

Homework1: Regression and Classification

Mingyi Xue*

April 30, 2018

1 Problem1 Regression

The loss function is defined as follows,

$$\omega^* = \arg \min_{\omega} \left\{ \frac{1}{n} \sum_{i=1}^n (x_i^T \omega - y_i)^2 + \frac{\lambda}{2} \|\omega\|^2 \right\} \quad (1)$$

Derive the gradient of equation(1),

$$\frac{\partial f(\omega)}{\partial \omega} = \frac{2}{n} X^T (X\omega - Y) + \lambda \omega \quad (2)$$

Algorithm 1 Gradient Descent with Fixed Step Size

Input: η : step size, ϵ : stopping condition, ω_0 : initial solution , $iter$: number of iteration

Output: ω

```
 $\omega \leftarrow \omega_0$ 
 $r_0 \leftarrow \|\nabla f(\omega)\|$ 
for  $i = 0 \rightarrow iter - 1$  do
   $g = \nabla f(\omega)$ 
  if  $\|g\| \leq \epsilon r_0$  then
    break
  end if
   $\omega \leftarrow \omega_0 - \eta g$ 
end for
```

Problem 1.1

pretreatment

We apply the gradient descent algorithm to dataset “cpusmall”, normalizing dataset so that the value of loss function and gradient won’t exceed the maximum of real number in python.

parameters

“cpusmall” is a dataset of 8192 sample and 12 features. ω_0 is set to a column of 12 zeros in each outer iteration for a specific stepsize. Other parameters are as follows,

*GSP student from Nanjing University

Table 1: default parameters

name	value
λ	1
ϵ	0.001
# of iterations	5000

results

The final output records the value of loss function $f(\omega)$ in the last iteration, whether ω has converged or not within 5000 iterations and the plot of reduction in the value of loss function. The results are as follows,

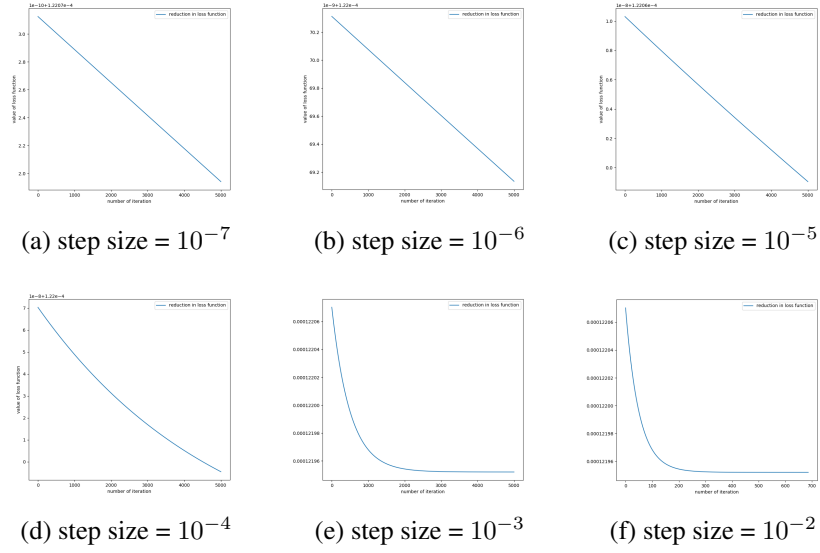


Figure 1: results of different step size

Table 2: results of different step size

step size η	final value of $f(\omega)$	converge
10^{-7}	1.22070×10^{-4}	not yet
10^{-6}	1.22069×10^{-4}	not yet
10^{-5}	1.22059×10^{-4}	not yet
10^{-4}	1.21996×10^{-4}	not yet
10^{-3}	1.21952×10^{-4}	not yet
10^{-2}	1.21952×10^{-4}	yes

Table 3: three best step size and its weight vector

	10^{-4}	10^{-3}	10^{-2}
ω_1	7.49×10^{-5}	7.45×10^{-5}	2.96×10^{-5}
ω_2	9.02×10^{-5}	8.97×10^{-5}	3.56×10^{-5}
ω_3	1.84×10^{-4}	1.83×10^{-4}	7.27×10^{-5}
ω_4	1.61×10^{-4}	1.60×10^{-4}	6.35×10^{-5}
ω_5	1.52×10^{-4}	1.51×10^{-4}	6.00×10^{-5}
ω_6	1.29×10^{-4}	1.28×10^{-4}	5.08×10^{-5}
ω_7	9.91×10^{-5}	9.85×10^{-5}	3.91×10^{-5}
ω_8	1.38×10^{-4}	1.37×10^{-4}	5.44×10^{-5}
ω_9	1.21×10^{-4}	1.21×10^{-4}	4.79×10^{-5}
ω_{10}	4.23×10^{-6}	4.20×10^{-6}	1.68×10^{-6}
ω_{11}	1.49×10^{-4}	1.48×10^{-4}	5.89×10^{-5}
ω_{12}	2.38×10^{-4}	2.36×10^{-4}	9.36×10^{-5}

conclusion

- When the step size is too small, the speed of reduction in loss function will be slow, and ω more iterations and time are needed to converge.
- When the step size is very large, ω may diverge. If we do not normalize the dataset at the beginning, $10^{-7} \sim 10^{-2}$ will be too large for the program and a step size smaller than 10^{-13} is doable.
- As a result, choosing the appropriate step size or learning rate is of significantly importance in machine learning. 10^{-3} and 10^{-2} are acceptable step sizes for this problem.

Problem 1.2

pretreatment

Normalize the dataset “cpusmall” as Problem 1.1.

Apply cross validation on the dataset. Zip matrix X with Y as a large matrix *data*, shuffle *data* and split it evenly into 5 folders, treating 4 folders of which as training set and the rest one as test set.

Run gradient descent algorithm on training set and calculate MSE of the test set.

parameters

Table 4: default parameters

name	value
λ	1
step size η	10^{-2}
# of folders	5
# of iterations	1000

$$\begin{aligned}
MSE &= \frac{1}{n_{test}} \sum_{i=1}^{n_{test}} \{(\mathbf{x}_i^T \boldsymbol{\omega} - y_i)^2\} \\
&= \frac{\|\mathbf{X}\boldsymbol{\omega} - \mathbf{Y}\|^2}{n_{test}}
\end{aligned} \tag{3}$$

results

Table 5: MSE

test folder	value of MSE
1	1.21159×10^{-4}
2	1.21787×10^{-4}
3	1.22295×10^{-4}
4	1.21572×10^{-4}
5	1.22356×10^{-4}
<i>mean</i>	1.21834×10^{-4}

Problem 1.3

parameters

Table 6: default parameters

name	value
ϵ	0.001
# of iterations	2000

results

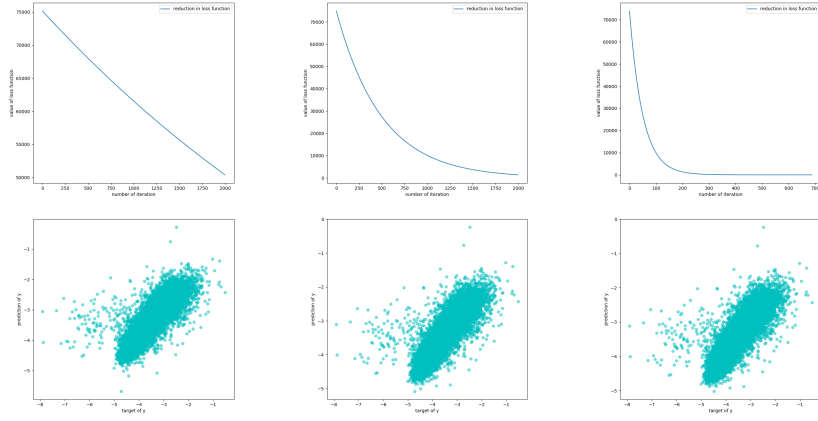
The final output records the value of test MSE, whether $\boldsymbol{\omega}$ has converged or not within 2000 iterations, the plot of reduction in the value of loss function and the scatter plot of estimated \hat{y} vs. target y .

Results are as follows,

Table 7: results of different step size

stepsize	value of test MSE	final value of $f(\boldsymbol{\omega})$	converge
10^{-4}	0.156556	50392.03	not yet
10^{-3}	0.150899	1370.50	not yet
10^{-2*}	0.150888	0.7328	yes
10^{-1}	—	—	diverge

* This step size is best.



(a) step size = 10^{-4}

(b) step size = 10^{-3}

(c) step size = 10^{-2}

2 Problem2 Classification

Problem 2.1

The loss function for logistic regression is defined as follows,

$$\omega^* = \arg \min_{\omega} \left\{ \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \omega^T x_i}) + \frac{\lambda}{2} \|\omega\|^2 \right\} \quad (4)$$

Derive the gradient of equation(4),

$$\frac{\partial f(\omega)}{\partial \omega} = \frac{1}{n} \sum_{i=1}^n \frac{-y_i}{1 + e^{-y_i \omega^T x_i}} x_i + \lambda \omega \quad (5)$$

In python, we can use numpy.ndarray broadcast to simplify calculation and save time,

$$\begin{aligned} K &= np.array(X @ \omega) * np.array(Y) \\ M &= -np.array(Y) * \frac{1}{1 + e^K} \\ \frac{\partial f(\omega)}{\partial \omega} &= \frac{1}{n} X^T M + \lambda \omega \end{aligned} \quad (6)$$

Problem 2.2

The initial ω_0 is randomly normal distributed.

Table 8: default parameters

dataset	accuracy
training set	69.92%
test set	68.28%