# Coordinate descent

**Qimin Chen**
A53284263
qic003@ucsd.edu

## 1 Coordinate descent method

### 1.1 Description

After computing the derivative of logistic loss with respect to $w_i$, pick the $k$ largest absolute value of $dw_i$ where $i \in \{1, 2, 3, ..., 13\}$ (for problem a and b $k = 1$), update the corresponding value of $w_i$ by subtracting its $dw_i \times \alpha$ where $\alpha$ is the step size and stop iterations when the difference between current loss and last loss is smaller than $10^{-4}$.

### 1.2 Pseudocode

---
**Algorithm 1** Coordinate descent

---
1: normalize the dataset and split data into training set and test set
2: randomly initialize weights $w$
3: compute initial loss $l_{prev}$
4: set current loss $l_{curr} = \infty$
5: **while** $l_{curr} - l_{prev} > 10^{-4}$ **do**
6:     compute $\frac{dE}{dw} = (Y_{pred} - Y_{true})X$
7:     sort and find index of $K$-largest $dw$
8:     **for** k=0,1,...,K **do**
9:         $w[k]$ -= $\alpha \times dw[k]$
10:     **end for**
11:     $l_{prev} = l_{curr}$
12:     recompute the current loss $l_{curr}$
13: **end while**

---

## 2 Convergence

The trainging loss stops decreasing and remains roughly the same for some epoche when reducing the learning rate then the model converges to the optimal loss.

## 3 Experimental results

Figure 1 shows the training loss on my coordinate descent(red line), random coordinate descent(yellow line) and logistic regression solver from scikit-learn(blue line) over 1000 iterations. It can be seen from Figure 1 that my coordinate descent performs better than random one because random coordinate descent merely picks coordinates at random and update corresponding weights which causes unstability on training phase.
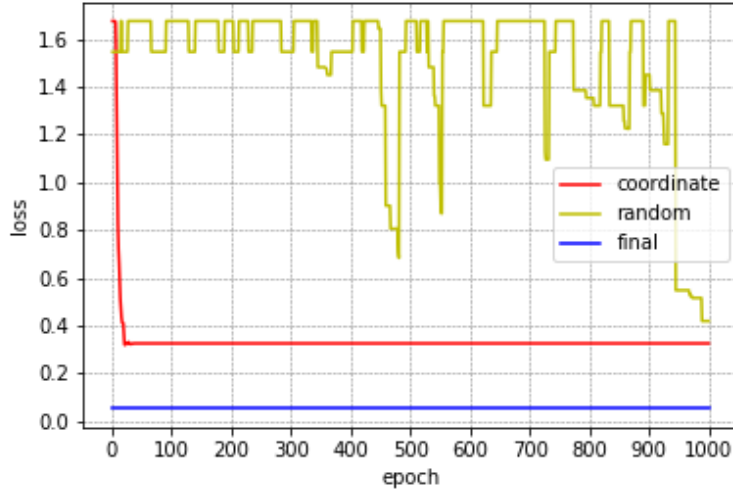
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107



Figure 1: training loss on my coordinate descent, random coordinate descent and logistic regression solver over 1000 iterations.

## 4 Critical evaluation

I think there is scope for further improvement of my coordinate descent scheme. Since I only update $k$ largest value of $dw$, it might not be appropriate to update $w$ by just subtracting its corresponding $dw$. One way to update the $k$ largest weights is to calculate Hessian matrix:

$$H_{\alpha\beta} = -\sum_t \sigma'(w \cdot x_t)x_{\beta t}x_{\alpha t}$$

$$w \leftarrow w - H^{-1}dw$$

where $\alpha, \beta \in \{1, 2, 3, ..., 13\}$ and $\sigma(z)$ is the sigmoid function.

## 5 Sparse coordinate descent

As the pseudocode described above, pick the $k$ largest absolute value of $dw_i$ where $i \in \{1, 2, 3, ..., 13\}$, update the corresponding value of $w_i$ by subtracting its $dw_i \times \alpha$

Table 1 shows the loss values for different values of $k$ on my coordinate descent and random descent over 1000 iterations. I do not think my method for sparse coordinate descent always find the best $k$-sparse solution when $L(\cdot)$ is convex.

2

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

| Experiment | My coordinate descent | Random coordinate descent |
|:---:|:---:|:---:|
| k = 1 | 0.3256 | 0.4191 |
| k = 3 | 0.2606 | 0.3868 |
| k = 5 | 0.3242 | 0.3837 |
| k = 7 | 0.3201 | 0.5803 |
| k = 9 | 0.3002 | 0.4190 |
| k = 11 | 0.2907 | 0.5480 |
| k = 13 | 0.2873 | 0.3868 |

Table 1: training loss for different values of $k$ on my coordinate descent and random descent over 1000 iterations.