
Use K-Means for prototype selection

Shuo Xu

Department of Computer Science and Engineering
University of California San Diego
s3xu@ucsd.edu

Abstract

In this project, we introduced a coordinate descent algorithm, which updates the variable vector one dimension at a time from the beginning to the last during one update. This method is proved to be more efficient and steady in terms of convergence.

1 Idea Description

Suppose our optimization problem is to minimize $L(w)$, where $w \in R^d$ and L is differentiable.

In my coordinate descent method, the update of weight vector w is performed one dimension at a time from w_1 to w_d , and each dimension is updated using their partial derivative.

Algorithm 1 My coordinate descent algorithm

```
1: procedure COORDINATEDDESCENT(epochs, learning_rate)
2:    $i \leftarrow 0$ 
3:   for  $i < \text{epochs}$  do
4:      $j \leftarrow 0$ 
5:     for  $j < d$  do
6:        $w\_grad = L(w).w\_grad()$  // a d-dimension vector
7:        $w[j] -= \text{learning\_rate} * w\_grad[j]$ 
8:        $j \leftarrow j + 1$ 
9:      $i \leftarrow i + 1$ 
```

2 Convergence

One trivial case is that $L(w)$ is convex and differentiable, the method would converge to the optimal loss. It's easy to find a counter case on those indifferentiable points where the weights cannot be updated.

3 Experimental Result

3.1 Data preprocess

The experiment is carried out on the first 2 classes of UCI wine dataset with 59 and 71 points each class.

Before training, all the data are normalized to $[0, 1]$ using following formula, where data is stored row wise.

$$data = \frac{data - \min(\text{train_data_set}, axis = 0)}{\max(\text{train_data_set}, axis = 0) - \min(\text{train_data_set}, axis = 0)}$$

The performance of baseline and k-means algorithm are compared in this section.

3.2 Performance

3.2.1 Standard Logistic Regression

We performed standard logistic regression over the dataset using sklearn, where $C = 10^{20}$ to reduce the effect of the penalty term. The models finally converged to

$$L^* = 9.992007221626415 \times 10^{-16}$$

Figure 1 is the training loss curve of my coordinate descent method and the random one with learning rate of 0.3. To be fair, the two methods start from the same initial weight via sharing the same random seed. It's obvious that our proposed algorithm not only performs as good as standard gradient descent in terms of final loss when converging, but is more steady and converges faster compared to the random approach.

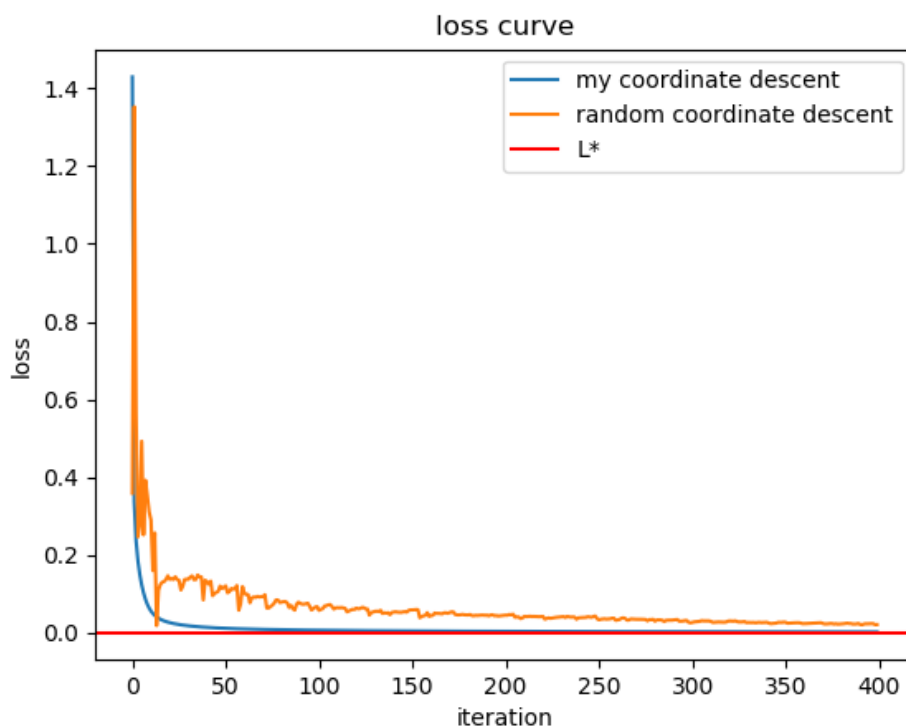


Figure 1: Training loss

4 Critical evaluation

Yes. I've got 2 ideas - greedy selection and batch selection.

Compared to the random approach, the method proposed only enforced a certain order when updating the weight. And we can further incorporate a greedy method. To be specific, after get the gradient of weight vector, we can select the dimension with the largest absolute value, which suggests that dimension would have the greatest change.

The idea of batch update is inspired by stochastic gradient descent and mini-batch gradient descent. Instead of update 1 weight at a time, we can update k dimension with the biggest absolute gradient.

Appendix of core methods implementation using Python

4.1 Base class of classifier

```
1 import numpy as np
2 from utils import onehot
3
4 class LogisticClassifier():
5     def __init__(self, X, y):
6         self.y_onehot = onehot(y)
7         self.X, self.y = np.column_stack((np.ones(len(X)), X)), y[:, np.
newaxis]
8         self.w = None
9
10    def logistic(self, s):
11        return np.array(1 / (1 + np.exp(-s)))
12
13    def predict(self, X):
14        probs = self.logistic(np.dot(np.column_stack((np.ones(len(X)), X)
), self.w))
15        prediction = np.zeros(probs.shape)
16        prediction[probs > 0.5] = 1
17        return np.squeeze(prediction)
18
19    def prob(self, X):
20        return self.logistic(np.dot(X, self.w))
21
22    def loss(self, X, y):
23        y_hat = self.logistic(np.dot(X, self.w))
24        return (- np.dot(y.T, np.log(y_hat)) / len(y_hat))[0][0]
25
26    def gradient(self, X, y):
27        y_hat = self.logistic(np.dot(X, self.w))
28        return np.sum((y_hat - y) * X, axis=0).reshape(-1, 1)
29
30    def accuracy(self, test_set=None):
31        if not test_set:
32            test_set = self.test_set
33        return np.sum(self.predict(test_set.X) == test_set.y) / len(
test_set.y)
```

4.2 Random coordinate Descent

```
1 class RandomCoordinateDescentClassifier(LogisticClassifier):
2     def __init__(self, X, y):
3         super().__init__(X, y)
4
5     def train(self, T=600, lr=0.3):
6
7         # initialization
8         train_X, train_y = self.X, self.y
9         self.train_losses = []
10        np.random.seed(11)
11        self.w = np.random.random((train_X.shape[1], train_y.shape[1]))
12        self.W = []
13
14        counter = 1
15        num_weights = len(self.w)
16        # gradient descent
17        for t in range(T):
18            i = np.random.randint(num_weights)
19            self.w[i] -= lr * self.gradient(train_X, train_y)[i]
```

```

20         self.W.append(self.w.tolist())
21
22         # compute losses on train dataset and holdout dataset
23         loss = self.loss(train_X, train_y)
24         self.train_losses.append(loss)
25
26         # save the parameters with best performance
27         print("loss {} = {}".format(counter, np.array(loss)))
28         counter += 1
29     self.w = np.array(min(self.W, key=lambda w: self.train_losses[self.W.
30                          index(w)]))
    return self.train_losses

```

4.3 My coordinate Descent

```

1 class MyCoordinateDescentClassifier(LogisticClassifier):
2     def __init__(self, X, y):
3         super().__init__(X, y)
4
5     def train(self, T=600, lr=0.3):
6
7         # initialization
8         train_X, train_y = self.X, self.y
9         self.train_losses = []
10        np.random.seed(11)
11        self.w = np.random.random((train_X.shape[1], train_y.shape[1]))
12        self.W = []
13
14        counter = 1
15        # gradient descent
16        for t in range(T):
17            for i in range(len(self.w)):
18                self.w[i] -= lr * self.gradient(train_X, train_y)[i]
19                self.W.append(self.w.tolist())
20
21            # compute losses on train dataset and holdout dataset
22            loss = self.loss(train_X, train_y)
23            self.train_losses.append(loss)
24
25            # save the parameters with best performance
26            print("loss {} = {}".format(counter, np.array(loss)))
27            counter += 1
28        self.w = self.W[np.argmin(self.train_losses)]
29        return self.train_losses

```