

Coordinate Descent vs. (Accelerated) Proximal Gradient for Lasso

Marcus Artner, Hannes Janker & Lukas Till Schawerda
Optimization for Data Science, University of Vienna

Abstract—First order methods like gradient descent are a crucial part for solving optimization problems. These methods, using the gradient of an objective function, play an essential role in a large number of machine learning and data science tasks. The basic framework of these algorithms is to approximate the problem by an easier model, solve the model and iteratively continuing this procedure. Applying these methods is often times much computationally cheaper than computing closed form solutions for the problems. In the following, two gradient methods are described and compared for the Lasso objective function.

I. INTRODUCTION

The aim of this project is to introduce and compare the two gradient methods coordinate descent and proximal gradient descent. These algorithms will be applied to the Lasso objective function. Firstly, both of the methods will be briefly explained. Moreover, a short introduction about the Lasso problem will be given. In the main part, the performance of the algorithms will be compared with implementations in Python. Besides that, it will be tried to achieve acceleration for both of the algorithms. Furthermore, the learning rate decreases with the epoch by a specified factor. In the final part, these convergence rates will be plotted and the performances will be compared.

II. THEORY ABOUT THE ALGORITHMS AND LASSO

Coordinate descent is an optimization algorithm that minimizes along the coordinate direction to minimize a function. Via a selection rule, the algorithm chooses one coordinate in each iteration and minimizes over the coordinate while all the other coordinates remain fixed. In other words, only one coordinate is modified per iteration. The algorithm basically contains two steps per iteration: select $i_k \in \{1, \dots, d\}$ and calculate

$$x_{k+1} = x_k + \gamma * e_{i_k},$$

where e_{i_k} is the i -th unit basis vector. There are two main variants. The first one chooses gradient-based step sizes:

$$x_{k+1} = x_k - \frac{1}{L} \nabla_{i_k} f(x_k) * e_{i_k}$$

The second one computes the exact coordinate minimization by solving the scalar problem:

$$\operatorname{argmin}_{\gamma \in R} f(x_k + \gamma * e_{i_k})$$

While the first one needs line search for the step sizes the second one is hyperparameter free. There are two technical

assumption that we make. We assume strong convexity and we assume coordinate-wise L -smoothness, i.e.

$$f(x + \gamma * e_i) \leq f(x) + \gamma \nabla_i f(x) + \frac{L}{2} \gamma^2$$

$$\forall x \in R^d, \forall \gamma \in R, \forall i \in [d]$$

There are several ways for selecting the coordinates. The most naive approach would be to simply cycle through all coordinates. A common practice is to select the coordinates uniformly at random ($i_k \in \{1, \dots, d\}$). Other ways of choosing the coordinates would be importance sampling and steepest coordinate descent.

Proximal gradient methods can be seen as a generalized form of projection which is used to solve non-differentiable convex optimization problems. Consider an objective function that can be composed as

$$f(x) = g(x) + h(x)$$

where g is a convex and differentiable function and h is a simple function which is convex but not necessarily differentiable. For the proximal gradient algorithm the proximal mapping for a function h and parameter α needs to be defined.

$$\operatorname{prox}_{\alpha h}(x) = \operatorname{argmin}_{y \in R} \{h(y) + \frac{1}{2\alpha} \|y - x\|^2\}$$

An iteration of the algorithm looks as following:

$$x_{k+1} = \operatorname{prox}_{\alpha h}(x_k - \alpha \nabla g(x_k))$$

The proximal mapping $\operatorname{prox}_{\alpha}(\cdot)$ can be computed analytically for many important functions h and it depends only on h , not on g . This means even if g is a complicated function, only its gradient needs to be computed. The computational cost of $\operatorname{prox}_{\alpha}(\cdot)$ depends of course on h . An important application is the Lasso objective which includes a non-smooth regularizer term.

The Lasso regression is a type of linear regression that uses shrinkage which means that the coefficients of the estimator are shrunk towards a central point. Lasso encourages simple, sparse models and is well-suited for models with high multicollinearity or for automatic model selection parts. The objective function which should be minimized is a quadratic and looks the following:

$$F(\beta) = \frac{1}{2} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

$\lambda \|\beta\|_1$ is called regularizer term. This kind of regularization can result in sparse models with few parameters and some coefficients can even become zero and will be eliminated from the model. How strong the so called penalty is, depends on the tuning parameter λ . The gradient of the objective function is given by

$$\frac{d}{d\beta} F(\beta) = X^T(Y - X\beta) + \lambda \text{sign}(\beta)$$

with the $\text{sign}(\cdot)$ -function being applied componentwise. Note that the gradient of the $\|\cdot\|_1$ -norm at the point 0 is thus set to 0.

In this setting, coordinate descent yields the updating scheme

$$\beta_{k+1}^j = \beta_k^j + \alpha X_{\cdot j}^T(Y - X\beta) + \lambda \text{sign}(\beta_k^j)$$

where β_k^j is the j -th component of β_k and $X_{\cdot j}$ is the j -th column of X . During each epoch, the order in which the components are updated is determined at random. If an acceleration parameter is specified, then Nesterov momentum is applied after each epoch.

Proximal gradient descent for the Lasso problem follows the updating scheme

$$\beta_{k+1} = \text{prox}_\alpha(\beta_k + \alpha X^T(Y - X\beta))$$

with

$$[\text{prox}_\alpha(\beta)]_j = \begin{cases} \beta^j + \lambda\alpha & \text{if } \beta^j < -\lambda\alpha \\ 0 & \text{if } -\lambda\alpha \leq \beta^j \leq \lambda\alpha \\ \beta^j - \lambda\alpha & \text{if } \lambda\alpha < \beta^j \end{cases}$$

III. IMPLEMENTATION

We implemented the updating schemes explained above in Python using the *numpy* library. All plots were creating using *matplotlib.pyplot*. The functions *coordinate_descent()* and *proximal_gradient_descent()* take the following inputs:

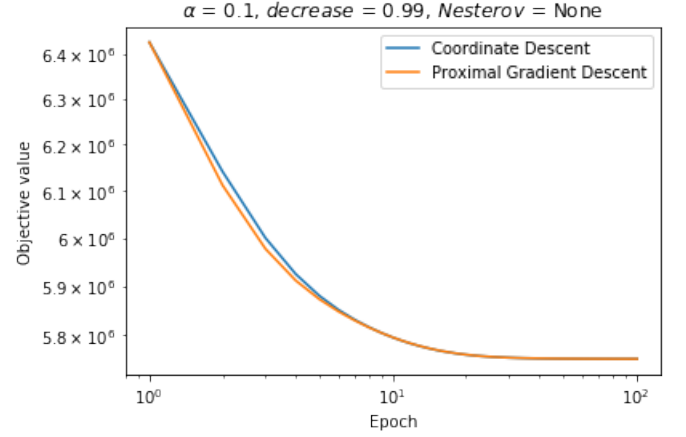
PARAMETER	INPUT
X	Data
y	Targets
beta_0	Starting point
alpha	Learning rate
num_epochs	Number of epochs
decr	Factor for decreasing learning rate
accelerate	Nesterov acceleration parameter

They output both the objective value over the course of learning and the final estimate for β . Auxiliary functions for computing intermediate results were also implemented, but are not explained in detail here. For further information one can consider the code.

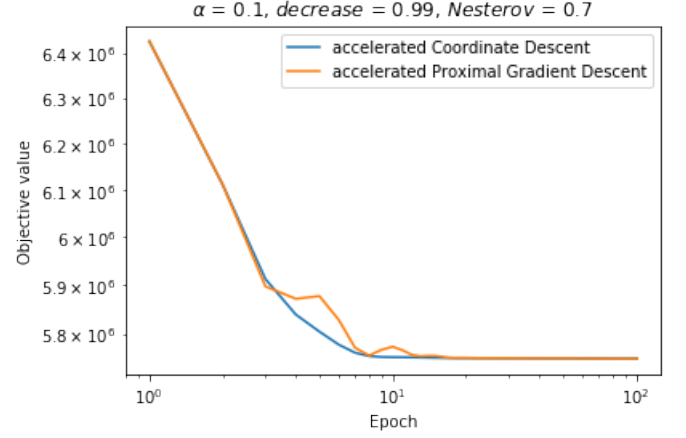
The *california_housing*-dataset from sklearn was used for evaluation. It has 20640 observations and 8 parameters.

IV. RESULTS

The initial weight vector β_0 was set to all zeros. First we compare coordinate descent and proximal gradient descent without acceleration.



We observe that Coordinate Descent, very similar to Proximal Gradient Descent converges rather fast after about 20 epochs. Neither algorithm performs significantly better than the other.



If we accelerate both algorithms with a Nesterov momentum parameter of 0.7, we observe slightly faster convergence over the not accelerated algorithms. Accelerated proximal gradient descent shows signs of overshooting in the first steps due to momentum.

Overall, both methods display excellent performance. This holds true even for various other data sets that we tried. Both methods are thus appropriate for solving the Lasso problem. We cannot conclude that one method is better than the other.

V. SUMMARY

To sum up, we implemented both, coordinate descent and proximal gradient descent for solving optimization problems. In particular, we applied the algorithms to the Lasso objective function using open source data sets like the *california_housing*-dataset. We were able to achieve a small acceleration by implementing Nesterov's momentum. Our

final conclusion is that none of the both algorithms can be considered better based on our findings.

VI. APPENDIX: IMPORTANT PARTS OF THE CODE

```
1 # Define Objective function
2
3 def objective(X,y,beta,lam):
4     obj = (1/2)*np.linalg.norm(y - X@beta, 2)**2 +
5         lam*np.linalg.norm(beta, 1)
6     return obj
7
8 ##### Coordinate Descent
9
10 def compute_coordinate_grad(X,y,beta,coordinate,lam)
11 :
12     """
13     INPUT
14     X...Data Matrix
15     Y...Targets
16     coordinate...direction which will be optimized
17     lam...lambda
18     """
19
20     # Dimensions
21     n,p = X.shape
22
23     # j-th column of X
24     x_j = X[:,coordinate]
25
26     # j-th element of beta
27     beta_j = beta[coordinate]
28
29     # Regularizer
30     reg = lam*np.sign(beta_j)
31
32     # Compute gradient
33     grad = -x_j@(y-X@beta) + reg
34
35     # Return results
36     return grad
37
38 def coordinate_descent(X, y, beta_0, lam, alpha,
39 num_epochs, decr = None, accelerate = None):
40     """
41     INPUT
42     X...Data Matrix
43     Y...Targets
44     beta_0...initial beta vector
45     alpha...learning rate
46     num_epochs...number of epochs
47     decr...Parameter to decrease learning rate in
48     each epoch
49     accelerate...Nesterov acceleration parameter
50     """
51
52     # Get Dimensions of data and prepare beta values
53     for acceleration
54
55     n,p = X.shape
56     beta = beta_0.copy()
57     beta_old = beta_0.copy()
58
59     # Initialize empty objective list
60
61     obj = []
62     obj.append(objective(X, y, beta_0, lam))
63
64     # Iterate over epochs
65
66     for epoch in range(num_epochs):
67
68         # Sample coordinates for current epoch
69
70         ind = np.random.permutation(p)
71
72         if accelerate is not None:
```

```

68         # Nesterov if accelerate is True
69
70         beta_new = beta.copy()
71         beta_temp = beta + accelerate * (beta -
72         beta_old)
73
74         for j in ind:
75
76             # Compute gradient at accelerated
77             point
78
79             grad = compute_coordinate_grad(X, y,
80             beta_temp, coordinate = j, lam = lam)
81
82             # update j-th component of beta
83
84             beta[j] = beta_temp[j] - alpha*grad
85             beta_old = beta_new.copy()
86
87         else:
88
89             #regular Coordinate descent
90
91             for j in ind:
92
93                 grad = compute_coordinate_grad(X, y,
94                 beta, coordinate = j, lam = lam)
95                 beta[j] -= alpha*grad
96
97             # append current objective
98
99             obj.append(objective(X,y,beta,lam))
100
101             # Adapt learning rate if decr = True
102
103             if decr is not None:
104
105                 alpha *= decr
106
107             return obj, beta
108
109 ##### Proximal Gradient Descent
110
111 # Soft Threshold Operator
112 def soft_thresh(beta, lam):
113     """
114     INPUT
115     beta...Vector beta
116     lam...lambda
117     """
118
119     res = np.zeros_like(beta)
120     p = len(beta)
121
122     # adjust beta according to formula
123
124     for i in range(p):
125         if beta[i] < -lam:
126             res[i] = beta[i] + lam
127         elif np.abs(beta[i]) < lam:
128             res[i] = 0
129         elif beta[i] > lam:
130             res[i] = beta[i] - lam
131     return res
132
133
134 def proximal_gradient_descent(X, y, beta_0, lam,
135     alpha, num_epochs, decr = None, accelerate =

```

```

136     INPUT
137     X...Data Matrix
138     Y...Targets
139     beta_0...initial beta vector
140     alpha...learning rate
141     num_epochs...number of epochs
142     decr...Parameter to decrease stepsize in each
143     epoch
144     accelerate...Nesterov acceleration parameter
145     """
146
147     n,p = X.shape
148     beta = beta_0.copy()
149     beta_old = beta_0.copy()
150
151     # initialize objective list
152
153     obj = []
154     obj.append(objective(X, y, beta_0, lam))
155
156     # Iterate over epochs
157
158     for epoch in range(num_epochs):
159
160         # Nesterov if accelerate = True
161
162         if accelerate is not None:
163
164             beta_new = beta.copy()
165             beta_temp = beta + accelerate * (beta -
166             beta_old)
167             content = beta_temp + alpha*X.T@(y -
168             X@beta)
169             beta = soft_thresh(content, lam*alpha)
170             beta_old = beta_new.copy()
171
172             # updating beta with soft threshold operator
173             according to formula
174
175         else:
176
177             content = beta + alpha*X.T@(y - X@beta)
178             beta = soft_thresh(content, lam*alpha)
179
180             # append current objective
181
182             obj.append(objective(X, y, beta, lam))
183
184             # adapt learnign rate if decr = True
185
186             if decr is not None:
187
188                 alpha *= decr
189
190     return obj, beta

```