

Algorithmic Infrastructure and Performance of Finding Optimal Conditional Sparse ℓ_p -norm Regression

John Hainline

Fall 2018

1 Introduction

1.1 Overview

The focus of this project was to implement a particular complex algorithm in a performant, robust, and scalable way. The algorithm was specifically designed to solve the *conditional sparse ℓ_p -norm regression problem*. I managed the requirements for this system using libraries contained in the Scala language, and scaled the system using the Google Kubernetes Engine.

The *conditional sparse ℓ_p -norm regression problem* can be motivated by first considering the following related task of trying to find a subset of data on which sparse ℓ_p -norm regression succeeds. This occurs where there is a *linear rule* with a small number of terms that *predict a response variable* accurately in the ℓ_p sense (where ℓ_2 , the squared error, is our primary focus).

Unfortunately any linear predictor we find may not be a good predictor of a new point. This is because we cannot be sure if the new point belongs in the subset that the linear predictor applies to. However, if we have the explicit description of the subset where our predictor works, we can tell whether or not this new point belongs to that subset, and so we know whether or not we can use the predictor. Moreover, an explicit subset description makes it possible to get a statistical guarantee that our prediction will, in expectation, be accurate with high probability.

A real world application of this would be discovering patterns of medical conditions that are likely to occur in a sub-population, but that rarely occur in the population as a whole. The vulnerabilities of a particular race, age group, sex, or combination of these, could be more easily ferreted out of a large and complex data set. Not only that, but an explicit description gives us an easily understood rule on which particular portions of the population this pattern of medical issues apply.

1.2 Problems

Unfortunately, searching arbitrary subsets of data is a combinatorial problem and quickly becomes intractable for even moderate data set size. However a method was proposed in an earlier work [Juba, 2017] to solve this problem using a loss function of the ℓ_p -norm and a conditional description using k -DNFs (described later).

The problem was defined as follows.

We are given a matrix $Y \in \mathbb{R}^{n \times m} = (\vec{y}_1, \vec{y}_2, \dots, \vec{y}_n)^T$ containing n vectors of size m and a prediction target $Z \in \mathbb{R}^n$. We construct a boolean matrix $X \in \mathbb{B}^{n \times p} = (\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n)^T$ containing n vectors of size p . We derive X from Y such that our \vec{x}_i partitions \vec{y}_i into a p dimensional space. How this partitioning is done is up to the user, but tends to be by forming many boolean "buckets" for the m variables in \vec{y}_i . We define c as a k -DNF conditional¹ that, when applied to X , gives us a sub-population. So $c(X)$ is a "filter" on Y , reducing our full population of n to some smaller number. On this sub-population we can find a linear predictor for Z defined as $\langle \vec{a}, Y \rangle$. So the coefficients \vec{a} define a sparse regression line on Y which ideally predicts Z with small loss.

We restrict ourselves to k -DNFs for our conditional because if we could find conditions with arbitrary conjunctions this would yield PAC-learning algorithms for general DNFs. Juba's earlier work [Juba, 2017] makes clear that it is unlikely any algorithm can hope to find conditions of this kind. So we restrict ourselves to k -DNFs, which are a natural and powerful class of boolean representations that come with useful theoretical guarantees.

1.3 Basic Solution

Our implementation used the Scala language. We dealt with the scale of this problem by designing our code to run on the Google Kubernetes Engine (GKE). We synchronized many machines using the Akka Cluster library, which provides a powerful, managed, peer-to-peer clustering system. To ensure a robust solution we used Akka's Reactive Streams implementation which has very flexible error-handling controls. We also attached this system to a MySQL database to enable saving of all our results, which provided us with safety should the system ever fail during its calculations.

¹ k -Disjunctive Normal Form (k -DNF) is an OR of ANDs of at most k literals, where a literal is either a Boolean attribute or the negation of a Boolean attribute. See https://en.wikipedia.org/wiki/Disjunctive_normal_form.

2 Algorithm

2.1 Finding Conditional

At a high level, the conditional c is an accumulation of some number of boolean "terms" that are OR'ed together. We use 2-DNFs in our experiments, which means each "term" is an AND of 2 boolean variables. A valid 2-DNF could look something like

$$c = (x_1 \wedge \neg x_2) \vee (x_1 \wedge \neg x_3) \vee (x_2 \wedge x_4) \quad (1)$$

or in *set* notation

$$c = (x_1 \cap \neg x_2) \cup (x_1 \cap \neg x_3) \cup (x_2 \cap x_4) \quad (2)$$

Recall matrix $X \in \mathbb{B}^{n \times p}$, describes $Y \in \mathbb{R}^{n \times m}$. So the *set* notation of c , when applied $c(X)$, describes a sub-population of Y . Note that in our 2-DNF there are a total of $\binom{p}{2}$ possible $(x_i \cap x_j)$ terms. In fact there actually are $4 \cdot \binom{p}{2}$, since we must account for $(x_i \cap x_j)$, $(\neg x_i \cap x_j)$, $(x_i \cap \neg x_j)$, and $(\neg x_i \cap \neg x_j)$.

We go about finding our condition c by searching across all possible regression lines \vec{a} defined by our points Y and for each we construct the best k -DNF possible. The k -DNF with the lowest error overall is returned.

To find a regression line \vec{a} on Y we need 2 points and 2 dimensions of $Y \in \mathbb{R}^{n \times m}$. So we need to select 2 of n points and 2 of m dimensions. Unfortunately there is no easy way of picking which 2 dimensions to use, so we will need to check all $\binom{m}{2}$ combinations. Likewise we have no idea which points may give us an ideal c and so we check all $\binom{n}{2}$ combinations.

2.2 Red-Blue Set Cover

Now that we have a regression line \vec{a} , we need to determine the best k -DNF $c(X)$ that defines a sub-population of points in Y that are near the regression line. This c can then be used to predict Z . Keeping in mind our earlier *set* notation, we can see that each $(x_i \cap x_j)$ term in our k -DNF becomes a *set* of included points in Y . Now we need to find which terms to include in c . To accomplish this, we employ a derivation of the Red-Blue Set Cover algorithm [Carr et al., 2000] to find the best terms for c that fit \vec{a} .

The Red-Blue Set Cover problem is defined as follows.

Given a finite set of "red" elements R , and a finite set of "blue" elements B and a family $S \subseteq 2^{R \cup B}$, the red-blue set cover problem is to find a subfamily $C \subseteq S$ which covers all blue elements, but which covers the minimum possible number of red elements [Carr et al., 2000].

In our case we use a greedy algorithm to find an approximate solution to this Red-Blue Set Cover problem. We consider the red-blue spectrum of a point to be its distance from \vec{a} . This is represented with a single value of increasing "redness" for every point, and we try to minimize the amount of "redness" introduced while maximizing the number of points included.

To start applying our solution to Red-Blue Set Cover, we first compute a "redness" value for every point $\vec{y}_i \in Y$. We then construct the total "redness" for each $(x_j \cap x_k)$ term and sort, successively taking the least "red" terms until reaching a balance between "redness" and point count.

Unfortunately, this will ultimately call our Set Cover algorithm $\binom{n}{2} \times \binom{p}{2}$ times. For even a moderately small data set with $n = 1000, p = 6$ this becomes $\binom{1000}{2} \binom{6}{2} = 499,500 \times 15 = 7,492,500$ Set Cover calls. The scaling issues are obvious.

2.3 Scaling

To deal with this we can simply sub-sample our Y . For many real world distributions, random sampling (assuming the sub-sample is not extremely small) will still find our k -DNF with high likelihood. The actual statistics is not broken down in this paper, but we very rarely failed when sub-sampling 200 out of 10,000 points. This sub-sampling is one of the **single biggest factors** in speeding up our calculations as it can reduce the number of *SetCover* calls by orders of magnitude, while only marginally reducing our probability of missing a hidden k -DNF.

See algorithms 1 and 2 for a pseudo-code implementation.

3 Implementation

We decided to write the system in the Scala programming language. This was due to language familiarity, and for the *Akka framework*, which facilitates computation across clusters of PCs. We specifically used the Akka Cluster², Akka Actor³, and Akka Stream⁴ libraries.

²<https://doc.akka.io/docs/akka/current/common/cluster.html#intro>

³<https://doc.akka.io/docs/akka/current/actors.html#introduction>

⁴<https://doc.akka.io/docs/akka/current/stream/stream-introduction.html>

3.1 Akka Library

3.1.1 Akka Actor

Akka *Actors* provide a high level abstraction for creating concurrent and distributed systems. The basic design is extremely simple. An *Actor* is simply an object that sends and receives messages using *mailboxes* (preventing concurrent access). *Actors* are stateful, but do not alter external state, so locking and other concurrency mechanisms are unnecessary. There are numerous other features of *Actors*, but at a high level this is what is important.

3.1.2 Akka Cluster

The Akka *Cluster* library provides a peer to peer based cluster membership service. It is fault tolerant and decentralized. We use this to establish an initial connection between servers. The Akka *Cluster* library also provides an easy mechanism for passing events or *Actor* messages across any/all members in the cluster, easily bridging any network boundaries.

3.1.3 Akka Streams

Akka *Streams* follow the "Reactive Streams" initiative. This standard tries to provide asynchronous stream processing with non-blocking back pressure. When trying to solve our k -DNF discovery problem, we utilize this API to provide a cross-machine work pulling strategy for running through all possible sparse regression lines \vec{a} in Y .

3.2 Architecture

The application is designed to be run either *locally*, or on the *Google Kubernetes Engine* (GKE). When running locally, events follow this pattern:

1. User sends a workload to the system. This names the data set the system will run, and provides relevant variables (such as μ and the subset size).
2. The ManageActor (an Akka Actor) receives the workload and:
 - (a) Loads (X, Y, Z) from the database.
 - (b) Prepares database tables for saving this workload's results.
 - (c) Creates Akka Stream objects, specifically BalanceHub and MergeHub, to handle work flowing through the system.
 - (d) Publishes the workload to the event manager.

3. The ComputeActor receives the workload and:
 - (a) Loads (X, Y, Z) from the database.
 - (b) Instantiates a number of SetCover objects (Akka Stream objects) that each take in work and give a result. These operate *in parallel*.
4. The BalanceHub, once a SetCover instance has been attached, begins sending work through the system.

For a visual breakdown of this process, see figure 1. By having the Akka Actors create Akka Streams independently, we can dynamically create and destroy SetCover instances, scaling up or down the number of CPUs we use.

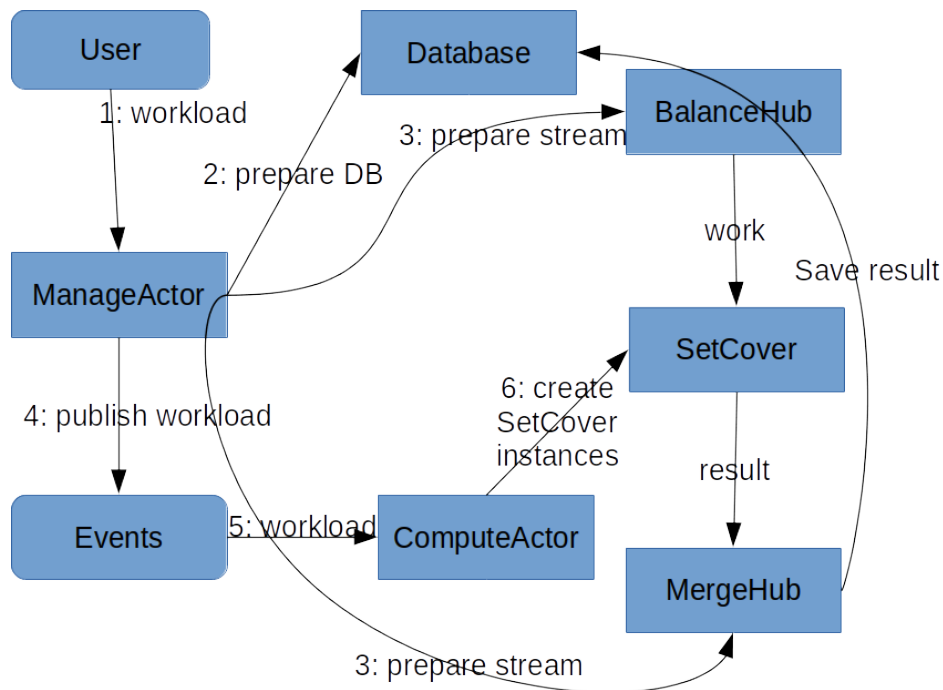


Figure 1: (1) The user sends a workload to the system. The ManageActor receives the workload and: (2) Prepares database tables for the workload. (3) Creates Akka Stream objects, specifically BalanceHub and MergeHub, to handle flowing work through the system. (4) Publishes the workload to the event manager. The event manager (5) sends the workload to the ComputeActor. Finally, the ComputeActor (6) instantiates some number of SetCover objects that each take in work and give a result.

For running in GKE, it should be noted that both Akka Actors as well as Akka Streams are capable of passing messages to separate machines on the network. Simply by instantiating ComputeActors on other machines, we get message passing over the network (rather than locally) which allows an entire cluster of machines to participate.

4 Experimental Evaluation

The code itself can be found at <https://github.com/johnhainline/sclr>. The latest design is able to run a data set with $\mathbf{X} \in \mathbb{B}^{10000 \times 50}$, $\mathbf{Y} \in \mathbb{R}^{10000 \times 20}$ (so $n = 10000$, $m = 20$, $p = 50$) and sub-sampling of 500 points. This meant $terms.size = 4 \cdot \binom{50}{2} = 4900$ and *SetCover* ran $\binom{20}{2} \cdot \binom{500}{2} = 27,702,500$ times.

We successfully ran 26 machines on Google Kubernetes Engine (GKE), each with 8 CPU cores. The entire run took roughly five days. We had a single instance of the *ManageActor* and Database running on 1 machine, while the other 25 all ran a single instance of the *ComputeActor*. When a workload reached the *ComputeActor* on each machine, they instantiated 8 *SetCover* instances apiece. The *ManageActor* used the *BalanceHub* to send work to the *SetCover* instances as fast as they could process data. The *ManageActor* used the *MergeHub* to collect all results from the *SetCover* instances and saved them in the database in batches. This gave us 200 CPUs each running *SetCover* in parallel at near 100% utilization.

4.1 Optimizations

4.1.1 Sub-sampling

Sub-sampling was our single greatest optimization. By reducing n to some small number s using random sampling we often run multiple orders of magnitude fewer *SetCover* calls than we would otherwise. This reduction, for many real world point distributions, only marginally reduces our probability of recovering an embedded k -DNF (assuming one exists in the data set).

4.1.2 Set Implementations

Originally we used a standard *Set* implementation for the *terms* variable in *SetCover*. As we had sets that were thousands of items long, it quickly became apparent that a *BitSet* implementation was necessary. This simple change increased performance roughly by a factor of 10.

4.1.3 Embedded Native Code

At one point we experimented with using Java Native Interface (JNI) to embed native C++ code, going so far as to implement *SetCover* in C++ and create bindings from Scala. While we did make this work, the inherent overhead in calling a C++ function from the JVM was so high that it was not performant. If we had used a different batching strategy (such that the C++ code was called once for a large group of *SetCover* calls) this may have become a viable strategy. Unfortunately this change would have required a rather drastic refactor of Scala code and so we pursued other methods instead.

4.1.4 Batching

The Akka Streams implementation allows easy batching of messages as they travel. We began batching messages that had to travel across the network (i.e. to/from the *SetCover* instances). We also batched our database save calls. These changes gave us a roughly $10\times$ speed up.

4.2 Recovery

4.2.1 Database

During a workload, every run of *SetCover* saves any resulting discovered k -DNFs to the database. Each record also contains all the relevant variables passed to *SetCover* (specifically i, j, k, l). With this information it is possible to know how far along we were in our computation and potentially resume after a catastrophic system shutdown (assuming the database still exists).

4.2.2 Sub-sampling

Because we sub-sample Y , attempts at resuming a failed workload could generate a completely different sub-sample and be unreproducible. We avoid this by allowing a random number generator seed to be passed in with the workload, which is used when sub-sampling Y .

4.3 Limitations

This algorithm's limitations are primarily derived from the non-parallelizable portions of the code. Specifically, sending information across the network as well as saving to the database. The latest version of code has reduced these issues using batching.

While we are sure there is a limit to the number of parallel *ComputeActor* machines (and therefore *SetCover* instances) we can work with, we have not reached it yet. As our largest run of the program took advantage of **200** CPUs simultaneously, we were afraid to bring a great deal more machines online due to the large cost in GKE. Earlier versions of our algorithm could only handle parallelization factors of **30** (before an improved *SetCover* implementation) and then **100** (before batching). Based on our experience at this point, it seems likely we could scale to somewhere around **300** – **400** CPUs before needing to further optimize.

4.4 Future Work

Running on a **10,000** point data set took about five days with roughly **200** CPUs. We still wish to increase the allowable data set size even further. There is still likely another order of magnitude speed up attainable in *SetCover*. However moving away from GKE and instead writing the system to easily run on a local computer cluster would probably yield greater dividends. Using thousands of CPUs on Kubernetes is of course possible, but local resources will certainly be far cheaper.

Algorithm 1: This algorithm breaks down the relevant search space of $X \in \mathbb{B}^{n \times p}$ and $Y \in \mathbb{R}^{n \times m}$, giving *SetCover* the parameters it needs. Note that *terms* must include the \neg versions of our $\binom{p}{2}$ possible combinations of X , which is why we have $4 \cdot \binom{p}{2}$ combinations.

```

1 def ComputeBestKDNF( $X, Y, Z, r$ ):
2    $X, Y = \text{SubsamplePoints}(X, Y, r)$ 
3    $terms = (x_i \cap x_j)$  for all  $4 \cdot \binom{p}{2}$  combinations of  $X$ 
4    $best\_terms = \emptyset$ 
5    $best\_error = \infty$ 
6   for  $(i, j) \in \binom{r}{2}$  points do
7     for  $(k, l) \in \binom{m}{2}$  dimensions do
8        $t, e = \text{SetCover}(X, Y, Z, terms, i, j, k, l)$ 
9       if  $e < best\_error$  then
10          $best\_terms = t$ 
11          $best\_error = e$ 
12       end
13     end
14   end
15   return  $best\_terms, best\_error$ 

```

Algorithm 2: The SetCover algorithm searches for the best k -DNF (lowest error in terms of "redness") for a given pair of data points i, j and dimensions k, l in $Y \in \mathbb{R}^{n \times m}$. The variable $terms$ holds every set generated by our $\binom{p}{2}$ terms $(x_i \cap x_j)$ in $X \in \mathbb{B}^{n \times p}$. μ dictates how large of a sub-population we are looking for.

```

1 def SetCover( $X, Y, Z, terms, i, j, k, l$ ):
2    $\vec{a} = \text{CalculateRegressionLine}(\vec{y}[k]_i, \vec{y}[l]_j)$ 
3    $r\_points = \text{CalculateRednessForPoints}(Y, \vec{a})$ 
4    $r\_terms = \text{CalculateRednessForTerms}(terms, r\_points)$ 
5    $sort = \text{SortByAverageRedness}(r\_terms)$ 
6    $best\_terms = \emptyset$ 
7    $best\_error = \infty$ 
8    $cur\_terms = \emptyset$ 
9   while ( $cur\_terms.size < \mu \cdot n$ ) do
10     $next = \text{NextUnusedTermInSort}(sort, cur\_terms)$ 
11     $cur\_terms.add(next)$ 
12  end
13   $cur\_error = \text{AverageRednessPerPoint}(cur\_terms)$ 
14  if  $cur\_error < best\_error$  then
15     $best\_error = cur\_error$ 
16     $best\_terms = cur\_terms$ 
17  end
18  return  $best\_terms, best\_error$ 

```

References

- [Carr et al., 2000] Carr, R. D., Doddi, S., Konjevod, G., and Marathe, M. (2000). On the red-blue set cover problem. In *Proceedings of the Eleventh Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '00, pages 345–353, Philadelphia, PA, USA. Society for Industrial and Applied Mathematics.
- [Juba, 2017] Juba, B. (2017). Conditional sparse linear regression. In *LIPICs-Leibniz International Proceedings in Informatics*, volume 67. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik. Proc. 8th ITCS.