
Unified Shapley Framework to Explain Prediction Drift

Aalok Shanbhag^{*1} Avijit Ghosh^{*1,2} Josh Rubin^{*1}

Abstract

Predictions are the currency of a machine learning model, and to understand the model’s behavior over segments of a dataset, or over time, is an important problem in machine learning research and practice. There currently is no systematic framework to understand this drift in prediction distributions over time or between two semantically meaningful slices of data, in terms of the input features and points. We propose GroupShapley and GroupIG (Integrated Gradients), as axiomatically justified methods to tackle this problem. In doing so, we re-frame all current feature/data importance measures based on the Shapley value as essentially problems of distributional comparisons, and unify them under a common umbrella. We axiomatize certain desirable properties of distributional difference, and study the implications of choosing them empirically.

1. Introduction

Fundamental Machine Learning theory is structured around the assumption that the training and test data belong to the same distribution. While this is a reasonable assumption for learning theory, for practical usage this is difficult to ascertain. In deployed models, the incoming stream of data might start to be significantly different than the static dataset that the model was trained on – a phenomenon named *Concept drift* (Gama et al., 2014; Wang & Abraham, 2015). Formally, concept or data drift is defined as the scenario where the distribution of the data $P(X)$, the label $P(Y)$ or the concept $P(Y|X)$ changes as compared to the training data that the model has seen. Previous drift detection methods have either focused on the overall error rate (Gama et al., 2004), or some other combination of the confusion matrix (Wang et al., 2013), either of which require the prediction labels, which is not guaranteed for a

machine learning model in production. Other work suggests using prediction drift as a proxy for concept drift in such cases (Žliobaite, 2010). It informs of the change in the model’s prediction distribution, and this information may be of importance to the practitioner, even if the model’s accuracy is not impacted. For example, a lending company may have a quarterly target of loans to be disbursed, which may be achieved in a month, if the prediction distribution of real world applicants differs from training data, thereby causing issues in business planning. A systematic method is thus needed for studying prediction drift and attributing it to a) the features of the model and b) the individual data points that constitute the distributional samples that are compared. We frame the question as follows: *Has the empirical distribution of inputs to the model drifted in a way that affects model behavior? If so, which features and which points in the sample have caused this shift?*

For this attribution to features and data, we focus on Shapley value based methods (Štrumbelj & Kononenko, 2014; Lundberg & Lee, 2017; Sundararajan & Najmi, 2020; Datta et al., 2016). We include Integrated Gradients (Sundararajan et al., 2017) in this broad family as it is equivalent to the Aumann-Shapley cost sharing method. Here, we adapt the Shapley framework, in the context of machine learning, for the following task: given two data samples of the same shape, and a function D which computes some metric of distributional difference on the predictions made on the given datasets by a model F , attribute the output of D to each point of the target dataset, and to each feature. By using the Shapley framework, we automatically inherit the Shapley axioms that have certain desirable properties, which we discuss in section 4.2.

Currently, there is no consensus on which of the many distributional difference metrics should be used to for calculating prediction drift, with previous work using measures like Jensen Shannon divergence (Pinto et al., 2019), Kolmogorov–Smirnov test (dos Reis et al., 2016) or the Wasserstein-1 distance (Miroshnikov et al., 2020). A comparative analysis of these methods is presented in Section 7. We demonstrate an axiomatic framework to choose the most appropriate distributional distance metric depending on the use case.

Contributions Our key contributions in this paper are:

^{*}Equal contribution ¹Fiddler Labs, Palo Alto, USA
²Northeastern University, USA. Correspondence to: Aalok Shanbhag <aalokshanbhag@gmail.com>.

- Establishing an axiomatic framework for calculating and explaining prediction drift using Shapley values and IG
- Extending the framework to explain arbitrary “groups” i.e. data and features together, thereby unifying several existing explanation methods
- Applying the Shapley values formulation to a function of distributional difference
- Axiomatization of measures of distributional difference
- Empirical analysis of the implications of choosing a particular metric of distributional difference to measure prediction drift, over a few handcrafted examples

2. Related Work

Concept Drift The problem of concept drift in machine learning has been extensively studied in the literature – spanning both sudden/instantaneous drift (Nishida et al., 2008) or slow/gradual drift (Stanley, 2003). Furthermore, the literature makes a distinction between “true” concept drift, and the “virtual” concept drift that happens due to a change in data distribution, essentially a sampling issue (Salganicoff, 1997).

In the literature, popular methods for detecting concept drift are AD-WIN (Bifet & Gavalda, 2007) and the Page-Hinkley Test (Gama et al., 2013), with both these methods assuming that labels are available for analysis, which is infeasible in a scenario where the model is deployed and constantly making predictions on new data.

Prediction Drift as a proxy In the absence of instantly available labels, other methods devolve to a measurement of the drift of the distributions of predictions as an ad-hoc method to detect concept drift. Methods using this approach include work by Pinto et al. (2019), Ghanta et al. (2019), dos Reis et al. (2016) and Žliobaite (2010). All these methods utilize different metrics to measure the difference in the distributions of the predictions of the new data points against a standard distribution.

Model explanations Methods to describe the contribution of input features towards the final value of the prediction have gained considerable interest in the present, both from researchers and practitioners. One class of these methods effectively utilize Shapley value (Shapley, 1953), a popular concept in game theory to measure the contribution of each feature. Another method, Integrated Gradients (Sundararajan et al., 2017), utilizes the path between the input and a baseline for each feature to measure the attribution of each feature on that path, as a specialized case of the Shapley-like cost sharing method

(commonly referred to as Shap). These methods have grown in popularity for quantifying the impact of features on a prediction at an instance level (Štrumbelj & Kononenko, 2014; Datta et al., 2016; Lundberg & Lee, 2017), on the loss at a global level (Covert et al., 2020), as also for quantifying the contribution of individual data points to a model’s performance (Ghorbani et al., 2020). They have also been proposed to understand feature importance for measures of fairness (Begley et al., 2020; Miroshnikov et al., 2020).

3. Terminology

In this section, we lay out the terminology and notations that we use throughout the paper.

Model function – the machine learning model function F , which takes a vector of features of shape n and returns an output(s). We limit the analysis to feature vectors instead of more general feature tensors, to avoid complications in notation. This does not mean however that the theory is particular to only models that accept single dimensional vectors, and can be extended quite easily.

Similarly, for the sake of simplicity and without loss of generality, for models which output a vector of values, we analyze only one output at a time. For example, classification models output a vector of length equal to the number of classes, of which there is a particular class of interest which we wish to analyze. $F : R^n \rightarrow R$. Akin to machine learning models with batch predict, F is also able to accept a batch input of shape (m, n) and return m outputs.

Sample (of points) – in the context of a model, a sample of m feature vectors of shape n , the complete sample hence being of shape (m, n) . The sample could be a single point ($m = 1$) or multiple ($m > 1$). It could be chosen randomly from a distribution, or could be chosen intentionally as per requirements, e.g. points corresponding to men over the age of 50 from New York or the feature vector corresponding to ID “x” in a database of customers of an online retail store.

Explicand – the input sample of shape (m, n) for which we want to explain the predictions, with respect to a particular model function.

Baseline – a sample with the same dimensions as the explicand, against which the explicand is explained. All Shapley value based methods have a baseline, though it may not be obvious due to being implicit in the formulation (Sundararajan & Najmi, 2020; Lundberg & Lee, 2017). The explanation is dependent on the choice of baseline, and various papers (Merrick

& Taly, 2020) have proposed certain choices of baselines, or ways to select one.

Value Function – the set function $v(s), v : 2^n \rightarrow R$ that is used in the Shapley value formulation to obtain the attribution of each player. Here n is the number of features, and 2^n refers to all possible combinations of feature presence (or absence).

Drift – a measure of distributional difference, commonly used in context of time dependence, but we use it in a general sense.

Distributional drift function – a function that given two samples (as defined above), returns a value characterizing the difference between them. We restrict ourselves to analyzing distributional differences over 1-D samples.

$$D(s_1, s_2) : R^n \times R^n \rightarrow R$$

Groups – combinations of the $n \times m$ feature–data-point components belonging to the explicand. These groups play the role of “players” in co-operative game theory for the purpose of Shapley and IG attributions of the drift value. Groups can be defined semantically, for example - males and females, and can be formed as combinations both in the feature and data-point dimensions. The Shapley value is calculated on the marginals of the resulting $m' \times n'$ groups as players that enter the coalition, over all such possible permutations.

4. Axioms

4.1. Axioms for attributions

From (Sundararajan & Najmi, 2020; Friedman & Moulin, 1999), we have the following desirable properties for attribution methods. In Section 6 we will formulate GroupShapley and GroupIG such that they are inherited. We state them here in terms of the group formulation for convenience. Reasons for their desirability are expanded on in the Appendix.

1. **Dummy** - this axiom states that a group that doesn't contribute to the game payout should get zero attribution.
2. **Efficiency** - the sum of the attributions over all groups is equal to the difference of the model function's output at the explicand and the baseline.
3. **Linearity** - the attributions of the linear combination of the two model functions, are the same linear combination of the attributions of the model functions, taken one at a time.

4. **Symmetry** - for model functions that are symmetric for two groups i and j , and the groups have the same value in both the explicand and baseline i.e. $g_i = g_j$ and $g'_i = g'_j$, the attributions to both the groups should be the same.
5. **Affine Scale Invariance** - requires the attributions to be invariant under the same affine transformation of both the model functions, and the groups.
6. **Demand Monotonicity** - for a model function that is monotonic for a group, the attribution of the group should only increase if the value of the group increases.
7. **Proportionality** - if the model function can be expressed as an additive sum of the input groups, and the baseline is zero, the attributions to each group are proportional to the group value.

4.2. Axioms for distributional drift functions

Miroshnikov et al. (2020) propose some desirable properties for a distributional drift function:

1. It should be **continuous** with respect to the change in the geometry of the distributions.
2. It should be **non-invariant** with respect to monotone transformations of the distributions.

Since our focus is on the distributional samples, and not the distributions themselves, we restate these properties for distributional drift measures D for two 1-D samples.

1. **Sensitivity** - the drift function should be continuous with respect to changes in the individual points in the samples. For example, given two 1-D samples S_1 and S_2 , if we change the value of any point in either, the function output should change.
2. **Differentiability** - the drift function should be differentiable with respect to the individual points in the samples - this is a stronger version of the continuity axiom
3. **Symmetry** - the drift function of two samples S_1 and S_2 should be symmetric i.e. $D(S_1, S_2) = D(S_2, S_1)$
4. **Identity of Indiscernibles** - the drift is zero if and only if both samples are the same $D(S_1, S_1) = 0$ and $D(S_1, S_2) \neq 0$ if $S_1 \neq S_2$
5. **Directionality** - the drift is signed based on the sample order $D(S_1, S_2) = -D(S_2, S_1)$. A metric cannot satisfy both Symmetry and Directionality, unless it's always zero

5. Prediction Drift

We define prediction drift as the change in the distribution of the predictions of a model between two semantically meaningful slices of data.

The need for studying prediction drift to answer the question raised above arises due to the following reasons:

1. Detecting drift in the distribution of individual features may not be sufficient. For instance, it could be that the predictions may drift despite no drift in any of the individual feature distributions. This is because the joint distributions of the features may have drift.

Reference Distribution				Target Distribution			
x	y	z	f(x, y, z)	x	y	z	f(x, y, z)
1.0	1.0	1.0	3.0	3.0	1.0	2.0	9.0
2.0	2.0	2.0	8.0	1.0	2.0	3.0	8.0
3.0	3.0	3.0	15.0	2.0	3.0	1.0	6.0

Table 1. The model function is $f(x, y, z) = xz + y + z$. The x , y and z distributions are unchanged at the univariate level, but the multivariate distribution has changed, so has the prediction distribution.

2. Furthermore, drift in individual features may not always lead to drift in predictions. This could, for instance, happen if the drifting feature is unimportant to the model.
3. Finally, detecting drift in the prediction distributions may not be sufficient either. For instance, while the predictions distributions may remain the same, it could still be that the input feature distributions have changed in a meaningful way that affects how the model reasons. Such a drift is still worth noting. For instance, the camera that feeds a face detection model could rotate over time, due to hinge failure. A robust model will be able to handle the distortion of the image for a while before it fails. The prediction distribution will not change initially, but the feature attributions over the pixels regions will change, which can serve as an early warning system.

We focus our attention on problems 1 and 2, leaving 3 for future work. To answer the aforementioned question, we rely on the following steps:

- Measure prediction drift for the model given two slices of data
- Attribute the drift to meaningful groups in the data.

Possible meaningful groups could be features of the model, n-tile buckets of predictions, or rule-based slices such as males vs females. We need to be careful to ensure that the the number of observations in each slice is proportionally similar for each sample, to avoid statistical anomalies seen in Simpson’s paradox. (Wagner, 1982)

Practically, for calculating the prediction drift given two data samples of unequal and/or large size, we suggest a bootstrapping approach. We sample from the two empirical distributions for a given number of repetitions and calculate the expected value of the prediction drift and the attributions and obtain statistical confidence bounds.

6. Group Shapley and Group IG Formulation

6.1. The Shapley value

Model function is $F : R^n \rightarrow R$

The Shapley value of a player i , playing an n-player coalitional game with a payout function v is defined as

$$\varphi_i(v) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|! (n - |S| - 1)!}{n!} (v(S \cup \{i\}) - v(S)) \quad (1)$$

6.2. Baseline Shapley

Baseline Shapley (Sundararajan & Najmi, 2020) or BShap, takes a function f , an explicand x and a baseline x' .

The value or payout function is $v(S) = f(x_S; x'_{N \setminus S})$

Here, the absence of a feature is modeled using the corresponding baseline value. BShap is equivalent to the Shapley-Shubik cost sharing method and satisfies the following axioms: Dummy, Linearity, Affine Scale Invariance, Demand Monotonicity, and Symmetry.

6.3. Integrated Gradients

The Integrated Gradients formulation is

$$IG_i(x, x', f) := (x_i - x'_i) \int_{\alpha=0}^1 \frac{\partial F(x' + \alpha \times (x - x'))}{\partial x_i} d\alpha \quad (2)$$

Integrated gradients is equivalent to the Aumann-Shapley cost sharing method for continuous functions.

Integrated Gradients satisfies the following axioms: Dummy, Linearity, Affine Scale Invariance, Proportionality, and Symmetry

6.4. Drift Group Shapley

We define Drift Group Shapley, or GroupShapley, as being parametrized by the following choices:

1. Choice of the explicand S_1 of shape (m, n)
2. Choice of the baseline S_2 of the same shape as the explicand
3. A model function $F : R^n \rightarrow R$
4. Additional functions, the chain of which we call G , which return two real valued outputs of equal shape for both the explicand and the sample
5. Choice of a distributional difference function D , that takes two equal shaped outputs of the function $G \circ F$ and returns a real valued output

The group formulation is:

$$\sum_{n=1}^G \varphi_i = D(G(F(S_1)), F(S_2)) \quad (3)$$

In GroupShapley, we explain the drift between the $G \circ F$ output of the explicand and the baseline. The number of players is equal to the number of groups times the number of features. The number of groups is the number of sub-divisions across rows. If the whole sample is one group, the features are the only players. If we have a row as it's own group, we end up with number of rows \times number of features groups to which we attribute the payout. To be precise, we are attributing the drift score to each group in the explicand, where a group is a cross section consisting of at least one row and at most all rows, and at least one feature or at most n features.

To simulate for missingness of a player, we replace the group of interest, with it's aligned counterpart from the reference dataset, similar to notion of the baseline in BShap or IG.

We now propose to frame every existing Shapley formulation as a prediction drift between some aspect of the model's behavior at the explicand and the baseline. We re-frame the two questions as:

1. Has the empirical distribution of inputs to the model drifted in a way that affects model behavior? becomes *Is there a difference in groups between the explicand and the baseline that affects some aspect of model behavior?*
2. If so, which features and which points in the sample have caused this shift? becomes *If so, which groups have caused it?*

We list the following existing methods which we attempt to bring under a common umbrella:

1. (Merrick & Taly, 2020) unifies BShap/KernelSHAP/QII, noting that the KernelSHAP (CES) and QII (RBShap) can be derived by taking the expectation of BShap over particular distributions, namely the input distribution for KernelSHAP and the joint marginal for QII. The approach in (Štrumbelj & Kononenko, 2014) is equivalent to kernelSHAP (Sundararajan & Najmi, 2020)

Therefore, we can consider KernelSHAP and QII to be the following case of GroupShapley: Explicand is of shape $(1, n)$, broadcast to (m, n) where m is the size of the background sample over which the expectation is calculated. The groups are the n features and the drift function is the expected value difference.

2. SAGE (Covert et al., 2020) is a global explanation method, where the aim is to attribute the loss of the model to the features, by suggesting that a feature whose removal increases the loss is more important. The loss is computed over a data sample of shape (m, n) . They propose using the conditional distribution as in CES in theory, but in practice use the marginal, as in RBShap. This is equivalent to GroupShapley on n groups, broadcasting the row dimension to $m \times j$ where j is the size of the background baseline sample the applicable distribution. The drift function is the expected value difference.
3. Distributional Shapley (Ghorbani et al., 2020) aims to find the value of a data point, given a model and an evaluation metric. There is no inherent concept of a baseline here, though we could trivially add a set of random data as the baseline. We can design G so as to make the $G \circ F$ of the baseline to be zero. The drift function is the expected value difference between the accuracy on the explicand and the artificially created zero value accuracy of the baseline. We note that it may be more instructive to introduce the notion of a baseline here, so as to ground the value of a datum in more definite terms. For example, is the data from source A more informative than source B.
4. In (Miroshnikov et al., 2020), they propose using Shapley values to explain the Wasserstein-1 distance between two prediction samples, each belonging to a class of a protected attribute like Gender, Race and so on. This is directly analogous to our scheme.

6.5. Drift Group Integrated Gradients

We define Drift Group IG, or GroupIG, as being parametrized by the following choices:

1. Choice of the explicand S_1 of shape (m, n)
2. Choice of the baseline S_2 of the same shape as the explicand
3. A model function $F : R^n \rightarrow R$ that is end-to-end differentiable with respect to the inputs
4. Additional functions, the chain of which we call G , which return two real valued outputs of equal shape for both the explicand and the sample. G has to be differentiable in terms of the individual samples
5. Choice of a distributional difference function D , that takes two equal shaped outputs of the function $G \circ F$ and returns a real valued output. Again, D has to be differentiable in terms of the original input samples

In GroupIG, we go from the baseline sample to the explicand in a straight line path. We can thus say that IG is a particular case of Drift Group IG, where $m = 1$, G is the identity function and the distributional difference function is the expected value difference. If we are using the Wasserstein-1 distance for a single input, we re-frame the function as the absolute distance between the prediction at the input and the baseline prediction.

7. Distributional Distance Metrics

We now discuss the properties of some of the widely used distance metrics for distances between two 1-D samples R_1, R_2 of data of length n .

7.1. Wasserstein-1 Distance

The Wasserstein-1 distance, also called the Earth Mover's distance or Mallows distance, is a well known metric from optimal transport theory, and widely used in statistics and machine learning. The mathematical properties which aid its suitability for our task are discussed below, building on prior work (Kolouri et al., 2018; Miroshnikov et al., 2020; Jiang et al., 2020)

For the case of two 1-D samples, which is the case we are focusing on, the W_p distance is the L_p norm of the sorted samples. The Wasserstein-1 distance is the special case where p is 1.

$W_p = (\frac{1}{n} \sum_{i=1}^n |R_1(i) - R_2(i)|^p)^{\frac{1}{p}}$. Hence for $p = 1$, it reduces to the mean of the L1 norm. (Levina & Bickel, 2001)

The W_1 distance for empirical samples satisfies the following distributional axioms: Sensitivity, Differentiability, Symmetry, and the Identity of Indiscernibles. (Proofs in Appendix)

7.2. Expected value difference

Expected value difference, can be understood simply as the difference in the Expected Value of two distributions. Given two samples, it's the difference in the mean. This is a very intuitive concept, and is the simplest measurement of distributional difference, corresponding to a change in the first order moment.

$$E(R_1, R_2) = \frac{1}{n} \sum_{i=1}^n (R_1(i) - R_2(i)).$$

The Expected value distance for empirical samples satisfies the following distributional axioms: Sensitivity, Differentiability, and Directionality but not the Identity of Indiscernibles. (Proofs in Appendix)

7.3. Jensen Shannon Divergence

The Jensen Shannon Divergence (JSD) given two probability distributions P and Q is defined as

$$JSD(P, Q) = \frac{1}{2}(D(P||M) + D(Q||M))$$

where $M = \frac{1}{2}(P + Q)$ and D is the Kullback-Liebler divergence.

While it is difficult to analyze JSD's behavior given empirical samples, we can see that it does not satisfy Sensitivity and Directionality. (Proofs in Appendix)

7.4. Kolmogorov-Smirnov Test Statistic for Two Samples

This is actually a test to determine if two empirical probability distributions differ, and yields a distance that is used as measure of distributional difference. (dos Reis et al., 2016)

The KS statistic distance is defined as

$$D(P, Q) = \sup_x |P(x) - Q(x)|$$

where $P(x)$ and $Q(x)$ are the empirical Cumulative Distribution Functions (CDF) of R_1 and R_2 and \sup is the supremum.

The KS Statistic satisfies only the Symmetry and the Identity of Indiscernibles axiom. (Miroshnikov et al., 2020)

8. The concept of Alignment

Given the need of a baseline in the Shapley value and IG formulations, it is natural to ask what is the right baseline, given that the attributions will differ with the choice of baseline. This is one of the most important questions in explainability. (Sundararajan et al., 2017) recommends choosing a baseline where the model's prediction is neutral. (Merrick & Taly, 2020) argues for contrastive explanations, with justification from norm theory (Kahneman & Miller,

Unified Shapley Framework to Explain Prediction Drift

Function(x, y, z)	Explicand [x, y, z]	Baseline [x, y, z]	Exp. value Difference	W ₁ Distance	Shapley	W ₁ Shapley	IG	W ₁ IG
xy	[1, 2, 3]	[0, 0, 0]	2.0	2.0	[1. 1. 0.]	[1. 1. 0.]	[1. 1. 0.]	[1. 1. 0.]
$x - y$	[1, 2, 3]	[0, 0, 0]	-1.0	1.0	[1. -2. 0.]	[0. 1. 0.]	[1. -2. 0.]	[-1. 2. 0.]
$x + y - z$	[1, 2, 3]	[0, 0, 0]	0.0	0.0	[1. 2. -3.]	[0. 0. 0.]	[1. 2. -3.]	[0. 0. 0.]
$xy - z^2$	[1, 2, 3]	[0, 0, 0]	-7.0	7.0	[1. 1. -9.]	[-0.33 -0.33 7.67]	[1., 1., -9., 0.]	[-1., -1., 9.]
$\min(x, y)$	[1, 2, 3]	[0, 0, 0]	1.0	1.0	[0.5 0.5 0.]	[0.5 0.5 0.]	[1 0 0.]	[1 0 0.]
$\text{abs}(x - y)$	[1, 2, 3]	[0, 0, 0]	1.0	1.0	[0. 1. 0.]	[0. 1. 0.]	[-1 2 0]	[-1 2 0]

Table 2. BShap and IG Attributions for functions using Expected Value Difference and Wasserstein-1. Note the sparser attributions using the Wasserstein distance

1986).

In GroupShapley and GroupIG when using the W_1 drift function, we take the counterpart in the other sample as baseline, when both samples are aligned by their sorted prediction values. The W_1 distance is based on the concept of optimal transport, and hence, the intuition extends naturally to the flow from the attributions, which make up the prediction from one distribution to the other.

For other drift metrics, there may not be a natural reason to align in any particular way. But the alignment of the W_1 distance still can be justified as comparing the most similar points in the two samples, if the prediction of model is viewed as a task specific dimensionality reduction. Fliptest (Black et al., 2020) uses a similar thought process for assessing individual fairness by creating counterfactuals via optimal transport.

The alternative, where no choice needs to be made, is to take the expectation over all possible alignments.

9. Analysis

We now look at some practical examples of how the choice of drift function impacts the explanations.

9.1. Simple Experiments

We analyze BShap and IG for a few functions in Table 2, using both the expected value difference and the W_1 distance. These are functions of three variables x , y , and z , the baseline for all is $[0, 0, 0]$, and the explicand is $[1, 2, 3]$. For the function $x - y$, we see that the W_1 attributions are different for different for both BShap and IG. It seems that the W_1 drift function gives sparser attributions for BShap, by compressing the attributions for the features that act in opposite direction to the eventual predicted value. For instance, for $x - y$, x is 1 and y is 2, so the prediction is -1, and the W_1 distance from the baseline prediction is 1. The method gives all the attribution to y , as it has the sign of the prediction. We can see this behavior for $x + y - z$ and $xy - z^2$ as well. This is reminiscent of how the L1 norm sparsifies coefficients in ridge regression, but we make no claims of there being any analogy between the two.

There is no reason to always prefer the explanation of one over the other, both can be justified in their own way and are a matter of choice, similar to how choosing a baseline is a choice depending on the question one is looking to answer.

9.2. Case Study

We now present a simple case study, to demonstrate how this might work in practice, by constructing a synthetic dataset. This allows us to inject known and controlled drifts in order to evaluate the effectiveness of various methods at finding them.

We create a dataset of the following features:

1. Location - {‘Springfield’, ‘Centerville’} - 70:30
2. Education - {‘GRAD’, ‘POST_GRAD’} - 80:20
3. Experience - years - (0, 50) - normally distributed
4. Engineer Type - {‘Software’, ‘Hardware’} - 85:15
5. Relevant Experience - years - (0, 50) - normally distributed

and ensure that $\text{experience} \geq \text{relevant experience}$.

The model predicts an individual’s salary from the features above, using the following formula:

$$\begin{aligned} \text{Salary} = & 50,000 + 20,000 \times \text{location} + 20,000 \\ & \times \text{education} + 5,000 \times \text{relevant_experience} \\ & + 100 \times \text{experience} + 10,000 \times \text{engineer_type} \end{aligned}$$

2000 events are created for each of three days. On the second day, a plausible data pipeline bug is introduced, whereby the location feature has the value “springfield” rather than “Springfield”. Because of this, all locations are identified as ‘Centerville’, which leads to an average salary drop for day two— a prediction drift. We now would like to attribute this the offending features. Figure 1 shows the drift over time measured by the various drift methods previously discussed.

In Figure 2, we calculate GroupShapley attributions over the fifteen feature-day combinations, and see that the job location feature gets the most attribution, as we would expect.

Additionally, we compare our approach to that of (Pinto et al., 2019), which measures drift using Jensen Shannon divergence and trains a Gradient Boosted Tree Classifier to identify the drift. The feature importances of the classifier are used to identify the cause. In the scenario described, it correctly gives the most attribution to the location feature. But if we introduce another spurious drift, of an unimportant feature like experience, the GBDT method selects the wrong feature. They do suggest a technique to remove time-trended features, but if the other feature also spikes in the same interval, that fix will not help either as seen in Figure 3.

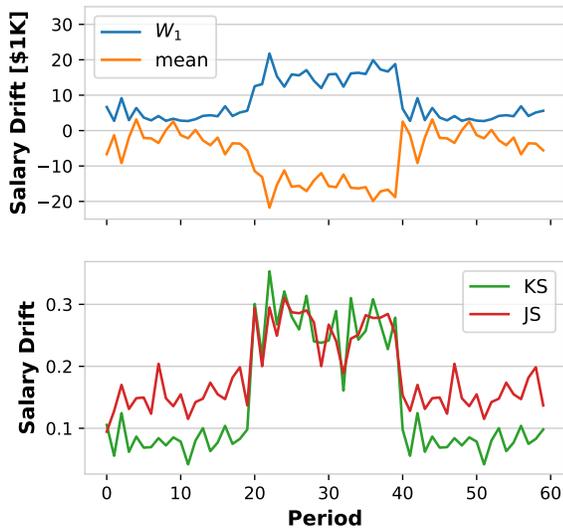


Figure 1. Comparison of four distributional drift functions. Wasserstein and means/expected value difference [upper] preserve the units of the predicted quantity and may provide a more intuitive scale. The Kolmogorov-Smirnov Statistic and Jensen Shannon Divergence [lower] are dimensionless and the scale reflects the degree of absolute distributional overlap. The central 20 periods have a data integrity error intentionally introduced which causes some applicants to have their location misinterpreted.

10. Conclusion and Future Work

We study the problem of prediction drift and attributing it, and propose it as a general framework of explainability, unifying several methods. We axiomatize certain desirable properties of distributional difference metrics, also demonstrating that explanation methods can be parameterized by the choice of this metric.

A more detailed study of the theoretical implications of choosing one distance metric over another for explanations is left for future work. Additionally, GroupShapley can be computationally expensive, and approximation schemes for faster calculations could be a future area of exploration.

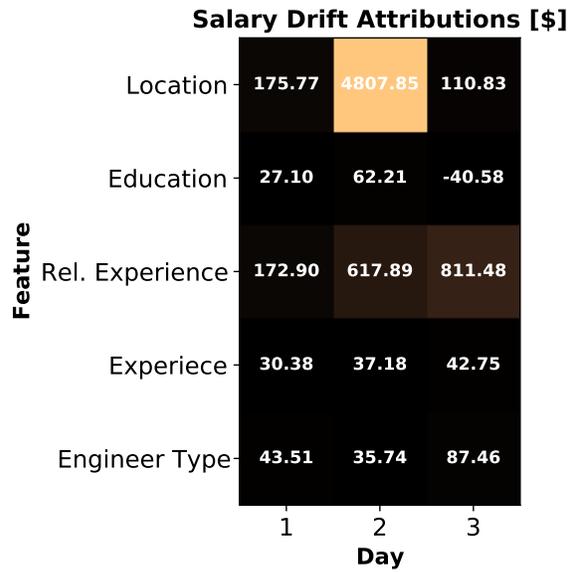


Figure 2. W_1 measures a drift of \$7022.46 over the complete dataset. By forming groups of [features]×[days], GroupShapley unambiguously identifies the source of drift as the “Location” feature on the second day.

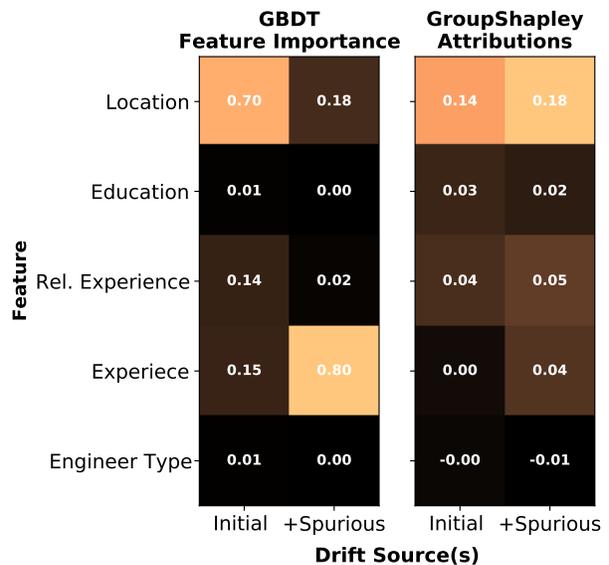


Figure 3. Comparison of tree-based feature importance method from (Pinto et al., 2019) and GroupShapley. Both methods initially identify the correct source of drift; but when an additional correlated feature drift is added, the GBDT method assigns it most of the importance, despite its minimal effect on the model output.

11. Appendix

11.1. Axioms

We will now go over the reasons for the desirability of the axioms:

1. **Dummy** - We do not want to credit a group/feature that makes no contribution to the model prediction.
2. **Efficiency** - This ensures a complete accounting of difference in the model's prediction between the explicand and the baseline.
3. **Linearity** - This property helps in avoiding counter-intuitive behavior when analyzing attributions of linear functions.
4. **Symmetry** - The purpose of this axiom is self-evident, if two groups contribute equally they should receive the same attribution.
5. **Affine Scale Invariance** - The justification for this is based on the idea that the units of measurement of individual features may not be comparable to each other, and secondly, within themselves, may not be canonical. For example, units of weight like pounds or kilograms are not more or less justified than the other, and the conversion to the other should not lead to a decrease in attribution. (Friedman & Moulin, 1999)
6. **Demand Monotonicity** - For a function that is monotonic with respect to a group, if the group value increases while all else is held constant, the function's value will increase. It is natural to want the attribution to the group to increase as compared to the previous scenario.
7. **Proportionality** - This ensures that the attributions to groups are proportional to their contribution in the additive sum of the group values. Let's look at a heat generation scenario. If there are three current sources, each supplying the same amount of current. The heat generated is proportional to the square of the current. The attribution to each should be one-third, compared to the zero baseline. Now if we combine two of the current sources, the attribution of the third should ideally remain the same.

11.2. Proofs for Drift Metrics satisfying Axioms

Wasserstein-1 Distance

Given two samples R_1 and R_2 of length n and sorted by value, the W_1 distance can be computed $W_1 = \frac{1}{n} \sum_{i=1}^n |R_1(i) - R_2(i)|$.

The W_1 distance for empirical samples satisfies the following distributional axioms:

Proofs:

1. **Sensitivity** - This is trivial to see, given that each point of the sample contributes to the overall sum.
2. **Differentiability** - The function is piece-wise differentiable, except at zero for each absolute difference.
3. **Symmetry** - The formula is symmetric in R_1 and R_2 .
4. **Identity of Indiscernibles** - The W_1 distance can be zero only if every element-pair in the two samples cancels each other out.

Expected value difference

Given two samples R_1 and R_2 of length n , the Expected value distance is $E(R_1, R_2) \approx \frac{1}{n} \sum_{i=1}^n (R_1(i) - R_2(i))$.

Proofs:

1. **Sensitivity** - Each point of the sample contributes to the overall sum.
2. **Differentiability** - One can see that the function is differentiable everywhere.
3. **Directionality** - The sign changes when the sample order is flipped.
4. **Identity of Indiscernibles** - This can be proved by a counter example. If there is a sample that only has values 1, and the other has equal number of zeros and twos. The two means will be equal and will cancel out, even though the two samples are not the same.

Jensen Shannon Divergence

The Jensen Shannon Divergence (JSD) given two probability distributions P and Q is defined as $JSD(P, Q) = \frac{1}{2}(D(P||M) + D(Q||M))$ where $M = \frac{1}{2}(P + Q)$ and D is the Kullback-Liebler divergence.

While it is difficult to analyze JSD's behavior given empirical samples, we can see that it does not satisfy Sensitivity and Directionality.

Proofs:

1. **Sensitivity** - This can be proved by a counter example. If there are two distributions that don't intersect anywhere, the JSD is one. Now if we translate the second distribution while ensuring there is no intersection, the JSD is still 1.

2. **Directionality** - JSD is symmetric to the change in the sample order.

Kolmogorov-Smirnov Test Statistic for Two Samples

For two distributions R_1 and R_2 , the KS statistic distance is $D(P, Q) = \sup_x |P(x) - Q(x)|$ where P and Q are the empirical Cumulative Distribution Functions (CDF) of R_1 and R_2

We can see from the definition that the KS Statistic satisfies Symmetry and the Identity of Indiscernibles axiom. For the other proofs please refer to ([Miroshnikov et al., 2020](#))

References

- Begley, T., Schwedes, T., Frye, C., and Feige, I. Explainability for fair machine learning. *arXiv preprint arXiv:2010.07389*, 2020.
- Bifet, A. and Gavalda, R. Learning from time-changing data with adaptive windowing. In *Proceedings of the 2007 SIAM international conference on data mining*, pp. 443–448. SIAM, 2007.
- Black, E., Yeom, S., and Fredrikson, M. Fliptest: fairness testing via optimal transport. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pp. 111–121, 2020.
- Covert, I., Lundberg, S., and Lee, S.-I. Understanding global feature contributions with additive importance measures. *Advances in Neural Information Processing Systems*, 33, 2020.
- Datta, A., Sen, S., and Zick, Y. Algorithmic transparency via quantitative input influence: Theory and experiments with learning systems. In *2016 IEEE symposium on security and privacy (SP)*, pp. 598–617. IEEE, 2016.
- dos Reis, D. M., Flach, P., Matwin, S., and Batista, G. Fast unsupervised online drift detection using incremental kolmogorov-smirnov test. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1545–1554, 2016.
- Friedman, E. and Moulin, H. Three methods to share joint costs or surplus. *Journal of economic Theory*, 87(2): 275–312, 1999.
- Gama, J., Medas, P., Castillo, G., and Rodrigues, P. Learning with drift detection. In *Brazilian symposium on artificial intelligence*, pp. 286–295. Springer, 2004.
- Gama, J., Sebastiao, R., and Rodrigues, P. P. On evaluating stream learning algorithms. *Machine learning*, 90(3): 317–346, 2013.
- Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., and Bouchachia, A. A survey on concept drift adaptation. *ACM computing surveys (CSUR)*, 46(4):1–37, 2014.
- Ghanta, S., Subramanian, S., Khermosh, L., Sundararaman, S., Shah, H., Goldberg, Y., Roselli, D. S., and Talagala, N. ML health: Fitness tracking for production models. *CoRR*, abs/1902.02808, 2019. URL <http://arxiv.org/abs/1902.02808>.
- Ghorbani, A., Kim, M., and Zou, J. A distributional framework for data valuation. In *International Conference on Machine Learning*, pp. 3535–3544. PMLR, 2020.
- Jiang, R., Pacchiano, A., Stepleton, T., Jiang, H., and Chiappa, S. Wasserstein fair classification. In *Uncertainty in Artificial Intelligence*, pp. 862–872. PMLR, 2020.
- Kahneman, D. and Miller, D. T. Norm theory: Comparing reality to its alternatives. *Psychological review*, 93(2): 136, 1986.
- Kolouri, S., Pope, P. E., Martin, C. E., and Rohde, G. K. Sliced-wasserstein autoencoder: An embarrassingly simple generative model. *arXiv preprint arXiv:1804.01947*, 2018.
- Levina, E. and Bickel, P. The earth mover’s distance is the mallows distance: Some insights from statistics. In *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, volume 2, pp. 251–256. IEEE, 2001.
- Lundberg, S. and Lee, S.-I. A unified approach to interpreting model predictions. *arXiv preprint arXiv:1705.07874*, 2017.
- Merrick, L. and Taly, A. The explanation game: Explaining machine learning models using shapley values. In *International Cross-Domain Conference for Machine Learning and Knowledge Extraction*, pp. 17–38. Springer, 2020.
- Miroshnikov, A., Kotsiopoulos, K., Franks, R., and Kannan, A. R. Wasserstein-based fairness interpretability framework for machine learning models. *arXiv preprint arXiv:2011.03156*, 2020.
- Nishida, K., Shimada, S., Ishikawa, S., and Yamauchi, K. Detecting sudden concept drift with knowledge of human behavior. In *2008 IEEE International Conference on Systems, Man and Cybernetics*, pp. 3261–3267, 2008. doi: 10.1109/ICSMC.2008.4811799.
- Pinto, F., Sampaio, M. O., and Bizarro, P. Automatic model monitoring for data streams. *arXiv preprint arXiv:1908.04240*, 2019.
- Salganicoff, M. Tolerating concept and sampling shift in lazy learning using prediction error context switching. In *Lazy learning*, pp. 133–155. Springer, 1997.
- Shapley, L. S. A value for n-person games, contributions to the theory of games ii (aw tucker and hw kuhn, eds.), 1953.
- Stanley, K. O. Learning concept drift with a committee of decision trees. *Informe técnico: UT-AI-TR-03-302, Department of Computer Sciences, University of Texas at Austin, USA*, 2003.

- Štrumbelj, E. and Kononenko, I. Explaining prediction models and individual predictions with feature contributions. *Knowledge and information systems*, 41 (3):647–665, 2014.
- Sundararajan, M. and Najmi, A. The many shapley values for model explanation. In *International Conference on Machine Learning*, pp. 9269–9278. PMLR, 2020.
- Sundararajan, M., Taly, A., and Yan, Q. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, pp. 3319–3328. PMLR, 2017.
- Wagner, C. H. Simpson’s paradox in real life. *The American Statistician*, 36(1):46–48, 1982. ISSN 00031305. URL <http://www.jstor.org/stable/2684093>.
- Wang, H. and Abraham, Z. Concept drift detection for streaming data. In *2015 international joint conference on neural networks (IJCNN)*, pp. 1–9. IEEE, 2015.
- Wang, S., Minku, L. L., Ghezzi, D., Caltabiano, D., Tino, P., and Yao, X. Concept drift detection for online class imbalance learning. In *The 2013 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–10. IEEE, 2013.
- Žliobaite, I. Change with delayed labeling: When is it detectable? In *2010 IEEE International Conference on Data Mining Workshops*, pp. 843–850. IEEE, 2010.