
“Why did the Model Fail?”: Attributing Model Performance Changes to Distribution Shifts

Haoran Zhang^{*1} Harvineet Singh^{*2} Shalmali Joshi³

Abstract

Performance of machine learning models may differ significantly in novel environments compared to during training due to shifts in the underlying data distribution. Attributing performance changes to specific data shifts is critical for identifying sources of model failures and designing stable models. In this work, we design a novel method for attributing performance differences between environments to shifts in the underlying causal mechanisms. To this end, we construct a cooperative game where the contribution of each mechanism is quantified as their Shapley value. We demonstrate the ability of our method to identify sources of spurious correlation and attribute performance drop to shifts in label and/or feature distributions on synthetic and real-world datasets.

1. Introduction

Machine learning models are widely deployed in dynamic environments ranging from recommendation systems to personalized clinical care. Such environments are prone to dataset shifts, which may lead to serious degradations in model performance (Guo et al., 2022; Chirra et al., 2018; Koh et al., 2021; Geirhos et al., 2020; Nestor et al., 2019). Importantly, the shifts are hard to anticipate and reduce the ability of model developers to design reliable systems.

When the performance of a model *does* degrade during deployment, it is crucial for the model developer to know *how* the dataset has shifted to cause this change. Cognizant of this information, the model developer can then take mitigating actions such as additional data collection, data augmentation, and model retraining (Ashmore et al., 2021; Zenke et al., 2017; Subbaswamy et al., 2019).

^{*}Equal contribution ¹Massachusetts Institute of Technology ²New York University ³Harvard University. Correspondence to: Haoran Zhang <haoranz@mit.edu>.

In this work, we present a method to attribute model performance changes to shifts in a given set of distributions. Dataset shifts can occur in various marginal or conditional distributions comprising variables involved in the model. Given that many distributions can change simultaneously, we define the effect of changing any set of distributions and use the concept of Shapley values (Roth, 1988) to attribute the change to individual distributions.

We build on a recent line of work that defines distribution shifts as interventions on causal mechanisms (Pearl & Bareinboim, 2011; Subbaswamy et al., 2019; 2021; Budhathoki et al., 2021; Thams et al., 2022). Most relevant is Budhathoki et al. (2021) which attributes a shift between two joint distributions to the causal mechanisms denoted by a sub-distributions (i.e. factorization of the joint distribution). We differ by attributing a change in *model performance* to sub-distributions. Note that each shifted sub-distributions may have different influences on model performance. In this work, we demonstrate that explaining performance discrepancy requires us to develop specialized methods.

We make the following contributions:

- We formalize the problem of attributing model performance changes due to distribution shifts.
- We propose a principled approach based on Shapley values for the attribution, and show that it satisfies several desirable properties.
- We validate the correctness of our method on synthetic and real-world datasets.

2. Related work

There has been extensive work that tests whether the data distribution has shifted (e.g. ones evaluated in Rabanser et al. (2019)). Past work has proposed to identify sub-distributions (factors constituting the joint distribution as determined by a generative model for the data) that comprise the shift between two joint distributions and order them by their contribution to the shift (Budhathoki et al., 2021). The method proposed in Budhathoki et al. (2021) does not attribute the contribution of each sub-distribution to model performance. Even a small change in some factors may have a large effect on model performance. Conversely, many shifts in

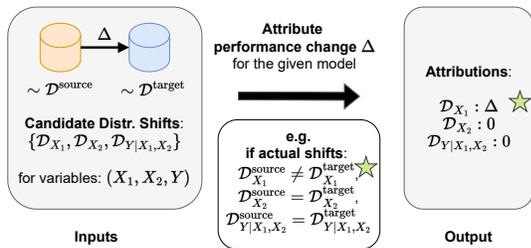


Figure 1. Inputs and outputs for attribution. The goal is to attribute the model’s performance change Δ between source and target distributions. In this example, out of the three candidate distributions only the marginal distribution of X_1 changes. Thus, the method attributes all of the performance change to X_1 .

sub-distributions may lead to no model performance degradation at all. Thereby, the attribution in Budhathoki et al. (2021) may yield multiple shifting sub-distributions that the model developer has to filter out to identify the relevant ones (see Property 2.2 below).

Shapley value-based attribution has recently become popular for interpreting model predictions (Štrumbelj & Kononenko, 2014; Lundberg & Lee, 2017; Wang et al., 2021). However, in most prior work, Shapley values have been leveraged for attributing changes in specific model predictions to variables (Sundararajan & Najmi, 2020). Challenges to appropriately interpreting such attributions and desirable properties thereof have been extensively discussed in Janzing et al. (2020); Kumar et al. (2021). Recently, Wu et al. (2021) decompose performance change to changes in marginal distributions using Shapley value framework. In contrast, we advance the use of Shapley values for interpreting model performance changes, particularly by attributing discrepancy to any sub-distribution.

Finally, recent work aims to find subsets of the dataset that have significantly worse (or better) performance (d’Eon et al., 2021; Eyuboglu et al., 2022). The main difference in our work is the data representations used for attribution. These works chose to identify subsets of data that are relevant to performance change whereas we find sub-distributions represented by causal mechanisms.

3. Preliminaries

Notation Consider a learning setup where we have some system variables denoted by V . We are given a fixed model f and a loss function $\ell : v, f \mapsto \ell(v, f) \in \mathbb{R}$ which assigns a real value to the model evaluated at a specific setting of the variables. For example, in the case of supervised learning, $V = (X, Y)$ comprises the features and labels, the model f maps features into the label space, and a loss function such as the squared error $\ell((x, y), f) := (y - f(x))^2$ is used to evaluate the model. Assume that the loss function is

computed separately for each data point. Then, performance of the model is summarized by the average of the losses,

$$\text{Perf}(\mathcal{D}) := \mathbb{E}_{v \sim \mathcal{D}}[\ell(v, f)]$$

on data from a given distribution \mathcal{D} . We use $\mathcal{D}^{\text{source}}$ to denote the distribution of data for the source domain and $\mathcal{D}^{\text{target}}$ for the target domain. Subscripts on \mathcal{D} refer to the distribution of specific variables. For example, \mathcal{D}_X is the distribution of features, and $\mathcal{D}_{Y|X}$ is the conditional distribution of labels given features.

Sources of performance change Model performance can change between development and deployment of models for many reasons (Jacobs & Wallach, 2021). We restrict our attention to a narrow yet important source of performance change which is the shift in the distribution of input features or labels. Since we have combinatorially many shifts that can be defined on subsets of $V = (X, Y)$, we leverage the knowledge of causal mechanisms in the form of a causal graph to identify potential shifts to consider (Pearl, 2009).

3.1. Problem setup

Suppose we are given a *candidate set* of distributions $\mathcal{C}_{\mathcal{D}}$ that may account for the shift from source $\mathcal{D}^{\text{source}}$ to target $\mathcal{D}^{\text{target}}$ distributions: $\text{Perf}(\mathcal{D}^{\text{target}}) - \text{Perf}(\mathcal{D}^{\text{source}})$. **Our goal is to attribute this change to each candidate distribution in the candidate set $\mathcal{C}_{\mathcal{D}}$.** For our method, we assume access to the model f , and samples from $\mathcal{D}^{\text{source}}$ as well as $\mathcal{D}^{\text{target}}$ (see Figure 1).

3.2. Choice of candidate set and causal mechanisms

Candidate distributions can be defined in multiple ways. For instance, it can be the set of marginal distributions on each system variable, $\mathcal{C}_{\mathcal{D}} = \{\mathcal{D}_{X_1}, \mathcal{D}_{X_2}, \dots\}$, or distribution of each variable after conditioning on the rest, $\mathcal{C}_{\mathcal{D}} = \{\mathcal{D}_{X_1|V \setminus X_1}, \mathcal{D}_{X_2|V \setminus X_2}, \dots\}$.

Motivated by Budhathoki et al. (2021), we propose to use the causal mechanisms constituting the structural causal graph for the system as our candidate set. A causal graph specifies a particular factorization of the joint distribution into a set of distributions (alternatively called causal mechanisms) (Pearl, 2009). That is $\mathcal{D}_V = \prod_{X_i \in V} \mathcal{D}_{X_i|\text{parent}(X_i)}$ where $\text{parent}(X_i)$ are the variables that have a directed edge to X_i in the causal graph. These distributions are assumed to be *independent*, i.e. an intervention to change one of the distributions does not change any other distribution in the factorization. We also assume that the causal graph is *sufficient* (Spirtes et al., 2000). Thus, the candidate set is

$$\mathcal{C}_{\mathcal{D}} = \{\mathcal{D}_{X_1|\text{parent}(X_1)}, \dots, \mathcal{D}_{X_i|\text{parent}(X_i)}, \dots\}_{i=1, \dots, |V|}$$

Advantages of using causal mechanisms. This choice of candidate set has three main advantages. First, it is *inter-*

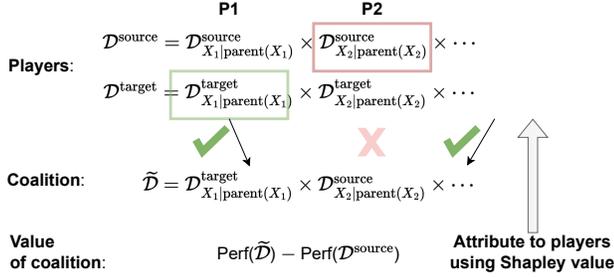


Figure 2. **Sketch of the game theoretic attribution method.** Each causal mechanism is a player that, if present in the coalition, changes to the target distribution and, if absent, remains fixed at the source distribution. This defines the distribution of the resulting coalition $\tilde{\mathcal{D}}$. Performance on $\tilde{\mathcal{D}}$ is estimated using importance sampling from training data samples. After computing values for each possible coalition, Shapley value (Eq. 2) gives the attribution to each player. Thus, we estimate the performance change under all possible ways to shift the mechanisms from source to target and use these to distribute the total performance change among the individual mechanisms.

pretable since the candidate shifts are specified by domain experts who constructed the causal graph. Second, it is *actionable* since identifying the causal mechanisms most responsible for performance change can inform training methods for handling distribution shifts (Subbaswamy et al., 2019). Third, it will lead to *succinct* attributions due to the independence property. Consider the case where only one conditional distribution $\mathcal{D}(X_i|\text{parent}(X_i))$ changes across domains. This will result in a change in distributions of all descendants of X_i (due to the factorization given above). In this case, a candidate set defined by all marginals is not succinct, as one would attribute performance changes to all marginals of descendants of X_i . Instead, focusing on our candidate set determined by the causal mechanism will isolate the appropriate conditional distribution.

4. Method

We motivate a game theoretic formulation for attributing performance changes to distributions over variable subsets.

4.1. Game theoretic attribution

Consider the following attribution game. The set of *players* in this game are the candidate distributions. A *coalition* of any subset of players determines the distributions that are allowed to shift, keeping the rest fixed. The *value* for the coalition is the model performance change between the resulting distribution for the coalition and the train distribution. See Figure 2 for an overview of the method.

Consider a coalition that consists of two factors in the candi-

date set, that is $\{\mathcal{D}_{X_1|\text{parent}(X_1)}, \mathcal{D}_{X_2|\text{parent}(X_2)}\} = \tilde{\mathcal{C}} \subset \mathcal{C}_{\mathcal{D}}$.

The resulting distribution for the coalition $\tilde{\mathcal{C}}$ is

$$\tilde{\mathcal{D}} = \mathcal{D}_{X_1|\text{parent}(X_1)}^{\text{target}} \mathcal{D}_{X_2|\text{parent}(X_2)}^{\text{target}} \prod_{i \in \mathcal{V} \setminus X_1, X_2} \mathcal{D}_{X_i|\text{parent}(X_i)}^{\text{source}}$$

The value of the coalition $\tilde{\mathcal{C}}$ with distribution $\tilde{\mathcal{D}}$ is given by

$$\text{Val}(\tilde{\mathcal{C}}) := \text{Perf}(\tilde{\mathcal{D}}) - \text{Perf}(\mathcal{D}^{\text{source}}) \quad (1)$$

Shapley value for this game gives the attribution of each player $d \in \mathcal{C}_{\mathcal{D}}$ as¹

$$\text{Attr}(d) = \sum_{\tilde{\mathcal{C}} \subset \mathcal{C}_{\mathcal{D}} \setminus \{d\}} \binom{|\mathcal{C}_{\mathcal{D}}| - 1}{|\tilde{\mathcal{C}}|} (\text{Val}(\tilde{\mathcal{C}} \cup \{d\}) - \text{Val}(\tilde{\mathcal{C}})) \quad (2)$$

Thus, to compute our attributions, we need estimates of model performance under $\tilde{\mathcal{D}}$. Note that we only have model performance estimates under $\mathcal{D}^{\text{source}}$ and $\mathcal{D}^{\text{target}}$ but not for any arbitrary coalition where a subset of the distributions have shifted. To compute this, we propose to use importance sampling.

4.2. Estimating attribution using importance sampling

Importance sampling allows us to re-weight the samples drawn from a given distribution, which can be $\mathcal{D}^{\text{source}}$ or $\mathcal{D}^{\text{target}}$, to simulate expectations for a desired distribution, which is the candidate $\tilde{\mathcal{D}}$ in our case. Thus, we re-write the value as

$$\begin{aligned} \text{Val}(\tilde{\mathcal{C}}) &= \text{Perf}(\tilde{\mathcal{D}}) - \text{Perf}(\mathcal{D}^{\text{source}}) \\ &= \mathbb{E}_{v \sim \tilde{\mathcal{D}}}[\ell(v, f)] - \mathbb{E}_{v \sim \mathcal{D}^{\text{source}}}[\ell(v, f)] \\ &= \mathbb{E}_{v \sim \mathcal{D}^{\text{source}}} \left[\frac{\tilde{\mathcal{D}}(v)}{\mathcal{D}^{\text{source}}(v)} \ell(v, f) \right] - \mathbb{E}_{v \sim \mathcal{D}^{\text{source}}}[\ell(v, f)] \end{aligned} \quad (3)$$

The importance weights are themselves a product of ratios of source and target distributions corresponding to the causal mechanisms in $\mathcal{C}_{\mathcal{D}}$ as follows:

$$w_{\tilde{\mathcal{C}}}(v) := \frac{\tilde{\mathcal{D}}(v)}{\mathcal{D}^{\text{source}}(v)} = \prod_{d \in \tilde{\mathcal{C}}} \frac{\mathcal{D}_d^{\text{target}}(v)}{\mathcal{D}_d^{\text{source}}(v)} =: \prod_{d \in \tilde{\mathcal{C}}} w_d(v) \quad (4)$$

There are multiple ways to estimate the importance weights $w_d(v)$, which are ratio of densities, in the literature (Sugiyama et al., 2012).

Computing Importance Weights Here, we use a simple approach for density ratio estimation via probabilistic classifiers as described in Sugiyama et al. (2012, Section 2.2), based on training probabilistic classifiers.

¹Here, we use exact Shapley value computation, though approximations can be made for larger candidate sets (Castro et al., 2009; Lundberg & Lee, 2017) for reduced computational effort.

Let D be a binary random variable, such that when $D = 1$, $X \sim \mathcal{D}^{\text{target}}(X)$, and when $D = 0$, $X \sim \mathcal{D}^{\text{source}}(X)$. Suppose $d = \mathcal{D}_{X_i|\text{parent}(X_i)}$, then

$$w_d = \frac{\mathbb{P}(D = 1|X_i)}{\mathbb{P}(D = 0|X_i)} \cdot \frac{\mathbb{P}(D = 1|\text{parent}(X_i))}{\mathbb{P}(D = 0|\text{parent}(X_i))}$$

Where each term is computed using a probabilistic classifier trained to discriminate data points from $\mathcal{D}^{\text{source}}$ and $\mathcal{D}^{\text{target}}$ from the concatenated dataset. In total, we need to learn $\mathcal{O}(|\mathcal{C}_{\mathcal{D}}|)$ models for computing all importance weights.

4.3. Properties of Our Method

Under perfect computation of importance weights, the Shapley values resulting from the performance-change game have the following desirable properties, which follow directly from properties of Shapley values (Winter, 2002).

Property 1. (Efficiency) $\sum_{d \in \mathcal{C}_{\mathcal{D}}} \text{Attr}(d) = \text{Val}(\mathcal{C}_{\mathcal{D}}) = \text{Perf}(\mathcal{D}^{\text{target}}) - \text{Perf}(\mathcal{D}^{\text{source}})$

By definition, we know that the sum of Shapley values equal the value of the all-player coalition. Thus, we distribute the total performance change due to the shift from source to target distribution to the shifts in causal mechanisms in the candidate set.

Property 2.1. (Null Player) $\mathcal{D}_d^{\text{source}} = \mathcal{D}_d^{\text{target}} \implies \text{Attr}(d) = 0$.

Property 2.2. (Relevance) Consider a mechanism d . If $\text{Perf}(\tilde{\mathcal{C}} \cup \{d\}) = \text{Perf}(\tilde{\mathcal{C}})$ for all $\tilde{\mathcal{C}} \subseteq \mathcal{C}_{\mathcal{D}} \setminus d$, then $\text{Attr}(d) = 0$.

We can verify that our method gives zero attribution to distributions that do not shift between the source and target, and distribution shifts which do not impact model performance. First, we observe that in both cases, $\text{Val}(\tilde{\mathcal{D}}) = \text{Val}(\tilde{\mathcal{D}} \cup \{d\})$. For Property 2.1, this is because $\tilde{\mathcal{D}} = \tilde{\mathcal{D}} \cup \{d\}$ for any $\tilde{\mathcal{D}} \subseteq \mathcal{C}_{\mathcal{D}}$ since the factor corresponding to d remains the same between source and target even when it is allowed to change as part of the coalition. For Property 2.2, this is clear from Eq. 3. By definition of Shapley value in Eq. 2, $\text{Attr}(d) = 0$.

Thus, the method attributes the overall performance change only to distributions that actually change in a way that affects the specified performance metric. The contribution of each distribution is computed by considering how much they impact the performance if they are made to change in different combinations alongside the other distributions.

5. Experiments



Figure 5. Causal graphs for (a) Synthetic and (b) Real-world data.

5.1. Synthetic datasets

We generate a synthetic binary classification dataset with three features according to the following data generating process, corresponding to the causal graph shown in Figure 5a. Here, $\xi_p : \{0, 1\} \rightarrow \{0, 1\}$ is a function that randomly flips the input with probability p .

$$G \sim \text{Ber}(0.5), \quad Y = \xi_q(G), \quad \tilde{Y} = \xi_{0.25}(Y) \\ X_1 = \mathcal{N}(\omega\tilde{Y}, 1), \quad X_2 = \mathcal{N}(\tilde{Y}, 1), \quad X_3 = \mathcal{N}(\tilde{Y} + \mu G, 1)$$

Where q, ω and μ are parameters of the data generating process. In the source environment, we set $q = 0.9, \omega = 1$ and $\mu = 5$. We generate 20,000 samples using these parameters, and train logistic regression (LR) and XGBoost (XGB, (Chen & Guestrin, 2016)) models on (X_1, X_2, X_3) to predict Y , using 3-fold cross-validation for model selection. We explore three data settings for the target environment:

- Conditional Label Shift: Vary $q \in [0, 1]$. Keep ω and μ at their source values. Only $P(Y|G)$ changes.
- Conditional Covariate Shift: Vary $\mu \in [0, 5]$. Keep q and ω at their source values. Only $P(X_3|G, Y)$ changes across domains.
- Combined Shift: Set $\omega = 0$. Vary $q \in [0, 1]$. Keep ω at its source value. Both $P(X_1|Y)$ and $P(Y|G)$ change across domains, but their specific contribution to model performance degradation is not known exactly.

We use our method to explain performance changes in AUROC and Brier score for each model on target environments generated within each setting (with $n = 20,000$), computing density ratios using XGB models. Note that the causal graph shown in Figure 5a implies five potential distributional shifts: $\mathcal{C}_{\mathcal{D}} = \{\mathcal{D}_G, \mathcal{D}_{Y|G}, \mathcal{D}_{X_1|Y}, \mathcal{D}_{X_2|Y}, \mathcal{D}_{X_3|G,Y}\}$.

5.2. Real-world datasets

For some real-world datasets, the causal graph in Figure 5b can model the spurious correlation between X and Y (here, due to the common cause G). We test our method on two such binary classification datasets where $G \in \{0, 1\}$.

- Waterbirds** (Wah et al., 2011). G is the background (water or land), Y is the type of bird (waterbird or landbird).

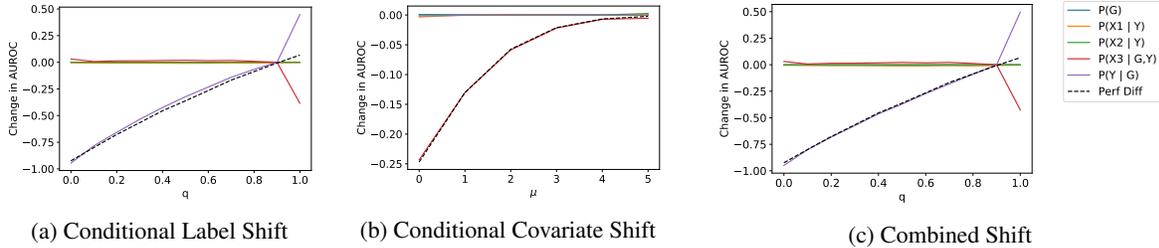


Figure 3. AUROC differences attributed by our model to five potential distributional shifts on the synthetic dataset for the LR model. We observe that the overall change (Perf Diff) is attributed to the true shift in all of the three cases. Rest of the shifts have zero attributed value.

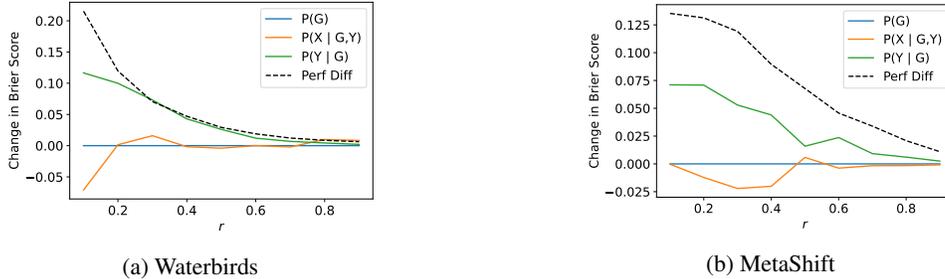


Figure 4. Brier score differences attributed by our model to three potential distribution shifts on two real-world datasets for the LR model.

- **MetaShift** (Liang & Zou, 2022). G is the background (indoor or outdoor), and Y is the type of animal.

In all cases, X is a vector representing the image. Here, we use static embeddings computed from a Imagenet pretrained ResNet-18 (He et al., 2016), and so $X \in \mathbb{R}^{512}$.

For the source environment, we split each dataset into a 75% train, 25% test set. We then use random upsampling on each set to ensure that $\mathbb{P}(G = 0) = \mathbb{P}(G = 1) = 0.5$. We train LR and XGB models to predict Y from X using 3-fold cross validation. In the target environments, we use the same train/test splits as the source environment, and vary a parameter $r \in (0, 1]$. For each target environment, we randomly downsample the two majority groups (across $G \times Y$) to fraction r , and randomly upsample the two minority groups to size $1/r$. Note that for $r = 1$, the source and target environments are identical. We again use random upsampling after this process to ensure that $\mathbb{P}(G)$ does not change. We use our method to explain the performance difference between the source test set and target test set, using AUROC and Brier score as metrics.

6. Results

In Figure 3, we show the output of our method across the three settings described in Section 5.1, with LR as the model of interest and AUROC as the metric. We show similar results for XGB and Brier score in Appendix A. We find that our method attributes all of the performance changes to the correct ground truth shifts, both when there is a sin-

gle shift (Settings 1 and 2) and when there are multiple shifts (Setting 3). In the case of Setting 3, we find that our method attributes all of the performance drop to a shift in $P(Y|G)$, which makes sense as the model relies only on the spurious information (G inferred from X_3) in the source environment.

In Figure 4, we show the results for the real world datasets, with LR as the model of interest and Brier score as the metric. We show similar results for XGB and AUROC in Appendix B. We find that our model largely attributes performance drops to the correct ($P(Y|G)$) distribution. However, some shift is still attributed to the $P(X|G, Y)$ distribution, which does not change. This is as a result of inaccuracies in the estimation of importance weights (e.g. due to overfitting or miscalibration), which is an area of future improvement for our method.

7. Conclusions

We propose a method to attribute changes in performance of a model deployed on a different distribution. Attribution is naturally understood as an interventional concept. We combine the notion of interventions on causal mechanisms with the Shapley value framework to provide a useful decomposition of performance changes. Improvements to the method include relaxing the assumption of sufficiency, improving the importance weighting procedure, and extending the experiments to additional settings such as unsupervised learning and reinforcement learning.

References

- Ashmore, R., Calinescu, R., and Paterson, C. Assuring the machine learning lifecycle: Desiderata, methods, and challenges. *ACM Computing Surveys (CSUR)*, 54(5): 1–39, 2021.
- Budhathoki, K., Janzing, D., Bloebaum, P., and Ng, H. Why did the distribution change? In Banerjee, A. and Fukumizu, K. (eds.), *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pp. 1666–1674. PMLR, 13–15 Apr 2021. URL <https://proceedings.mlr.press/v130/budhathoki21a.html>.
- Castro, J., Gómez, D., and Tejada, J. Polynomial calculation of the shapley value based on sampling. *Computers & Operations Research*, 36(5):1726–1730, 2009.
- Chen, T. and Guestrin, C. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794, 2016.
- Chirra, P., Leo, P., Yim, M., Bloch, B. N., Rastinehad, A. R., Purysko, A., Rosen, M., Madabhushi, A., and Viswanath, S. Empirical evaluation of cross-site reproducibility in radiomic features for characterizing prostate mri. In *Medical Imaging 2018: Computer-Aided Diagnosis*, volume 10575, pp. 105750B. International Society for Optics and Photonics, 2018.
- d’Eon, G., d’Eon, J., Wright, J. R., and Leyton-Brown, K. The spotlight: A general method for discovering systematic errors in deep learning models, 2021. URL <https://arxiv.org/abs/2107.00758>.
- Eyuboglu, S., Varma, M., Saab, K. K., Delbrouck, J.-B., Lee-Messer, C., Dunnmon, J., Zou, J., and Re, C. Domino: Discovering systematic errors with cross-modal embeddings. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=FPCMqjI0jXN>.
- Geirhos, R., Jacobsen, J.-H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., and Wichmann, F. A. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- Guo, L. L., Pfohl, S. R., Fries, J., Johnson, A. E., Posada, J., Aftandilian, C., Shah, N., and Sung, L. Evaluation of domain generalization and adaptation on improving model robustness to temporal dataset shift in clinical medicine. *Scientific reports*, 12(1):1–10, 2022.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Jacobs, A. Z. and Wallach, H. Measurement and fairness. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pp. 375–385, 2021.
- Janzing, D., Minorics, L., and Blöbaum, P. Feature relevance quantification in explainable ai: A causal problem. In *International Conference on artificial intelligence and statistics*, pp. 2907–2916. PMLR, 2020.
- Koh, P. W., Sagawa, S., Marklund, H., Xie, S. M., Zhang, M., Balsubramani, A., Hu, W., Yasunaga, M., Phillips, R. L., Gao, I., et al. Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*, pp. 5637–5664. PMLR, 2021.
- Kumar, I., Scheidegger, C., Venkatasubramanian, S., and Friedler, S. Shapley residuals: Quantifying the limits of the shapley value for explanations. *Advances in Neural Information Processing Systems*, 34, 2021.
- Liang, W. and Zou, J. Metashift: A dataset of datasets for evaluating contextual distribution shifts and training conflicts. *arXiv preprint arXiv:2202.06523*, 2022.
- Lundberg, S. M. and Lee, S.-I. A unified approach to interpreting model predictions. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf>.
- Nestor, B., McDermott, M. B., Boag, W., Berner, G., Naumann, T., Hughes, M. C., Goldenberg, A., and Ghassemi, M. Feature robustness in non-stationary health records: caveats to deployable model performance in common clinical machine learning tasks. In *Machine Learning for Healthcare Conference*, pp. 381–405. PMLR, 2019.
- Pearl, J. *Causality*. Cambridge university press, 2009.
- Pearl, J. and Bareinboim, E. Transportability of causal and statistical relations: A formal approach. In *Twenty-Fifth AAAI Conference on Artificial Intelligence*, 2011.
- Rabanser, S., Gunnemann, S., and Lipton, Z. Failing loudly: An empirical study of methods for detecting dataset shift. In Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/>

- 846c260d715e5b854ffad5f70a516c88-Paper.pdf.
- Roth, A. E. *The Shapley value: essays in honor of Lloyd S. Shapley*. Cambridge University Press, 1988.
- Spirtes, P., Glymour, C. N., Scheines, R., and Heckerman, D. *Causation, prediction, and search*. MIT press, 2000.
- Subbaswamy, A., Schulam, P., and Saria, S. Preventing failures due to dataset shift: Learning predictive models that transport. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 3118–3127, 2019.
- Subbaswamy, A., Adams, R., and Saria, S. Evaluating model robustness and stability to dataset shift. In Banerjee, A. and Fukumizu, K. (eds.), *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pp. 2611–2619. PMLR, 13–15 Apr 2021. URL <http://proceedings.mlr.press/v130/subbaswamy21a.html>.
- Sugiyama, M., Suzuki, T., and Kanamori, T. Density-ratio matching under the bregman divergence: a unified framework of density-ratio estimation. *Annals of the Institute of Statistical Mathematics*, 64(5):1009–1044, 2012.
- Sundararajan, M. and Najmi, A. The many shapley values for model explanation. In *International conference on machine learning*, pp. 9269–9278. PMLR, 2020.
- Thams, N., Oberst, M., and Sontag, D. Evaluating robustness to dataset shift via parametric robustness sets. *arXiv preprint arXiv:2205.15947*, 2022.
- Štrumbelj, E. and Kononenko, I. Explaining prediction models and individual predictions with feature contributions. *Knowl. Inf. Syst.*, 41(3):647–665, dec 2014. ISSN 0219-1377. doi: 10.1007/s10115-013-0679-x. URL <https://doi.org/10.1007/s10115-013-0679-x>.
- Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. The caltech-ucsd birds-200-2011 dataset. 2011.
- Wang, J., Wiens, J., and Lundberg, S. Shapley flow: A graph-based approach to interpreting model predictions. In Banerjee, A. and Fukumizu, K. (eds.), *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pp. 721–729. PMLR, 13–15 Apr 2021. URL <https://proceedings.mlr.press/v130/wang21b.html>.
- Winter, E. The shapley value. *Handbook of game theory with economic applications*, 3:2025–2054, 2002.
- Wu, E., Wu, K., and Zou, J. Explaining medical ai performance disparities across sites with confounder shapley value analysis, 2021. URL <https://arxiv.org/abs/2111.08168>.
- Zenke, F., Poole, B., and Ganguli, S. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*, pp. 3987–3995. PMLR, 2017.

A. Additional Experimental Results On Synthetic Data

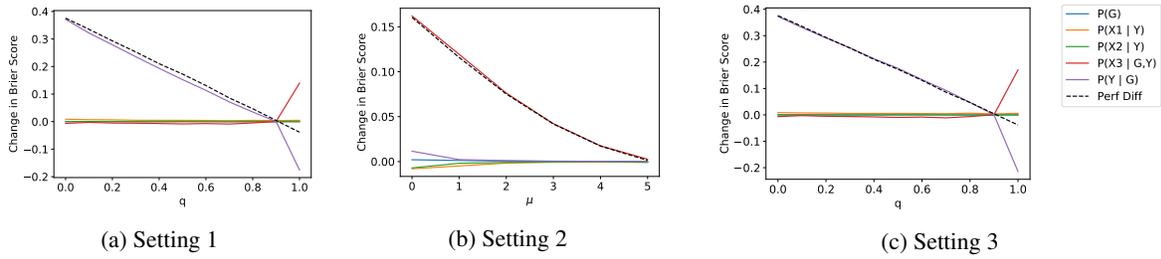


Figure A.1. Brier score differences attributed by our model to five potential distributional shifts on the synthetic dataset for the LR model.

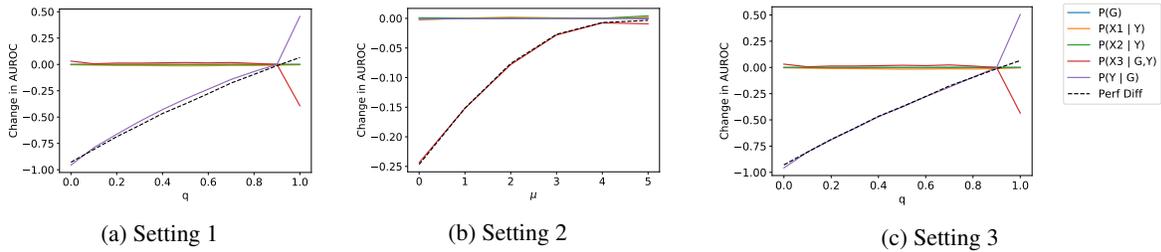


Figure A.2. AUROC differences attributed by our model to five potential distributional shifts on the synthetic dataset for the XGB model.

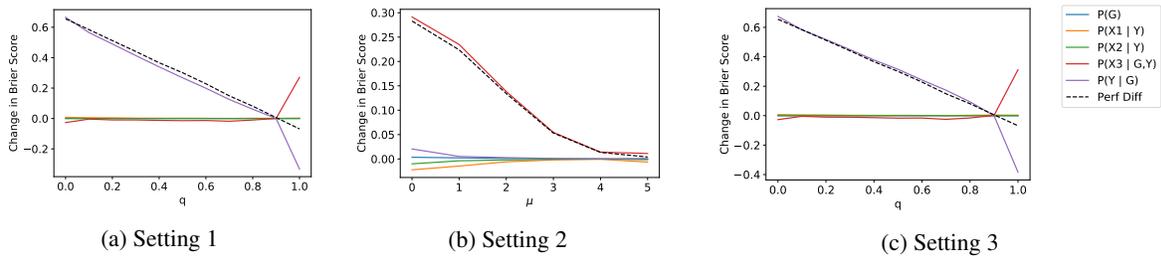
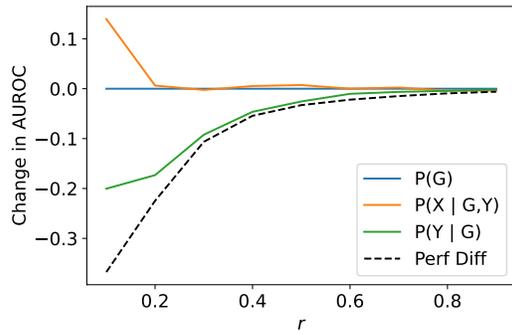
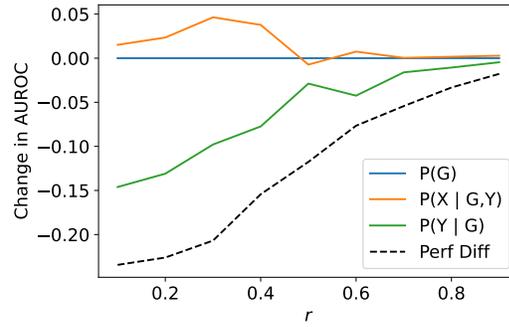


Figure A.3. Brier score differences attributed by our model to five potential distributional shifts on the synthetic dataset for the XGB model.

B. Experimental Results On Real-World Data

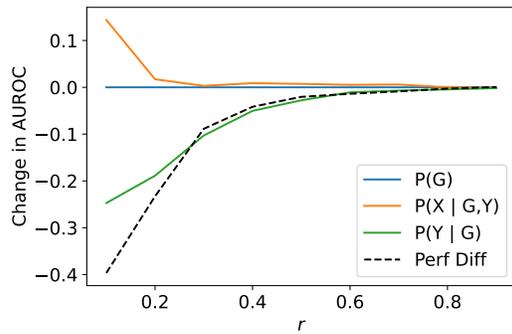


(a) Waterbirds

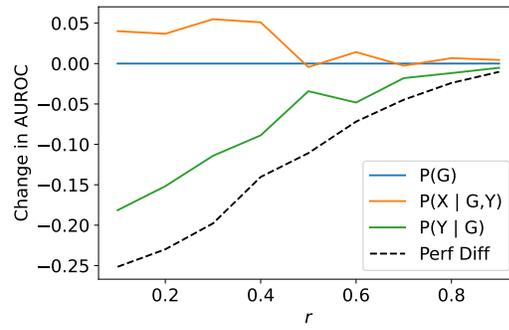


(b) MetaShift

Figure B.4. AUROC differences attributed by our model to three potential distributional shifts on two real-world datasets for the LR model.



(a) Waterbirds



(b) MetaShift

Figure B.5. AUROC differences attributed by our model to three potential distributional shifts on two real-world datasets for the XGB model.

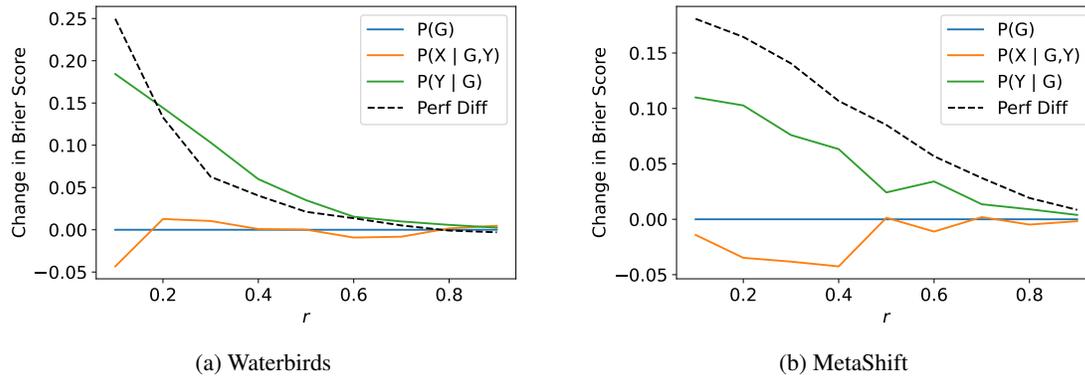


Figure B.6. Brier score differences attributed by our model to three potential distributional shifts on two real-world datasets for the XGB model.