

Debugging Machine Learning Tasks arXiv:1603.07292v1 [cs.LG] 23 Mar 2016 Aleksandar Chakraborty, Aditya Nori, Sriram Rajamani, Shayak Sen, and Deepak Vijaykumar
A University of Colorado, Boulder B Microsoft Research C Carnegie Mellon University D IBM Research March 24, 2016 Abstract Unlike traditional programs (such as operating systems or word processors) which have large amounts of code, machine learning tasks use programs with relatively small amounts of code (written in machine learning libraries), but voluminous amounts of data. Just like developers of traditional programs debug errors in their code, developers of machine learning tasks debug and fix errors in their data. However, algorithms and tools for debugging and fixing errors in data are less common, when compared to their counterparts for detecting and fixing errors in code. In this paper, we consider classification tasks where errors in training data lead to misclassifications in test points, and propose an automated method to find the root causes of such misclassifications. Our root cause analysis is based on Pearl’s theory of causation, and uses Pearl’s PS (Probability of Sufficiency) as a scoring metric. Our implementation, Psi, encodes the computation of PS as a probabilistic program, and uses recent work on probabilistic programs and transformations on probabilistic programs (along with gray-box models of machine learning algorithms) to efficiently compute PS. Psi is able to identify root causes of data errors in interesting data sets.

1 Introduction

Machine learning techniques are used to perform data-driven decision-making in a large number of diverse areas including image processing, medical diagnosis, credit decisions, insurance decisions, email spam detection, speech recognition, natural language processing, robotics, information retrieval and online advertising. Over time, these techniques have been honed and tuned, and are now at a stage where machine learning libraries [1, 2] are used as black-boxes by programmers with little or no expertise in the details of the machine learning algorithms themselves. The black-box nature of the reuse, however, has an unfortunate downside. Current implementations of machine learning techniques provide little insight into why a particular decision was made. Because of this absence of transparency, debugging the outputs of a machine learning algorithm has become incredibly hard. Most programmers who implement machine learning use libraries to build models from voluminous training data, and then use these models to perform predictions. These machine learning libraries often employ complex, stochastic, or approximate, search and optimization algorithms that search for an optimal model for a given training data set. The model is then applied to a set of unseen test samples in the hope of satisfactory generalization. When generalization fails, i.e., an incorrect result is produced for a test input, it is often difficult to debug the cause of the failure. Such failures can arise due to several reasons. Common causes for failure include bugs in the implementation of the machine learning algorithm, incorrect choice of features, incorrect setting of parameters (such as degree of the polynomial for regression or number of layers in a neural network) when invoking the machine learning library, and noise in the training set. Over time, bugs in implementation of machine learning algorithms get detected and fixed. There is a lot of work in feature selection [3], and parameter choices can be made by systematically building models for various parameter values and choosing the model with the best validation score [4]. However, since training data is typically voluminous, errors in training data are common and notoriously difficult to debug. This suggests a new class of debugging problems where programs (machine learning classifiers) are learnt from data and bugs in a program are now the result of faults in the data. In this paper, we focus on debugging machine learning tasks in the presence of errors in training data. Specifically, we consider classification tasks, which are typically implemented using algorithms such as logistic regression [5] and boosted decision trees [6]. Suppose we train a classifier on training data (which has errors), and the classifier produces incorrect results for one or more test points. We desire to produce an automated procedure to identify the root cause of this failure. That is, we would like to identify a subset of training points that influences the classification for these test points the most. Therefore, correcting mistakes in these training points is most likely to fix the incorrect results. Our algorithm for identifying root causes is inspired by the structural equations framework of causation, as formulated by Judea Pearl [7, 8]. We think of each of the training data points as possible causes of the misclassification in the test data set, and calculate for each such training point, a score corresponding to how likely it is that the current label for that point is the cause for the misclassification of the test data set. A simple measure of the score of a training point can be obtained by merely flipping the label of the training point and observing if the flip improves the results of the classifier on test points. However, such a simple measure does not work when errors exist in several training points, and several training points together cause the incorrect results in the test points. Thus, the score we calculate for each training point t considers 2^k alternate counterfactual worlds, where training points are labeled with several possible values (other than the value in the training data), and sums up the probability that flipping the label of t causes the misclassification error in the test data, among all such alternate worlds. In Pearl’s framework, such a score is called the probability of sufficiency or PS for short. One of the main difficulties in calculating the probability of sufficiency is that the classifier (or model) needs to be relearned for alternate worlds. Each of these model computing steps (also called as training steps) is expensive. We use a gray-box view of the machine learning library, and profile key intermediate values (that are hand-picked for each machine learning algorithm) during the initial training phase. Using these values, we build a gray-box abstraction of the training process by which the model for a new training set (which is obtained by flipping certain number of training labels) can be obtained efficiently without the need to perform complete (and expensive) retraining. Finally, we are able to amortize the cost of computing the PS score by sharing common work across the computation for different training points. In order to carry out these optimizations, we model the PS computation as a probabilistic program [9]. Probabilistic programs allow us to represent all of the above optimizations such as using gray-box models, using instrumented values from actual training runs, and sharing work across multiple PS computations as program transformations. We are also able to leverage recent progress in efficient inference of probabilistic programs to scale the computation of PS scores to large data sets. We have implemented our root cause detection algorithm in a tool Psi. Psi currently works with two popular classifiers: (1) logistic regression, and (2) boosted decision trees. For these classifiers, Psi runs a production quality implementation of the techniques, profiles specific values and builds an abstract gray-box model of the classifier, which avoids expensive re-training. Armed with this gray-box model, Psi performs scalable inference to compute the PS values for all points in the training set. Psi is able to identify root causes of misclassifications in several interesting data sets. In summary, the main contributions of this paper are as follows: We propose using the structural equations framework of causality, and specifically Pearl’s PS score to compute root causes of failures in machine learning algorithms. We model the PS computation as a probabilistic program, and this enables us to leverage efficient techniques developed to perform inference on probabilistic programs to calculate PS scores. We build gray-box models of the machine learning techniques by profiling actual training runs of the library, and using profiled values to build abstract models of the training process. We amortize work across PS computations of different training points. Probabilistic programs allow us to carry out these optimizations and reason about them as program transformations. We have built a tool Psi implementing the approach for logistic regression.

3 Training Phase

Training Set ML Training Algorithm **Test Set** Classifier Class Labels $\{1, 1\}$

Evaluation Phase

Figure 1: A two-stage design flow of a machine learning task: training phase in which the ML algorithm A is applied to training set S to learn classifier h , and evaluation phase to judge the quality of h on test set T . and boosted decision trees. Psi is able to identify root causes of misclassifications in several interesting data sets. We hypothesize that this approach can be generalized to other machine learning tasks as well.

2 Overview

We motivate our approach through the experience of Alice, a typical developer who uses machine learning.

2.1 Typical Scenario

Alice is not a machine learning expert, but needs to write a classifier for images of vehicles and animals. Mallory is a machine learning expert who built a classification library using state-of-the-art machine learning techniques. Alice decides to use Mallory’s library, and since machine learning libraries are driven by data, she carefully collects some amount of training data $\{x_i\}_i$ with images of cats, dogs, elephants trucks, cars, buses etc., with labels $y_i = 1$ or $y_i = 0$, stating whether an image is that of a vehicle or an animal respectively. She partitions it into a training set $S = \{(x_i, y_i)\}_i$ $M_i = 1$, and a test set T , and picks out her favorite ML algorithm, logistic regression, to learn a binary classifier that separates vehicles from animals. Alice runs Mallory’s