 = m$, and $o_k = 1$) with intermediary models ($f_i(x) = (h_i \circ \cdots \circ h_1)(x)$), DeepSHAP attributions are computed as:

$$\psi^k = \hat{\phi}(h_k, x^e, x^b) \tag{2}$$

$$\psi^i = \hat{\phi}(h_i, x^e, x^b)\big(\psi^{i+1} \oslash (f_i(x^e) - f_i(x^b))\big), \ i \in 1, \cdots, k-1. \tag{3}$$

---

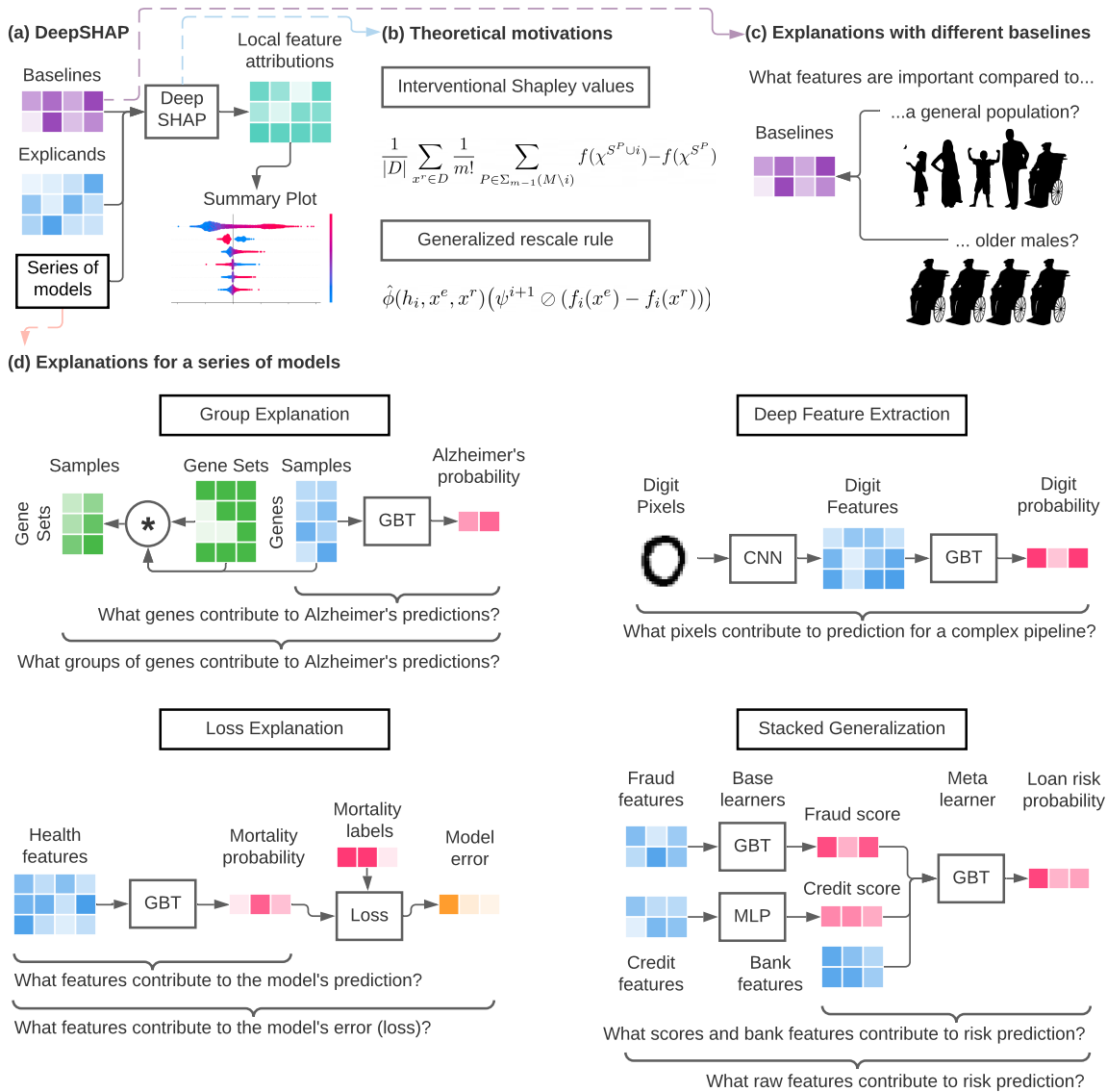[2][22] show an analogous result for an "input distribution" under an observational conditional expectation lift.

**(a) DeepSHAP**

Baselines

Explicands

Series of models

Deep SHAP

Local feature attributions

Summary Plot

**(b) Theoretical motivations**

Interventional Shapley values

$$\frac{1}{|D|}\sum_{x^r \in D}\frac{1}{m!}\sum_{P \in \Sigma_{m-1}(M \setminus i)} f(\chi^{S^P \cup i}) - f(\chi^{S^P})$$

Generalized rescale rule

$$\hat{\phi}(h_i, x^e, x^r)\big(\psi^{i+1} \oslash (f_i(x^e) - f_i(x^r))\big)$$

**(c) Explanations with different baselines**

What features are important compared to...

Baselines

...a general population?

... older males?

**(d) Explanations for a series of models**

Group Explanation

Samples | Gene Sets | Samples | Alzheimer's probability

Gene Sets * Genes → GBT

What genes contribute to Alzheimer's predictions?

What groups of genes contribute to Alzheimer's predictions?

Deep Feature Extraction

Digit Pixels | Digit Features | Digit probability

CNN → GBT

What pixels contribute to prediction for a complex pipeline?

Loss Explanation

Health features | Mortality probability | Mortality labels | Model error

GBT → Loss

What features contribute to the model's prediction?

What features contribute to the model's error (loss)?

Stacked Generalization

Fraud features | Base learners | Fraud score | Meta learner | Loan risk probability

GBT

Credit score

MLP → GBT

Credit features | Bank features

What scores and bank features contribute to risk prediction?

What raw features contribute to risk prediction?

Figure 1: **DeepSHAP estimates Shapley value feature attributions to explain a series of models using a baseline distribution.** (a) Local feature attributions with DeepSHAP require explicands (samples being explained), a baseline distribution (samples being compared to), and a model that is comprised of a series of models. They can be visualized to understand model behavior (Appendix Section A.3). (b) Theoretical motivation behind DeepSHAP (Methods Sections 6.1 and 6.4). (c) The baseline distribution is an important, but often overlooked, parameter that changes the scientific question implicit in the local feature attributions we obtain. (d) Explaining a series of models enables us to explain groups of features, model loss, and complex pipelines of models (deep feature extraction and stacked generalization). Experimental setups are described in Appendix Section A.2.

We use Hadamard division to denote an element-wise division of $\vec{a}$ by $\vec{b}$ that accommodates zero division, where, if the denominator $b_i$ is 0, we set $a_i/b_i$ to 0. The attributions $\hat{\phi}$ for a particular model in the stack are computed utilizing DeepLIFT with the rescale rule for deep models [18], interventional TreeSHAP for tree models [21], or exactly for linear models. Each intermediate attribution $\psi^i$ serves as feature attribution that

satisfies efficiency for $h_i$'s input features, where the attribution in the raw feature space is given by $\psi^1$. This approach takes inspiration from the chain rule applied specifically to deep networks in [18], that we extend to more general classes of models.

# 3    Incorporating a baseline distribution

We now use DeepSHAP to explain deep models with different choices of baseline distributions to empirically evaluate our theoretical connections to interventional conditional expectations.

## 3.1    Baseline distributions avoid bias



Figure 2: **Using a single all-black baseline image (DeepLIFT) leads to biased attributions compared to attributions with a randomly sampled baseline distribution (DeepSHAP).** The image is the explicand. The attribution plots are the sum of the absolute value of the feature attributions for the three channels of the input image. The pixel distribution is the distribution of pixels in terms of their grayscale values. The attribution distribution is the amount of attribution mass upon a group of pixels binned by their grayscale values.

We show that single baseline attributions are biased in a CNN that achieves 75.56% test accuracy (hyperparameters in Appendix Section A.2.1) in the CIFAR10 data set [31]. We aim to demonstrate that single baselines can lead to bias in explanations by comparing attributions using either a single baseline (an all-black image) as in DeepLIFT or a random set of 1000 baselines (random training images) as in DeepSHAP. Although the black pixels in the image are qualitatively important, using a single baseline leads to biased attributions with little attribution mass for black pixels (Figure 2). In comparison, averaging over multiple baselines leads to qualitatively more sensible attributions. Quantitatively, we show that despite the prevalence of darker pixels (pixel distribution plots in Figure 2), single baseline attributions are biased to give them low attribution, whereas averaging over many baselines more sensibly assigns a large amount of credit to dark pixels (attribution distribution plots in Figure 2). To generalize this finding beyond DeepSHAP, we replicate this bias for IME and IG, two popular feature attribution methods that similarly rely on baseline distributions (Appendix Section A.4.1).

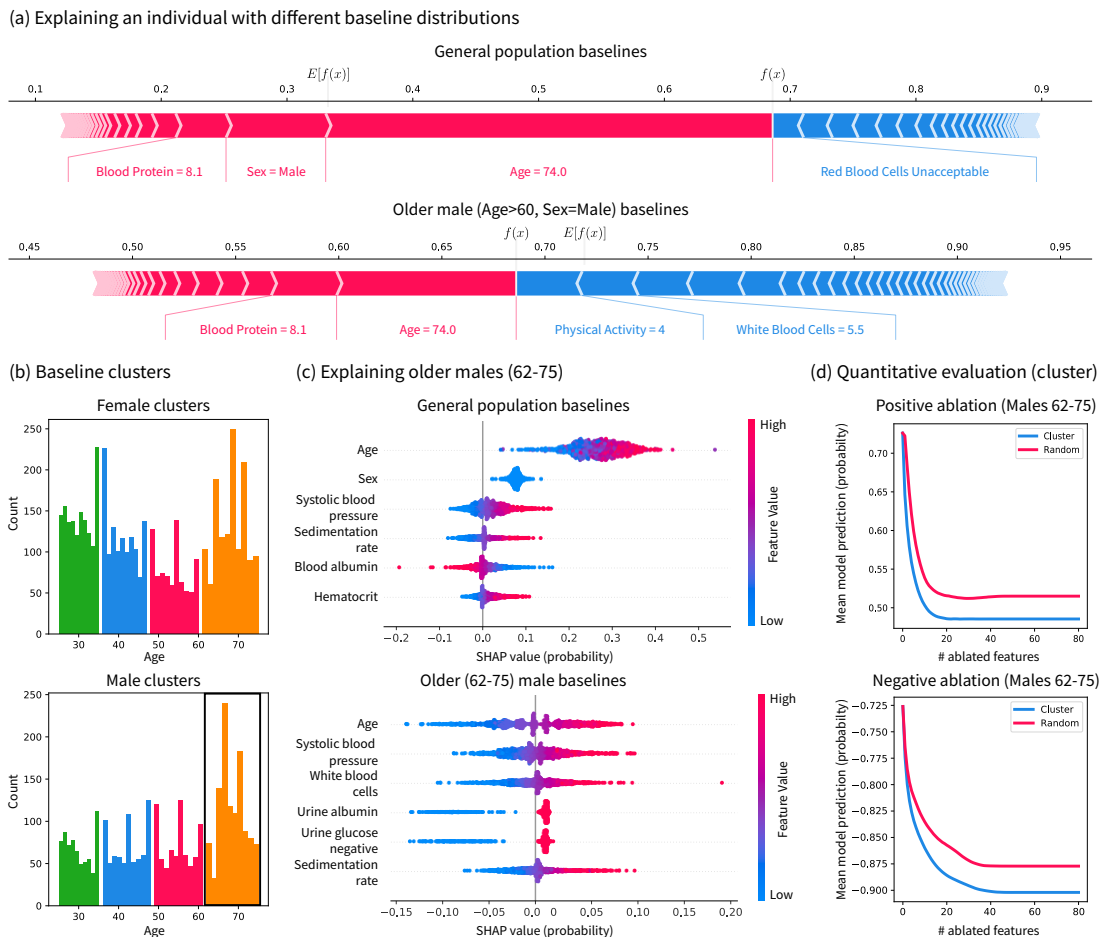## 3.2 Natural scientific questions with baseline distributions



Figure 3: **The baseline distribution is an important parameter for model explanation.** (a) Explaining an older male explicand with both a general population baseline distribution and an older male baseline distribution. (b) Automatically finding baseline distributions using 8-means clustering on age and sex. (c) Explaining the older male subpopulation (62-75 years old) with either a general population baseline or an older male baseline. (d) Quantitative evaluation of the feature attributions via positive and negative ablation tests where we mask with the mean of the older male subpopulation. Note that (b) shows summary plots (Appendix Section A.3.3) and (c) shows dependence plots (Appendix Section A.3.2).

To demonstrate the importance of baseline distributions as a parameter, we explain an MLP (hyper-parameters in Appendix Section A.2.2) with 0.872 ROC AUC for predicting fifteen year mortality in the NHANES I data set. We use DeepSHAP to explain an explicand relative to a baseline distribution drawn uniformly from all samples (Figure 3a (top)). This explanation places substantial emphasis on age and sex because it compares the explicand to a population that includes many younger/female individuals. However, in practice epidemiologists are unlikely to compare a 74-year old male to the general population. Therefore, we can manually select a baseline distribution of older males to reveal novel insights, as in Figure 3a (bottom). The impact of gender is gone because we compare only to males, and the impact of age is lower because we compare only to older individuals. Furthermore, the impact of physical activity is much higher because being physically active is more healthy compared to older individuals, who are less active than the general population. This example illustrates that the baseline distribution is an important parameter for feature attributions.

To provide a more principled approach to choosing the baseline distribution parameter, we propose

k-means clustering to select a baseline distribution (detail in Methods Section 6.3). Previous work analyzed clustering in the attribution space or contrasting to negatively/positively labelled samples [22]. In Figure 3b, we show clusters according to age and gender. Then, we explain many older male explicands using either a general population or an older male population baseline distribution (Figure 3c). When we compare to the older male baselines, the importance of age is centered around zero, sex is no longer important, and the importance orderings of remaining features change. Further, the inquiry we make changes from "What features are important for older males relative to a general population?" to "What features are important for older males relative to other older males?". To quantitatively evaluate whether our attributions answer the second inquiry, we can ablate features in order of their positive/negative importance by masking with the mean of the older male baseline distribution (Figure 3d, (Methods Section 6.8)). In both plots, lower curves indicate attributions that better estimated positive and negative importance. *For both tests, attributions with a baseline distribution chosen by k-means clustering substantially outperforms a baseline distribution drawn from the general population.*

We find that our clustering-based approach to selecting a baseline distribution has a number of advantages. Our recommendation is to choose baseline distributions by clustering according to non-modifiable, yet meaningful, features like age and gender. This yields explanations that answer questions relative to inherently interpretable subpopulations (e.g., older males). The first advantage is that choosing baseline distributions in this way decreases variance in the features that determined the clusters and subsequently reduces their importance to the model. This is desirable for age and gender because individuals typically cannot modify their age or gender in order to reduce their mortality risk. Second, this approach could potentially reduce model evaluation on off-manifold samples when computing Shapley values [32, 33] by considering only baselines within a reasonable subpopulation. The final advantage is that the flexibility of choosing a baseline distribution allows feature attributions to answer natural contrastive scientific questions [22] that improve model comprehensibility, as in Figure 3c.

# 4    Explaining a series of models

DeepSHAP (and DeepLIFT) have been shown to be very fast and performant explanation methods for explaining deep models [18, 34, 35]. In this section, we instead focus on evaluating our extension of DeepSHAP to accommodate a series of mixed models (trees, neural networks, and linear models) and address four impactful applications.

## 4.1    Group attributions identify meaningful gene sets

We explain two MLPs trained to predict Alzheimer's disease status and breast cancer tumor stage from gene expression data with test ROC AUC of 0.959 and 0.932, respectively. We aim to demonstrate that our approach to propagating attributions to groups contributes to model interpretability by validating our discoveries with scientific literature. Gene expression data is often extremely high dimensional; as such, solutions such as gene set enrichment analysis (GSEA) are widely used [36]. In contrast, we aim to attribute importance to gene sets while maintaining efficiency by proposing a *group rescale rule* (Methods Section 6.7). This rule sums attributions for genes belonging to each group and then normalizes according to excess attribution mass due to multiple groups containing the same gene. It generalizes to arbitrary groups of features beyond gene sets, such as categories of epidemiological features (e.g., laboratory measurements, demographic measurements, etc.).

We can validate several key genes identified by DeepSHAP. For Alzheimer's disease, the overexpression of SERPINA3 has been closely tied to prion diseases [37], and UBTD2 has been connected to frontotemporal dementia – a neurodenerative disorder [38]. For breast cancer tumor stage, UBE2C was positively correlated with tumor size and histological grade [39]. In addition to understanding gene importance, understanding higher level importance can be obtained using gene sets, i.e., groups of genes defined by biological pathways or co-expression. We obtain gene set attributions by grouping genes according to a curated gene set.[3]: the KEGG (Kyoto Encyclopedia of Genes and Genomes) pathway database (additional gene set attributions in Appendix Section A.4.3)

---

[3]Curated gene sets available here: https://www.gsea-msigdb.org/gsea/msigdb/collections.jsp#C2
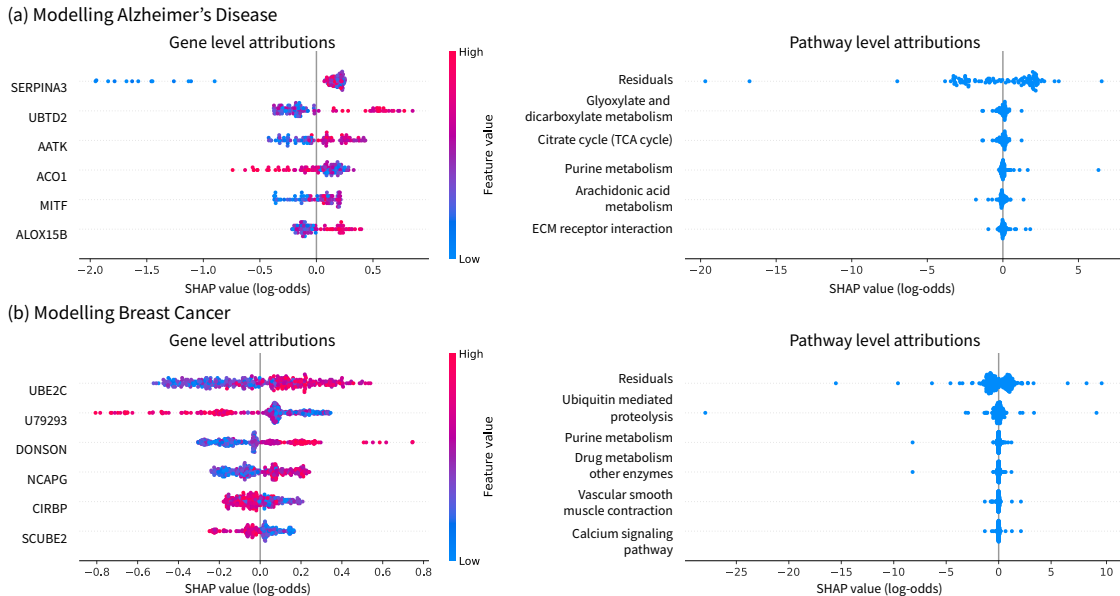
Figure 4: **Propagating attributions to gene sets enables higher level understanding.** (a) Gene and gene set attributions for predicting Alzheimer's disease using gene expression data. (b) Gene and gene set attributions for predicting breast cancer tumor stage using gene expression data. Residuals in the gene set attributions summarize contributions for genes that are not present in any gene set and describes variations in output not described by the pathways we analyzed. Note that (a) and (b) show summary plots (Appendix Section A.3.3).

Next, we verify important gene sets identified by DeepSHAP. For Alzheimer's disease, the glyoxylate and dicarboxylate metabolism pathway was independently identified based on metabolic biomarkers [40]; several studies have demonstrated aberrations in the TCA cycle in Alzheimer's disease brain [41]; and alterations of purine-related metabolites are known to occur in early stages of Alzheimer's disease [42]. For breast cancer, many relevant proteins are involved in ubiquitin-proteasome pathways [43] and purine metabolism was identified as a major metabolic pathway differentiating a highly metastatic breast cancer cell line from a slightly metastatic one [44]. *Identifying these phenotypically relevant biological pathways demonstrates that our group rescale rule identifies important pathways.*

## 4.2 Loss attributions provide insights to model behavior

We examine an NHANES (1999-2014) mortality prediction GBT model (0.868 test set ROC AUC) to show how explaining the model's loss (loss explanations) provides important insights different from insights revealed by explaining the model's output (output explanations). DeepSHAP lets us explain transformations of the model's output. For instance, we can explain a binary classification model in terms of its log-odds predictions, its probability predictions (often easier for non-technical collaborators to understand; see Appendix Section A.4.2), or its loss computed based on the prediction. Here, we focus on local feature attributions that explain per-sample loss.[4]

We train our model on the first five release cycles of the NHANES data (1999-2008) and evaluate it on a test set of the last three release cycles (2009-2014) (Figure 3a). As a motivating example, we simulate a covariate shift in the weight variable by re-coding it to be measured in pounds, rather than kilograms, in release cycles 7 and 8 (Figure 3b). Then, we ask, "Can we identify the impact of the covariate shift with feature attributions?" Comparing the train and test output attributions, release cycles 7 and 8 are skewed, but they mimic the same general shape of the training set attributions. If we did not color by release cycles, it might be difficult to identify the covariate shift. In contrast, for loss attributions with positive labels,

---

[4]This is analogous to model monitoring in [21], which is in fact enabled via the generalized rescale rule we present here.
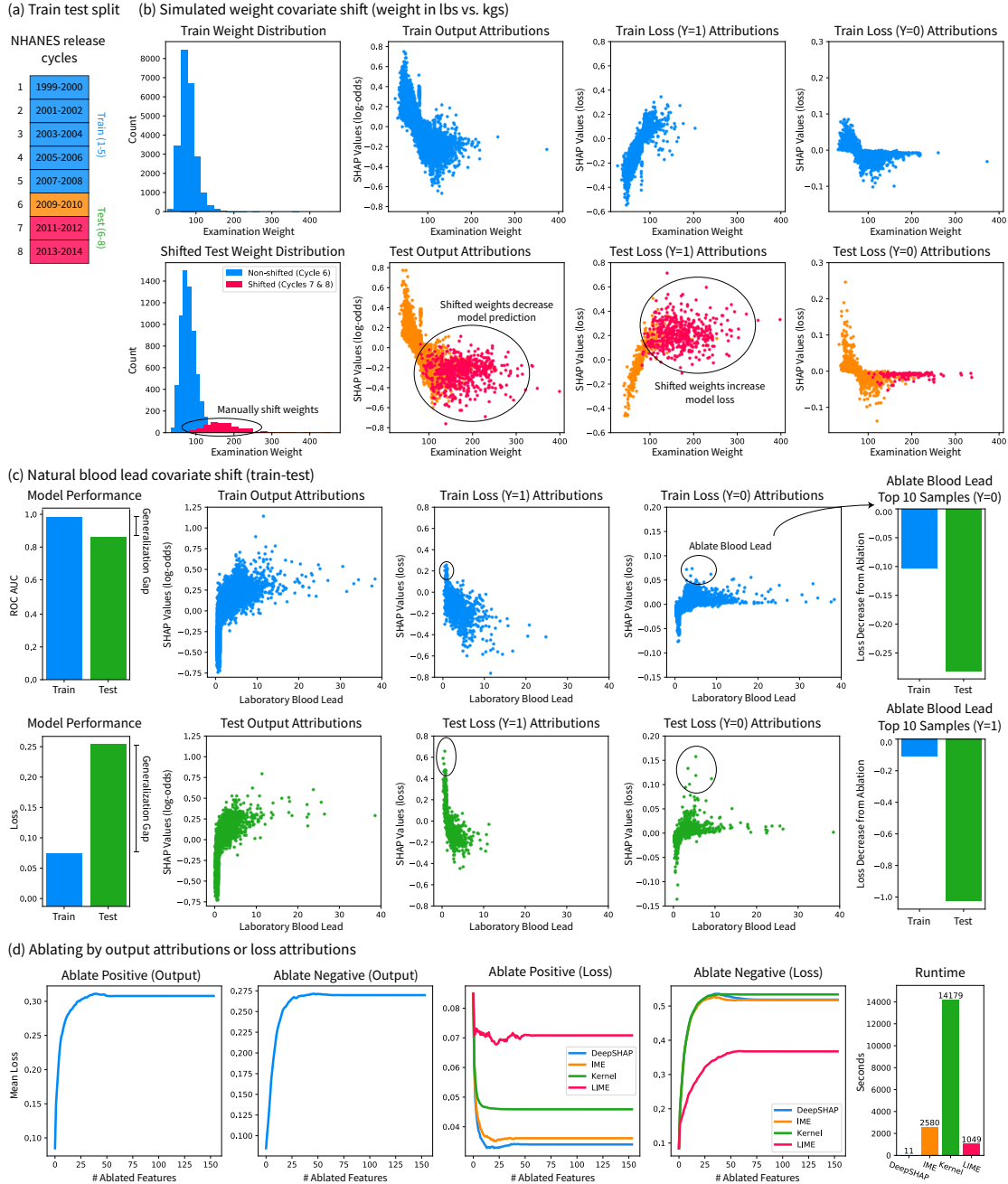
Figure 5: **Explanations of the model's loss rather than the model's prediction yields new insights.**
(a) We train on the first five cycles of NHANES (1999-2008) and test on the last three cycles (2009-2014). (b)
We identify a simulated covariate shift in cycles 7-8 (2011-2014) by examining loss attributions. (c) Under
a natural covariate shift, we identify and quantitatively validate test samples for which blood lead greatly
increases the loss in comparison to training samples. (d) We ablate output attributions (DeepSHAP) and
loss attributions (DeepSHAP, IME, KernelSHAP, and LIME) to show their respective impacts on model loss.
We compare only to model-agnostic methods for loss attributions because explaining model loss requires
explaining a series of models. Note that (b) and (c) show dependence plots (Appendix Section A.3.2).

we can identify that the falsely increased weight leads to many misclassified samples where the loss weight
attribution exceeds the expected loss. Although such debugging is powerful, it is not perfect. Note that in

the negatively labelled samples, we cannot clearly identify the covariate shift because higher weights are protective and lead to more confident negative mortality prediction.

Next, we examine the natural generalization gap induced by covariate shift over time, which shows a dramatically different loss in the train and test sets (Figure 5c). We can see that output attributions are similarly shaped between the train and test distributions; however, the loss attributions in the test set are much higher than in the training set. We can quantitatively verify that negative blood lead affects model performance more in the test set by ablating blood lead for the top 10 samples in the train and test sets according to their loss distributions. From this, we can see that blood lead constitutes a substantial covariate shift in the model's loss and helps explain the observed generalization gap.

As an extension of the quantitative evaluation in Figure 5c, we can visualize the impact on the model's loss of ablating by output attributions compared to ablating by loss attributions (Figure 5d). This ablation test (Methods Section 6.8) asks "What features are important to the model's performance (loss)?" Ablating the positive and negative attributions both increase the mean model loss by hiding features central to making predictions. However, ablating by the negative loss attribution directly increases the loss far more drastically than ablating by the output. More so, ablating positive loss attributions clearly decreases the mean loss, which is not achievable by output attribution ablation. Finally, we compare loss attributions computed using either a model-agnostic approach or DeepSHAP. In this setting, DeepSHAP is two orders of magnitude faster than model-agnostic approaches (IME, KernelSHAP, and LIME) while showing extremely competitive positive loss ablation performance and the best negative loss ablation performance.
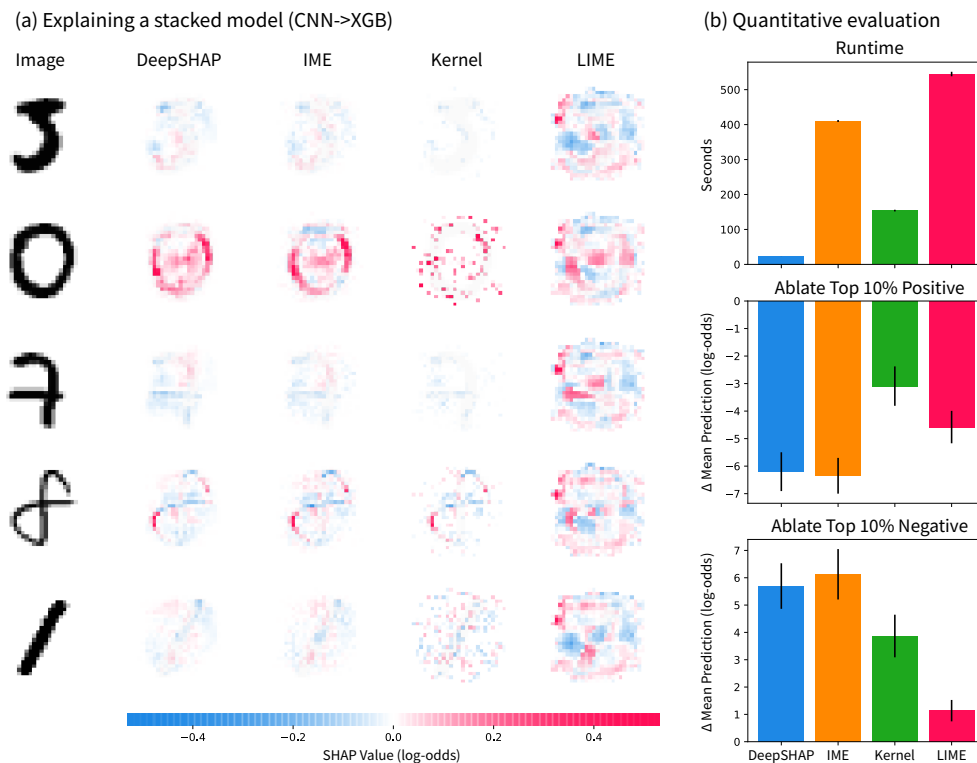
## 4.3    Explaining deep image feature extractors



Figure 6: **Explaining a series of models comprised of a convolutional neural network feature extractor and a gradient boosted tree classifier.** (a) Explanations from DeepSHAP and state of the art model-agnostic approaches. Each model-agnostic approach has a "number of samples" parameter which we set to 100,000. (b) Quantitative evaluation of approaches, including runtime and ablation of the top 10% of positive and negative features.

We compare DeepSHAP explanations to a number of model-agnostic explanations for a series of two models: a CNN feature extractor fed into a GBT model that classifies MNIST zeros with 0.998 test accuracy. In this example, non-linear transformations of the original feature space improve performance of the downstream model (Appendix Section A.4.4) but make model-specific attributions impossible. Qualitatively, we can see that DeepSHAP and IME are similar, whereas KernelSHAP is similar for certain explicands but not others[5] (Figure 6a). Finally, LIME's attributions show the shape of the original digit, but there is a consistent attribution mass around the surrounding parts of the digit. Qualitatively, we observe that the DeepSHAP attributions are sensible. The pixels that constitute the zero digit and the absence of pixels in the center of the zero are important for a positive zero classification.[6]

In terms of quantitative evaluations, we report the runtime and performance of the different approaches in Figure 6b. We see that DeepSHAP is an order of magnitude faster than model-agnostic approaches, with KernelSHAP being the second fastest. Then, we ablate the top 10% of important positive or negative pixels to see how the model's prediction changes. If we ablate positive pixels, we would expect the model's predictions to drop, and vice versa for negative pixels; doing both showed that DeepSHAP outperforms KernelSHAP and LIME, and performs comparably to IME at greatly reduced computational cost.

## 4.4   Explaining distributed proprietary models

We evaluate DeepSHAP explanations for a **consumer scoring** example that feeds a simulated GBT fraud score model and a simulated MLP credit score model into a GBT bank model, which classifies good risk performance (0.681 test ROCAUC) (Figure 7). Consumer scores (e.g., credit scores, fraud scores, health risk scores, etc.) describe individual behavior with predictive models [12]. A vast industry of data brokers generates consumer scores based on a plethora of consumer data. For instance, a single data broker in a 2014 FTC study had 3000 data segments on nearly every consumer in the United States, and another broker added three billion new records to its databases each month [45].[7] Unfortunately, explaining the models that use consumer scores can obscure important features. For instance, explaining the bank model in Figure 7a will tell us that fraud and credit scores are important (in Figure 7c), but these scores are inherently opaque to consumers [12]. The truly important features may instead be those that these scores use. A better solution might be model-agnostic methods that explain the entire pipeline at once. However, the model-agnostic approaches require access to all models. In Figure 7a, a single institution would have to obtain access to fraud, credit, and bank models to use the standard model-agnostic approaches (Figure 7b (left)). This may be fundamentally impractical because each of these models is proprietary. This opacity is concerning given the growing desire for transparency in artificial intelligence [12, 45, 46].

DeepSHAP naturally addresses this obstacle by enabling attributions to the original features without forcing companies to share their proprietary models *if* each institution in the pipeline agrees to work together and has a consistent set of baselines. Furthermore, DeepSHAP can combine any other efficiency-satisfying feature attribution method in an analogous way (e.g., integrated/expected gradients [19]). Altogether, DeepSHAP constitutes an effective way to "glue" together explanations across distributed models in industry. In particular, in Figure 7a, the lending institution can explain its bank model in terms of bank features and fraud and credit scores. The bank then sends fraud and credit score attributions to their respective companies, who can use them to generate DeepSHAP attributions to the original fraud and credit features. The fraud and credit institutions then send the attributions back to the bank, which can provide explanations in terms of the original, more interpretable features to their applicants (Figure 7d).

We first quantitatively verify that the DeepSHAP attributions for this pipeline are comparable to the model agnostic approaches in Figure 7b. We once again see that DeepSHAP attributions are competitive with the best performing attributions methods for ablating the top 5 most important positive or negative features. Furthermore, we see that DeepSHAP is several orders of magnitude faster than the best performing ablation methods (KernelSHAP and IME) and an order of magnitude faster and much more performant than LIME.
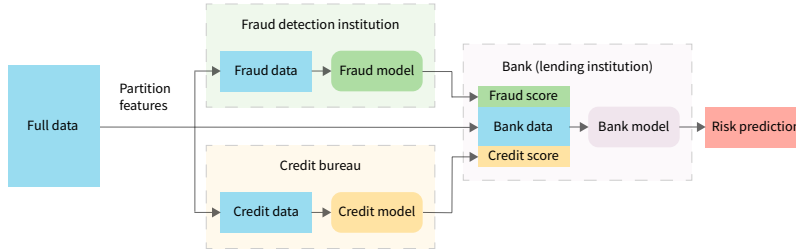
We can qualitatively verify the attributions in Figures 7c-d. In Figure 7c, we find that the fraud and credit scores are extremely important to the final prediction. In addition, bank features including low revolving

---

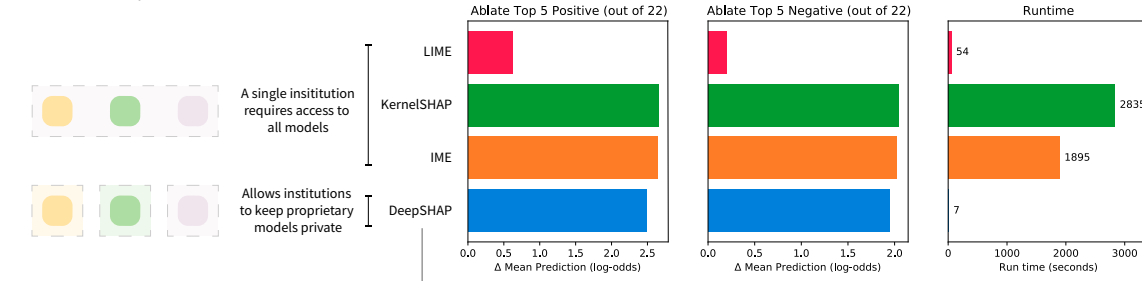[5]One potential reason is the regularization in the default settings of the package.

[6]The importance of the absence of pixels in the center of the zero is revealed because we use a baseline distribution; and it would not be revealed with an all-white baseline image.

[7]Furthermore, even the HELOC data set that we use initially had an "ExternalRiskEstimate" feature that we removed for clarity.
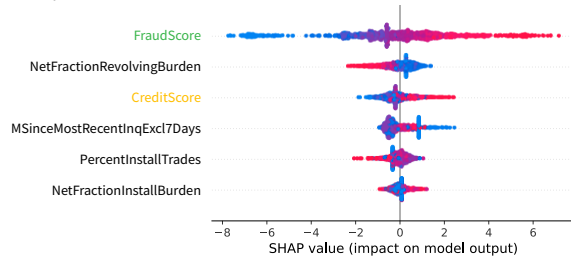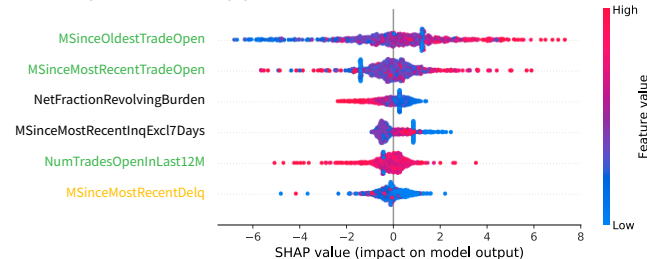
Figure 7: **Explaining a stacked generalization pipeline of models for the HELOC data set (details in Appendix Section A.1.7)** (a) A simulated model pipeline in the financial services industry. We partition the original set of features into fraud, credit, and bank features. We train a model to predict risk using fraud data and a model to predict risk using credit data. Then, we use the outputs of the fraud and credit models as scores alongside additional bank features to predict the final customer risk. (b) Ablation tests (ablating top 5 positive/negative features out of a total 22 features) comparing model agnostic approaches (LIME, KernelSHAP, IME with 100 samples), which require access to all models in the pipeline, and DeepSHAP, which allows institutions to keep their proprietary models private. (c) Summary plot of the top six features the bank model uses to predict risk (TreeSHAP). (d) Summary plot of the top six features the entire pipeline uses to explain risk (DeepSHAP). The green features originate from the fraud data, and the yellow features from the credit data. We explain 1000 randomly samples explicands using 100 randomly sampled baselines for all attribution methods. Note that (c) and (d) show summary plots (Appendix Section A.3.3).

balance divided by credit limit ("NetFractionRevolvingBurden") and low number of months since inquisitions ("MSinceMostRecentInqExcl7Days") are congruously important to good risk performance. Then, in Figure 7d we use the generalized rescale rule to obtain attributions in the original feature space. Doing so uncovers important variables hidden by the fraud and credit scores. In particular, we see that the fraud score heavily relied on a high number of months since the applicants's oldest trade ("MSinceOldestTradeOpen"), and the credit score relied on a low number of months since recent delinquency ("MSinceMostRecentDelq") in order to identify applicants that likely had good risk performance. Importantly, the pipeline we analyze in Figure 7a also constitutes a stacked generalization ensemble, which we analyze more generally in Appendix Section

A.4.5.

# 5  Discussion

In this manuscript, we presented examples where explaining a series of models is critical. Series of models are prevalent in a variety of applications (health, finance, environmental science, etc.), where understanding model behavior contributes important insights. Furthermore, having a fast approach to explain these complex pipelines may be a major desiderata for a diagnostic tool to debug ML models.

The practical applications we focus on in this paper include gene set attribution, where the number of features far surpasses the number of samples. In this case, we provide a rule that aggregates group attributions to higher level groups of features while maintaining efficiency. Second, we demonstrate the utility of explaining transformations of a model's default output (Appendix Section A.4.2). Explaining the probability output rather than the log-odds output of a logistic model yields more naturally interpretable feature attributions. Furthermore, explaining the loss of a logistic model enables debugging model performance and identification of covariate shift. A third application is neural network feature extraction, where pipelines may include transformations of the original features fed into a different model. In this setting we demonstrate the computational tractability of DeepSHAP compared to model-agnostic approaches. Finally, because our approach propagates feature attributions through a series of models while satisfying efficiency at each step (Methods Section 6.5), the intermediary attributions at each part of the network can be interpreted as well. We use this to understand the importance of both consumer scores and the original features used by the consumer scores.

In consumer scoring, distributed proprietary models (i.e., models that exist in different institutions) have historically been an obstacle to transparency. This lack of transparency is particularly concerning given the prevalence of consumer scores, with some data brokers having thousands of data segments on nearly every American consumer [45]. In addition, many new consumer scores fall outside the scope of previous regulations (e.g., the Fair Credit Reporting Act and the Equal Credit Opportunity Act) [12]. In fact, these new consumer scores that depend on features correlated with protected factors (e.g., race) can reintroduce discrimination hidden behind proprietary models, which is an issue that has historically been a concern in credit scores (the oldest existing example of a consumer score) [12]. DeepSHAP naturally enables feature attributions in this setting and takes a significant and practical step towards increasing the transparency of consumer scores and provides a tool to help safeguard against hidden discrimination.

It should be noted that we focus specifically on evaluating DeepSHAP for a series of mixed model types. Previous work evaluates the rescale rule for explaining deep models, specifically. The original presentation of the rescale rule [18] demonstrates its applicability to deep networks in explaining digit classification and regulatory DNA classification. Schwab & Karlen shows that for explaining deep networks, DeepSHAP is a very fast yet performant approach in terms of an ablation test for explaining MNIST and CIFAR images.[8] Finally, Sixt *et al.* shows that many modified back propagation feature attribution techniques are independent of the parameters of later layers, with the exception of DeepLIFT. This particularly significant finding suggests that compared to most fast back propagation-based deep feature attribution approaches, DeepSHAP (which relies on the same rescale rule as DeepLIFT) is not ignorant of later layers in the network.

Although DeepSHAP works very well for explaining a series of mixed model types in practice, an inherent limitation is that it is not guaranteed to satisfy the desirable axioms (e.g., implementation invariance) that other feature attribution approaches satisfy (assuming exact solutions to their intractable problem formulations) [15, 17, 19]. This suggests that DeepSHAP may be more appropriate for model debugging or for identifying scientific insights that warrant deeper investigation, particularly in settings where models or the input dimension is huge and tractability is a major concern. However, for applications where high-stakes decision making is important, it may be more appropriate to run axiomatic approaches to completion or use interpretable models [47]. Furthermore, in many real world circumstances, such as distributed proprietary models based on credit risk scores, exact axiomatic approaches and interpretable models are not feasible. In these cases DeepSHAP represents a promising direction that allows multiple agents to collaboratively build explanations while maintaining separation of model ownership.

---

[8]Although their approach, CXPlain, is comparably fast at attribution time, it has the added cost of training a separate explanation model.

# 6 Methods

## 6.1 The Shapley value

The Shapley value is a solution concept for allocating credit among players $(M = 1, \cdots, m)$ in an $m$-person game. The game is fully described by a set function $v(S) : \mathcal{P}(S) \to \mathbb{R}^1$ that maps the power set of players $S \subseteq M$ to a scalar value. The Shapley value for player $i$ is the average marginal contribution of that player for all possible permutations of remaining players:

$$\phi_i(v) = \frac{1}{m!} \sum_{P \in \Sigma_{m-1}(M)} (v(S^P \cup i) - v(S^P)). \tag{4}$$

We denote the finite symmetric group $\Sigma_{n-1}(M)$, which is the set of all possible permutations, and $S^P$ to be the set of players before player $i$ in the permutation $P$. The Shapley value is a provably unique solution under a set of axioms (Appendix Section A.5.1). One axiom that we focus on in this paper is *efficiency*:

$$\sum_{i=1}^{m} \phi_i(v) = v(M) - v(\emptyset). \tag{5}$$

**Adapting the Shapley value for feature attribution of ML models**

Unfortunately, the Shapley value cannot assign credit for an ML model $(f(x) : \mathbb{R}^m \to \mathbb{R}^1)$ directly because most models require inputs with values for every feature, rather than a subset of features. Accordingly, feature attribution approaches based on the Shapley value define a new set function $v(S)$ that is a "lift" of the original model [48]. In this paper, we focus on local feature attributions,[9] which describe a model's behavior for a single sample, called an explicand $(x^e)$. A "lift" is defined as:

$$\mu(f, x^e, S) : \mathbb{R}^m \times 2^m \to \mathbb{R}^1. \tag{6}$$

One common lift is the *observational conditional expectation*, where the lift is the conditional expectation of the model's output holding features in $S$ fixed to $x_S^e$ and $X$ is a multivariate random variable with joint distribution $D$:

$$\mu_D^{obs}(f, x^e, S) = \mathbb{E}_D[f(X)|X_S = x_S^e]. \tag{7}$$

Another common lift is the *interventional conditional expectation* with a flat causal graph, where we "intervene" on features by breaking the dependence between features in $X_S$ and the remaining features using the causal inference *do*-operator [30]:

$$\mu_D^{int}(f, x^e, S) = \mathbb{E}_D[f(X)|do(X_S = x_S^e)]. \tag{8}$$

Both approaches have tradeoffs that have been described elsewhere [14, 22, 32, 33, 49]. Here, we focus on the interventional approach because it is most closely related to DeepSHAP and does not require estimating the joint density of $X$.

The Shapley values computed for any lift will satisfy efficiency in terms of the lift $\mu$. However, for the interventional and observational lift described above, the Shapley value will also satisfy efficiency in terms of the model's prediction:

$$\sum_i \phi_i^{\mu_D}(f, x^e) = f(x^e) - \mathbb{E}_D[f(X)]. \tag{9}$$

This means that attributions can naturally be understood to be in the scale of the model's predictions (e.g., log-odds or probability for binary classification).

---

[9]As opposed to global feature attributions, which measure feature importance of a model across an entire data set.

## 6.2 Interventional Shapley values baseline distribution

We can define a single baseline lift

$$\mu_{x^b}^{int}(f, x^e, S) = \mathbb{E}_{\{x^b\}}[f(X)|do(X_S = x_S^e)] = \chi^S, \tag{10}$$

where $\chi^S$ is a spliced sample and $\chi_i^S = x_i^e$ if $i \in S$, else $\chi_i^S = x_i^b$.

Then, we can decompose the Shapley value $\phi_i(f, x^e)$ for the interventional conditional expectation lift (eq. 8) (henceforth referred to as the interventional Shapley value) into an average of Shapley values with single baseline lifts (proof in Appendix Section A.5.2):

$$\phi_i(f, x^e, D) = \frac{1}{|D|} \sum_{x^b \in D} \underbrace{\frac{1}{m!} \sum_{P \in \Sigma_{m-1}(M)} f(\chi^{S^P \cup i}) - f(\chi^{S^P})}_{\text{Shapley value for single baseline lift}} \tag{11}$$

$$= \frac{1}{|D|} \sum_{x^b \in D} \phi_i(f, x^e, x^b). \tag{12}$$

Here, $D$ is an empirical distribution with equal probability for each sample in a baseline data set. An analogous result exists for the observational conditional distribution lift using an input distribution[22].[10]

In the original DeepLIFT paper, [18] recommend two heuristic approaches to define baseline distributions: (1) choosing a sensible single baseline and (2) averaging over multiple baselines. In addition, DeepSHAP, as previously described in [17], created attributions with respect to a single baseline equal to the expected value of the inputs. In this paper, we show that from the perspective of Shapley values with an interventional conditional expectation lift, averaging over feature attributions computed with single baselines drawn from an empirical distribution is the correct approach. One exception to this are linear models, where taking the average as the baseline is equivalent to averaging over many single baseline feature attributions [49]. Interventional Shapley values computed with a single baseline satisfy efficiency in terms of the model's prediction:

$$\sum_i \phi_i(f, x^e, x^b) = f(x^e) - f(x^b). \tag{13}$$

## 6.3 Selecting a baseline distribution

As in the previous section, we define a baseline distribution $D$ over which we compute Shapley values with single baseline lifts. This baseline distribution is naturally chosen to be a distribution over the training data $X^{train}$, where each sample $x^j \in \mathbb{R}^m$ has equal probability. The interpretation of this distribution is that the explicand is compared to each baseline in $D$. This means that the interventional Shapley values implicitly create attributions that explain the model's output relative to a baseline distribution.

Although the entire training distribution is a natural and interpretable choice of baseline distribution, it may be desirable to use others. To automate the process of choosing such an interpretable baseline distribution, we turn to unsupervised clustering. We utilize k-means clustering on a reduced version of the training data ($\hat{X}^{train}$) comprised of $\hat{x}^j = [x_i^j \forall i \in M_r]$ with a reduced set of features ($M_r$). The output of the k-means clustering are clusters $C_1, \cdots, C_k$ with means $\mu_1, \cdots, \mu_k$ that minimize the following objective on the reduced training data:

$$\underset{C_1, \cdots, C_k}{\arg\min} \sum_{i=1}^{k} \sum_{\hat{x} \in C_i} ||\hat{x} - \mu_i||^2. \tag{14}$$

Then, the cluster selected as a baseline distribution explaining an explicand $x^e$ is chosen based on:

$$\underset{i}{\arg\min} ||\hat{x}^e - \mu_i||^2. \tag{15}$$

---

[10]The attributions for these single baseline games are also analogous to baseline Shapley in [14].

## 6.4 A generalized rescale rule to explain a series of models

We define a *generalized rescale rule* to explain an arbitrary series of models that propagates approximate Shapley values with an interventional conditional expectation lift for each model in the series.[11] To describe the approach, we define a *series of models* to be a composition of functions $f_k(x) = (h_k \circ \cdots \circ h_1)(x)$, and we define intermediary models $f_i(x) = (h_i \circ \cdots \circ h_1)(x)$, $i = 1, \cdots, k$. We define the domain and codomain of each model in the series as $h_i(x) : \mathbb{R}^{m_i} \to \mathbb{R}^{o_i}$. Then, we can define the propagation for a single baseline[12] recursively:

$$\psi^k = \hat{\phi}(h_k, x^e, x^b) \tag{16}$$

$$\psi^i = \hat{\phi}(h_i, x^e, x^b)\big(\psi^{i+1} \oslash (f_i(x^e) - f_i(x^b))\big), \ i \in 1, \cdots, k-1. \tag{17}$$

We use Hadamard division to denote an element-wise division of $\vec{a}$ by $\vec{b}$ that accommodates zero division, where if the denominator $b_i$ is 0, we set $a_i/b_i$ to 0. Additionally, $\hat{\phi}$ are an appropriate feature attribution technique that approximates interventional Shapley values while crucially satisfying efficiency for the model $h_i$ it is explaining. In this paper, we utilize DeepLIFT (rescale) for deep models, TreeSHAP for tree models, and exact interventional Shapley values for linear models. We define efficiency as $\hat{1}_{1 \times m_i}\hat{\phi}(h_i, x^e, x^b) = f_i(x^e) - f_i(x^b)$ where $\hat{1}_{a \times b}$ is a matrix of ones with shape $a \times b$ and the approximate Shapley value functions $\hat{\phi}$ return matrices in $\mathbb{R}^{(m_i \times o_i)}$. The final attributions in the original feature space are:

$$\phi_i(f_k, x^e, x^b) = \psi_i^1. \tag{18}$$

Furthermore, this approach yields intermediate attributions that serve as meaningful feature attributions. In particular, $\psi^i$ can be interpreted as the importance of the inputs to the model $(h_k \circ \cdots \circ h_i)$, where the new explicand and baseline are $(h_{i-1} \circ \cdots \circ h_1)(x^e)$ and $(h_{i-1} \circ \cdots \circ h_1)(x^b)$, respectively. This approach takes inspiration from the chain rule applied specifically for deep networks in [18], but we extend it to more general classes of models.

## 6.5 Efficiency for intermediate attributions

As one might expect, each intermediate attribution $\psi^i$ satisfies efficiency:

**Theorem 1.** Each attribution $\psi^i \in \mathbb{R}^m, \forall i \in 1, \cdots, k$ satisfies efficiency and sums up to $f_k(x^e) - f_k(x^b)$.

*Proof.* We will prove by induction that

$$\hat{1}_{1 \times m_i}\psi^i = f_k(x^e) - f_k(x^b), \forall i \in 1, \cdots, k. \tag{19}$$

For simplicity of notation, denote $\hat{\phi}^i = \hat{\phi}(h^i, x^e, x^b)$.
*Assumption:* Each $\hat{\phi}$ satisfies efficiency

$$\hat{1}_{1 \times m_i}\hat{\phi}^i = f_i(x^e) - f_i(x^b). \tag{20}$$

*Base Case:* By our assumption,

$$\hat{1}_{1 \times m_k}\psi^k = f_k(x^e) - f_k(x^b). \tag{21}$$

*Induction Step:*

$$\psi^i = \hat{\phi}\big(\psi^{i+1} \oslash (f_i(x^e) - f_i(x^b))\big) \tag{22}$$

$$\hat{1}_{1 \times m_i}\psi^i = \hat{1}_{1 \times m_i}\hat{\phi}\big(\psi^{i+1} \oslash (f_i(x^e) - f_i(x^b))\big) \tag{23}$$

$$= (f_i(x^e) - f_i(x^b))\big(\psi^{i+1} \oslash (f_i(x^e) - f_i(x^b))\big) \tag{24}$$

$$= \hat{1}_{1 \times o_i}\psi^{i+1} \tag{25}$$

$$= \hat{1}_{1 \times m_{i+1}}\psi^{i+1} \tag{26}$$

$$= f_k(x^e) - f_k(x^b). \tag{27}$$

---

[11]This generalized chain rule will also generalize to any feature attribution method that satisfies efficiency.
[12]The baseline distribution attribution $\phi_i(f, x^e, D)$ is then simply the average across many of these single baseline attributions $\phi_i(f, x^e, x^b)$.

*Conclusion:* By the principle of induction, each intermediate attribution satisfies efficiency (eq. 19).

<div style="text-align: right">□</div>

Then, because the interventional Shapley value with a baseline distribution is the average of many single baseline attributions, it satisfies a related notion of efficiency:

$$\sum_i \phi_i(f_k, x^e) = \sum_i \sum_{x^b \in D} \phi_i(f_k, x^e, x^b) \tag{28}$$

$$= \sum_{x^b \in D} \sum_i \phi_i(f_k, x^e, x^b) \tag{29}$$

$$= \sum_{x^b \in D} f_k(x^e) - f_k(x^b) \tag{30}$$

$$= f_k(x^e) - \frac{1}{|D|} \sum_{x^b \in D} f_k(x^b). \tag{31}$$

This can be naturally interpreted as the difference between the explicand's prediction and the expected value of the function across the baseline distribution.

An additional property of the generalized rescale rule is that although it is an approximation to the interventional Shapley values in the general case, if every model in the composition is linear ($h_i(x) = \beta x$), then this propagation exactly yields the interventional Shapley values (Appendix Section A.5.3).

## 6.6   Connecting DeepLIFT's rules to the Shapley values

Now we can connect the Shapley values to DeepLIFT's Rescale and RevealCancel rules. Both rules aim to satisfy an efficiency axiom (what they call *summation to delta*) and can be connected to an interventional conditional expectation lift with a single baseline (as in Section 6.3).

In fact, multi-layer perceptrons are a special case where the models in the series are non-linearities applied to linear functions. We first represent deep models as a composition of functions $(h_1 \circ \cdots \circ h_k)(x)$. The Rescale and RevealCancel rules canonically apply to a specific class of function: $h_i(x) = (f \circ g)(x)$, where $f$ is a non-linear function and $g$ is a linear function parameterized by $\beta \in \mathbb{R}^m$. We can interpret both rules as an approximation to interventional Shapley values based on the following definition.

**Definition 6.1** (k-partition approximation)**.** A k-partition approximation to the Shapley values splits the features in $x \in \mathbb{R}^m$ into $K$ disjoint sets. Then, it exactly computes the Shapley value for each set and propagates it linearly to each component of the set.

The Rescale rule can be described as a 1-partition approximation to the Interventional Shapley values for $h_i(x)$, while the RevealCancel rule can be described as a 2-partition approximation that splits according to whether $\beta_i x_i > t$, where the threshold $t = 0$. This k-partition approximation lets us consider alternative variants of the Rescale and RevealCancel rules that incur exponentially larger costs in terms of $K$ and for different choices of thresholds.

## 6.7   Explaining groups of input features

Here, we further generalize the Rescale rule to support groupings of features in the input space. Having such a method can be particularly useful when explaining models with very large numbers of features that are more understandable in higher level groups. One natural example is gene expression data, where the numbers of features is often extremely large.

We introduce a *group rescale rule* that facilitates higher level understanding of feature attributions. We can define a set of groups $G_1, \cdots, G_o$ whose members are the input features $x_i$. If each group is disjoint and covers the full set of features, then a natural group attribution that satisfies efficiency is the sum:

$$\phi_{G_j}^0(f, x^e) = \sum_{i \in G_j} \phi_i(f, x^e). \tag{32}$$

If the groups are not disjoint or do not cover all input features, then the above attributions do not satisfy efficiency. To address this, we define a residual group $G_R$ that covers all input features not covered by the remaining groups. Then, the new attributions are a rescaled version of eq. 32

$$\phi_{G_j}(f, x^e) = \phi^0_{G_j}(f, x^e) \times \frac{\sum \phi_{G_j}(f, x^e)}{\sum \phi_i(f, x^e)}. \tag{33}$$

We can naturally extend this approach to accommodate non-uniform weighting of group elements, although we do not experiment with this in our paper.

## 6.8 Ablation Tests

We evaluate our feature attribution methods with *ablation tests* [21, 50]. In particular, we rely on a simple yet intuitive ablation test. For a matrix of explicands $X^e \in \mathbb{R}^{n_e, m}$, we can get attributions $\phi(f, X^e) \in \mathbb{R}^{n_e, m}$. The ablation test is defined by three parameters: (1) the feature ordering, (2) an imputation sample $x^b \in \mathbb{R}^m$, and (3) an evaluation metric. Then, the ablation test replaces features one at a time with the baseline's feature value based on the feature attributions to assess the impact on the evaluation metric. We can iteratively define the ablation test based on modified versions of the original explicands:

$$X^{e,0} = X^e \tag{34}$$

$$X^{e,k} = X^e \odot I_k(\phi) + x^b \odot (1 - I_k(\phi)), \forall k \in 1, \cdots, m. \tag{35}$$

Note that $x^b := [\ \underbrace{x^b \cdots x^b}_{n_e \text{ elements}}\ ]^T$ and $I_k(\phi) := I_k(\phi(f, X^e)) = \arg\max_{k, axis=1}(\phi(f, X^e))$, where $\arg\max_{k, axis=1}(G)$ returns an indicator matrix of the same size as $G$ and 1 indicates that the element was in the maximum $k$ elements across a particular axis.

Then, the ablation test measures the mean model output (e.g., the predicted log-odds, predicted probability, the loss, etc.) if we ablate $k$ features to be the average over the predictions for each ablated explicand:

$$\frac{1}{n_e} \sum_{i \in 1, \cdots, n_e} f(X_i^{e,k}). \tag{36}$$

Note that for our ablation tests we focus on either the positive or the negative elements of $\phi$, since the expected change in model output is clear if we ablate only by positive or negative attributions.

Ablation tests are a natural approach to test whether feature attributions are correct for a set of explicands. For feature attributions that explain the predicted log-odds, a natural choice of model output for the ablation test is the mean of the log-odds predictions. Then, as we ablate increasing numbers of features, we expect to see the model's output change. When we ablate the most positive features, the mean model output should decrease substantially. As we ablate additional features, the mean model output should still decrease, but less drastically so. This implies that, for positive ablations, lower curves imply attributions that better described the model's behavior. In contrast, for negative ablations

# 7 Acknowledgements

# 8 Competing Interests

The authors declare that there are no competing interests.

# 9 Data availability

The NHANES I, NHANES 1999-2014, CIFAR, and MNIST data sets are all publicly available. The HELOC data set can be obtained by accepting the data set usage license: (`https://community.fico.com/s/explainable-machine-learning-challenge?tabset-3158a=a4c37`). Metabric data access is restricted and requires getting an approval through Sage Bionetworks Synapse website: `https://www.synapse.org/#!Synapse:syn1688369` and `https://www.synapse.org/#!Synapse:syn1688370`. ROSMAP data access is restricted and requires getting an approval through Sage Bionetworks Synapse website: `https://www.synapse.org/#!Synapse:syn3219045` and is available as part of the AD Knowledge Portal [51].

# 10 Code availability

The code for the experiments is available here: `https://github.com/suinleelab/DeepSHAP`.

# 11 Author Contribution

H.C. contributed to study design, data analysis, and manuscript preparation. S.M.L. contributed to data analysis and manuscript preparation. S.L. contributed to study design, and manuscript preparation.

# References

1. Wang, S.-Q., Yang, J. & Chou, K.-C. Using stacked generalization to predict membrane protein types based on pseudo-amino acid composition. *Journal of Theoretical Biology* **242,** 941–946 (2006).

2. Healey, S. P. *et al.* Mapping forest change using stacked generalization: An ensemble approach. *Remote Sensing of Environment* **204,** 717–728 (2018).

3. Bhatt, S. *et al.* Improved prediction accuracy for disease risk mapping using Gaussian process stacked generalization. *Journal of The Royal Society Interface* **14,** 20170520 (2017).

4. Doumpos, M. & Zopounidis, C. Model combination for credit risk assessment: A stacked generalization approach. *Annals of Operations Research* **151,** 289–306 (2007).

5. *Otto Group Product Classification Challenge* `https://www.kaggle.com/c/otto-group-product-classification-challenge/discussion/14335`.

6. Wolpert, D. H. Stacked generalization. *Neural Networks* **5,** 241–259 (1992).

7. Guo, H. & Gelfand, S. B. *Classification trees with neural network feature extraction* in *Proceedings 1992 IEEE Computer Society Conference on Computer Vision and Pattern Recognition* (1992), 183–184.

8. Chen, Y., Jiang, H., Li, C., Jia, X. & Ghamisi, P. Deep feature extraction and classification of hyperspectral images based on convolutional neural networks. *IEEE Transactions on Geoscience and Remote Sensing* **54,** 6232–6251 (2016).

9. Xu, Y. *et al. Deep learning of feature representation with multiple instance learning for medical image analysis* in *2014 IEEE international conference on acoustics, speech and signal processing (ICASSP)* (2014), 1626–1630.

10. Liang, H., Sun, X., Sun, Y. & Gao, Y. Text feature extraction based on deep learning: a review. *EURASIP Journal on Wireless Communications and Networking* **2017,** 1–12 (2017).

11. Jahankhani, P., Kodogiannis, V. & Revett, K. *EEG signal classification using wavelet feature extraction and neural networks* in *IEEE John Vincent Atanasoff 2006 International Symposium on Modern Computing (JVA'06)* (2006), 120–124.

12. Dixon, P. & Gellman, R. *The scoring of America* in *World Privacy Forum* (2014).

13. Fay, B. *Credit Scoring: FICO, VantageScore; Other Models* Nov. 2020. `https://www.debt.org/credit/report/scoring-models/`.

14. Sundararajan, M. & Najmi, A. *The many Shapley values for model explanation* in *Proceedings of the International Conference on Machine Learning* (2020), 513–523.

15. Strumbelj, E. & Kononenko, I. An efficient explanation of individual classifications using game theory. *The Journal of Machine Learning Research* **11,** 1–18 (2010).

16. Ribeiro, M. T., Singh, S. & Guestrin, C. *"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), 1135–1144.

17. Lundberg, S. M. & Lee, S.-I. *A unified approach to interpreting model predictions* in *Advances in Neural Information Processing Systems* (2017), 4765–4774.

18. Shrikumar, A., Greenside, P. & Kundaje, A. *Learning important features through propagating activation differences* in *Proceedings of the 34th International Conference on Machine Learning-Volume 70* (2017), 3145–3153.

19. Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. *arXiv preprint arXiv:1703.01365* (2017).

20. Breiman, L., Friedman, J., Stone, C. J. & Olshen, R. A. *Classification and regression trees* (CRC press, 1984).

21. Lundberg, S. M. *et al.* Explainable AI for Trees: From Local Explanations to Global Understanding. *CoRR* **abs/1905.04610.** arXiv: 1905.04610. `http://arxiv.org/abs/1905.04610` (2018).

22. Merrick, L. & Taly, A. *The Explanation Game: Explaining Machine Learning Models Using Shapley Values* in *International Cross-Domain Conference for Machine Learning and Knowledge Extraction* (2020), 17–38.

23. A Bennett, D., A Schneider, J., Arvanitakis, Z. & S Wilson, R. Overview and findings from the religious orders study. *Current Alzheimer Research* **9,** 628–645 (2012).

24. Bennett, D. A. *et al.* Religious orders study and rush memory and aging project. *Journal of Alzheimer's Disease* **64,** S161–S189 (2018).

25. Curtis, C. *et al.* The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature* **486,** 346–352 (2012).

26. Pereira, B. *et al.* The somatic mutation profiles of 2,433 breast cancers refine their genomic and transcriptomic landscapes. *Nature Communications* **7,** 1–16 (2016).

27. Cox, C. S. *Plan and operation of the NHANES I Epidemiologic Followup Study, 1992* **35** (National Ctr for Health Statistics, 1998).

28. LeCun, Y. The MNIST database of handwritten digits. *http://yann. lecun. com/exdb/mnist/* (1998).

29. *FICO, Xml challenge* `https://community.fico.com/s/explainable-machine-learning-challenge`. [Online; accessed 01-June-2021].

30. Janzing, D., Minorics, L. & Blöbaum, P. Feature relevance quantification in explainable AI: A causality problem. *arXiv preprint arXiv:1910.13413* (2019).

31. Krizhevsky, A., Hinton, G., *et al. Learning multiple layers of features from tiny images* tech. rep. (Citeseer, 2009).

32. Kumar, I. E., Venkatasubramanian, S., Scheidegger, C. & Friedler, S. Problems with Shapley-value-based explanations as feature importance measures. *arXiv preprint arXiv:2002.11097* (2020).

33. Frye, C., de Mijolla, D., Cowton, L., Stanley, M. & Feige, I. Shapley-based explainability on the data manifold. *arXiv preprint arXiv:2006.01272* (2020).

34. Schwab, P. & Karlen, W. *CXPlain: Causal explanations for model interpretation under uncertainty* in *Advances in Neural Information Processing Systems* (2019), 10220–10230.

35. Sixt, L., Granz, M. & Landgraf, T. When Explanations Lie: Why Many Modified BP Attributions Fail. *arXiv,* arXiv–1912 (2019).

36. Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences* **102,** 15545–15550 (2005).

37. Vanni, S. *et al.* Differential overexpression of SERPINA3 in human prion diseases. *Scientific Reports* **7,** 1–13 (2017).

38. Taskesen, E. *et al.* Susceptible genes and disease mechanisms identified in frontotemporal dementia and frontotemporal dementia with Amyotrophic Lateral Sclerosis by DNA-methylation and GWAS. *Scientific Reports* **7,** 1–16 (2017).

39. Mo, C.-h. *et al.* The clinicopathological significance of UBE2C in breast cancer: a study based on immunohistochemistry, microarray and RNA-sequencing data. *Cancer Cell International* **17,** 83 (2017).

40. Yu, J. *et al.* High-throughput metabolomics for discovering potential metabolite biomarkers and metabolic mechanism from the APPswe/PS1dE9 transgenic model of Alzheimer's disease. *Journal of Proteome Research* **16,** 3219–3228 (2017).

41. Atamna, H. & Frey II, W. H. Mechanisms of mitochondrial dysfunction and energy deficiency in Alzheimer's disease. *Mitochondrion* **7,** 297–310 (2007).

42. Alonso-Andres, P., Albasanz, J. L., Ferrer, I. & Martin, M. Purine-related metabolites and their converting enzymes are altered in frontal, parietal and temporal cortex at early stages of Alzheimer's disease pathology. *Brain Pathology* **28,** 933–946 (2018).

43. Ohta, T. & Fukuda, M. Ubiquitin and breast cancer. *Oncogene* **23,** 2079–2088 (2004).

44. Kim, H.-Y. *et al.* Comparative metabolic and lipidomic profiling of human breast cancer cells with different metastatic potentials. *Oncotarget* **7,** 67111 (2016).

45. Schmitz, A. J. Secret consumer scores and segmentations: Separating haves from have-nots. *Mich. St. L. Rev.,* 1411 (2014).

46. Goodman, B. & Flaxman, S. European Union regulations on algorithmic decision-making and a "right to explanation". *AI Magazine* **38,** 50–57 (2017).

47. Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence* **1,** 206–215 (2019).

48. Merrill, J., Ward, G., Kamkar, S., Budzik, J. & Merrill, D. Generalized Integrated Gradients: A practical method for explaining diverse ensembles. *arXiv preprint arXiv:1909.01869* (2019).

49. Chen, H., Janizek, J. D., Lundberg, S. & Lee, S.-I. True to the Model or True to the Data? *arXiv preprint arXiv:2006.16234* (2020).

50. Hooker, S., Erhan, D., Kindermans, P.-J. & Kim, B. *A benchmark for interpretability methods in deep neural networks* in *Advances in Neural Information Processing Systems* (2019), 9737–9748.

51. Greenwood, A. K. *et al.* The AD Knowledge Portal: A Repository for Multi-Omic Data on Alzheimer's Disease and Aging. *Current Protocols in Human Genetics* **108,** e105 (2020).

# A   Appendix

## A.1   Data Sets

### A.1.1   NHANES I

The National Health and Nutrition Examination Survey (NHANES) I [27] is a national longitudinal study conducted on a random sample of individuals from the United States. NHANES I investigates a number of demographics and socioeconomic variables. We utilize the NHANES I Epidemiologic Follow-up Study (NHEFS) which is designed to investigated the relationships between clinical, nutritional, and behavioral factors originally assessed in NHANES I. The NHEFS study comprised a series of follow up studies that trace the cohort (all persons 25-74 years of age who completed a medical examination in NHANES I) and measure additional variables as well as collect death certificates.

### A.1.2   NHANES 1999-2014

The National Health and Nutrition Examination Survey (NHANES) continually collects information on subsamples of the civilian noninstitutionalized US population in two-year cycles. We collected the data from these cycles from 1999-2014 yielding a total of eight release cycles. The surveys collect a variety of laboratory, questionnaire, examination, and demographic data. In particular, the features collected do not match across cycles, so we only utilize variables that are consistently collected across cycles.

### A.1.3   ROSMAP Alzheimer's Gene Expression

Gene expression data collected from the Religious Orders Study (ROS) and Memory and Aging Project (MAP) [23, 24]. ROS is a longitudinal cohort study of aging and Alzheimer's disease run by Rush University enrolling individuals from religious communities for longitudinal clinical analysis and brain donation. MAP is a longitudinal epidemiologic cohort study of common chronic conditions of aging run by Rush University that aims to complement the ROS study by enrolling individuals with wider life experiences and socioeconomic status. Both studies aim to study aging and risk of Alzheimer's disease. We utilize gene expression data collected using ChIP-seq and predict Alzheimer's disease status of the corresponding patients.

### A.1.4   METABRIC Breast Cancer Gene Expression

Gene expression data collected from Molecular Taxonomy of Breast Cancer International Consortium (METABRIC) [25, 26]. METABRIC analyzes genomic and transcriptomic information from a set of 995

breast cancer tumors. Although the original analyses of the transcriptomic information from the tumors are used for a variety of analyses, we purely utilize the transcriptomic information to predict tumor status.

### A.1.5   CIFAR

The CIFAR10 data set consists of $32 \times 32$ color images with 10 possible classes that are a labeled subset of the 80 million tiny images data set [31]. The mutually classes include airplanes, automobiles, birds, cats, deers, dogs, frogs, horses, ships, and trucks. In particular, the images were collected by colleagues at MIT and NYU with natural images collected on a number of search engines.

### A.1.6   MNIST

The MNIST database consists of $28 \times 28$ black and white handwritten digits [28]. The digits are size-normalized and centered in a fixed-size image. There are ten possible classes that correspond to the digits $0, \cdots, 9$.

### A.1.7   HELOC

The Home Equity Line of Credit (HELOC) data set [29] is an anonymized data set of HELOC applications from real homeowners. A HELOC is a line of credit offered by a bank as a percentage of home equity. The outcome is whether the applicant will repay their HELOC account within two years. Financial institutions use predictions of loan repayment to decide whether applicants qualify for a line of credit. The data set was released as part of a FICO xML Challenge (`https://community.fico.com/s/explainable-machine-learning-challenge`) and can be obtained under appropriate agreement to a data set usage license.

## A.2   Experimental Setup

### A.2.1   CIFAR Multiple Baseline

**Model Hyperparameters:** The model explained is a CNN with the following sequence of layers: a convolutional layer with 32 filters of shape 3 by 3 with a ReLU activation, a convolutional layer with 32 filters of shape 3 by 3 with a ReLU activation, a max pooling layer with of size 2 by 2, a dropout layer with 0.25 probability, a convolutional layer with 64 filters of shape 3 by 3 with a ReLU activation, a convolutional layer with 64 filters of shape 3 by 3 with a ReLU activation, a max pooling layer with of size 2 by 2, a dropout layer with 0.25 probability, a dense layer with 512 nodes with ReLU activation, a dropout layer with probability 0.5, a dense output layer with softmax activation. RMSprop with a learning rate of 0.0001 and decay of $1 \times 10^{-6}$ is used to optimize the network for a categorical cross entropy loss over 100 epochs with batch sizes of 32. The test accuracy achieved by the model is 75.56%.

**Experimental setup:** In Figure 2 we explain three explicands with black objects: a plane, a horse, and an ostrich. For the single baseline attributions we utilize DeepSHAP with a single black image as the baseline. For multiple baselines we utilize DeepSHAP with a baseline distribution of 1000 randomly sampled images from the training data set. The feature attribution plots take the local feature attributions for the softmax output corresponding to the true label. For simple visualization we take the absolute value of the attributions and average across channels to get a grayscale image which we plot after normalizing the attribution values between zero and one. The pixel distributions are the number of pixels in the gray scale version of the explicand image that fell within ten equally sized gray scale value bins. The attribution distribution is the sum of the attribution mass for the pixels in the original image that correspond to each gray scale value bin.

### A.2.2   NHANES Multiple Baseline

**Model Hyperparameters:** The model we explain is an MLP with four hidden layers with 100 nodes each. The hidden layers have ReLU activation functions and dropout layers in between. The final output node is a sigmoid activation trained to minimize binary cross entropy loss to optimize mortality classification. RMSprop with a learning rate of 0.001 is used to optimize the network over 50 epochs with batch sizes of 128. The test ROC achieved by the model is .872.

**Experimental setup:** In Figure 3a, the explicand is a randomly chosen older male individual from the NHANES data set. In the top force plot we show DeepSHAP attributions for the explicand with a

baseline distribution of 1000 randomly chosen samples from the training set. In the bottom force plot we show DeepSHAP attributions for the same explicand with a baseline distribution of 1000 randomly chosen samples from older (>60 years old) males from the training set. In Figure 3b, the clusters are obtained by k-means clustering (k=8) the training data with only two features: age and sex. In Figure 3c, the explicands are the older male cluster in the training data (n=1137). We show two summary plots where the top are DeepSHAP attributions with a baseline distribution of 1000 randomly chosen samples from the training set and the bottom uses the older male cluster as a baseline distribution. In Figure 3d, we perform an ablation test that ablates all explicands in the older male clusters according to either their most positive or most negative local feature attributions. When ablating we impute by the mean feature value in the older male cluster. Then we evaluate the model's prediction across all of the explicands after ablating features one at a time.

### A.2.3   Gene set explanations

**Model hyperparameters:** We train two GBToost classifiers (an implementation of gradient boosting trees) to predict our binary phenotypes (Alzheimer's and breast cancer tumor stage) based on transcriptomic data. The classifiers are trained with a learning rate of 0.3, a max tree depth of 6, and automatic heuristic tree construction. We train with a validation set and 10 early stopping rounds. For Alzheimer's classification we achieve a test ROC AUC of 0.959 and for breast cancer tumor stage classification we achieve a test ROC AUC of 0.932.

    **Experimental setup:** The feature attributions for the tree model are obtained using Interventional Tree Explainer [21]. These attributions correspond to the importance of each gene to the log odds of the output phenotypes. In order to explain these attributions in terms of groups we utilize our group rescale rule to propagate the gene attributions to pathway attributions. For Alzheimer's we fix the baseline distribution to be the training data set and for breast cancer which has more samples we fix a baseline distribution of 100 random samples from the training set for breast cancer.

### A.2.4   NHANES Loss explanations

**Model hyperparameters:** We train an GBToost classifier (an implementation of gradient boosting trees) to predict our mortality based on epidemiological features. The classifier is trained with a learning rate of 0.3, a max tree depth of 6, and automatic heuristic tree construction. We train with a validation set and 10 early stopping rounds. For the weight-shifted test set we achieve an ROC AUC of 0.860 and for the non-shifted test set we achieve an ROC AUC of 0.868.

    **Experimental setup:** In Figure 5b, we generate output feature attributions using Interventional Tree Explainer and use DeepSHAP's generalized rescale rule to explain the loss in addition to the output. The loss and output attributions are explained with respect to the same baseline distribution of 1000 random samples from the training set. The loss attributions for positive labelled explicands and negative labelled explicands are very different, leading us to plot them as separate dependence plots. In Figure 5c, we generate output and loss feature attributions as before. For the ablation, we do a simplified univariate ablation where we impute the blood loss to the mean of the baseline distribution for samples selected based on largest loss attributions. In Figure 5d, we perform an ablation test that ablates 1000 explicands from the training set according to either their loss or output local feature attributions. When ablating we impute by the mean value of a given feature in the explicands. Then we evaluate the model's prediction across all of the explicands after ablating features one at a time.

### A.2.5   MNIST Feature Extraction

**Model hyperparameters:** We train a CNN model to classify all digits in MNIST. The CNN model consists of a convolutional layer with 32 filters of size 3 by 3 with ReLU activation, a max pooling layer with pools of size 2 by 2, a convolutional layer with 64 filters of size 3 by 3 with ReLU activation, a max pooling layer with pools of size 2 by 2, a dense layer with 100 nodes and ReLU activation, and the dense output layer with 10 nodes and softmax activation. We utilize categorical cross-entropy loss, an Adam optimizer with learning rate 0.001, and train for 10 epochs. Then, in order to utilize the model to to extract higher level features from raw MNIST images, we remove the final output layer. The GBToost model we train to predict zeros

using the MNIST features has a max tree depth of 5, a learning rate of 0.5, and a binary logistic objective. This model achieves a test accuracy of 0.998 for predicting zeros.

**Experimental setup:** In this experiment, we train a CNN model and use it to extract features that are fed into an GBT model. In Figure 6a, we show the feature attributions for DeepSHAP and three model-agnostic approaches. Each model-agnostic approach uses a number of samples which is set to 100,000. All models utilize the same baseline distribution of 100 random images to explain the five images we selected. In Figure 6b we report the runtime of these feature attribution approaches, and the ablation of the top 10% of features. In order to ablate the top ten positive (or negative) features, we simply select the pixels with the largest positive (or negative) attribution in the five explicands and impute them with the mean pixels across the baseline distribution. We obtain confidence intervals by repeating this 20 times for different randomly selected sets of five explicands, where we enforce that at least one zero occurs within the five explicands, because it is the class of interest.

### A.2.6 HELOC Stacked Generalization

**Model hyperparameters:** In this experiment we train two base-learners. One base learner is a GBT classifier that represents a fraud detection model which utilizes the following features: "MSinceOldestTradeOpen", "MSinceMostRecentTradeOpen", and "NumTradesOpeninLast12M". Although this classifer represents a fraud detection model, we train it to predict risk using a learning rate of 0.1, 100 estimators, and a max tree depth of 3. The other base learner is that represents a credit scoring model which utilizes the following features: "AverageM-InFile", "NumSatisfactoryTrades", "NumTrades60Ever2DerogPubRec", "NumTrades90Ever2DerogPubRec", "PercentTradesNeverDelq", "MSinceMostRecentDelq", "MaxDelq2PublicRecLast12M", "MaxDelqEver", and "NumTotalTrades". We train the base learner to predict risk using an MLP consisting of two hidden layers with 100 nodes and ReLU activations and an output layer consisting of a single dense node with sigmoid activation. The binary cross-entropy loss function is optimized using stochastic gradient descent and a learning rate of 0.005. The meta learner is a GBT classifier that represents a bank risk prediction model which utilizes the remaining HELOC features in addition to the outputs of the two base learners. The meta learner uses the following hyperparameters: learning rate of 0.1, 100 estimators, and a max tree depth of 3.

**Experimental setup:** In Figure 7a, we first train a GBT and MLP base-learner on disjoint subsets of features from the training data. Then we generate scores for the training data and append it to the remaining features. The remaining features and consumer scores are used to train a final GBT model. Finally, we evaluate the final GBT on a held out test data set. In Figure 7b, we create explanations for the meta model using interventional Tree Explainer for the GBT. Then in 7c, we use the generalized rescale rule to propagate the attributions back through the base-learners (GBT and MLP) to obtain attributions in the original feature space.

### A.2.7 NHANES Stacked Generalization

**Model hyperparameters:** In this experiment we train five base-learners - MLPs. The MLPs consist of two hidden layers with 100 nodes and ReLU activations. The output layer is a single dense node with sigmoid activation. The binary cross-entropy loss function is optimized via stochastic gradient descent with a learning rate of 0.005. Then we train a two meta-models that use the outputs of the MLPs as inputs. The first is a logistic regression model with an L2 penalty and regularization strength of 1. The second is a gradient boosted trees classifier with a learning rate of 0.1, 100 estimators, and a max tree depth of 3.

**Experimental setup:** In Figure 13a, we first train five MLP base-learners on training data. Then we embed held out validation data using the predictions of the five MLP base-learners. This embedded validation data is used to train the logistic regression and gradient boosting trees models. Finally, all models are evaluated on a held out test data set. In Figure 13b, we create meta-level explanations using interventional Shapley value attributions for the linear models (average voting and logistic regression) [49], and interventional Tree Explainer for the GBT. Then in 13c, we use the generalized rescale rule to propagate the attributions back through the base-learner MLPs to obtain attributions in the original feature space.

## A.3 Feature attribution plots

In this section we describe a number of plotting techniques for conveying information about local feature attributions. These plots were first introduced in [21].

### A.3.1 Force plots

Force plots show the feature attributions for a single explicand in terms of how they drive the model's prediction for the explicand away from the average model prediction across the baseline distribution. The width of the bars indicate the feature attribution value with red indicating a positive affect and blue indicating a negative one. The features corresponding to the largest bars are below with their actual values for the explicand.

### A.3.2 Dependence plots

Dependence plots show the feature attributions for many explicands for a single feature. Every point corresponds to a single explicand where the x-axis is the value of the feature and the y-axis is the the feature attribution value. The coloring of the points often denotes the value of a separate feature.

### A.3.3 Summary plots

Summary plots show the feature attributions for many explicands and multiple features. Summary plots stack multiple subplots plots for each individual feature. For the feature plots, every point corresponds to a single explicand where the x-axis is the feature attribution value and the y-axis is vertical dispersion representing the frequency of samples with a particular feature attribution value. Finally, the color of each point represents the normalized feature value, with red representing a high value and blue representing a low one. Intermediary feature values are interpolations between red and blue.

## A.4 Results

### A.4.1 Additional CIFAR bias examples

We present additional examples of bias for IME and integrated gradients in Figures 8 and 9.



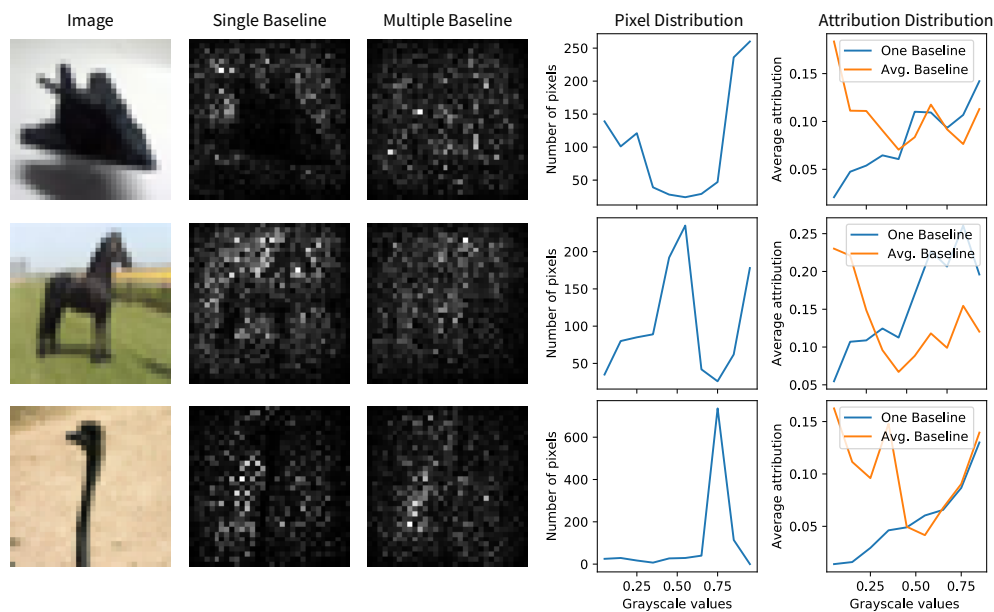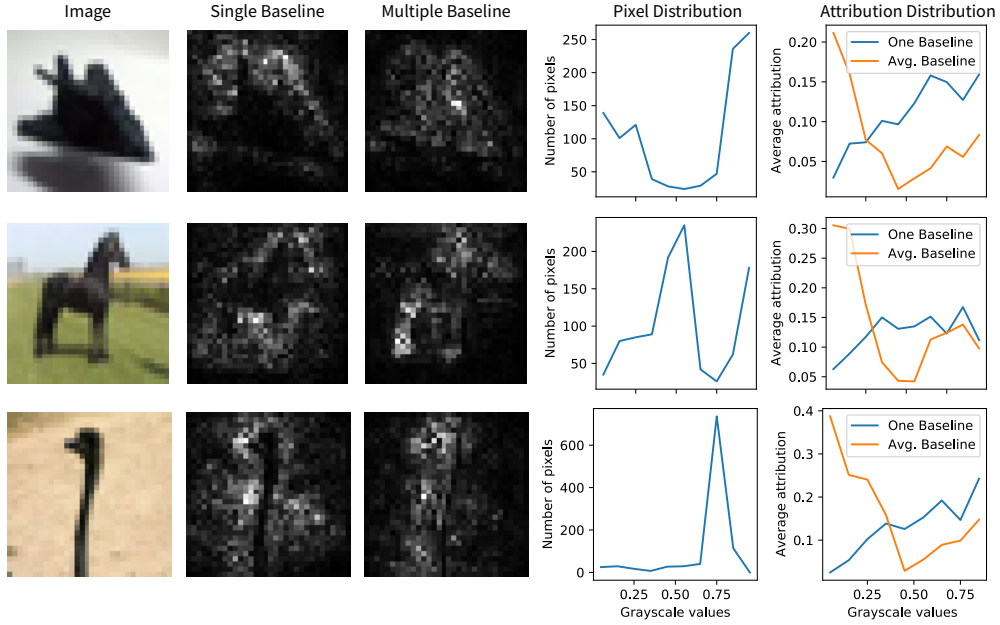Figure 8: Demonstrating bias of a single baseline for IME.

Figure 9: Demonstrating bias of a single baseline for integrated/expected gradients.

## A.4.2 Probability vs. log-odds explanations

In Figure 10 we illustrate the difference between explanations in log-odds versus probability space using attributions obtained from rescaling the log-odds explanations provided by TreeSHAP.
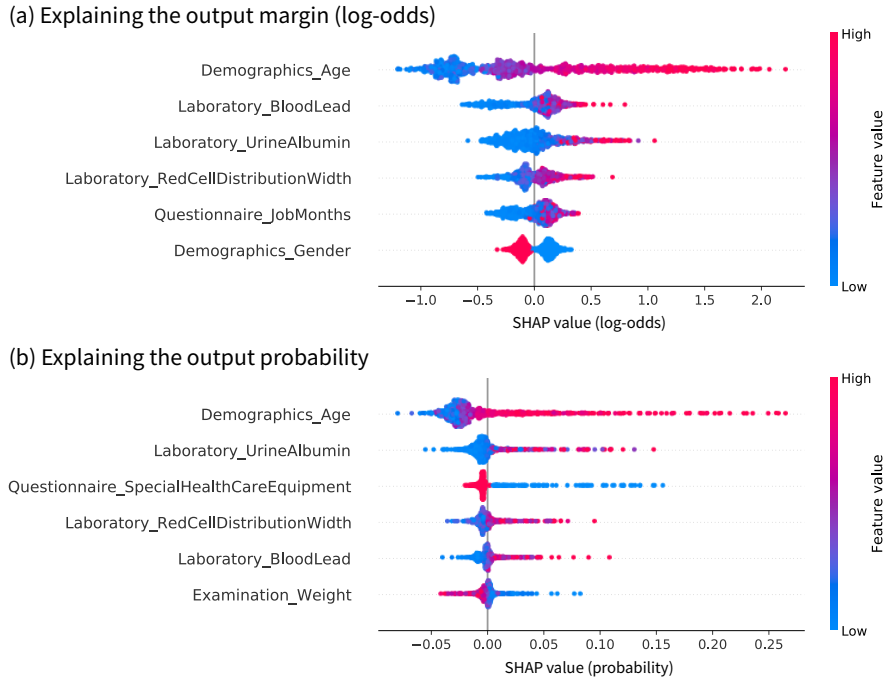


Figure 10: The summary plot for the log-odds model output differs to the summary plot for the probability output in terms of ordering of important features. This is to be expected because of the non-linear mapping between log-odds and probability. Often times, it can be useful to communicate scientific findings in terms of the probability output of the model, although the log-odds output is also natural as it is the output margin.

### A.4.3 Additional gene sets

We present attributions aggregated by the Reactome canonical pathway gene set and the Biological Process gene ontology gene set in Figure 11.

(a) Alzheimer's



(a) Breast cancer tumor stage



Figure 11: Additional gene set attributions for the Reactome canonical pathway gene set and the Biological Process gene ontology gene set. Analogous to the attributions in Figure 4

### A.4.4 Improved predictive performance of feature extraction

In Figure 12 we demonstrate the efficacy of deep feature extraction fed into a tree model for MNIST.

### A.4.5 Stacked generalization

We compare five bagged MLP base-learners (feature attributions in Figures 14-18) and three meta-learners (average voting, logistic regression, and gradient boosting trees) that use the base-learners' predictions as features for NHANES (1999-2014) mortality prediction with performance in Figure 13a. We see that average voting outperforms any individual MLP and is improved upon by a non-uniform weighting scheme (logistic regression). Finally, stacked generalization with a gradient boosted tree meta-model outperforms both linear approaches.

Since our framework enables attributions that satisfy efficiency at each layer, we obtain the importance each meta model assigns to each base-learner (Figure 13b), which is much harder to do for model-agnostic methods because it will require separately estimating the importance for each layer. Although the average voting scheme assigns equal importance to each base model, each MLP's predictions are different, leading to the different shapes in the summary plots. In contrast, the logistic regression model downweights MLP0 and MLP3 and primarily relies on MLP2 and MLP4 which achieved the highest performance. The gradient boosting tree model uses the base-learners in a non-linear fashion. For MLP4, high predictions actually decreases the overall prediction of the meta-learner. These meta-level explanations reveal novel insights that explanations in the original feature space would not. Finally, we can also propagate the meta-level explanations back to the original input space and verify that most models give similarly reasonable feature attributions in Figure 7c.
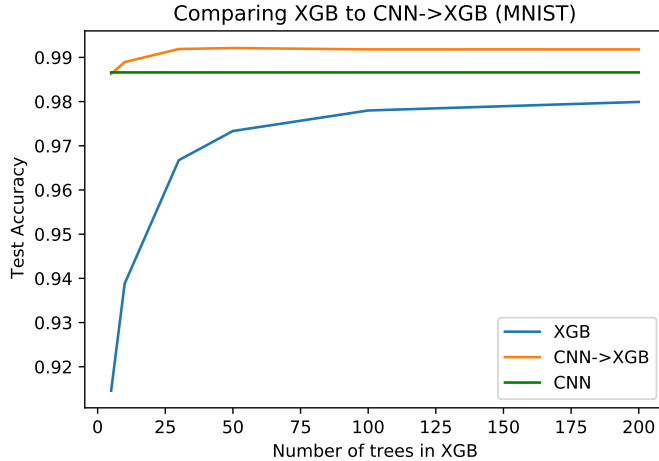
Figure 12: Investigating the utility of CNN feature extraction in MNIST. We compare a CNN on the raw digits to an GBT model trained on the raw digits to an GBT model trained to classify digits on the basis of the features extracted by the trained CNN. We vary the number of estimators in GBToost to investigate how well underparamterized trees classify digits with different features. Overall, using GBT models with the features extracted from the CNN yield much higher accuracy for trees with the same number of estimators.

## A.5   Methods

### A.5.1   Shapley value axioms

The Shapley values satisfy a number of desirable properties in terms of the set function $v$. It is uniquely defined by three axioms:

- **Efficiency:** The sum of the Shapley values for each player equals the value of the game with the set of all players (the grand coalition):

$$\sum_{i=1}^{m} \phi_i(v) = v(M) - v(\emptyset) \tag{37}$$

- **Monotonicity:** If a player $i$ always increases game $v_1$'s value more than they would company $v_2$ for all possible remaining sets of players, then $i$'s attribution for $v_1$ should be greater than or equal to their attribution in $v_2$:

$$v_1(S \cup i) - v_1(S) \geq v_2(S \cup i) - v_2(S) \forall S \subseteq N \setminus i \implies \phi_i(v_1) \geq \phi_2(v_2) \tag{38}$$

- **Missingness:** Employees $i$ that don't help or hurt the company's profit must have no attribution:

$$v(S \cup i) = v(S) \forall S \subseteq N \setminus i \implies \phi_i(v) = 0 \tag{39}$$

While the above three axioms determine the Shapley values as a unique solution concept for credit allocation, the Shapley values have a number of additional desirable properties:

- **Symmetry:** If two players have the same marginal impact for all subsets, then they should have the same Shapley value:

$$v(S \cup i) = v(S \cup j) \forall S \subseteq N \setminus i, j \implies \phi_i(v) = \phi_j(v) \tag{40}$$

- **Linearity:** The Shapley values for a linear combination of games is equal to the linear combination of Shapley values for each game:
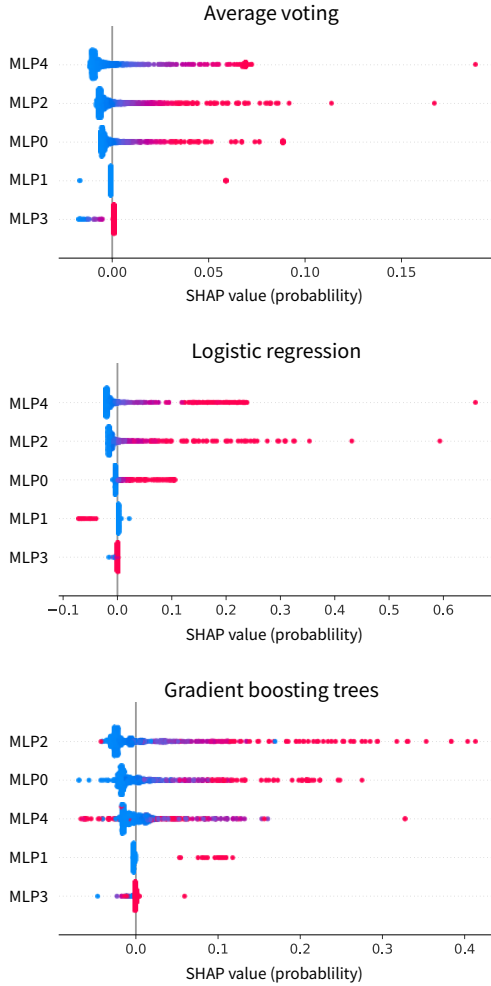
$$\phi_i(v_1 + v_2) = \phi_i(v_1) + \phi_i(v_2) \tag{41}$$

and

$$\phi_i(av) = a\phi_i(v) \tag{42}$$

28

(a) Model performance

| Model | MLP0 | MLP1 | MLP2 | MLP3 | MLP4 | AV | LR | GBT |
|-------|------|------|------|------|------|-----|-----|-----|
| ROC AUC | 0.8207 | 0.5292 | 0.8330 | 0.5042 | 0.8331 | 0.8405 | 0.8425 | **0.8444** |

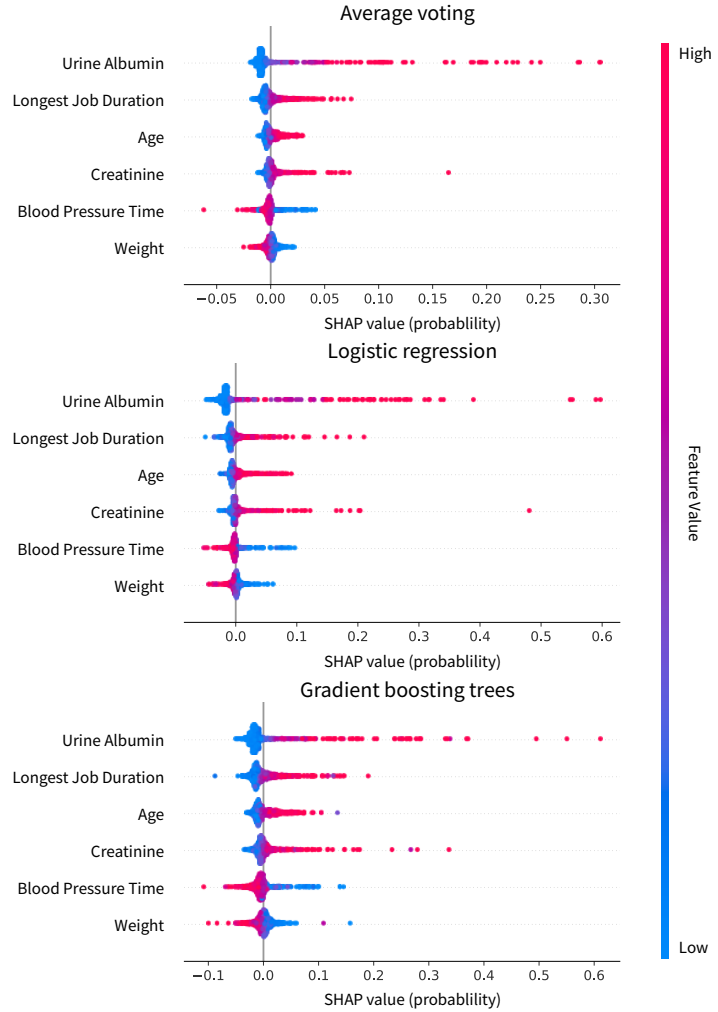(b) Meta-level explanations

(c) Raw feature explanations

Figure 13: **Explaining stacked generalization by looking at meta-level and raw feature explanations.** (a) The test set performance of the five MLP models and three meta-models that use make predictions based on the MLP models' predictions. (b) Intermediary explanations for the meta-models that assign credit based on which MLP was important to the meta-model's prediction. (c) Raw feature explanations obtained by propagating the credit for each meta-model in (b) to the original feature space. (b) and (c) show summary plots (Appendix Section A.3.3).

### A.5.2 Baseline distribution proof for interventional Shapley values

*Proof.* Define $D$ to be the data distribution, $N$ to be the set of all features, and $f$ to be the model being explained. Additionally, define $\mathcal{X}(x, x', S)$ to return a sample where the features in $S$ are taken from $x$ and the remaining features from $x'$. Define $C$ to be all combinations of the set $N \setminus \{i\}$ and $P$ to be all
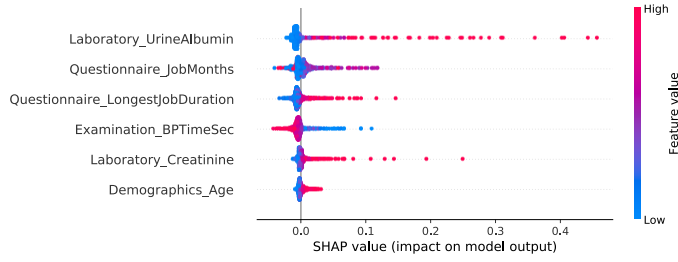
Figure 14: Feature attributions for base learner MLP0.
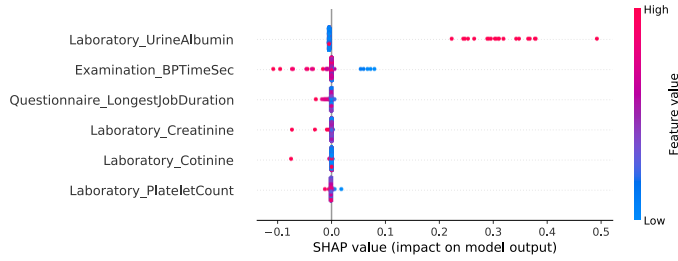


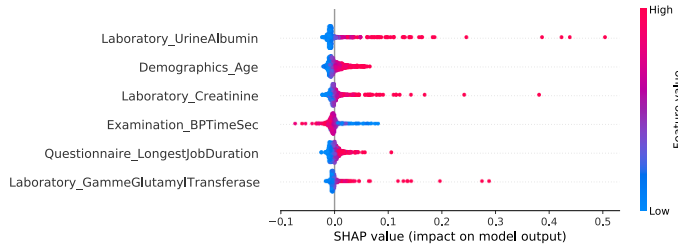Figure 15: Feature attributions for base learner MLP1.



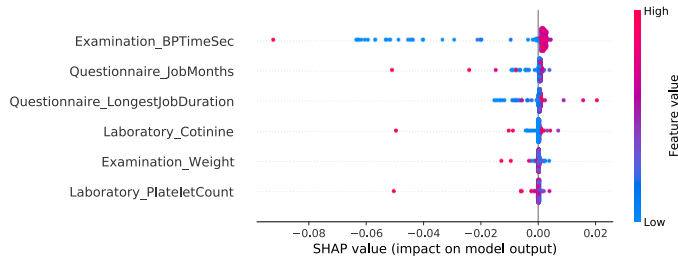Figure 16: Feature attributions for base learner MLP2.



Figure 17: Feature attributions for base learner MLP3.

permutations of $N \setminus \{i\}$. Starting with the definition of SHAP values for a single feature: $\phi_i(x)$

$$= \sum_{S \in C} W(|S|, |N|)(\mathbb{E}_D[f(X)|x_{S \cup \{i\}}] - \mathbb{E}_D[f(X)|x_S])$$

$$= \frac{1}{|P|} \sum_{S \subseteq P} \mathbb{E}_{\mathcal{D}}[f(x)|\mathrm{do}(x_{S \cup \{i\}})] - \mathbb{E}_{\mathcal{D}}[\mathrm{do}(f(x)|x_S)]$$

$$= \frac{1}{|P|} \sum_{S \subseteq P} \frac{1}{|D|} \sum_{x' \in D} f(\mathcal{X}(x, x', S \cup \{i\})) - f(\mathcal{X}(x, x', S))$$

$$= \frac{1}{|D|} \sum_{x' \in D} \underbrace{\frac{1}{|P|} \sum_{S \subseteq P} f(\mathcal{X}(x, x', S \cup \{i\})) - f(\mathcal{X}(x, x', S))}_{\text{single baseline SHAP value}}$$
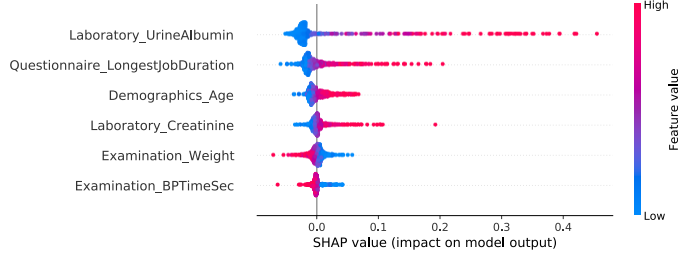
30

Figure 18: Feature attributions for base learner MLP4.

where the second step depends on an interventional conditional expectation [30] which is very close to Random Baseline Shapley in [14]). $\qquad\square$

### A.5.3 Generalized rescale rule is exact for linear models

We define a series of models composed of linear functions: $f_k(x) = B^k \cdots B^2 B^1 x$ where $B^i \in \mathbb{R}^{o_i \times m_i}$, $m_1 = m$, and $o_k = 1$. If we define $\hat{\phi}$ to return Interventional Shapley values for linear models ($\phi(f, x^e, x^b) = \beta(x^e - x^b)$ where $x^e$ and $x^b$ are the inputs to the linear model and $f(x) = \beta x$ [49]). Then, the generalized rescale rule gives:

$$\psi^k = B^k(f_{k-1}(x^e) - f_{k-1}(x^b)) \tag{43}$$

$$\psi^i = B^i(f_{i-1}(x^e) - f_{i-1}(x^b))\big(\psi^{i+1} \oslash (f_i(x^e) - f_i(x^b))\big), \ i \in 1, \cdots, k-1 \tag{44}$$

Therefore,

$$\phi_i(f_k, x^e, x^b) = B^k \cdots B^2 B^1(x^e - x^b) \tag{45}$$

This coincides with the interventional Shapley values for $f_k(x)$ since the composition of linear models is linear.

### A.5.4 Explaining ensembles of models

Explaining ensembles of models is straightforward for Shapley values, because of the linearity property (Appendix Section A.5.1). In particular, bagged and boosted ensembles are linear functions of individual models:

$$f(x) = \beta_1 f_1(x) + \cdots \beta_k f_k(x), \tag{46}$$

where bagged ensembles have $\beta_i = \frac{1}{k}$ and boosted ensembles have $\beta_i = 1$. In order to explain these models, it suffices to explain the ensemble model $f$ with:

$$\phi_i(f) = \beta_1 \phi_i(f_1) + \cdots \beta_k \phi_k(f_k) \tag{47}$$

Furthermore, explaining linear meta-models for stacked ensembles is also encompassed by the linearity property of Shapley values. In contrast, in order to explain stacked ensembles with mixed model types as in Section 4.4, we employ our generalized rescale rule.