

# Logistic Regression Via Coordinate Descent

Chunlin Chen

Department of Electrical and Computer Engineering  
University of California, San Diego  
chc126@ucsd.edu

## Abstract

In convex optimization problems, gradient descent (GD) and stochastic gradient descent (SGD) have been proven effective for differentiable objective functions. In this paper, we pay attention to a different approach, coordinate descent (CD), which updates one certain coordinate of the variable vector at each iteration, and compare the performances under different coordinate selection methods on a classic logistic regression problem.

## 1 Coordinate Descent Method

We consider a standard unconstrained optimization problem:

$$\min_{\omega \in \mathbb{R}^d} L(\omega) \quad (1)$$

Suppose the objective function  $L(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$  is continuous and differentiable everywhere, then a general coordinate descent algorithm can be described as follows (Wright, 2015):

---

### Algorithm 1 General Coordinate Descent Pipeline

---

- 1: Initialization:  $k \leftarrow 0, \omega_0 \in \mathbb{R}^d$
  - 2: **repeat**
  - 3:   Choose index  $i_k \in \{1, 2, \dots, d\}$
  - 4:    $\omega^{k+1} \leftarrow \omega^k - \alpha_k [\nabla L(\omega)]_{i_k} \mathbf{e}_{i_k}$
  - 5:    $k \leftarrow k + 1$
  - 6: **until** convergence
- 

We use  $[\nabla L(\omega)]_i$  to denote the  $i$ th component of the gradient  $\nabla L(\omega)$ , and  $\mathbf{e}_i$  denotes the vector with a 1 in the  $i$ th coordinate and 0's elsewhere. Note that the update of the value of the selected coordinates is just applying gradient descent with respect to scalars. We will focus more on the selection of coordinates to be updated. The trivial method is to choose coordinates uniformly at random, i.e. random-feature coordinate descent. In the following we will discuss two other methods we implemented.

## 1.1 Cyclic Coordinate Descent

One natural way of coordinate selection is to choose coordinates in a cyclic ordering:

$$i_k = (k + 1) \bmod d, k = 0, 1, 2, \dots \quad (2)$$

Therefore the cyclic coordinate descent algorithm can be described as follows:

---

### Algorithm 2 Cyclic Coordinate Descent

---

- 1: Initialization:  $k \leftarrow 0, \omega_0 \in \mathbb{R}^d$
  - 2: **repeat**
  - 3:    $i_k = (k + 1) \bmod d$
  - 4:    $\omega^{k+1} \leftarrow \omega^k - \alpha_k [\nabla L(\omega)]_{i_k} \mathbf{e}_{i_k}$
  - 5:    $k \leftarrow k + 1$
  - 6: **until** convergence
- 

## 1.2 Max Gradient Magnitude Coordinate Descent

Another way of coordinate selection is to choose the coordinate that has the steepest direction of gradient to update, i.e. the coordinate that has the largest absolute value of gradient, which will empirically accelerate the convergence rate:

$$i_k = \underset{i}{\operatorname{argmax}} |[\nabla L(\omega)]_i| \quad (3)$$

---

### Algorithm 3 Max Gradient Magnitude Coordinate Descent

---

- 1: Initialization:  $k \leftarrow 0, \omega_0 \in \mathbb{R}^d$
  - 2: **repeat**
  - 3:    $i_k = \underset{i}{\operatorname{argmax}} |[\nabla L(\omega)]_i|$
  - 4:    $\omega^{k+1} \leftarrow \omega^k - \alpha_k [\nabla L(\omega)]_{i_k} \mathbf{e}_{i_k}$
  - 5:    $k \leftarrow k + 1$
  - 6: **until** convergence
- 

We call this method the max gradient magnitude coordinate descent.

### 1.3 Convergence

Both of our methods will undoubtedly converge to the optimal loss as long as given sufficient iteration steps, since our methods are essentially applying general gradient descent to a selected vector coordinate every iteration instead of updating the whole, therefore the value of the objective function will decrease after each iteration and will eventually converge to a global optimum when the function is convex. Different coordinate selection methods determine the direction of gradient descent and only effect the convergence rates.

## 2 Objective Function for Logistic Regression

To simplify the derivation, the bias  $b$  is already absorbed into  $\omega$  in following discussion. For binary labels  $y \in \{0, 1\}$ , we define the model:

$$\Pr(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\omega^\top \mathbf{x}}} \quad (4)$$

$$\Pr(y = 0|\mathbf{x}) = 1 - \Pr(y = 1|\mathbf{x}) = \frac{1}{1 + e^{\omega^\top \mathbf{x}}} \quad (5)$$

Above equations can be rewritten as:

$$\Pr(y|\mathbf{x}) = \left( \frac{1}{1 + e^{-\omega^\top \mathbf{x}}} \right)^y \left( \frac{1}{1 + e^{\omega^\top \mathbf{x}}} \right)^{1-y} \quad (6)$$

Given data  $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}) \in \mathbb{R}^d \times \{0, 1\}$ , the loss function for logistic regression is:

$$L(\omega) = - \sum_{i=1}^n \ln \Pr(y^{(i)}|\mathbf{x}^{(i)}) = \sum_{i=1}^n y^{(i)} \ln(1 + e^{-\omega^\top \mathbf{x}^{(i)}}) + (1 - y^{(i)}) \ln(1 + e^{\omega^\top \mathbf{x}^{(i)}}) \quad (7)$$

Take the gradient of  $L(\omega)$  with respect to  $\omega$ :

$$\nabla L(\omega) = \sum_{i=1}^n \left( \frac{1}{1 + e^{-\omega^\top \mathbf{x}^{(i)}}} - y^{(i)} \right) \mathbf{x}^{(i)} \quad (8)$$

## 3 Experimental Results

We test our method on the `wine` dataset. First we run a standard logistic regression solver from `scikit learn` to obtain the optimal loss  $L^* = 2.37 \times 10^{-6}$ . Then, we run the three methods, random, cyclic, and max gradient magnitude, and record the loss curves. The result are shown in Figure 1.

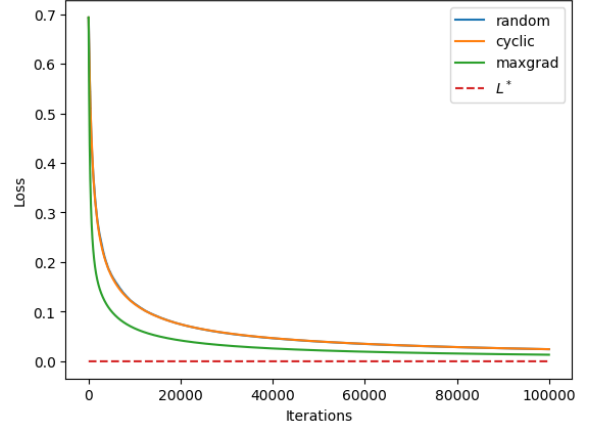


Figure 1: Loss Curves of Different Coordinate Descent Methods.

## 4 Critical Evaluation

From Figure 1, we can find that the losses of all three methods decrease with iteration step and eventually asymptote to  $L^*$ . On top of that, we can see that the loss curve of the max gradient magnitude method has the fastest convergence rate, which indicates the significance of coordinate selection in its effect on convergence rate. The loss curves of cyclic ordering and uniformly random selection are nearly coincident, showing that as the number of iterations are large enough, these two selection methods are essentially traversals of all coordinates and demonstrate relatively low convergence rate.

There exist much scope for improvement of our methods. Note that our loss function does not include the regularization term, which fails to penalize the complexity of  $\omega$  and could lower the convergence rate. For further improvement, we can introduce regularization into our objective function. Moreover, our methods suppose the objective function must be differentiable everywhere, which is not applicable to some common regularization techniques, e.g. the  $l_1$  or Lasso regularization. In such cases, subgradient methods may be applied to solve this problem.

## References

Stephen J Wright. 2015. Coordinate descent algorithms. *Mathematical programming*, 151(1):3–34.

## A Code

```
1 from sklearn.datasets import load_wine
2 from sklearn.linear_model import
   LogisticRegression
```

200	3	from sklearn.metrics import log_loss	250
201	4	from sklearn.preprocessing import	251
		StandardScaler	
202	5	import numpy as np	252
203	6	import matplotlib.pyplot as plt	253
	7		
204	8	data = load_wine()	254
205	9	XTrain, yTrain = data.data[:130], data.	255
		target[:130]	
206	10	XTrain = StandardScaler().fit_transform(	256
		XTrain)	
207			257
208	11		258
209	12	standard_lr = LogisticRegression(penalty	259
		=None, solver='newton-cg').fit(	
210		XTrain, yTrain)	260
211	13	loss_opt = log_loss(yTrain, standard_lr.	261
		predict_proba(XTrain))	
212	14		262
213	15	def sigmoid(x):	263
	16	return 1.0 / (1 + np.exp(-x))	
214	17		264
215	18	def CD(selection, lr=1e-2, iters=100000)	265
		:	
216	19	X = np.insert(XTrain, 0, 1, axis=1)	266
217	20	y = yTrain[:, np.newaxis]	267
218	21	w = np.zeros((X.shape[1], 1))	268
	22	losses = []	
219	23	loss_0 = log_loss(y, sigmoid(X @ w))	269
220	24	losses.append(loss_0)	270
221	25	print('Initial Loss: {}'.format(loss_0	271
		))	
222	26	iter = 0	272
223	27	while iter < iters:	273
	28	grad = np.sum((sigmoid(X @ w) - y) *	274
224		X, axis=0) / X.shape[0]	
225	29	if selection == 'cyclic':	275
	30	index = iter % w.shape[0]	
226	31	elif selection == 'maxgrad':	276
227	32	index = np.argmax(np.abs(grad))	277
	33	elif selection == 'random':	
228	34	index = np.random.randint(0, 14)	278
229	35	w[index] -= lr * grad[index]	279
230	36	loss = log_loss(y, sigmoid(X @ w))	280
	37	losses.append(loss)	
231	38	if (iter+1) % 10000 == 0:	281
232	39	print('Iteration: {}, Loss: {}'.	282
		format(iter+1, loss))	
233	40	iter += 1	283
234	41	return losses	284
	42		
235	43	loss_cyclic = CD(selection='cyclic')	285
236	44	loss_maxgrad = CD(selection='maxgrad')	286
237	45	loss_rand = CD(selection='random')	287
	46		
238	47	step = np.arange(0, len(loss_rand))	288
239	48	loss_std = loss_opt * np.ones(len(	289
		loss_rand))	
240	49	plt.xlabel('Iterations')	290
241	50	plt.ylabel('Loss')	291
242	51	plt.plot(step[:100000], loss_rand	292
		[:100000], label='random')	
243	52	plt.plot(step[:100000], loss_cyclic	293
		[:100000], label='cyclic')	
244	53	plt.plot(step[:100000], loss_maxgrad	294
		[:100000], label='maxgrad')	
245	54	plt.plot(step[:100000], loss_std	295
246		[:100000], linestyle='--', label='\$L	296
247		^*\$')	297
248	55	plt.legend()	298
249	56	plt.show()	299