
Technical Report: A Stratification Approach to Partial Dependence for Codependent Variables

Terence Parr
University of San Francisco
parrt@cs.usfca.edu

James D. Wilson
University of San Francisco
jdwilson4@usfca.edu

Abstract

1 Introduction

Partial dependence, the isolated effect of a specific variable or variables on the response variable, y , is important to researchers and practitioners in many disparate fields such as medicine, business, and the social sciences. For example, in precision medicine, physicians are interested in the relationship between an individual’s demographic or clinical features and their susceptibility to illness. Business analysts at a car manufacturer might need to know how changes in their supply chain affect defect rates. Climate scientists are interested in how different atmospheric carbon levels affect temperature.

For an explanatory matrix, \mathbf{X} , with a single variable, x_1 , a plot of the y against x_1 visualizes the marginal effect of feature x_1 on y exactly. Given two or more features, one can similarly plot the marginal effects of each feature separately, however, the analysis is complicated by the interactions of the variables. Variable interactions, codependencies between features, result in marginal plots that do not isolate the specific contribution of a feature of interest to the target. For example, a marginal plot of sex (male/female) against body weight would likely show that, on average, men are heavier than women. While true, men are also taller than women on average, which likely accounts for most of the difference in average weight. It is unlikely that two “identical” people, differing only in sex, would be appreciably different in weight.

Rather than looking directly at the data, there are several partial dependence techniques that interrogate fitted models provided by the user: Friedman’s original partial dependence (FPD) Friedman [2000], Individual Conditional Expectations (ICE) Goldstein et al. [2015], Accumulated Local Effects (ALE) Apley [2016], and most recently SHAP Lundberg and Lee [2017]. Model-based techniques dominate partial dependence research because interpreting the output of a fitted model has several advantages. Models have a tendency to smooth over noise. Models act like analysis preprocessing steps, potentially reducing the computational burden on model-based partial dependence techniques; e.g., ALE is $O(n)$ for the n records of \mathbf{X} . Partial dependence techniques that interrogate models also provide insight into models themselves; i.e., how variables affect model behavior. It is also true that, in many cases, a predictive model is the primary goal so creating a suitable model is not an extra burden.

peer through the lens of a fitted model. As far as we can tell this is the only way it has been solved.

, although the techniques vary in their accuracy in the presence of codependent features

there are however several disadvantages to using a model. Many analysts do not need a predictive model nor would they know how to choose, tune, and assess a model. Could also be the case that a technique is not available in the desired deployment environment.

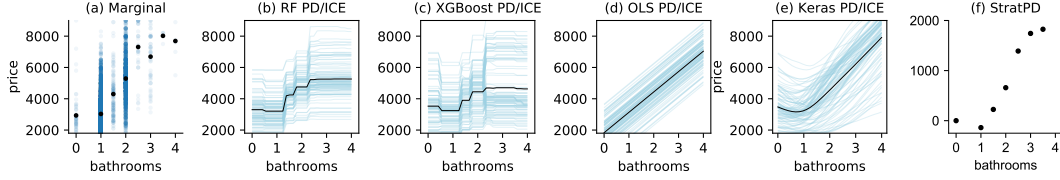


Figure 1: Plots of bathrooms versus rent price using New York City apartment rent data. (a) marginal plot, (b) PD/ICE plot derived from random forest, (c) PD/ICE plot derived from gradient boosted machine, and (d) PD/ICE plot derived from ordinary least squares regression; sample size is 10,000 observations of ~50k. The PD/ICE plots are different for the same data set, depending on the chosen user model. `X['bathrooms'].unique()` shows `(array([0. , 1. , 1.5, 2. , 2.5, 3. , 3.5, 4.]), array([54, 8151, 140, 1539, 39, 67, 3, 7]))`. STRATPD has missing last value, not enough data.

For the same data applying the same technique but using different models, we get different answers, which calls into question the validity of the curves. If we include the cost of cross validating grid search, particularly for deep learning, partial dependence can be expensive. can be used as an EDA tool to help with choosing a model: a nonparametric technique could also inform which machine learning model to use if a model is desired.

The techniques differ in algorithm simplicity, performance, and ability to isolate codependent variables. a nonparametric technique could also inform which machine learning model to use if a model is desired.

key is "all else being equal", which implies you don't want curves affected by other variables. Interaction plots are also very useful, such as ICE, but here our goal is the pure partial dependence curve. In the future, we hope to consider extracting interaction between variables like SHAP.

we introduce an ideal definition of partial dependence that does not rely on predictions from a fitted model based upon partial derivatives and then estimate partial derivatives nonparametrically to get partial dependence. The technique seems to isolate variables well and has linear behavior for numeric variables and mildly quadratic behavior for categorical variables in practice. The theoretical complexity is $O(n^2)$ like FPD.

SHAP is mean centered FPD for independent variables, proof in supplemental material.

state up front it only gets pure partial dependence, no interaction and has quadratic theoretical complexity, but it has the advantage that it doesn't require a fitted model. Sometimes there is an advantage to a model, smoothing etc. But, in many cases lack of model increases the accessibility of the tool to analysts and could prevent nonexpert machine learning practitioners from interpretation errors from poorly fit or tuned models.

2 Partial dependence without model predictions

Definition 1 The *idealized partial dependence* of y on feature x_j for smooth generator function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ evaluated at $x_j = z$ is the cumulative sum up to z :

$$PD_j(z) = \int_{\min(x_j)}^z \frac{\partial y}{\partial x_j} dx_j \quad (1)$$

$PD_j(z)$ is the value contributed to y by x_j at $x_j = z$ and $PD_j(\min(x_j)) = 0$. The advantages of this partial dependence definition are that it does not depend on predictions from a fitted model and is insensitive to collinear or otherwise codependent features, unlike the Friedman's original definition that he points out is less accurate for codependent data sets. We will denote Friedman's as FPD_j to distinguish it from this ideal, PD_j .

For example, consider quadratic equation $y = x_1^2 + x_2 + 100$ as a generator of data in $[0, 3]$. The partial derivatives are $\frac{\partial y}{\partial x_1} = 2x_1$ and $\frac{\partial y}{\partial x_2} = 1$, giving $PD_1 = x_1^2$ and $PD_2 = x_2$.

The obvious disadvantage of this feature impact definition is that function f , from which PD_j is derived, is unknown in practice, so symbolically computing the partial derivatives is not possible. But, if we could compute accurate partial dependence curves by some other method, then this definition would still represent a viable means to obtain feature impacts.

STRATPD stratifies a data set into groups of observations that are similar, except in the variable of interest, x_j , through the use of a single decision tree. Any fluctuation of the response variable within a group (decision tree leaf) is likely due to x_j . The β_1 coefficient of a simple local linear regression fit to the (x_j, y) values within a group provides an estimate of $\frac{\partial y}{\partial x_j}$ in that group's x_j range. Averaging the partial derivative estimates across all such groups yields the overall $\frac{\partial y}{\partial x_j}$ partial derivative approximation. The cumulative sum of the estimated partial derivative yields the partial dependence curve.

3 Existing work

FPD

ICE

ALE

SHAP

4 Experimental results

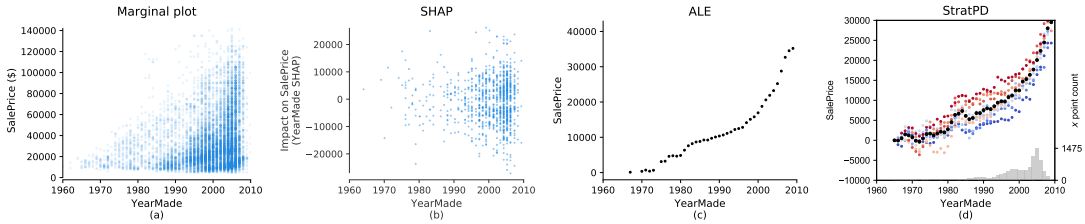


Figure 2: (a) Marginal plot of bulldozer **YearMade** versus **SalePrice** using subsample of 20k observations, (b) partial dependence drawn by SHAP interrogating an RF with 40 trees and explaining 1000 values with 100 observations as background data, (c) STRATPD partial dependence.

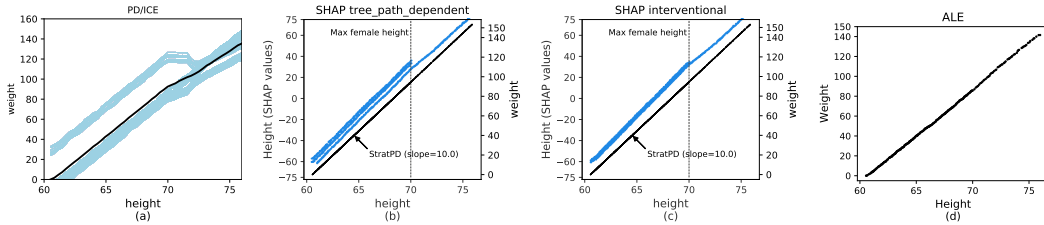


Figure 3: SHAP partial dependence plots of response body weight on feature **height** using 2000 synthetic observations from Equation (??). SHAP interrogated an RF with 40 trees and explained all 2000 samples; the interventional case used 100 observations as background data.

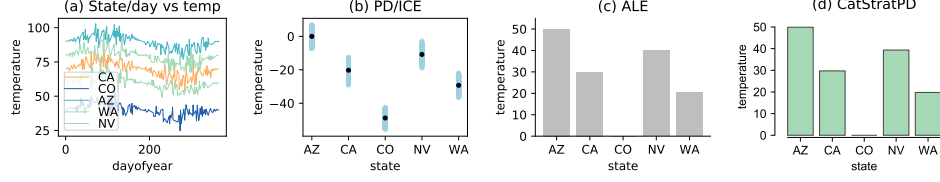


Figure 4: foo.

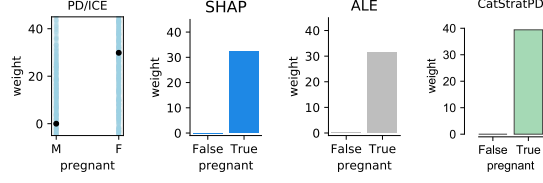


Figure 5: foo.

5 Algorithms

StratPD

```

Fit tree regressor to all but x_c with hyper parameter min_slopes_per_x
For each leaf:
    y bar = Group leaf samples by x_c, computing average y per unique x_c
    dx = discrete difference between adjacent unique x_c
    dy = discrete difference between adjacent average y bar
    add (x[i], x[i+1], dy[i]/dx[i]) for each unique x_c to list D

for each x in unique x_c from X:
    slopes = [slope for (a, b, slope) in D if x >= a and x < b]
    count[x] = |slopes|
    dydx[x] = mean(slopes)

Drop slope estimates computed using fewer than min_slopes_per_x values
pdx = discrete difference between adjacent unique x_c
pdy = cumulative sum of dydx * pdx
return pdx, [0]+pdy // insert 0 for pdx[0] since sum contributed from beyond left is 0

```

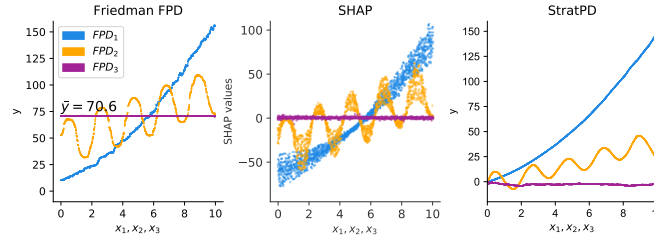


Figure 6: $y = x_1^2 + x_1x_2 + 5x_1\sin(3x_2) + 10$ where $x_1, x_2, x_3 \sim U(0, 10)$ and x_3 does not affect y . No noise added.

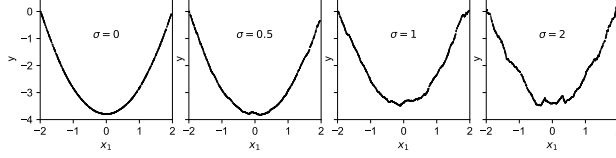


Figure 7: $y = x_1^2 + x_1 + 10 + N(0, \sigma)$ where $x_1, x_2 \sim U(-2, 2)$ and $\sigma \in [0, 0.5, 1, 2]$.

Algorithm: *StratPD*

Input: $\mathbf{X}, \mathbf{y}, c, \text{min_samples_leaf}, \text{min_slopes_per_x}$

Output: pdx, pdy : Unique x_c , partial dependence values across x_c

T := Decision tree regressor fit to $(\mathbf{X}_{\bar{c}}, \mathbf{y})$ with hyper-parameter: min_samples_leaf

for each leaf $l \in T$ **do**

$(\mathbf{x}_l, \mathbf{y}_l) = \{(x_c^{(i)}, y^{(i)})\}_{i \in l}$ // Get leaf samples

$\mathbf{ux} := \text{unique}(\mathbf{x}_l)$

$\bar{\mathbf{y}} :=$ Group leaf records $(\mathbf{x}_l, \mathbf{y}_l)$ by value of \mathbf{x}_l , computing \bar{y} per unique value

$\mathbf{dx} := \mathbf{ux}^{(i+1)} - \mathbf{ux}_{i=1..|\mathbf{ux}|-1}^{(i)}$ // Discrete difference

$\mathbf{dy} := \bar{\mathbf{y}}^{(i+1)} - \bar{\mathbf{y}}_{i=1..|\mathbf{ux}|-1}^{(i)}$

 Add tuples $(\mathbf{ux}^{(i)}, \mathbf{ux}^{(i+1)}, \mathbf{dy}^{(i)}/\mathbf{dx}^{(i)})_{i=1..|\mathbf{ux}|-1}$ to list \mathbf{d}

end

$\mathbf{ux} := \text{unique}(\{x_c^{(i)}\}_{i=1..n})$

for each $x \in \mathbf{ux}$ // Counts slopes, compute average slope per unique x_c value

do

$\text{slopes} := [\text{slope for } (a, b, \text{slope}) \in \mathbf{d} \text{ if } x \geq a \text{ and } x < b]$

$c_x := |\text{slopes}|$

$\text{dydx}_x := \text{slopes}$

end

$\text{dydx} := \text{dydx}[c \geq \text{min_slopes_per_x}]$ // Drop slope estimates computed from too few

$\mathbf{ux} := \mathbf{ux}[c \geq \text{min_slopes_per_x}]$

$\text{pdx} := \mathbf{ux}^{(i+1)} - \mathbf{ux}_{i=1..|\mathbf{ux}|-1}^{(i)}$

$\text{pdy} := [0] + \text{cumulative_sum}(\text{dydx} * \text{pdx})$ // integrate, inserting 0 for leftmost x_c

return pdx, pdy

CatStratPD

Fit tree regressor to all but x_c with hyper parameter min_slopes_per_x

For each leaf:

$\bar{\mathbf{y}}$ = Group leaf samples by categories of x_c , computing average \bar{y} per unique category x_c

 Compute unique categories and counts per category

refcat is randomly chosen category from x_c

For each unique category x **in leaf:**

$\text{delta}[\text{cat}, \text{leaf}] = \text{Subtract } \bar{\mathbf{y}} \text{ for refcat from all } \bar{\mathbf{y}}$ ($\text{refcat delta will be 0}$)

end

Let $\text{Avg}[\text{cat}]$ be vector with running sum mapping category to count

$\text{work} = \text{set of leaf indexes}$

while more work and something changed and less than max iterations:

for each leaf in leaves:

if cat in $\text{delta}[:, \text{leaf}]$ intersects with Avg :

$j = \text{random category in intersection}$

 adjust $\text{delta}[:, \text{leaf}]$ to be relative to j so $\text{delta}[j, \text{leaf}] == 0$ then add $\text{Avg}[j]$ so comparable

 merge into Avg

$\text{work} -=$ all j merged this iteration

Algorithm: *CatStratPD***Input:** $\mathbf{X}, \mathbf{y}, c, \text{min_samples_leaf}$ **Output:** $\Delta^{(k)}$ = category k 's effect on y where $\text{mean}(\Delta^{(k)}) = 0$
 $n^{(k)}$ = number of supported observations per category k T := Decision tree regressor fit to $(\mathbf{X}_{\bar{c}}, \mathbf{y})$ with hyper-parameter: min_samples_leaf // Get average y delta relative to random ref category for each sample in each leafLet $\Delta_{x,l}$ be dictionary mapping (category, leaf) to delta from ref categoryLet $\text{Count}_{x,l}$ be dictionary mapping (category, leaf) to count**for each leaf** $l \in T$ **do** $(\mathbf{x}_l, \mathbf{y}_l) = \{(x_c^{(i)}, y^{(i)})\}_{i \in l}$ // Get leaf samples $\mathbf{ux}, \mathbf{cx} := \text{unique}(\mathbf{x}_l)$ // Get unique categories, counts from leaf samples $\bar{\mathbf{y}} :=$ Group leaf records $(\mathbf{x}_l, \mathbf{y}_l)$ by categories of \mathbf{x}_l , computing \bar{y} per unique category $\text{refcat}_l :=$ random category from \mathbf{y} **for each** $x \in \mathbf{ux}$ **do** $\text{Count}_{x,l} := \mathbf{cx}_x$ $\Delta_{x,l} := \bar{\mathbf{y}} - y[\text{refcat}_l]$ **end****end** $\text{work} := 1 \dots |\text{uniq_refcats}|$ Let Avg_x be vector with running sum mapping category to count**while** $\text{len}(\text{work}) > 0$ **and** $\text{len}(\text{completed}) > 0$ **and** $\text{iteration} \leq \text{max_iter}$ **do****end**

References

- Jerome H. Friedman. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29:1189–1232, 2000.
- Alex Goldstein, Adam Kapelner, Justin Bleich, and Emil Pitkin. Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *Journal of Computational and Graphical Statistics*, 24(1):44–65, 2015. doi: 10.1080/10618600.2014.907095. URL <https://doi.org/10.1080/10618600.2014.907095>.
- Daniel W Apley. Visualizing the effects of predictor variables in black box supervised learning models. *arXiv preprint arXiv:1612.08468*, 2016.
- Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4765–4774. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>.