
OPTIMAL EXPLANATIONS OF LINEAR MODELS

A PREPRINT

Dimitris Bersimas

Sloan School of Management
Massachusetts Institute of Technology
Cambridge, MA 02139
dbertsim@mit.edu

Arthur Delarue

Operations Research Center
Massachusetts Institute of Technology
Cambridge, MA 02139
adelarue@mit.edu

Patrick Jaillet

Dep. of Electrical Engineering and Computer Science
Massachusetts Institute of Technology
Cambridge, MA 02139
jaillet@mit.edu

Sebastien Martin

Operations Research Center
Massachusetts Institute of Technology
Cambridge, MA 02139
92sebastien@gmail.com

July 11, 2019

ABSTRACT

When predictive models are used to support complex and important decisions, the ability to explain a model’s reasoning can increase trust, expose hidden biases, and reduce vulnerability to adversarial attacks. However, attempts at interpreting models are often ad hoc and application-specific, and the concept of interpretability itself is not well-defined. We propose a general optimization framework to create explanations for linear models. Our methodology decomposes a linear model into a sequence of models of increasing complexity using coordinate updates on the coefficients. Computing this decomposition optimally is a difficult optimization problem for which we propose exact algorithms and scalable heuristics. By solving this problem, we can derive a parametrized family of interpretability metrics for linear models that generalizes typical proxies, and study the tradeoff between interpretability and predictive accuracy.

1 Introduction

As machine learning models influence a growing fraction of everyday life, from cancer diagnoses to parole decisions to loan applications [1, 2, 3], individuals often want to understand the reasons for the decisions that affect them [4]. Model interpretability is of significant interest to the machine learning community [5], even though the lack of a well-defined concept of interpretability [6] means researchers often focus on proxies (e.g. sparsity or coefficient integrality in linear models).

Building “simple” models by optimizing interpretability proxies is a helpful but incomplete way to enhance interpretability. In some cases, practitioners prefer more complex models (e.g. deeper decision trees) because they can be explained and justified in a more compelling way [7]. Recent EU legislation [8] guarantees citizens a “right to explanation,” not a right to be affected only by sparse models. But explaining a model in simple terms is a major challenge in interpretable machine learning, because it is typically an ad hoc, audience-specific process. Formalizing the process of model explanation can yield ways to create models that are easier to explain, and rigorously quantify the tradeoff between interpretability and accuracy.

In this paper, we focus on linear models, and explore ways to decompose them into a sequence of interpretable coordinate updates. We propose a general optimization framework to measure and optimize the interpretability of these sequences. We then discuss how to create linear models with better explanations, leading to a natural set of interpretability metrics, and show that we can generalize various aspects of linear model interpretability. In particular,

- Section 2 introduces *coordinate paths* and motivates their use to explain linear models.
- Section 3 presents a set of metrics to evaluate the interpretability of coordinate paths and extends them into interpretability metrics for models. This allows us to study the *price of interpretability*, i.e., the Pareto front between accuracy and interpretability. We show that our metrics are consistent with existing approaches and exhibit desirable properties.
- Section 4 presents both optimal and scalable algorithms to compute coordinate paths and interpretable models.
- Section 5 discusses various practical uses of our framework and other extensions.

1.1 Related work

Many interpretable machine learning approaches involve optimizing some characteristics of the model as proxies for interpretability. Examples include sparsity for linear models [9], number of splits for decision trees [10], number of subspace features for case-based reasoning [11], or depth for rule lists [12, 13]. Some approaches optimize these proxies directly, while others fit auxiliary simple models to more complex black-box models [14, 15, 16, 17, 18, 19].

In the specific case of linear models, the typical interpretability proxy of sparsity (small number of nonzero coefficients) has been a topic of extensive study over the past twenty years [9]. Sparse regression models can be trained using heuristics such as LASSO [20], stagewise regression [21] or least-angle regression [22], or using scalable mixed-integer approaches [23]. More recently, another factor of interpretability in linear models has involved imposing integrality on the coefficients [24, 25], which allows to think of the output as tallying up points from each feature into a final score.

Training low-complexity models often affects predictive accuracy, and the tradeoff between the two can be difficult to quantify [26]. Similarly, the limitations of an *ex post* explanation relative to the original black box model can be difficult to explain to users [27]. And it is not clear that practitioners always find models that optimize these proxies more interpretable [7]. Recent landmark works [28, 6, 27] have argued that any study of interpretability must include input from human users. The framework we propose is both human-driven and mathematically rigorous, as users can define their own understanding of interpretability and quantify the resulting tradeoff with accuracy.

2 A Sequential View of Model Construction

Given a dataset with feature matrix $X \in \mathbb{R}^{n \times d}$ and labels $y \in \mathbb{R}^n$, a linear model is a vector of coefficients $\beta \in \mathbb{R}^d$, associated with a cost $c(\beta)$ that measures how well it fits the data, such as the mean-squared error $c(\beta) = (1/n)\|X\beta - y\|^2$ (potentially augmented with a regularization term for out-of-sample error).

We will motivate our approach to explaining linear models with a toy example. The goal is to predict a child’s age y_A given height X_H and weight X_W . The normalized features X_H and X_W have correlation $\rho = 0.9$ and are both positively correlated with the normalized target. Solving the ordinary least squares problem yields optimal coefficients $\beta^* = (2.12, -0.94)$:

$$y_A = 2.12 \cdot X_H - 0.94 \cdot X_W + \varepsilon, \quad (1)$$

with ε the error term. The mean squared error (MSE) of β^* is $c(\beta^*) = \mathbb{E}[\varepsilon^2] = 0.25$.

2.1 Coordinate paths

We propose a framework to construct an explanation of β^* by decomposing the model into a sequence of interpretable building blocks. In particular, we consider sequences of linear models leading to β^* where each model is obtained by changing one coefficient from the preceding model. We choose these *coordinate steps* because they correspond to the natural idea of adding a feature or updating an existing coefficient, and we will show they have interesting properties. We discuss other potential steps in Section 5. Table 1 shows 3 possible decompositions of β^* into coordinate steps. This is a natural way to decompose β^* , notice for example that decomposition 1a corresponds to introducing the model coefficient by coefficient.

We refer to these sequences of models as *coordinate paths*. Formally, we define a coordinate path of length K as a sequence of K models $\beta = (\beta_1, \dots, \beta_K)$ such that $\beta_k \in \mathcal{S}(\beta_{k-1})$ for all $1 \leq k \leq K$, where $\beta_0 = 0$, and $\mathcal{S}(\beta)$ is the set of linear models that are one coordinate step away from β , i.e., $\mathcal{S}(\beta) = \{\theta \in \mathbb{R}^d : \|\beta - \theta\|_0 \leq 1\}$. \mathcal{P}_K is the set of all coordinate paths of length K , and $\mathcal{P} = \cup_{K=1}^{\infty} \mathcal{P}_K$ is the set of all finite coordinate paths. An *explanation* of a model β is a coordinate path $\beta \in \mathcal{P}$ such that the last model is β . $\mathcal{P}_K(\beta)$ is the set of explanations of β of length K (potentially empty), and $\mathcal{P}(\beta)$ is the set of all possible explanations of β (typically infinite).

The examples in Table 1 are all explanations of β^* , and it is natural to ask which is the most useful or interpretable one. Formally speaking, it is of interest to define an *interpretability loss* $\mathcal{L}(\cdot)$ on the space of coordinate paths \mathcal{P} , such that

Table 1: Three decompositions of β^* into a sequence of coordinate steps.

$\beta = (\beta_H, \beta_W)$	$c(\beta)$	$\beta = (\beta_H, \beta_W)$	$c(\beta)$	$\beta = (\beta_H, \beta_W)$	$c(\beta)$
$\beta_0 = (0, 0)$	2.04	$\beta'_0 = (0, 0)$	2.04	$\hat{\beta}_0 = (0, 0)$	2.04
$\beta_1 = (2.12, 0)$	1.13	$\beta'_1 = (0, -0.94)$	4.74	$\hat{\beta}_1 = (1.70, 0)$	0.60
$\beta_2 = (2.12, -0.94)$	0.25	$\beta'_2 = (2.12, -0.94)$	0.25	$\hat{\beta}_2 = (1.70, -0.94)$	0.43
				$\hat{\beta}_3 = (2.12, -0.94)$	0.25

(a) Two steps, starting with X_H (b) Two steps, starting with X_W (c) Three steps, changing X_H twice

$\mathcal{L}(\beta) < \mathcal{L}(\beta')$ when β is more interpretable than β' . Then finding the best possible explanation for any model β can be written as the optimization problem

$$\min_{\beta \in \mathcal{P}(\beta)} \mathcal{L}(\beta).$$

For any path interpretability loss $\mathcal{L}(\cdot)$, it is then easy to consider the interpretability loss $\mathcal{L}(\beta)$ of a *model* β as the interpretability loss of the best explanation $\beta \in \mathcal{P}(\beta)$, i.e.

$$\mathcal{L}(\beta) = \begin{cases} \infty, & \text{if } \mathcal{P}(\beta) = \emptyset, \\ \min_{\beta \in \mathcal{P}(\beta)} \mathcal{L}(\beta), & \text{otherwise.} \end{cases} \quad (2)$$

How to select a path interpretability loss $\mathcal{L}(\beta)$? A natural choice is to consider that an explanation is better if it is shorter. Formally, we define the *path complexity* loss $\mathcal{L}_c(\beta) = |\beta|$, corresponding to the length of the coordinate path. For any model β , we can define the corresponding interpretability loss

$$\mathcal{L}_c(\beta) = \min_{\beta \in \mathcal{P}(\beta)} \mathcal{L}_c(\beta) = \min_{\beta \in \mathcal{P}(\beta)} |\beta|,$$

which we call *model complexity* (minimum number of coordinate steps required to reach β). Interestingly, for any model β , $\mathcal{L}_c(\beta)$ corresponds to the number of non-zero coefficients of β . The natural metric of coordinate path length thus recovers the usual interpretability proxy of model sparsity.

Consider the different coordinate paths in Table 1. If we use the interpretability loss \mathcal{L}_c , i.e., if we consider shorter paths (and thus sparser models) to be more interpretable, then β (path 1a) and β' (path 1b) are equally interpretable. Though both paths verify $c(\beta_2) = c(\beta'_2) = 0.25$, we notice that $c(\beta_1) = 1.13 < c(\beta'_1) = 4.74$. Indeed, β'_1 is a particularly inaccurate model, as weight is actually positively correlated with age. Since a coordinate path represents an explanation of the final model, the costs of intermediate models should play a role in quantifying the interpretability of a path; higher costs should be penalized. The path complexity loss $\mathcal{L}_c(\cdot)$ does not consider intermediate model costs at all, and therefore cannot capture this effect.

2.2 Incrementality

To explore alternatives to path complexity, consider the example of greedy coordinate paths. where the next model at each step is chosen by minimizing the cost $c(\cdot)$:

$$\beta_{k+1}^{\text{greedy}} \in \arg \min \left\{ c(\beta), \beta \in \mathcal{S}(\beta_k^{\text{greedy}}) \right\} \quad \forall k \geq 1. \quad (3)$$

This approach is appealing from an explanation standpoint, because we always select the coordinate step which most improves the model. However, many steps may be required to obtain an accurate model (slow convergence). Returning to toy example (1) and considering only paths of length 2, we compute the greedy coordinate path β^{greedy} by solving (3) twice [21]. Comparing to β from Table 1a, we have $c(\beta_1^{\text{greedy}}) = 0.42 < 1.13 = c(\beta_1)$, but $c(\beta_2^{\text{greedy}}) = 0.39 > 0.25 = c(\beta_2)$. The improvement of the first model comes at the expense of the second step.

Deciding which of the two paths β and β^{greedy} is more interpretable is a hard question. It highlights the tradeoff between the desirable incrementality of the greedy approach and the optimality of the second model. For paths of length 2, there is a continuum of solutions that trade off MSE in the first and second steps, shown in Figure 1. The next section introduces interpretability losses that formalize this tradeoff.

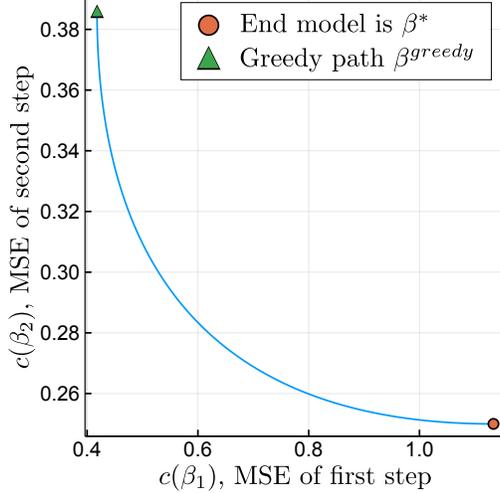


Figure 1: Tradeoff between the cost of the first and second models for a coordinate path of length 2 on problem (1).

3 Defining an Interpretability Loss

3.1 Coherent Interpretability Losses

In example (1), comparing the interpretability of β and β' is easy, because they have the same length and β has a better cost than β' at each step. In contrast, comparing the interpretability of β and β^{greedy} is not trivial, because β has a better final cost, but β^{greedy} has a better initial cost.

We can define the *cost sequence* of a coordinate path $\beta \in \mathcal{P}_K$ as the infinite sequence (c_1, c_2, \dots) such that $c_k = c(\beta_k)$ if $k \leq K$, and $c_k = 0$ otherwise. Then we call a path interpretability loss $\mathcal{L}(\cdot)$ *coherent* if the following conditions hold for any two paths $\beta, \beta' \in \mathcal{P}$ with cost sequences c and c' .

- (a) If $c = c'$, then $\mathcal{L}(\beta) = \mathcal{L}(\beta')$.
- (b) If $c_k \leq c'_k \forall k$, then $\mathcal{L}(\beta) \leq \mathcal{L}(\beta')$.

Condition (a) means that in our modeling framework, the interpretability of a path depends only on the sequence of costs along that path. Condition (b) formalizes the intuition that paths with fewer steps or better steps are more interpretable. Under any coherent interpretability loss $\mathcal{L}(\cdot)$ in toy example (1), β is more interpretable than β' , but β may be more or less interpretable than β^{greedy} depending on the specific choice of coherent interpretability loss.

In addition, consider a path $\beta \in \mathcal{P}_K$ and remove its last step to obtain a new path $\beta' \in \mathcal{P}_{K-1}$. This is equivalent to setting the K -th element of the cost sequence $c(\beta)$ to zero. Since $c(\cdot) \geq 0$, we have that $c(\beta') \leq c(\beta)$, which implies $\mathcal{L}(\beta') \leq \mathcal{L}(\beta)$. In other words, under a coherent interpretability loss, removing a step from a coordinate path can only make the path more interpretable. We also notice that the path complexity \mathcal{L}_c (sparsity) is a coherent path interpretability loss.

3.2 A Coherent Model Interpretability Loss

Condition (b) states that a path with at least as good a cost at each step as another path must be at least as interpretable. This notion of Pareto dominance suggests a natural path interpretability loss:

$$\mathcal{L}_\alpha(\beta) = \sum_{k=1}^{\infty} \alpha_k c(\beta)_k = \sum_{k=1}^{|\beta|} \alpha_k c(\beta_k).$$

In other words, the interpretability loss \mathcal{L}_α of a path β is the weighted sum of the costs of all steps in the path. This loss function is trivially coherent and extremely general. It is specified by the infinite sequence of parameters α , which specify the relative importance of the accuracy of each step in the model for the particular application at hand.

Defining a family of interpretability losses with infinitely many parameters allows for significant modeling flexibility, but it is also cumbersome and overly general. We therefore propose to select $\alpha_k = \gamma^k$ for all k , replacing the infinite

sequence of parameters $(\alpha_1, \alpha_2, \dots)$ with a single parameter $\gamma > 0$. In this case, following 2, we propose the following interpretability loss function on the space of models.

Definition 1 (Model interpretability). Given a model $\beta \in \mathbb{R}^d$, its interpretability loss $\mathcal{L}_\gamma(\beta)$ is given by

$$\mathcal{L}_\gamma(\beta) = \begin{cases} \infty, & \text{if } \mathcal{P}(\beta) = \emptyset, \\ \min_{\beta \in \mathcal{P}(\beta)} \mathcal{L}_\gamma(\beta) = \sum_{k=1}^{|\beta|} \gamma^k c(\beta_k), & \text{otherwise.} \end{cases} \quad (4)$$

By definition, \mathcal{L}_γ is a coherent interpretability loss. The parameter γ captures the tradeoff between favoring more incremental models or models with a low complexity, as formalized in Theorem 1.

Theorem 1 (Consistency of interpretability measure). *Assume that $c(\cdot)$ is bounded and nonnegative.*

(a) *Let $\beta^+, \beta^- \in \mathbb{R}^d$ with $\mathcal{L}_{\text{complexity}}(\beta^+) < \mathcal{L}_{\text{complexity}}(\beta^-)$, or $\mathcal{L}_{\text{complexity}}(\beta^+) = \mathcal{L}_{\text{complexity}}(\beta^-)$ and $c(\beta^+) < c(\beta^-)$.*

$$\lim_{\gamma \rightarrow \infty} \mathcal{L}_\gamma(\beta^-) - \mathcal{L}_\gamma(\beta^+) = +\infty. \quad (5)$$

(b) *Given models $\beta^+, \beta^- \in \mathbb{R}^d$, if there is $\beta^+ \in \mathcal{P}(\beta^+)$ such that $c(\beta^+) \leq c(\beta^-)$ for all $\beta^- \in \mathcal{P}(\beta^-)$, then*

$$\lim_{\gamma \rightarrow 0} \mathcal{L}_\gamma(\beta^-) - \mathcal{L}_\gamma(\beta^+) \geq 0. \quad (6)$$

Intuitively, in the limit $\gamma \rightarrow +\infty$, (a) states that the most interpretable models are the ones with minimal complexity, or minimal costs if their complexity is the same. (b) states that in the limit $\gamma \rightarrow 0$ the most interpretable models are the ones that can be constructed with greedy steps. All proofs are provided in the supplement.

3.3 The Price of Interpretability

Given the metric of interpretability defined above, we want to compute models that are Pareto-optimal with respect to $c(\cdot)$ and $\mathcal{L}_\gamma(\cdot)$ (more generally $\mathcal{L}_\alpha(\cdot)$). Computing these models can be challenging, as our definition of model interpretability requires to optimize over paths of any length. We can get around this if we can at least find the most interpretable path of a fixed length K , i.e.,

$$\min_{\beta \in \mathcal{P}_K} \mathcal{L}_\alpha(\beta) = \sum_{k=1}^K \alpha_k c(\beta_k) \quad (7)$$

Indeed, the following result shows that we can compute Pareto-optimal solutions by solving a sequence of optimization problems (7) for various K .

Proposition 1 (Price of interpretability). *Pareto-optimal models that minimize the interpretability loss \mathcal{L}_α and the cost $c(\cdot)$ can be computed by solving the following optimization problem:*

$$\min_{K \geq 0} \left(\min_{\beta \in \mathcal{P}_K} c(\beta_K) + \lambda \sum_{k=1}^K \alpha_k c(\beta_k) \right), \quad (8)$$

where $\lambda \in \mathcal{R}$ is a tradeoff parameter between cost and interpretability.

Notice that the inner minimization problem in (8) is simply problem (7) with appropriate modifications of the coefficients $(\alpha_1, \dots, \alpha_K)$. We can use this decomposition to compute the price of interpretability in the toy problem (1), with the interpretability loss \mathcal{L}_γ chosen such that $\gamma = 1$. Figure 2 shows all Pareto-optimal models with respect to cost (MSE) and interpretability loss.

By defining the general framework of coordinate paths and a natural family of coherent interpretability loss functions, we can understand exactly how much we gain or lose in terms of accuracy when we choose a more or less interpretable model. Our framework thus provides a principled way to answer a central question of the growing literature on interpretability in machine learning.

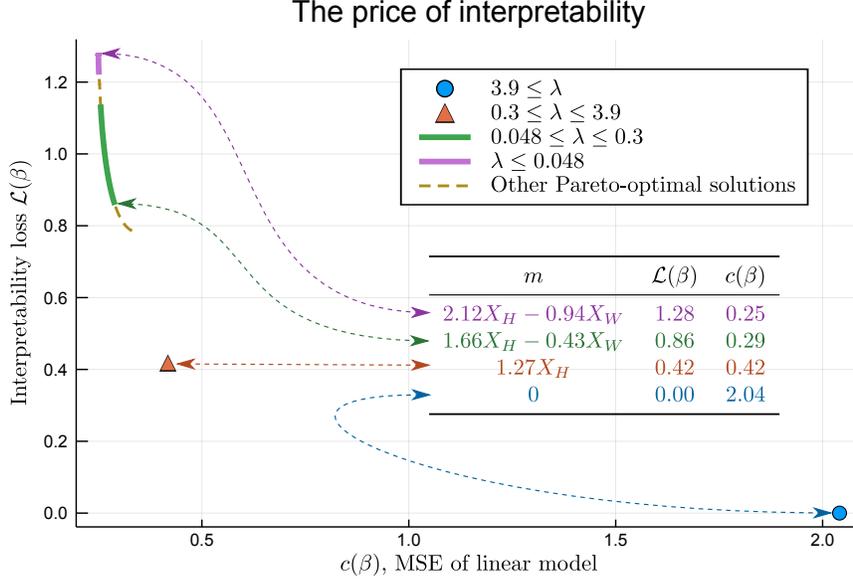


Figure 2: Pareto front between interpretability loss $\mathcal{L}(\beta) = \mathcal{L}_\gamma(\beta)$ (with $\gamma = 1$) and cost $c(\beta)$ on the toy OLS problem (1), computed by varying λ in (8). The dashed line represents Pareto-optimal solutions that cannot be computed by this weighted-sum method [29]. Note that the front is discontinuous. The number of steps in the corresponding optimal coordinate paths is respectively 0, 1, 2, 3 for the blue, orange, green and purple segments. The inset table describes several interesting Pareto-optimal models.

4 Computing the Price of Interpretability

4.1 Algorithms

Optimal. Given the step function $\mathcal{S}(\cdot)$ and the convex quadratic cost function $c(\cdot)$, problem (7) can be written as a convex integer optimization problem using special ordered sets of type 1 (SOS-1 constraints), and solved using Gurobi or CPLEX for small problems:

$$\min \sum_{k=1}^K c(\beta_k) \quad \text{s.t.} \quad \text{SOS-1}(\beta_{k+1} - \beta_k) \quad \forall 0 \leq k < K, \quad (9)$$

where β_0 designates the starting linear model.

Local improvement. In higher-dimensional settings, or when K grows large, the formulation above may no longer scale. Thus it is of interest to develop a fast heuristic for such instances.

A feasible solution β to problem (9) can be written as a vector of indices $i \in \{1, \dots, d\}^K$ and a vector of values $\delta \in \mathbb{R}^K$, such that for $0 \leq k < K$,

$$\beta_{k+1}^i = \begin{cases} \beta_k^i + \delta_k, & \text{if } i = i_k \\ \beta_k^i, & \text{if } i \neq i_k. \end{cases}$$

The vector of indices i encodes which coefficients are modified at each step, while the vector of values δ encodes the value of each modified coefficient. Thus problem (9) can be rewritten as

$$\min_i \min_\delta C(i, \delta) := \sum_{k=1}^K c \left(\beta_0 + \sum_{j=1}^k \delta_j e_{i_j} \right), \quad (10)$$

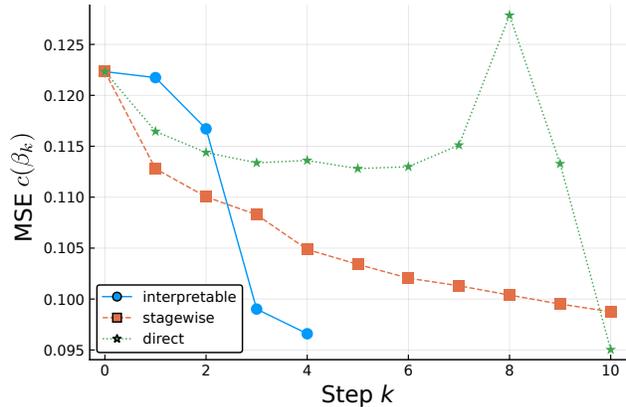
where e_i designates the i -th unit vector. The inner minimization problem is an “easy” convex quadratic optimization problem, while the outer minimization problem is a “hard” combinatorial optimization problem. We propose the following local improvement heuristic for the outer problem: given a current vector of indices i , we randomly sample one step κ in the coordinate path. Keeping all i_k constant for $k \neq \kappa$, we iterate through all d possible values of i_κ and

Table 2: Convergence of local improvement heuristics for different batch sizes q .

Method	Time to convergence (s)	Optimality gap (%)
Exact	5.078	0.00
Local improvement ($q = 1$)	0.004	0.02
Local improvement ($q = 2$)	0.019	0.00

	Feature				MSE
	MealPct	AvgInc	ELPct	ExpnStu	
β_0	-0.87	-	-	-	0.122
β_1		0.23	-	-	0.122
β_2	-0.59		-	-	0.117
β_3			-0.18	-	0.099
β_4				0.07	0.097

(a) Coordinate path from old model to new model. MealPct is the percentage of students qualifying for reduced-price lunch, AvgInc is the average income, ELPct is the percentage of English Learners, ExpnStu is the expenditure per student.



(b) Comparison between coordinate path and other approaches

Figure 3: Example of a Pareto-efficient coordinate path. On the left we see the benefits of each coefficient modification. On the right we compare the coordinate path with the forward stagewise path which greedily selects the best β_{k+1} given β_k , and the “direct” path, which adds the optimal least squares coefficients one by one. The direct method is only good when all the coefficients have been added, whereas the greedy approach is good at first but then does not converge. The coordinate path is willing to make some suboptimal steps in preparation for very cost-improving steps.

obtain d candidate vectors \hat{i} . For each candidate, we solve the inner minimization problem and keep the one with lowest cost. A general version of this algorithm, where we sample not one but q steps at each iteration, is provided in the supplement.

4.2 Results

Optimal vs heuristic. In order to empirically evaluate the local improvement heuristic, we run it with different batch sizes q on a small real dataset, with 100 rows and 6 features (after one-hot encoding of categorical features). The goal is to predict the perceived prestige (from a survey) of a job occupation given features about it, including education level, salary, etc.

Given this dataset, we first compute the optimal coordinate path of length $K = 10$. We then test our local improvement heuristic on the same dataset. Given the small size of the problem, in the complete formulation a provable global optimum is found by Gurobi in about 5 seconds. To be useful, we would like our local improvement heuristic to find a good solution significantly faster. We show convergence results of the heuristic for different values of the batch size parameter q in Table 2. For both batch sizes, the local improvement heuristic converges two orders of magnitude faster than Gurobi. With a batch size $q = 2$, the solution found is optimal.

Insights from a real dataset. We now explore the results of the presented approach on a dataset from the 1998-1999 California test score dataset. Each data point represents a school, and the variable of interest is the average standardized test score of students from that school. The ten continuous features and the target variables are centered and rescaled to have unit variance.

In our example, we assume that we already have a regression model available to predict the number of trips: it was trained using only the percentage of students qualifying for a reduced-price lunch. This model has an MSE of 0.122 (compared to an optimal MSE of 0.095). We would like to update this model in an interpretable way given the availability of all features in the dataset. In our framework, this corresponds to problem (9) where β_0 is no longer 0 but the available starting model.

	Season:			Weekday:		Weather:		MSE
	Ftemp	Day	Hum	Summer	Wind	Sunday	Clear	
β_1	0.56							0.293
β_2		0.55						0.146
β_3			-0.20					0.129
β_4				0.15				0.119
β_5					-0.14			0.110
β_6						-0.07		0.108
β_7							0.06	0.107

(a) Interpretable path for the bike-sharing dataset.

	Season:		Weather:		Season:			MSE
	Ftemp	Day	Spring	Clear	Hum	Summer	Wind	
β_1	0.65							0.289
β_2	0.55	0.54						0.146
β_3	0.50	0.52	-0.10					0.143
β_4	0.47	0.52	-0.12	0.18				0.128
β_5	0.49	0.51	-0.13	0.09	-0.15			0.121
β_6	0.50	0.55	-0.07	0.10	-0.14	0.13		0.114
β_7	0.49	0.54	-0.06	0.08	-0.18	0.14	-0.14	0.107

(b) Sequence of models with an increasing number of features obtained via LASSO.

Figure 4: Comparison of interpretable path and sequence of models of increasing sparsity selected by LASSO. Ftemp is the “feels like” temperature, Day is the number of days since data collection began, Hum is humidity, Wind is wind speed. Season, Weather and Weekday are categorical variables.

In Figure 3 we study one particular coordinate path on the Pareto front of interpretability and efficiency. The path (and the new model it leads to) is shown in Figure 3a. It can be obtained from the original model in just four steps. First we add the district average income with a positive coefficient, then we correct the coefficient for reduced-price lunch students to account for this new feature, and finally we add the percentage of English learners and the per-student spending. The final model has an MSE of 0.097 which is near-optimal. When we compare this path to other methods (see Figure 3b) we see that our interpretable formulation allows us to find a good tradeoff between a greedy formulation and a formulation that just sets the coefficients to their final values one by one.

5 A General Framework

5.1 Different Steps for Different Notions of Interpretability

So far, we have focused exclusively on paths in which linear models are constructed in a series of coordinate steps. However, the choice of what constitutes a step is ultimately a modeling choice which encodes what a user in a particular application considers a simple building block. Choosing a different step function $\mathcal{S}(\cdot)$ can lead to other notions of interpretability. For instance, choosing $\mathcal{S}_{\text{SLIM}}(\beta) = \{\theta : \|\beta - \theta\|_0 \leq 1, \|\beta - \theta\|_1 \in \mathbb{Z}\}$ imposes integer coordinate updates at each step. This is related to the notion of interpretability introduced by the score-based methods of Ustun and Rudin [25]. Another way to think about score-based methods is to choose $\mathcal{S}'_{\text{SLIM}}(\beta) = \{\theta : \|\beta - \theta\|_0 \leq 1, \|\beta - \theta\|_1 \in \{0, 1\}\}$, which imposes that each step adds one point to the scoring system. The fundamental idea of optimally decomposing models into a sequence of simple building blocks is general, and can be applied not only to more general linear models (e.g. ridge regression or logistic regression by suitably modifying the cost c) but also to other machine learning models in general (for example, a decision tree can be decomposed into a sequence of successive splits).

5.2 Human-in-the-loop Model Selection

Viewing a coordinate path as a nested sequence of models of increasing complexity can be useful in the context of human-in-the-loop analytics. Consider the problem of selecting a linear model by a human decision maker. For example, consider a city planner that would like to understand bike-sharing usage in Porto, by training a linear model using a dataset from the UCI ML repository [30], where each of the 731 data points represents a particular day (18 features about weather, time of year, etc.), and the variable of interest is the number of trips recorded by the bike-sharing system on that day. The decision-maker may prefer a sparse model, but may not know the exact desired level of sparsity. Given a discrete distribution on the choice of the level of sparsity $K : p_k = \mathbb{P}(\{k \text{ will be chosen by the decision maker}\})$, we can choose $\alpha_k = p_k$ and solve (7) to find paths β that minimize the expected cost $\mathbb{E}_k[c(\beta_k)]$.

In Table 4a, we show the path β obtained by assuming that the desired level of sparsity is uniformly random between 1 and 7. We can compare the result to a sequence of linear models obtained using LASSO to select an increasing set of features, shown in Table 4b. The expected costs are respectively 0.145 and 0.150, and the MSE is essentially the same as each step. The two sequences also use almost the same features, in a similar order. However, the coordinate path can be read much more easily: because only one coefficient changes at each step, the whole path can be described with $\Theta(K)$ parameters, while the path constructed using LASSO/sparse regression needs $\Theta(K^2)$ parameters.

Acknowledgements

Research funded in part by ONR grant N00014-18-1-2122.

References

- [1] Sendhil Mullainathan and Ziad Obermeyer. Does machine learning automate moral hazard and error? *American Economic Review*, 107(5):476–480, 2017.
- [2] Jon Kleinberg, Himabindu Lakkaraju, Jure Leskovec, Jens Ludwig, and Sendhil Mullainathan. Human decisions and machine predictions. *The quarterly journal of economics*, 133(1):237–293, 2017.
- [3] Richard Berk. An impact assessment of machine learning risk forecasts on parole board decisions and recidivism. *Journal of Experimental Criminology*, 13(2):193–216, 2017.
- [4] Berkeley J Dietvorst, Joseph P Simmons, and Cade Massey. Overcoming algorithm aversion: People will use imperfect algorithms if they can (even slightly) modify them. *Management Science*, 64(3):1155–1170, nov 2016.
- [5] Alex A. Freitas. Comprehensible classification models. *ACM SIGKDD Explorations Newsletter*, 15(1):1–10, 2014.
- [6] Zachary C. Lipton. The Mythos of Model Interpretability. *arXiv preprint arXiv:1606.03490*, 2016.
- [7] Nada Lavrač. Selected techniques for data mining in medicine. *Artificial intelligence in medicine*, 16(1):3–23, 1999.
- [8] Bryce Goodman and Seth Flaxman. European Union regulations on algorithmic decision-making and a "right to explanation". pages 1–9, 2016.
- [9] Trevor Hastie, Robert Tibshirani, and Martin Wainwright. *Statistical learning with sparsity: the lasso and generalizations*. CRC press, 2015.
- [10] Leo Breiman. *Classification and regression trees*. New York: Routledge, 1984.
- [11] Been Kim, Cynthia Rudin, and Julie Shah. The Bayesian Case Model: A Generative Approach for Case-Based Reasoning and Prototype Classification. In *Neural Information Processing Systems (NIPS) 2014*, 2014.
- [12] Benjamin Letham, Cynthia Rudin, Tyler H. McCormick, and David Madigan. Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model. *Annals of Applied Statistics*, 9(3):1350–1371, 2015.
- [13] Hongyu Yang, Cynthia Rudin, and Margo Seltzer. Scalable Bayesian Rule Lists. In *Proceedings of the 34th International Conference on Machine Learning*, 2017.
- [14] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “Why Should I Trust You?” Explaining the Predictions of Any Classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.
- [15] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*. Springer series in statistics New York, NY, USA:, 2001.
- [16] Anupam Datta, Shayak Sen, and Yair Zick. Algorithmic Transparency via Quantitative Input Influence :. In *2016 IEEE Symposium on Security and Privacy*, 2016.
- [17] Hamsa Bastani, Osbert Bastani, and Carolyn Kim. Interpreting Predictive Models for Human-in-the-Loop Analytics. *arXiv preprint arXiv:1705.08504*, pages 1–45, 2018.
- [18] Himabindu Lakkaraju, Ece Kamar, Rich Caruana, and Jure Leskovec. Interpretable & Explorable Approximations of Black Box Models. *FAT/ML*, jul 2017.
- [19] Cristian Bucilă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '06*, page 535, New York, New York, USA, 2006. ACM, ACM Press.
- [20] Robert J. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
- [21] Jonathan Taylor and Robert J. Tibshirani. Statistical learning and selective inference. *Proceedings of the National Academy of Sciences*, 112(25):7629–7634, jun 2015.
- [22] Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Least Angle Regression. *Annals of Statistics*, 32(2):407–499, apr 2004.
- [23] Dimitris Bertsimas, Angela King, and Rahul Mazumder. Best subset selection via a modern optimization lens. *Annals of Statistics*, 44(2):813–852, 2016.
- [24] Jongbin Jung, Connor Concannon, Ravi Shroff, Sharad Goel, and Daniel G Goldstein. Simple rules for complex decisions. feb 2017.
- [25] Berk Ustun and Cynthia Rudin. Supersparse linear integer models for optimized medical scoring systems. *Machine Learning*, 102(3):349–391, 2016.

- [26] Leo Breiman. Statistical modeling: The two cultures. *Statistical science*, 16(3):199–231, 2001.
- [27] Leilani H Gilpin, David Bau, Ben Z Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal. Explaining Explanations : An Approach to Evaluating Interpretability of Machine Learning. *arXiv preprint arXiv:1806.00069*, 2018.
- [28] Finale Doshi-Velez and Been Kim. Towards A Rigorous Science of Interpretable Machine Learning. *arXiv preprint arXiv:1702.08608*, (MI):1–13, 2017.
- [29] I. Y. Kim and O. L. De Weck. Adaptive weighted-sum method for bi-objective optimization: Pareto front generation. *Structural and Multidisciplinary Optimization*, 29(2):149–158, 2005.
- [30] Hadi Fanaee-T and Joao Gama. Event labeling combining ensemble detectors and background knowledge. *Progress in Artificial Intelligence*, 2(2-3):113–127, 2014.

A Proof of Theorem 1

Proof of part (a). As $c(\cdot)$ is bounded, we have $c_{\max} \in \mathbb{R}$ such that $0 < c(\cdot) \leq c_{\max}$.

Let $\mathbf{m}^+ \in \mathcal{P}(m^+)$ be a path of optimal length to the model m^+ , i.e., $|\mathbf{m}^+| = \mathcal{L}_c(m^+)$. Let $\mathbf{m}^- \in \mathcal{P}(m^-)$ be any path leading to m^- (not necessarily of optimal length). By assumption, we have $|\mathbf{m}^-| \geq |\mathbf{m}^+|$, and by definition of model interpretability, we have $\mathcal{L}_\gamma(m^+) \leq \mathcal{L}_\gamma(\mathbf{m}^+)$. Therefore we obtain:

$$\mathcal{L}_\gamma(\mathbf{m}^-) - \mathcal{L}_\gamma(m^+) \geq \mathcal{L}_\gamma(\mathbf{m}^-) - \mathcal{L}_\gamma(\mathbf{m}^+) \quad (11)$$

$$= \sum_{k=1}^{|\mathbf{m}^-|} \gamma^k c(m_k^-) - \sum_{k=1}^{|\mathbf{m}^+|} \gamma^k c(m_k^+) \quad (12)$$

$$= \gamma^{|\mathbf{m}^+|} \left(\sum_{k=1}^{|\mathbf{m}^+|-1} \frac{1}{\gamma^{|\mathbf{m}^+|-k}} (c(m_k^-) - c(m_k^+)) + (c(m_{|\mathbf{m}^+|}^-) - c(m_{|\mathbf{m}^+|}^+)) + \sum_{k=|\mathbf{m}^+|+1}^{|\mathbf{m}^-|} \gamma^{k-|\mathbf{m}^+|} c(m_k^-) \right) \quad (13)$$

$$\geq \gamma^{|\mathbf{m}^+|} \left(-c_{\max} \sum_{k=1}^{|\mathbf{m}^+|-1} \frac{1}{\gamma^{|\mathbf{m}^+|-k}} + (c(m_{|\mathbf{m}^+|}^-) - c(m^+)) + \sum_{k=|\mathbf{m}^+|+1}^{|\mathbf{m}^-|} \gamma^{k-|\mathbf{m}^+|} c(m_k^-) \right), \quad (14)$$

where (12) follows from the definition of model interpretability, (13) is just a development of the previous equation, and (14) just bounds the first sum and uses $m_{|\mathbf{m}^+|}^+ = m^+$ for the middle term.

If $\mathcal{L}_c(m^+) < \mathcal{L}_c(m^-)$, we have $|\mathbf{m}^+| < |\mathbf{m}^-|$, and therefore the last sum in (14) is not empty and for $\gamma \geq 1$ we can bound it:

$$\sum_{k=|\mathbf{m}^+|+1}^{|\mathbf{m}^-|} \gamma^{k-|\mathbf{m}^+|} c(m_k^-) \geq \gamma^{|\mathbf{m}^-|-|\mathbf{m}^+|} c(m_{|\mathbf{m}^-|}^-) \geq \gamma c(m^-). \quad (15)$$

Therefore, for $\gamma \geq 1$ we have:

$$\mathcal{L}_\gamma(\mathbf{m}^-) - \mathcal{L}_\gamma(m^+) \geq \gamma^{|\mathbf{m}^+|} \left(\gamma c(m^-) - c(m^+) - c_{\max} \sum_{k=1}^{|\mathbf{m}^+|-1} \frac{1}{\gamma^{|\mathbf{m}^+|-k}} \right). \quad (16)$$

This bound is valid for all the path \mathbf{m}^- leading to m^- , in particular the one with optimal interpretability loss, therefore we have (for $\gamma \geq 1$):

$$\mathcal{L}_\gamma(m^-) - \mathcal{L}_\gamma(m^+) \geq \gamma^{|\mathbf{m}^+|} \left(\gamma c(m^-) - c(m^+) - c_{\max} \sum_{k=1}^{|\mathbf{m}^+|-1} \frac{1}{\gamma^{|\mathbf{m}^+|-k}} \right). \quad (17)$$

which implies (as $c(m^-) > 0$):

$$\lim_{\gamma \rightarrow +\infty} \mathcal{L}_\gamma(m^-) - \mathcal{L}_\gamma(m^+) = +\infty \quad (18)$$

We now look at the case $\mathcal{L}_c(m^+) = \mathcal{L}_c(m^-)$ and $c(m^+) < c(m^-)$. For $\gamma \geq 1$, we can easily bound parts of equation (14):

$$c(m_{|\mathbf{m}^+|}^-) + \sum_{k=|\mathbf{m}^+|+1}^{|\mathbf{m}^-|} \gamma^{k-|\mathbf{m}^+|} c(m_k^-) \geq \gamma^{|\mathbf{m}^-|-|\mathbf{m}^+|} c(m_{|\mathbf{m}^-|}^-) \geq c(m^-). \quad (19)$$

Putting it back into (14), we obtain (for $\gamma \geq 1$)

$$\mathcal{L}_\gamma(\mathbf{m}^-) - \mathcal{L}_\gamma(m^+) \geq \gamma^{|\mathbf{m}^+|} \left((c(m^-) - c(m^+)) - c_{\max} \sum_{k=1}^{|\mathbf{m}^+|-1} \frac{1}{\gamma^{|\mathbf{m}^+|-k}} \right). \quad (20)$$

This bound is independent of the path \mathbf{m}^- leading to m^- , therefore we have

$$\mathcal{L}_\gamma(\mathbf{m}^-) - \mathcal{L}_\gamma(\mathbf{m}^+) \geq \gamma^{|\mathbf{m}^+|} \left((c(\mathbf{m}^-) - c(\mathbf{m}^+)) - c_{max} \sum_{k=1}^{|\mathbf{m}^+|-1} \frac{1}{\gamma^{|\mathbf{m}^+|-k}} \right) \xrightarrow{\gamma \rightarrow +\infty} +\infty, \quad (21)$$

which ends the proof. \square

Proof of part (b). Consider two paths $\mathbf{m}^+, \mathbf{m}^- \in \mathcal{P}$, such that $\mathbf{c}(\mathbf{m}^+) \preceq \mathbf{c}(\mathbf{m}^-)$. By definition of the lexicographic order, either the two paths are the same (in that case the theorem is trivial), or there exist $K \geq 1$ such that:

$$\begin{cases} \mathbf{c}(\mathbf{m}^+)_k = \mathbf{c}(\mathbf{m}^-)_k & \forall k < K \\ \mathbf{c}(\mathbf{m}^+)_K < \mathbf{c}(\mathbf{m}^-)_K. \end{cases}$$

We have:

$$\mathcal{L}_\gamma(\mathbf{m}^-) - \mathcal{L}_\gamma(\mathbf{m}^+) = \sum_{k=1}^{|\mathbf{m}^-|} \gamma^k c(\mathbf{m}^-)_k - \sum_{k=1}^{|\mathbf{m}^+|} \gamma^k c(\mathbf{m}^+)_k \quad (22)$$

$$= \sum_{k=1}^{\infty} \gamma^k \left(\mathbf{c}(\mathbf{m}^-)_k - \mathbf{c}(\mathbf{m}^+)_k \right) \quad (23)$$

$$= \sum_{k=1}^{K-1} \gamma^k \left(\mathbf{c}(\mathbf{m}^-)_k - \mathbf{c}(\mathbf{m}^+)_k \right) + \gamma^K \left(\mathbf{c}(\mathbf{m}^-)_K - \mathbf{c}(\mathbf{m}^+)_K \right) + \sum_{k=K+1}^{\infty} \gamma^k \left(\mathbf{c}(\mathbf{m}^-)_k - \mathbf{c}(\mathbf{m}^+)_k \right) \quad (24)$$

$$= \gamma^K \left(\mathbf{c}(\mathbf{m}^-)_K - \mathbf{c}(\mathbf{m}^+)_K + \sum_{k=K+1}^{\infty} \gamma^{k-K} \left(\mathbf{c}(\mathbf{m}^-)_k - \mathbf{c}(\mathbf{m}^+)_k \right) \right), \quad (25)$$

where (23) just applies the definition of the sequence \mathbf{c} , and (25) uses $\mathbf{c}(\mathbf{m}^+)_k = \mathbf{c}(\mathbf{m}^-)_k \quad \forall k < K$.

The term inside the parenthesis in (25) converges to $\mathbf{c}(\mathbf{m}^-)_K - \mathbf{c}(\mathbf{m}^+)_K > 0$ when $\gamma \rightarrow 0$, as the paths are finite. Therefore

$$\lim_{\gamma \rightarrow 0} \mathcal{L}_\gamma(\mathbf{m}^-) - \mathcal{L}_\gamma(\mathbf{m}^+) \geq 0, \quad (26)$$

which completes the proof. The very end of the theorem is an immediate consequence. \square

B Proof of Proposition 1

Proof. First, a solution of

$$\min_{\mathbf{m} \in \mathcal{M}} (c(\mathbf{m}) + \lambda \mathcal{L}(\mathbf{m}))$$

is Pareto optimal between the cost $c(\cdot)$ and the interpretability $\mathcal{L}_\alpha(\cdot)$ as it corresponds to the minimization of a weighted sum of the objectives. Furthermore, we can write

$$\begin{aligned} \min_{\mathbf{m} \in \mathcal{M}} (c(\mathbf{m}) + \lambda \mathcal{L}_\alpha(\mathbf{m})) &= \min_{\mathbf{m} \in \mathcal{M}} \left(c(\mathbf{m}) + \lambda \min_{\mathbf{m} \in \mathcal{P}(\mathbf{m})} \mathcal{L}_\alpha(\mathbf{m}) \right) = \min_{\mathbf{m} \in \mathcal{M}, \mathbf{m} \in \mathcal{P}(\mathbf{m})} (c(\mathbf{m}) + \lambda \mathcal{L}_\alpha(\mathbf{m})) \\ &= \min_{\mathbf{m} \in \mathcal{M}, K \geq 0, \mathbf{m} \in \mathcal{P}_K(\mathbf{m})} \left(c(\mathbf{m}_K) + \lambda \sum_{k=1}^K \alpha_k c(\mathbf{m}_k) \right) \\ &= \min_{K \geq 0, \mathbf{m} \in \mathcal{P}_K} \left(c(\mathbf{m}_K) + \lambda \sum_{k=1}^K \alpha_k c(\mathbf{m}_k) \right). \end{aligned}$$

\square

C Local improvement algorithm

Algorithm 1 Local improvement heuristic. Inputs: regression cost function $c(\cdot)$; starting vector of indices i^0 . Parameters: $q \in \mathbb{N}$ controls the size of the batch, $T \in \mathbb{N}$ controls the number of iterations.

```

1: function LOCALIMPROVEMENT( $c(\cdot)$ ,  $i^0$ ,  $q$ ,  $T$ )
2:   for  $1 \leq t \leq T$  do
3:      $i^* \leftarrow i^0$ 
4:      $\delta^* \leftarrow \arg \min_{\delta} C(i^0, \delta)$ 
5:      $C^* \leftarrow C(i^0, \delta^*)$ 
6:     Randomly select  $\mathcal{K} = \{\kappa_1, \dots, \kappa_q\} \subset \{1, \dots, K\}$  ▷ subset of cardinality  $q$ 
7:      $\hat{i} \leftarrow i^*$ 
8:      $\hat{\delta} \leftarrow \delta^*$ 
9:     for  $(f_1, \dots, f_q) \in \{1, \dots, d\}^q$  do
10:      for  $1 \leq p \leq q$  do
11:         $\hat{i}_{\kappa_p} = f_p$ 
12:         $\hat{\delta} \leftarrow \arg \min_{\delta} C(\hat{i}, \delta)$ 
13:        if  $C(\hat{i}, \hat{\delta}) < C^*$  then
14:           $C^* \leftarrow C(\hat{i}, \hat{\delta})$ 
15:           $i^* \leftarrow \hat{i}$ 
16:           $\delta^* \leftarrow \hat{\delta}$ 
17:   return  $i^*, \delta^*$ 

```
