

Final Project

Convex Clustering

220041040 Qiuyan YUAN

221041041 Yining ZHONG

117030035 Xinrong JIANG

116020277 Sixuan XIANG

220041032 Hao HAN

2021.12.28

—
MDS6106

Introduction to Optimization
—

Abstract

Convex clustering is an interesting topic worthy of study. It can solve some of the problems of traditional unsupervised clustering methods, such as poor model effect, the need to manually set some important parameters, and so on. In this project, by setting reasonable convex optimization goals, we can make each sample point not far from the cluster center to which it belongs, and at the same time make it possible to distinguish between different cluster centers. We try to choose different parameters to define our problem, use AGM, Newton CG, BFGS, L-BFGS, and other different methods to optimize this problem, and the difference between using different matrix norms for the penalty term of the loss function, using Mini batch SGD and other methods is studied. Some meaningful and interesting observation results are obtained.

1. Introduction

This report is going to investigate different optimization models and utilize minimization methodologies to solve convex clustering problems for large-scale unsupervised learning.

In some pattern recognition problems, the training data consists of a set of input vectors X without any corresponding target values. The goal in such unsupervised learning problems may be to discover groups of similar examples within the data, where it is called clustering.

Let $\mathbb{R}^{d \times n} \ni A = (a_1, a_2, \dots, a_n)$ be a given data matrix with n observations and d features.

The convex clustering model for these n observations solves the following convex optimization problem:

$$\min_{x \in \mathbb{R}^{d \times n}} \frac{1}{2} \sum_{i=1}^n \|x_i - a_i\|^2 + \lambda \sum_{i=1}^n \sum_{j=i+1}^n \|x_i - x_j\|_p \quad (1)$$

where $\lambda > 0$ is a regularization parameter and $\|\cdot\|_p$ denotes the p -norm. Here, $\|\cdot\|$ denotes the standard Euclidean norm (as usual). The p -norm above with $p \geq 1$ ensures the convexity of the model. Typically, p is chosen to be 1, 2, or ∞ . (Since this model is strongly convex, the optimal solution for a given positive λ is unique.)

After solving (1) and obtaining the optimal solution $X^* = (x_1^*, \dots, x_n^*)$, we assign a_i and a_j to the same cluster if and only $\|x_i^* - x_j^*\| \leq \varepsilon$ for a given tolerance $\varepsilon > 0$.

After several group discussions, we finish the final project together, and the details of the work division are as follows.

Participation	name
Data preparation, AGM, and Weighted AGM	Qiuyan YUAN
Newton-CG, and Weighted Newton-CG	Yining ZHONG
BFGS, and L-BFGS	Sixuan XIANG
SGD, and Mini-batch SGD	Hao HAN
Norm regularization	Xinrong JIANG

Table 1: Individual work

2. Data Preparation

- Data preparation
 - a) Self-generated Data Sets

We generated seven 2-dimension datasets with different shapes and sizes for the project. Figure 1 displays the distribution of them on a two-dimensional plane.

Dataset	n	classes	description
200_random_4	200	4	4 random center points, the spread of the data points is controlled by the selected variance.
2000_random_4	2000	4	4 random center points, the spread of the data points is controlled by the selected variance.

200_moon_2	200	2	Moon-shaped data sets that generated by <code>sklearn.datasets.make_moons()</code> .
200_circle_2	200	2	Ring-shaped data sets that generated by <code>sklearn.datasets.make_moons()</code> .
200_dot_3	200	3	3 random center points, the spread of the data points is controlled by the same variance.
2000_dot_3	2000	3	3 random center points, the spread of the data points is controlled by the same variance.
200_dot_6	200	6	6 random center points, the spread of the data points is controlled by the same variance.

Table 2: Self-generated Data Sets

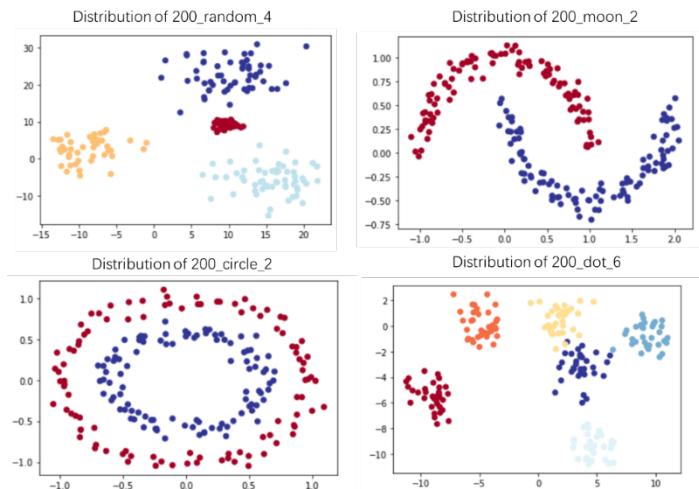


Figure 1: Distribution of 4 types of artificial data sets

b) Real Data Sets

All the provided datasets are used for the project. The basic information of the datasets is summarized in Table x. However, we do not have enough computing resources to run the algorithm under the large data set. We mainly use the wine data set.

Dataset	d	n	classes
wine	13	178	3
vowel	10	528	11
segment	19	2310	7
mint	784	60000	10

Table 3: Real Data Sets

3. Experiment

- AGM
 - Introduction to Accelerated Gradient Method

The Accelerated Gradient Method is based on acceleration techniques (extrapolation step) and Momentum. In this project, we apply the Accelerated Gradient Method (AGM) with fixed step size and basic extrapolation strategy.

Algorithm1. Accelerated Gradient Method (AGM)

Begin

Initialization: Choose a point $x_0 \in R^n$ and set $x^{-1} = x^0$

for $k = 0, 1, 2, \dots$ **do**

 select an extrapolation parameter β_k :

$$\beta_k = \frac{\theta_k(1 - \theta_{k-1})}{\theta_{k-1}} = \frac{\theta_k}{\theta_{k-1}} - \theta_k \text{ and } \theta_{-1} = \theta_0 = 1, \theta_k = t_k^{-1}$$

$$t_k = \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}, \beta_k = \frac{t_{k-1} - 1}{t_k}, t_{-1} = t_0 = 1$$

 compute the step $y^{k+1} = x^k + \beta_k(x^k - x^{k-1})$

 select step size $\alpha_k = 1/L, L = 1 + n\lambda/\delta$

 set $x^{k+1} = y^{k+1} - \alpha_k \nabla f(y^{k+1})$.

if $\|\nabla f(x^{k+1})\| \leq tol$ **do**

return x^{k+1}

end

By applying the accelerated gradient method to test datasets, we can find that for the datasets with large dimensions and sample size, AGM can't converge with the low tolerance level. The reason is that only based on the first-order information, our step update is of low efficiency. Even the AGM methods have relaxed the Armijo condition, the decrease of the objective function value can't be guaranteed.

■ Experiment

a) Clustering results on different datasets

 1) Cluster Result

Six artificial data sets are used to test the AGM algorithm. They are 200_random_4, 200_dot_3, 200_dot_6, 2000_random_4, 200_moon_2, 200_circle_2. According to the experimental results, the AGM algorithm performs well on small, scattered data sets, such as 200_random_4 and 200_dot_3. In a small data set with 6 clusters, 200_dot_6, the clustering results are reasonable but not completely classified into 6 classes. The left image of the figure below represents the original distribution of the artificial data set, and the right image represents the clustering result of the AGM.

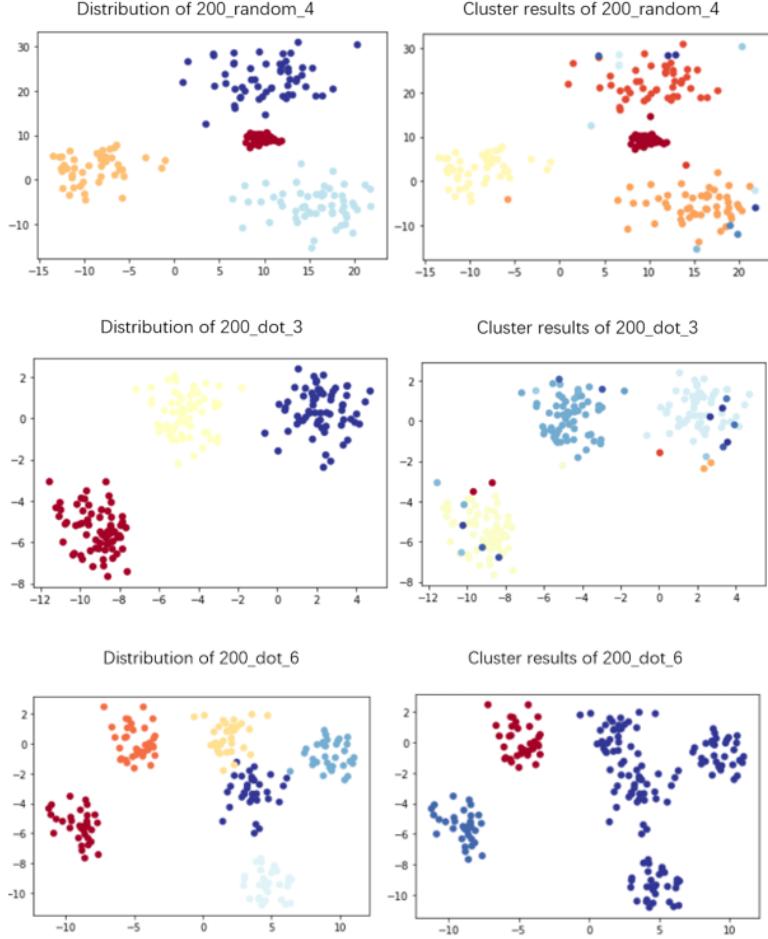


Figure 2: original distribution and cluster results

However, on large data sets, such as the 2000_random_4 data set which has 2-dimension, 2000 data points, and 4 clusters, the AGM algorithm cannot converge. After Iteration reaches 30000 and the running time exceeds 10 hours, the norm gradient is still greater than 1000, and effective clustering cannot be performed.

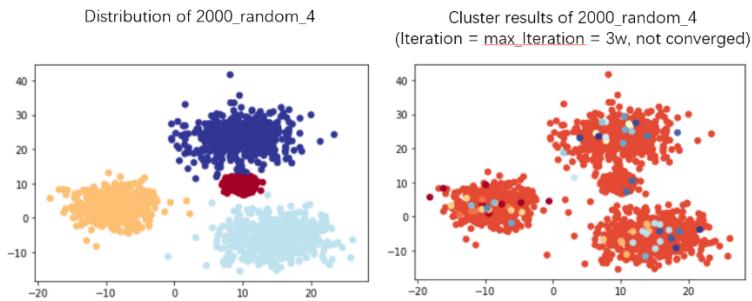


Figure 3: original distribution and cluster results of 2000_random_4

On data sets with special shapes, such as moon-shaped and circular data sets, although AGM converges, the clustering effect is not significant, and the two types of data cannot be distinguished.

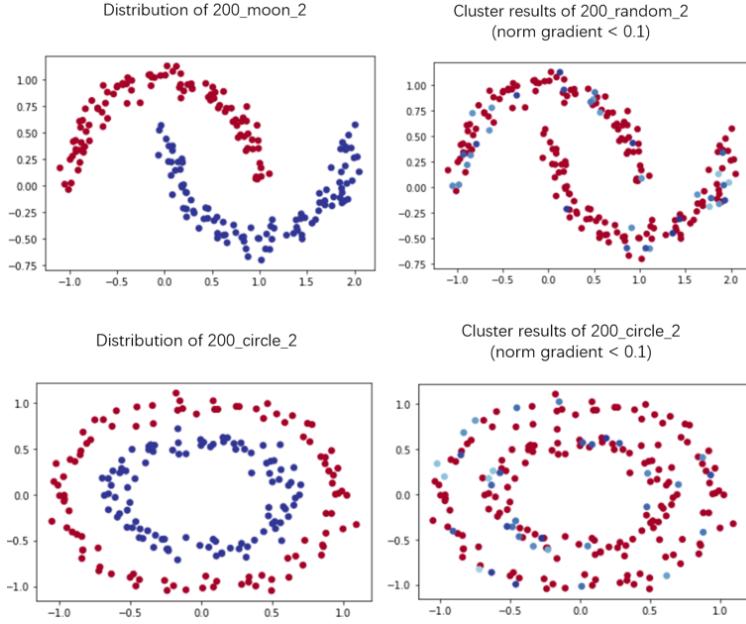


Figure 4: original distribution and cluster results

2) CPU Time and Convergence rate

AGM used on large datasets such as 2000_random_4 and circle shape datasets such as 200_circle_2 cannot converge even though the iteration reaches 30000. The CPU time for different data with 200 data points is around 10 minutes with AGM. It is too slow to apply to real data sets.

Dataset Name	Iteration Num	CPU Time	Norm Gradient
200_random_4	13486	358.22s	< 0.1
200_dot_3	13473	752.88s	< 0.1
200_dot_6	11530	635.88s	< 0.1
2000_random_4	30000	62955.85s	1171
200_moon_2	26651	587.12s	< 0.1
200_circle_2	30000	641.45s	1.04

Table 4: performance of AGM on different datasets

- Newton-CG

- Introduction to Newton-CG Method

The Newton-CG Method is a truncated Newton method, which is to approximate the solution of the Newton equation to get the idea of the inspection direction because solving the equation is time-consuming. We require an approximate solution to be found with a certain number of iterations.

In this method, we choose CG Method to approximate the Newton function $\nabla^2 f_k p_k + \nabla f_k = 0$

Algorithm2. Newton-CG Method

Given initial point x_0

for $k = 0, 1, 2, \dots$ **do**

Define tolerance $\epsilon_k = \min(0.5, \sqrt{||\nabla f_k||}) ||\nabla f_k||$;

Set $z_0 = 0, r_0 = \nabla f_k, d_0 = -r_0 = -\nabla f_k$;

for $j = 0, 1, 2 \dots$

if $d_j^T B_k d_j \leq 0$
if $j = 0$
return $p_k = -\nabla f_k;$
else
return $p_k = z_j;$

Set $\alpha_j = r_j^T r_j / d_j^T B_k d_j;$

Set $z_{j+1} = z_j + \alpha_j d_j;$
Set $r_{j+1} = r_j + \alpha_j B_k d_j;$
if $\|r_{j+1}\| < \epsilon_k$
return $p_k = z_{j+1};$

Set $\beta_{j+1} = r_{j+1}^T r_{j+1} / r_j^T r_j;$

Set $d_{j+1} = -r_{j+1} + \beta_{j+1} d_j;$
end (for)

Set $x_{k+1} = x_k + \alpha_k p_k,$ where α_k satisfies the Wolfe, Goldstein, or Armijo backtracking conditions (using $\alpha_k = 1$ if possible);

end

■ Calculate the gradient and hessian

$$\min_{X \in \mathbb{R}^{d \times n}} f_{\text{clust}}(X) := \frac{1}{2} \sum_{i=1}^n \|x_i - a_i\|^2 + \lambda \sum_{i=1}^n \sum_{j=i+1}^n \varphi_{\text{hub}}(x_i - x_j)$$

We can use a vector to represent X , so $X \in \mathbb{R}^{dn}, x_i \in \mathbb{R}^d, i = 1, 2, \dots, n, X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

Then the gradient of x_i is,

$$\frac{\partial f}{\partial x_i} = x_i - a_i + \lambda \sum_{j=1}^n \varphi'_{\text{hub}}(x_i - x_j), \quad \varphi'_{\text{hub}}(y) = \begin{cases} \frac{1}{\delta} y & |y| \leq \delta \\ \frac{y}{|y|} & |y| > \delta \end{cases}$$

The whole gradient can be represented like, $\frac{\partial f}{\partial X} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$

Then the hessian matrix is $\begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{bmatrix},$

where $\frac{\partial^2 f}{\partial x_i \partial x_i} = I + \lambda \sum_{j=1}^n \varphi''_{hub}(x_i - x_j)$ and $\frac{\partial^2 f}{\partial x_i \partial x_j} = -\lambda \varphi''_{hub}(x_i - x_j)$.

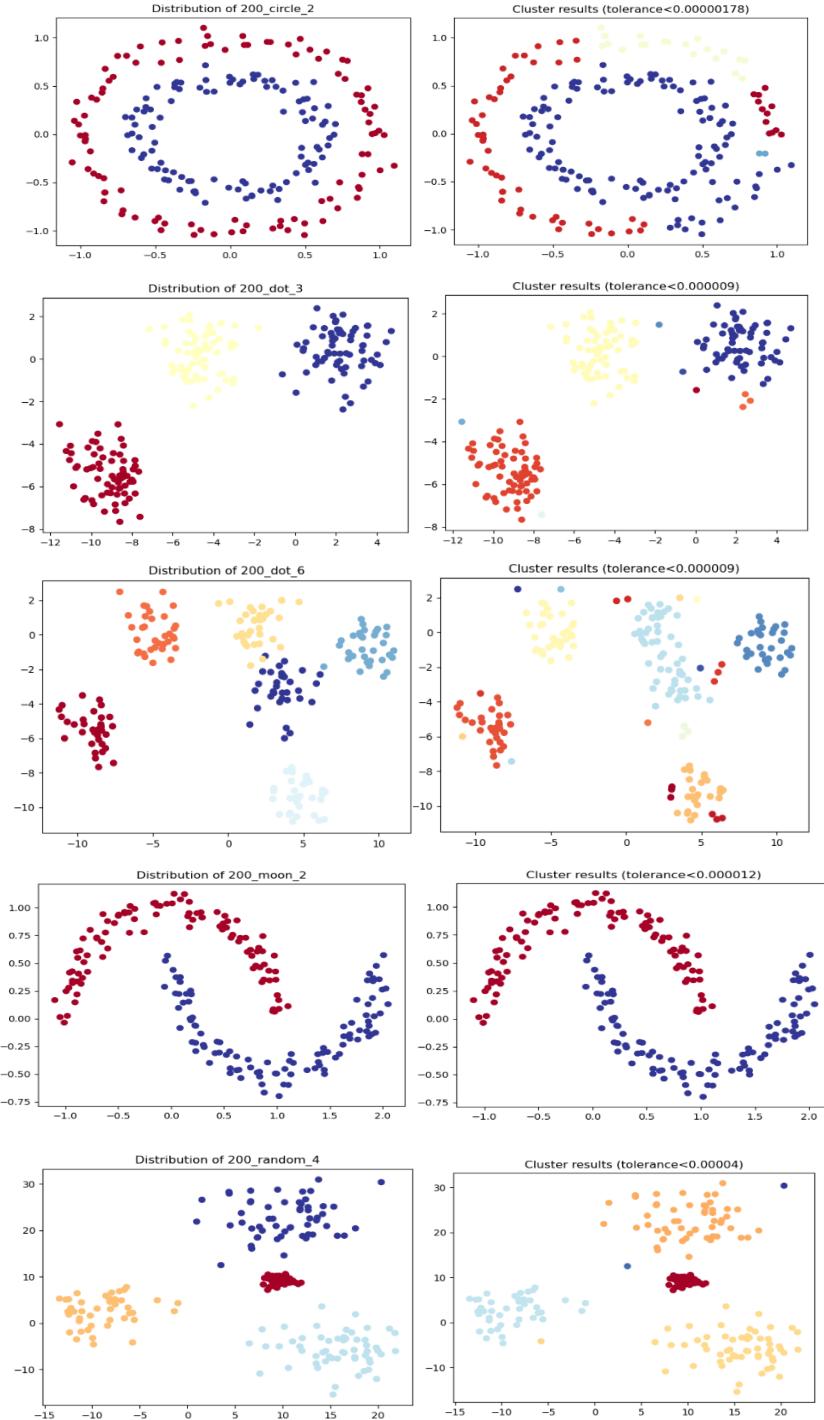
For which $\varphi''_{hub}(y) = \begin{cases} \frac{1}{\delta} I & \|y\| \leq \delta \\ \frac{I}{\|y\|} - \frac{y \cdot y^T}{\|y\|^3} & \|y\| > \delta \end{cases}$

According to these formulas, try to implement them in python.

■ Experiment

a) Clustering results on different datasets

1) Cluster Result



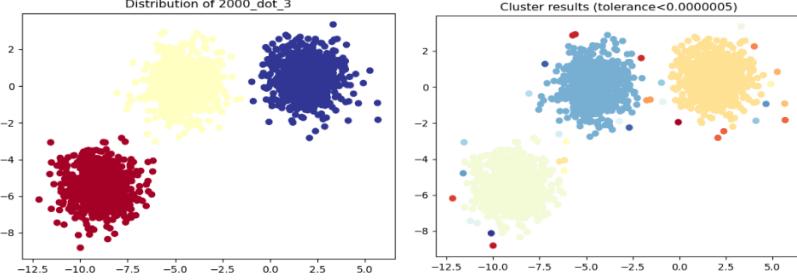


Figure 7: The clustering results

These are the results visualized after clustering in a different dataset. The pictures on the left are the distribution of the original data, and pictures on the right side are clustering results using Newton-CG Method. After experiments, we can find in the circle dataset, the result is not quite good although after too many tests trying on different tolerance. However, the result in the moon dataset is perfect in that all points are clustered in the right classes. In the other datasets, the solutions are also ok that except for some outliers, most of the classes are distinguished.

All these are visualizations of artificial data sets, but for real datasets, they can't be visualized like above. Then we use the convergence line and relative error to evaluate.

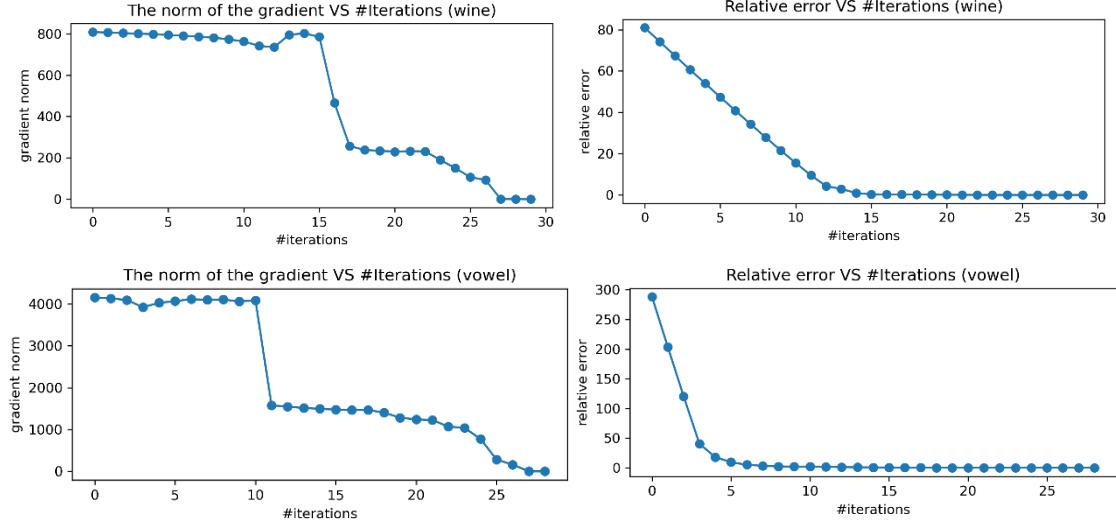


Figure 8: The relationship between Relative Error/Gradient Norm and #Iterations

From the convergence and relative error lines we can see that using Newton-CG, there is a vain and rapid decline during the iterations. Then it gradually converges to 0. After we can compare the CPU time and the number of iterations.

2) CPU Time and iteration number

We can compare the performance of Newton-CG on different data sets.

Dataset Name	Iteration Num	CPU Time	Norm Gradient
200_circle_2	15	44 s	<6.7e-5
200_dot_3	60	2 min 46 s	<2e-5
200_dot_6	76	3 min 20 s	<2e-5
200_moon_2	19	1 min	<5.7e-5
200_random_4	193	8 min 31 s	<7.9e-6
2000_dot_3	14	86 min 37 s	<2.2e-6

Table 5: Artificial datasets (no weight, tol=1e-3, lamda=0.5)

Dataset Name	d	n	Iteration Num	CPU Time	Norm Gradient	Tol
wine	13	178	29	1.63 min.	<1.7e-5	1e-3
vowel	10	528	28	20.19min	<0.0006	1e-3

Table 6: Real datasets (no weight, tol=1e-3 lamda=0.5)

We can see that although the number of iterations of Newton-CG is not too large, the size of the dataset can lead to a long time to run an iteration, like 2000_dot_3, which needs more than one hour to finish, so we only run one large size dataset.

b) The impact of regularization parameters λ

Set $\lambda = 0.05, 0.1, 0.2, 0.5, 0.8, 1, 3, 5, 7.5, 10$. And using the dataset “200_moon_2”, “200_dot_6” to run different λ and compare them.

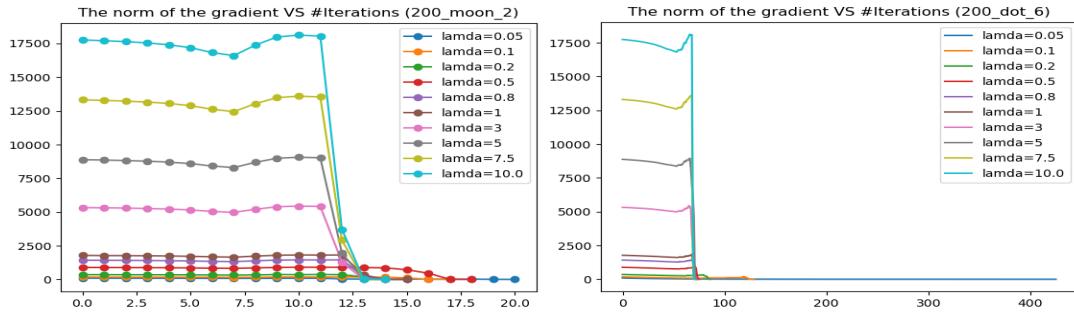


Figure 9: The relationship between Relative Error/Gradient Norm and #Iterations

λ	0.05	0.1	0.2	0.5	0.8	1	3	5	7.5	10
200_moon_2	21	18	16	19	16	16	15	15	15	15
200_dot_6	426	130	87	76	74	73	73	74	72	72

Table 7: The iterations of different λ

From the two plots of different datasets with different λ , we find that the bigger λ , the bigger the initial gradient and the more dramatic decline in convergence. And we compare the number of iterations then find that the slower λ is, the more iterations need. And this is also understandable, the smaller λ means that the smaller the difference in different directions, so it may need more attempts to find a better direction.

● Weighted Models

■ Introduction of weighted model

Now, we change the objective function, giving weights based on the input data A, so we also need to change the gradient and hessian. But in fact, we only need to add weight to the functions.

The gradient of x_i need to be changed like

$$\frac{\partial f}{\partial x_i} = x_i - a_i - \lambda \sum_{j=1}^{i-1} w_{ji} \varphi''_{hub}(x_j - x_i) + \lambda \sum_{j=i+1}^n w_{ij} \varphi''_{hub}(x_j - x_i).$$

Then the whole formula is,

$$\frac{\partial f}{\partial x_i} = x_i - a_i + \lambda \sum_{j=1}^{i-1} w_{ji} \varphi''_{hub}(x_i - x_j) + \lambda \sum_{j=i+1}^n w_{ij} \varphi''_{hub}(x_j - x_i)$$

The $\frac{\partial^2 f}{\partial x_i \partial x_i} = I + \lambda \sum_{j=1}^n w \varphi''_{hub}(x_i - x_j)$

And $\frac{\partial^2 f}{\partial x_i \partial x_j} = -\lambda w \varphi''_{hub}(x_i - x_j)$

The $w = \begin{cases} w_{ij} & j > i \\ w_{ji} & j < i \\ 0 & j = i \end{cases}$

■ AGM V.S. Weighted AGM

Use wine data set to test AGM and weighted AGM ($\lambda = 0.05$, $tol = 1e-1$). The norm gradient of AGM started to drop significantly after 6000 iterations, while that of the weighted methods drop steadily from the beginning. Although the number of iterations of the two methods is similar, the weighted method is faster with less CPU time on Wine. The weighted AGM can achieve a better result and converge faster.

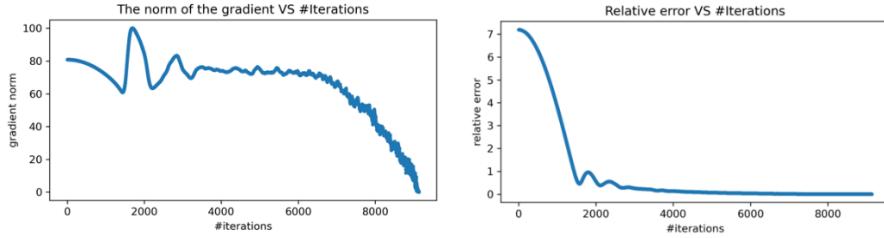


Figure 5: norm gradient and relative error of AGM

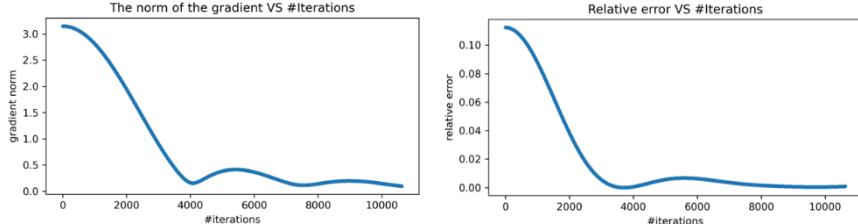


Figure 6: norm gradient and relative error of Weighted AGM

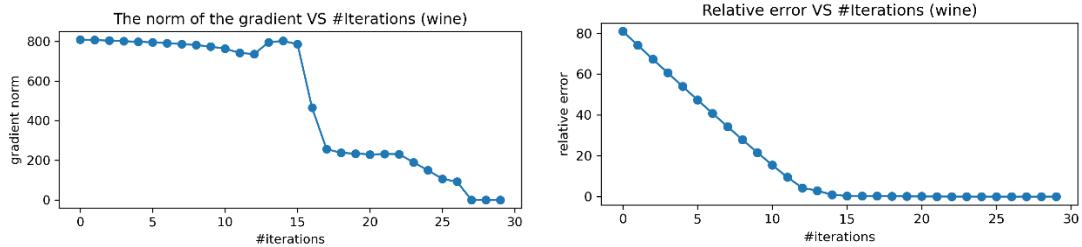
Method	Iteration Num	CPU Time	Norm Gradient
Weighted	10633	47.15 min	<0.1
No weighted	9136	169.01 min	<0.1

Table 8: compare weighted and no weighted AGM

■ Newton-CG V.S. Weighted Newton-CG

- 1) Clustering results (Wine (lamda=0.05, tol=1e-3, k=5))

No weighted



Weighted

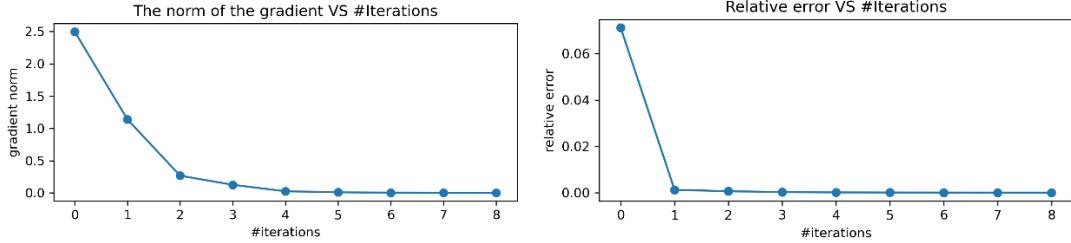


Figure 10: The relationship between Relative Error/Gradient Norm and #Iterations

Comparing the weighted model and no weighted model, we can find that the no weighted model needs more iteration times and the gradients change from a large number to 0, but for the weighted model, the gradient changes are not so violent because of adding weights.

2) CPU time

Method	Iteration Num	CPU Time	Norm Gradient
No weighted	43	2 min 48 s	<0.0003
Weighted	8	25 s	<0.0008

Table 9: The comparison between weighted model and no weighted model

From the CPU time and the number of iterations, we can also find that the weighted model can do better on the wine dataset.

● Weighted AGM V.S. Weighted Newton CG

Weighted Newton CG performs far better than weighted AGM, no matter on iteration num or on CPU time. AGM is only based on the first-order information, the step update is of low efficiency. Because of using the second-order derivative information, the number of iterations and CPU time are dramatically better than AGM.

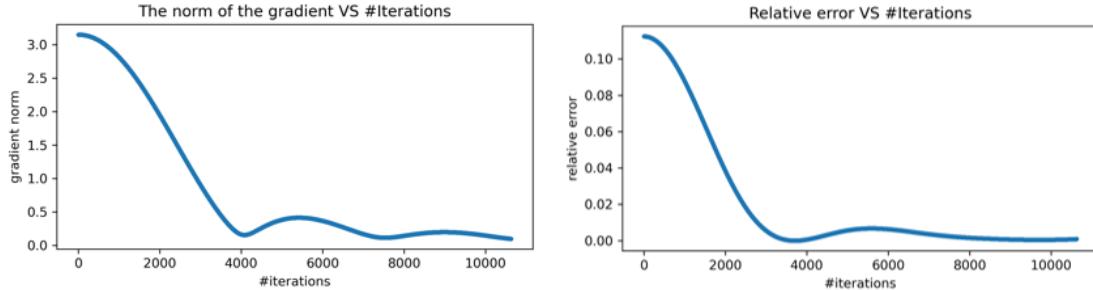


Figure 11: norm gradient and relative error of Weighted AGM

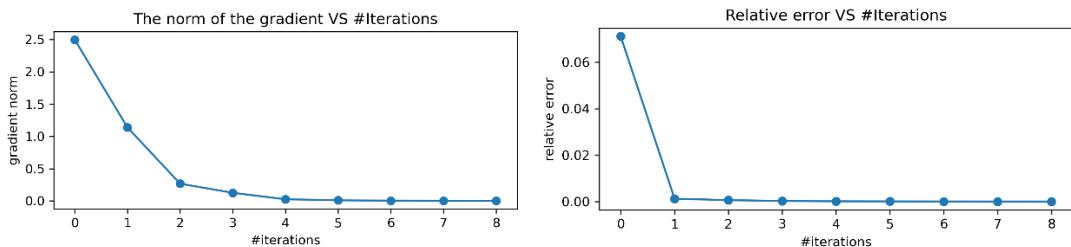


Figure 12: norm gradient and relative error of Weighted Newton CG

Method	Iteration Num	CPU Time	Norm Gradient
Weighted AGM	10633	47.15 min	<0.1
Weighted Newton CG	8	0.42min	<0.0008

Table 10: compare weighted AGM and weighted Newton CG

4. Extensions and Stochastic Optimization

- BFGS, L-BFGS

- Introduction

Newton's method is very computationally expensive. While the computation of the gradient scales as $O(n)$, the computation of the inverse Hessian scales as $O(n^3)$ (computing the Hessian scales as $O(n^2)$, inverting it as scales as $O(n^3)$). As the dimensions of our problem increase, the overhead in memory and time gets out of hand very quickly. It's clear at this point that the benefit of an increased convergence rate will be far outweighed by the large cost of the additional computation time. To address this issue, quasi-Newton methods are introduced, where we can have faster convergence than gradient descent, but a lower operational cost per iteration than Newton's method.

In quasi-Newton methods, instead of computing the actual Hessian, we just approximate it with a positive-definite matrix B , which is updated from iteration to iteration using information computed from previous steps. We immediately see that this scheme would yield a much less costly algorithm compared to Newton's method because instead of computing a large number of new quantities at each iteration, we're largely making use of previously computed quantities.

The specific update for B is given by the specific quasi-Newton method used.

Here, we will focus on one of the most popular methods, known as the BFGS method:

The Globalized BFGS-Method with Armijo step size

Begin

Initialization: $x^0 \in R^n$ and asymmetric pos. def. matrix $H_0 \in R^{n*n}$, and choose $\sigma, \gamma \in (0,1)$.

for $k = 0, 1, 2, \dots$ **do**

$d^k = -H_k \nabla f(x^k)$

Computing the stepsize α_k using Armijo line search

$x^{k+1} = x^k + \alpha_k d^k$

if $\|\nabla f(x^{k+1})\| \leq tol$ **do**

return x^{k+1}

else do

$s^k = x^{k+1} - x^k, \quad y^k = \nabla f(x^{k+1}) - \nabla f(x^k)$

if $(s^k)^T y^k < 0$ **do**

$H_{k+1} = H_k$

else do

$$H_{k+1} = H_k + \frac{(s^k - H_k y^k)(s^k)^T + s^k (s^k - H_k y^k)^T}{(s^k)^T y^k} - \frac{((s^k - H_k y^k)^T y^k)}{((s^k)^T y^k)^2}$$

end

■ Experiment

Initial settings:

tol	maxiter	lambda	delta	s	sigma	gamma
1e-3	1e5	0.5	0.001	1	0.5	0.1

Table 11: Global variables

a) Tests on Artificial Data Set

We consider two artificial data set: data_200_dot_3 and data_200_dot_6. Visualization of the BFGS

performance on these two data sets is shown below. We find that BFGS gives a satisfying classification for both problems. Significant decline of the gradient norm occurs within ten iterations and the relative error approach zero within five iterations.

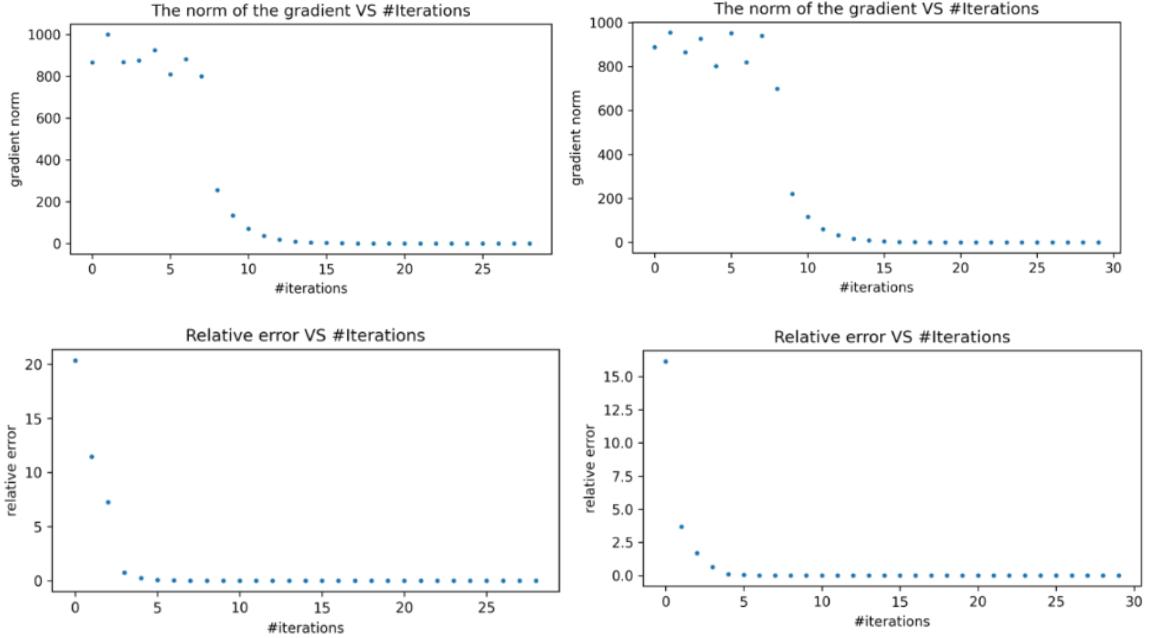


Figure 13: The relationship between Relative Error/Gradient Norm and #Iterations
(left: 200_dot_3; right: 200_dot_6)

b) Comparison with Newton-CG

Utilizing Newton-CG, running the experiments on the same data sets (data_200_dot_3 and data_200_dot_6) gives the following results.

As we expected, the BFGS method greatly reduces CPU time. In both tests, the running time of BFGS is less than one second, while Newton-CG needs two to three minutes. Moreover, the number of iterations of BFGS is also less than half of Newton-CG.

Dataset Name	Iteration Number		CPU Time		Norm Gradient	
	Newton-CG	BFGS	Newton-CG	BFGS	Newton-CG	BFGS
200_dot_3	60	28	2 min 46 s	0.13 s	< 2e-5	< 2e-5
200_dot_6	76	29	3 min 20 s	0.14 s	< 2e-5	< 2e-5

Table 12: Newton-CG vs. BFGS

So far, we have built the approximations B_k and H_k of the Hessian and inverse of the Hessian as full $n \times n$ matrices. However, if n is large, this might not be possible or too expensive. However, we do not need H_k explicitly, we are only interested in $d^k = -H_k \nabla f(x^k)$.

For a vector $v \in R^n$, it holds that:

$$H_{k+1}v = H_k \left[v - \frac{(s^k)^\top v}{(s^k)^\top y^k} \cdot y^k \right] + \frac{(s^k)^\top v}{(s^k)^\top y^k} s^k - \frac{(y^k)^\top H_k \left[v - \frac{(s^k)^\top v}{(s^k)^\top y^k} \cdot y^k \right]}{(s^k)^\top y^k} s^k$$

Hence, $H_{k+1}v$ can be computed recursively: $\beta_k = \frac{(s^k)^\top v}{(s^k)^\top y^k}$, $q^k = v - \beta_k y^k$, $p^k =$

$$H_k q^k, \quad H_{k+1}v = p^k + \left[\beta_k - \frac{(y^k)^\top p^k}{(s^k)^\top y^k} \right] s^k$$

The idea of the limited memory BFGS method is to store only the last m pairs.

$$\{y^{k-m}, y^{k-m+1}, \dots, y^k\} \quad \text{and} \quad \{s^{k-m}, s^{k-m+1}, \dots, s^k\}$$

and to use this information to build a much cheaper BFGS-type approximation of the Hessian.

c) Experiments

In this part, we are going to compare the BFGS method and L-BFGS method by testing the wine data set.

Both methods converge eventually. However, the L-BFGS method converges slower than the BFGS method. In the first thirty iterations, the gradient decline of L-BFGS method is inapparent. In the BFGS method, we find a significant decline within 7 iterations. It is a trade-off between accuracy and saving storage space. To save storage space, the L-BFGS only store and use the last m pairs of $\{y^{k-m}, y^{k-m+1}, \dots, y^k\}$ and $\{s^{k-m}, s^{k-m+1}, \dots, s^k\}$. Since the L-BFGS does not utilize full information, and this might be the reason that it didn't find the right way until the thirtieth iteration.

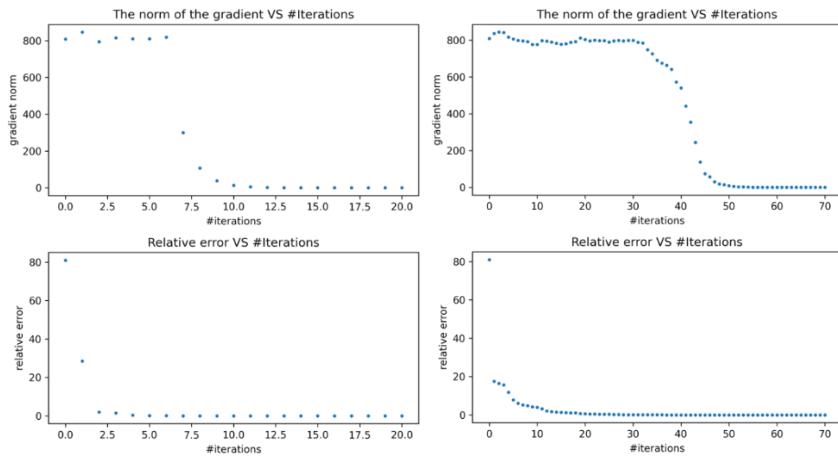


Figure 14: BFGS vs. L-BFGS
(left: BFGS; right: L-BFGS)

- SGD, Mini-Batch SGD

- Introduction

- a) Stochastic gradient descent[1]

Stochastic gradient descent (often abbreviated SGD) is an iterative method for optimizing an objective function with suitable smoothness properties. It can be regarded as a stochastic approximation of gradient descent optimization since it replaces the actual gradient (calculated from the entire data set) with an estimate thereof (calculated from a randomly selected subset of the data).

Specifically, we often use empirical risk $E_n(f)$ to measure the training set performance in machine learning.

$$E_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

To minimize the empirical risk, gradient descent (GD) is often used. Each iteration updates the weights w based on the gradient of $E_n(f_w)$.

$$w_{t+1} = w_t - \eta \frac{1}{n} \sum_{i=1}^n \nabla_w Q(z_i, w_t)$$

where η is an adequately chosen learning rate.

The stochastic gradient descent algorithm (SGD) is a drastic simplification. Instead of computing the gradient of $E_n(f_w)$ exactly, each iteration estimates this gradient based on a single randomly picked example z_t :

$$w_{t+1} = w_t - \eta_t \nabla_w Q(z_t, w_t)$$

The stochastic process $\{w_t, t = 1, \dots\}$ depends on the examples randomly picked at each iteration. Since the stochastic algorithm does not need to remember which examples were visited during the previous iterations, it can process examples on the fly in a deployed system.

b) Mini-Batch Stochastic Gradient Descent

During training, Mini-Batch Stochastic Gradient Descent (mini-batch SGD) processes a group of examples per iteration. For notational simplicity, assume that n is divisible by the number of mini-batches m . Then we partition the examples into m mini-batches, each of size $b = n/m$.

Given a random minibatch $I \subset \{1, \dots, n\}$ of size b , we can define the objective function on I as

$$Q_I(w) = \frac{1}{|I|} \sum_{i \in I} Q_i(w) \text{ we have } Q(w) = \mathbf{E}_I[Q_I(w)]$$

In the simple case that $\Omega = \mathbb{R}^d$ the mini-batch SGD employs the following stochastic update rule: at each iteration t , we pick mini-batch $I_t \subset \{1, \dots, n\}$ of size b at random and let

$$w_t = w_{t-1} - \eta_t \nabla Q_{I_t}(w)$$

Although b times more examples are processed in an iteration, the mini-batch training can converge much slower than that of standard SGD with the same number of processed examples. In practice, the convergence rate slows down dramatically in terms of the number of examples processed, when we use a large mini-batch size [2].

■ Experiment

a) Data

After preliminary experiments, the gradient descent algorithm takes a long time (perhaps more than 3 hours) to reach convergence under large data sets (the amount of data usually exceeds several thousand). We do not have enough computing resources to run the algorithm under the large data set, so we consider using the wine data set and the artificial data set (200 two-dimensional points, 6 cluster centers) generated in the data preparation stage as the experimental data sets.

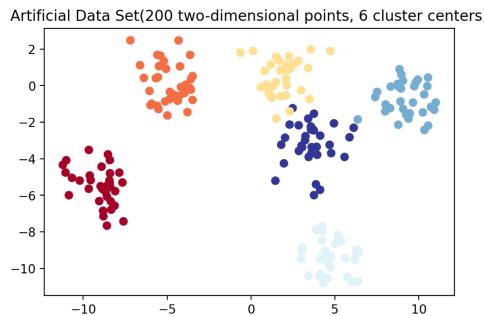


Figure 15: Display of true labels of the artificial data set

b) The impact of different batch sizes

In the Mini-Batch SGD algorithm, the batch size is a very important hyperparameter. We are here to discuss the influence of different batch sizes on convergence speed, time, convergence effect, and other factors.

Except for batch size, other hyperparameter settings of artificial data set and wine data set are as follows. We realized that convergence results usually require decreasing learning rates satisfying the conditions $\sum_t \eta_t^2 < \infty$ and $\sum_t \eta_t = \infty$ [2]. So we choose $\alpha = 1/k$ ($k = \#$ iterations) as our learning rate.

hyperparameter	artificial data set	wine data set
λ	0.1	0.1
δ	0.0001	0.0001
Learning rate (α)	1/k ($k = \#$ iterations)	0.1/k ($k = \#$ iterations)
Gradient tolerance (tol)	0.1	0.1
Max iteration number	30000	20000

Table 13: hyperparameter settings of artificial data set and wine data set

1) CPU Time

Iteration Num	BS=1(SGD)	BS=25	BS=50	BS=100	BS=200(GD)
10000	0.86s	11.79s	23.02s	45.71s	90.67s
20000	1.72s	23.56s	46.19s	91.59s	181.67s
30000	2.60s	35.35s	69.23s	137.30s	272.65s

Table 14: relationship between CPU Time and #iterations/batch size(manual data set)

Iteration Num	BS=1(SGD)	BS=50	BS=100	BS=178(GD)
10000	1.49s	38.15s	75.14s	132.47s
20000	2.95s	76.27s	150.23s	265.15s

Table 15: the relationship between CPU Time and #iterations/batch size (wine data set)

When the batch size is 1, Mini-Batch SGD is SGD. When the batch size is equal to the number of data in the data set, Mini-Batch SGD is GD. From the table, we can find that when the batch size is unchanged, the time taken for every 10,000 iterations is approximate. When the batch size is changed, the larger the batch size, the more CPU time is consumed, and the time consumed is approximately linearly related to the batch size.

2) Convergence rate

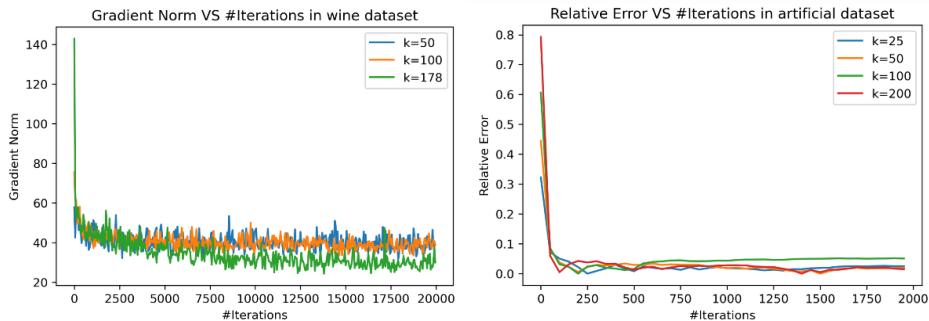


Figure 16: The relationship between Relative Error/Gradient Norm and #Iterations

We can see from the figure that, in general, the larger the batch size, the smaller the relative error, and the smaller the gradient norm after the same number of iterations. However, we must take into account that when the batch size is different, the time spent in each iteration is different. Therefore, the batch

size should be a certain number between [1, the amount of data in the data set], so that the Mini-Batch SGD algorithm can achieve the effect of GD in a relatively short time. If the batch size is too small, the convergence will be very unstable. If the batch size is too large, it will take a lot of time, especially on a large data set.

3) Experimental result

Because the data in the artificial data set is two-dimensional, we can display them on a two-dimensional plane. The clustering results corresponding to different k are selected as follows.

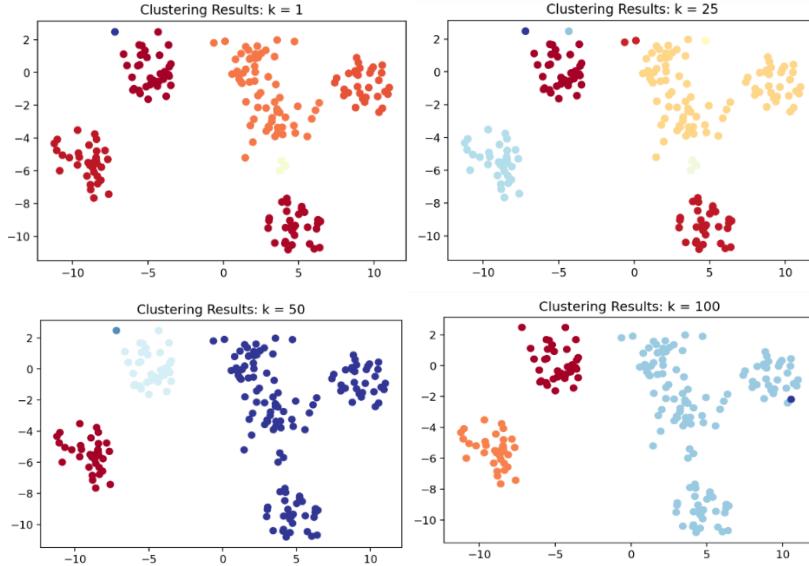


Figure 17: Clustering results of different k

It can be seen from the figure that under different batch sizes, the original clustering problem can achieve a reasonable clustering result, which can indicate that the difference in batch size does not affect the convergence of the function.

- Types of norm regularizations

- Introduction

Regularization is a very important technique in machine learning to prevent overfitting. The difference between L1 and L2 is that L2 is the sum of the square of the weights, while L1 is the sum of the weights. And the maximum norm is the maximum weight in the vector. These three types of norm regularizations in Huber function are as follows:

$$\|x\|_1 = \sum_{i=1}^N |x_i|, \|x\|_2 = (\sum_{i=1}^N |x_i|^2)^{\frac{1}{2}}, \|x\|_\infty = \max_i |x_i|, i = \operatorname{argmax}_i |x_i|$$

The differences between their properties can be promptly summarized as follows:

L2 regularization	L1 regularization
Computational efficiency due to having analytical solutions	Computational inefficient on non-sparse cases
Non-sparse outputs	Sparse outputs
No feature selection	Built-in feature selection

Table 16: the difference between L1 regularization and L2 regularization

From the computational efficiency perspective, L1-norm does not have an analytical solution, but L2-norm does. This allows the L2-norm solutions to be calculated computationally efficiently. However, L1-norm solutions have sparsity properties that allow them to be used along with sparse algorithms,

making the calculation more computationally efficient.

Sparsity refers to that only very few entries in a matrix (or vector) are non-zero. L1-norm has the property of producing many coefficients with zero values or very small values with few large coefficients.

Built-in feature selection is frequently mentioned as an important property of the L1-norm, which the L2-norm does not. This is a result of the L1-norm, which tends to produce sparse coefficients (explained below). L2-norm produces non-sparse coefficients, so it does not have this property.

Moreover, L2-norm has unique solutions while L1-norm does not. Similar to L1 norm regularization, the maximum norm may have several solutions. It treats a matrix as if it were a vector, with the entries written in a rectangular array instead of a line. In a convex problem, the corners maybe not be on the axis. So, the solutions will give lots of coefficients shrunk to have absolute value, but not sparse solutions.

Even though we put lots of effort into how to get a specific correct hessian of the Huber function for different norm regularization, we failed. Since the L1 norm is not differentiable, it is hard to get the Hessian matrix, and we can just mainly focus on AGM instead of Newton-CG. When calculating the gradient, a subgradient is utilized. The subgradient of the L1 norm is as follows:

$$\frac{\partial}{\partial \mathbf{x}} \|\mathbf{x}\|_1 = \text{sign}(\mathbf{x}) = \begin{cases} +1 & x_i > 0 \\ -1 & x_i < 0 \\ [-1,1] & x_i = 0 \end{cases}$$

Similarly, the subgradient of the maximum norm is as follows:

$$\frac{\partial \|\mathbf{x}\|_\infty}{\partial x} = \text{sign}(x_i) \cdot e_i, \quad \frac{\partial \|\mathbf{x}\|_\infty^2}{\partial x} = 2 \cdot x_i \cdot e_i, \quad i = \text{argmax } |x_i|$$

Then, the gradient of the Huber function in L1 norm and L infinite norm can be denoted respectively as follows:

$$\begin{aligned} \varphi'_{\text{hub}} &= \begin{cases} \frac{1}{\delta} \cdot \|y\| \cdot \text{sign}(y), & \|y\| \leq \delta \\ \text{sign}(y), & \|y\| > \delta \end{cases} \\ \varphi'_{\text{hub}} &= \begin{cases} \frac{1}{\delta} y_i \cdot e_i, & \|y\| \leq \delta, i = \text{argmax } |x_i| \\ \text{sign}(y_i) \cdot e_i, & \|y\| > \delta \end{cases} \end{aligned}$$

■ Experiment

We test AGM on data_200_dot_6 with 3 different norm types. Set lamda=0.05, tol=1e-1, max iteration = 30000. The results are shown in the following table. When we use L1 norm and L-inf norm, it cannot converge even reaching the max iteration. The run time is also longer when comparing the performance of other norms with L2 norm.

Norm type	Iteration Num	CPU Time	Norm Gradient
L1 norm	30000	754.01s	0.96
L2 norm	11530	239.04s	< 0.1
L-inf norm	30000	1249.03s	0.51

Table 17: performance of AGM on different types of norm regularizations w.r.t dataset 200_dot_6

Compared among different types of norm regularization, the results show that the ℓ_2 -($\|\cdot\|_2$) norm regularization method based on the Huber function has the fastest convergence. The performance in relative error has no highlight.

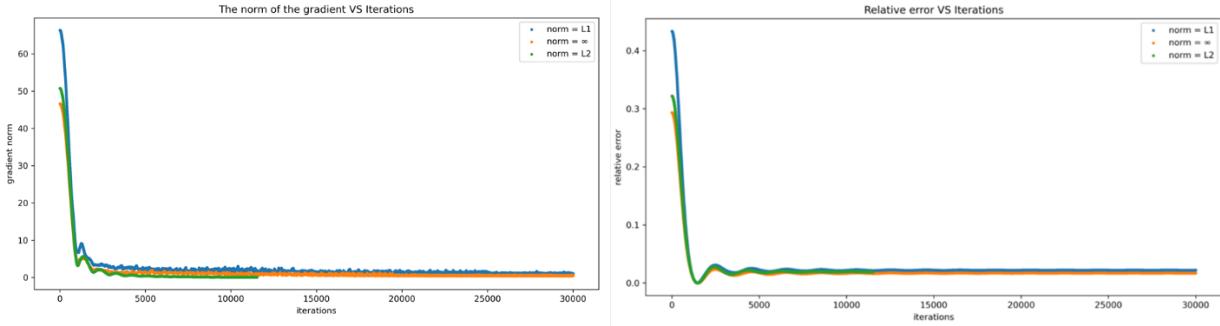


Figure 17: The Gradient Norm & relative error vs Iterations

5. Conclusion

From our experiments, we learned a lot about AGM, Newton CG, different optimization methods (BFGS/L-BFGS.....), matrix norms, mini-batch SGD, etc. AGM algorithm cannot perform well on large-scale or special shapes data sets. It is also really time-consuming with loose tolerance, which is only based on first-order information. Apply weighted methods to improve the AGM, it can achieve better results and better runtimes.

Regarding Newton-CG, we find that the number of iterations of the Newton method is significantly smaller than that of AGM, which of course is due to the use of the second-order derivative information. But in the experiment, how to speed up and store the Hessian matrix has become a very important issue. In the future, we can continue to try to store the hessian matrix in a dictionary and optimize the codes to achieve speed improvements.

In the BFGS method, instead of computing the actual Hessian, we just approximate it with a positive-definite matrix B. Therefore, it can achieve faster convergence than gradient descent, but the lower operational cost per iteration than Newton's method. The L-BFGS method aims to build a much cheaper BFGS-type approximation of the Hessian by storing only the last m pairs of $\{y^{k-m}, y^{k-m+1}, \dots, y^k\}$ and $\{s^{k-m}, s^{k-m+1}, \dots, s^k\}$. Although the L-BFGS is much cheaper, it takes more iterations and a longer time to converge.

Mini-batch SGD can be regarded as a stochastic approximation of gradient descent optimization since it replaces the actual gradient (calculated from the entire data set) with an estimate thereof (calculated from a randomly selected subset of the data). Generally, it can increase the calculation speed without greatly losing accuracy. This method is especially useful on large data sets. However, choosing a good batch size is also a question worth discussing.

6. References

References

- [1] Bottou, Léon. "Stochastic gradient descent tricks." Neural networks: Tricks of the trade. Springer, Berlin, Heidelberg, 2012. 421-436.
- [2] Li, Mu, et al. "Efficient mini-batch training for stochastic optimization." Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining. 2014.

7. Appendix

```
In [ ]:
#Load libraries
import pandas
import numpy as np
from sklearn.datasets import load_wine
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
import matplotlib.pyplot as plt
import random
import math
```

```
In [ ]:
import time
from sklearn.preprocessing import normalize
from scipy.io import loadmat
from scipy import sparse
from scipy.sparse import csc_matrix
import scipy
```

1.1 Input Data

```
In [ ]:
#generate random x0
def generate_random_mat(X):
    ??????????
    shape = X.shape
    csc_matrix((data, indices, indptr),shape)
    return

def generate_random_array(X):
    A=np.random.rand(X.shape[0],X.shape[1])
    return A
```

```
In [ ]:
#input dataset wine(mat)
from google.colab import drive
drive.mount('/content/drive')
```

```
In [ ]:
wine_data_mat = loadmat("/content/drive/MyDrive/wine_data.mat",mat_dtype=True)
wine_label_mat = loadmat("/content/drive/MyDrive/vowel_label.mat",mat_dtype=True)
wine_data=wine_data_mat['A']
wine_label=wine_label_mat['b']
wine_data
```

```
In [ ]:
type(wine_data)
```

```
In [ ]:
#input 200_dot_3
data_200_dot_3 = np.loadtxt('/content/drive/MyDrive/op_data/200_dot_3.txt')
x_200_dot_3 = data_200_dot_3[0]
y_200_dot_3 = data_200_dot_3[1]
c_200_dot_3 = data_200_dot_3[2]
data_200_dot_3 = np.stack((x_200_dot_3,y_200_dot_3), axis=1)
#normalization
#data_200_dot_3 = normalize(data_200_dot_3, axis=0, norm='max')
```

```
In [ ]:
```

```
A=csc_matrix((np.array([1,1]),(np.array([0,1]),np.array([0,1])),shape=(2,2))
print(A)
x0=np.concatenate(list(map(lambda i:A.getcol(i).toarray(),range(A.shape[1]))))
print(x0)
```

1.2 Data Generation

```
In [ ]: #Data Generation1: 多组随机分布, 可重叠
np.random.seed(2)
def data_generation(n,d,range,theta):
    arr_1 = np.random.uniform(range[0][0],range[0][1],(d,1))
    arr_2 = np.random.uniform(range[1][0],range[1][1],(d,1))
    data_reference = np.hstack((arr_1,arr_2))
    data_group = np.random.choice(d,n)
    data_generation = list()
    for i in data_group:
        a = [data_reference[i][0]+random.normalvariate(0,theta[i]),data_reference[i][1]]
        data_generation.append(a)
    return data_group, np.array(data_generation)

label_generation1, data_generation1 = data_generation(20000,4,[-10,35],[-20,20])
```

```
In [ ]: label_generation1[:10]
```

```
In [ ]: data_generation1[:10]
```

```
In [ ]: x=data_generation1[:,0]
y=data_generation1[:,1]
c=label_generation1

fig, ax = plt.subplots()
scatter = ax.scatter(x, y, c=c, cmap=plt.cm.RdYlBu)
plt.show()
```

```
In [ ]: from google.colab import drive
drive.mount('/content/drive')
```

```
In [ ]: np.savetxt('/content/drive/MyDrive/op_data/20000_random_4.txt',(x,y,c),fmt="%f %f %d")
```

```
In [ ]: np.loadtxt('/content/drive/MyDrive/op_data/20000_random_4.txt')
```

```
In [ ]: #Data Generation2: 月牙形
np.random.seed(2)
def data_generation(n,range,theta):
    random.seed(3)
    a = random.uniform(range[0][0]/3,range[0][1]/3)
    b = random.uniform(range[0][0]/2,range[0][1]/2)
    x1 = np.linspace(range[0][0], (range[0][0]+range[0][1])/2, int(n/2), dtype=np.float32)
    noise = np.random.normal(0,theta,x1.shape).astype(np.float32)
    y1 = -np.square(x1) / a + b + noise

    x2 = np.linspace((range[0][0]+range[0][1])/2, range[0][1], int(n/2), dtype=np.float32)
```

```

y2 = np.square(x2 - ((range[0][0]+range[0][1])/2)) / a - b + noise

group1 = np.hstack((x1,y1))
group2 = np.hstack((x2,y2))
data_generation = np.vstack((group1,group2))

data_group = np.hstack((np.zeros(int(n/2)),np.ones(int(n/2)))) 
return data_group, data_generation

label_generation2, data_generation2 = data_generation(600,[-40,40],20)

```

In []:

```

x=data_generation2[:,0]
y=data_generation2[:,1]
c=label_generation2

fig, ax = plt.subplots()
scatter = ax.scatter(x, y, c=c, cmap=plt.cm.RdYlBu)
plt.show()

```

In []:

```

from sklearn.datasets import make_moons
data_generation2, label_generation2 = make_moons(n_samples=200, noise=0.07, random_state=2)
x=data_generation2[:,0]
y=data_generation2[:,1]
c=label_generation2

fig, ax = plt.subplots()
scatter = ax.scatter(x, y, c=c, cmap=plt.cm.RdYlBu)
plt.show()

```

In []:

```
np.savetxt('/content/drive/MyDrive/op_data/20000_moon_2.txt',(x,y,c),fmt="%f")
```

In []:

```
np.loadtxt('/content/drive/MyDrive/op_data/20000_moon_2.txt')
```

In []:

```

#Data Generation3: 环形
np.random.seed(2)
from sklearn.datasets import make_circles
data_generation3, label_generation3 = make_circles(n_samples=200, noise=0.07, random_state=2)
x=data_generation3[:,0]
y=data_generation3[:,1]
c=label_generation3

fig, ax = plt.subplots()
scatter = ax.scatter(x, y, c=c, cmap=plt.cm.RdYlBu)
plt.show()

```

In []:

```
np.savetxt('/content/drive/MyDrive/op_data/20000_circle_2.txt',(x,y,c),fmt="%f")
```

In []:

```
data = np.loadtxt('/content/drive/MyDrive/op_data/20000_circle_2.txt')
data
```

In []:

```

x = data[0]
y = data[1]
label = data[2]

```

```
In [ ]: label
```

```
In [ ]: #Data Generation4: 点簇型
from sklearn.datasets import make_blobs
data_generation4, label_generation4 = make_blobs(n_samples=200, n_features=2,
x=data_generation4[:,0]
y=data_generation4[:,1]
c=label_generation4

fig, ax = plt.subplots()
scatter = ax.scatter(x, y, c=c, cmap=plt.cm.RdYlBu)
plt.show()
```

```
In [ ]: np.savetxt('/content/drive/MyDrive/op_data/200_dot_6.txt',(x,y,c),fmt="%f",de
```

```
In [ ]: np.loadtxt('/content/drive/MyDrive/op_data/200_dot_6.txt')
```

```
In [ ]: #数据集生成 https://blog.csdn.net/qq\_42730750/article/details/117286533
```

2.1 AGM

```
In [ ]:
from numpy.linalg import norm
def f(x,a,lambda1,delta):
    f_1 = np.sum((norm(x[i]-a[i]))**2 for i in range(len(x)))
    x_diff = list()
    f_2 = 0
    for i in range(len(a)):
        # for j in range(i+1,len(a)):
        #     x_diff.append(x[i]-x[j])
        xi_diff = x[i] - x[i+1:]
        xi_diff_norm = norm(xi_diff, ord=2, axis=1)
        mask = xi_diff_norm <= delta
        f_2 += np.sum(xi_diff_norm[mask] ** 2 / (2*delta))
        f_2 += np.sum(xi_diff_norm[~mask] - delta/2)

    #f_2 = phi(x_diff,delta)
    return 1/2*(f_1) + lambda1 * f_2

def g(x,a,lambda1,delta):
    grad_1 = x - a
    grad_2 = np.zeros(x.shape)
    for i in range(x.shape[0]):
        xi_diff = x[i] - x[i+1:]
        xi_diff_norm = norm(xi_diff, ord=2, axis=1)
        mask = xi_diff_norm <= delta

        grad_2[i] += lambda1 * np.sum(xi_diff[mask] / delta, axis=0)
        grad_2[i] += lambda1 * np.sum(xi_diff[~mask] / np.expand_dims(xi_diff_norm[~mask], axis=0), axis=0)

    return grad_1 + grad_2
```

```
In [ ]:
def function(x,A,lambda1,delta):
    ...
```

```

Parameters
-----
x : numpy.ndarray
    shape: (row*col,1)
    DESCRIPTION: The independent variable of function 'f_clust'.


Returns
-----
type: float
    DESCRIPTION: The solution of f_clust.

PS: A, row, col, lamd, delta all store as global variables.
'''

row=A.shape[0]
col=A.shape[1]
f1=[np.power(np.linalg.norm(x[i*row:(i+1)*row]-A[:, i]).reshape(-1, 1)),2)
f2=0
for i in range(col):
    for j in range(i+1,col):
        n=np.linalg.norm(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
        f2+=0.5/delta*np.power(n,2) if n<=delta else n-delta*0.5
return 0.5*sum(f1)+lamd*f2

def gradient(x,A,lambda):
'''


Parameters
-----
x : numpy.ndarray
    shape: (row*col,1)
    DESCRIPTION: The independent variable of gradient function.

Returns
-----
type: numpy.ndarray
    shape: (row*col,1)
    DESCRIPTION: The solution of gradient function.

'''

row=A.shape[0]
col=A.shape[1]
solution1=np.asarray(np.concatenate(list(map(lambda i:x[i*row:(i+1)*row]-x[(i+1)*row:(i+2)*row],range(0,col-1)))))
solution2=np.zeros([row*col,1])
for i in range(col):
    for j in range(col):
        n=np.linalg.norm(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
        solution2[i*row:(i+1)*row]+=(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
return solution1+lambda*solution2

```

In []: gradient(x,A,0.05,1e-4)

```

In [ ]:
def AGM(obj,grad,x,options):
    alpha = options['alpha']
    tol = options['tol']
    isprint = options['isprint']

    k = 0
    x_pre = x
    t = 1
    t_pre = 1

    gradient = grad(x,A,lambda1,delta)

```

```

x_lst = [x]
gradient_lst = [norm(gradient)]
while norm(gradient) >= tol and k < 30000:
    # accelerate
    beta = (t_pre - 1)/t
    y = x + beta * (x - x_pre)
    # update
    x_pre = x
    x = y - alpha * grad(y,A,lambda1,delta)
    t_pre = t
    t = 1/2 * (1 + np.sqrt(1+4*t**2))
    gradient = grad(x,A,lambda1,delta)
    x_lst.append(x)
    gradient_lst.append(norm(gradient))
    k = k + 1
    if isprint:
        print("Iteration:",k,"obj:",obj(x,A,lambda1,delta),"norm_gradient")
return x,x_lst,gradient_lst

```

In []:

```

def invoke_AGM(f,g,x_init,options):
    global lambda1
    global delta
    lambda1 = options['lambda']
    delta = options['delta']

    L = 1+x_init.shape[0]*lambda1/delta
    alpha = 1/L
    options["alpha"] = alpha
    solution, x_lst, gradient_lst = AGM(f,g,x_init,options)

    print("Stepsize:",alpha,"NumIteration:",len(x_lst)-1)
    return solution,x_lst,gradient_lst

```

In []:

```

def plot_convergence_figure(obj, x_lst, gradient_lst, options, x_init):
    lambda1 = options['lambda']
    alpha = options['alpha']

    assert len(x_lst) == len(gradient_lst)
    iteration_num = len(gradient_lst)

    #The norm of the gradient VS #iterations
    plt.figure(dpi = 300, figsize=(6, 6.5))
    ax1 = plt.subplot(2, 1, 1)
    scatter = ax1.scatter(range(iteration_num), gradient_lst, s = 5)
    plt.title("The norm of the gradient VS #Iterations")
    ax1.set_xlabel('#iterations')
    ax1.set_ylabel('gradient norm')

    #Relative error VS #Iterations
    ax2 = plt.subplot(2, 1, 2)
    best_function_value = obj(x_lst[-1],x_init,lambda1,delta)
    scatter = ax2.scatter(range(iteration_num), [abs(obj(each,x_init,lambda1,delta) - best_function_value)) for each in x_lst])
    plt.title("Relative error VS #Iterations")
    ax2.set_xlabel('#iterations')
    ax2.set_ylabel('relative error')

    plt.tight_layout()
    plt.show()

```

In []:

```
#test data_200_dot_3
```

```

T1 = time.time()
options = {
    "lambda":0.05,
    "delta":1e-4,
    "tol":1e-1,
    "isprint":True
}

A=data_200_dot_3
solution,x_lst,gradient_lst = invoke_AGM(f,g,A,options)

T2 = time.time()
run_time = (T2-T1)/60

```

In []: run_time

In []: plot_convergence_figure(f,x4_lst,gradient4_lst,options,data_200_dot_3)

In []:

```

#test data_200_dot_3 with random x0
T1 = time.time()
options = {
    "lambda":0.05,
    "delta":1e-4,
    "tol":1e-1,
    "isprint":True
}

A=generate_random_array(data_200_dot_3)
solution,x_lst,gradient_lst = invoke_AGM(f,g,A,options)

T2 = time.time()
run_time = (T2-T1)/60

```

In []:

```

#test data_200_dot_3
T1 = time.time()
options = {
    "lambda":0.05,
    "delta":1e-4,
    "tol":1e-1,
    "isprint":True
}

A=data_200_dot_3
solution,x_lst,gradient_lst = invoke_AGM(function,gradient,A,options)

T2 = time.time()
run_time = (T2-T1)/60

```

In []:

```

#test wine
T1 = time.time()
options = {
    "lambda":0.05,
    "delta":1e-4,
    "tol":1e-1,
    "isprint":True
}
A=wine_data

```

```
x0=np.concatenate(list(map(lambda i:A.getcol(i).toarray(),range(A.shape[1]))))
solution,x_lst,gradient_lst = invoke_AGM(function,gradient,x0,options)
T2 = time.time()
run_time = (T2-T1)/60
```

Weighted AGM

```
In [ ]:
def function(x,A,lambda):
    '''
    Parameters
    -----
    x : numpy.ndarray
        shape: (row*col,1)
        DESCRIPTION: The independent variable of function 'f_clust'.

    Returns
    -----
    type: float
        DESCRIPTION: The solution of f_clust.

    PS: A, row, col, lamd, delta all store as global variables.
    '''

    v=0.5
    row=A.shape[0]
    col=A.shape[1]
    f1=[np.power(np.linalg.norm(x[i*row:(i+1)*row]-A[:, i].reshape(-1, 1)),2)
    f2=0
    for i in range(col-1):
        for j in index_list[i]:
            n=np.linalg.norm(x[i*row:(i+1)*row]-x[(j+i)*row:(j+i+1)*row])
            w=np.exp((-v)*A_dist_list[i][j])
            f2+=w*0.5/delta*np.power(n,2) if n<=delta else w*(n-delta*0.5)
            #f2+=w*n
    return 0.5*sum(f1)+lambda*f2

def gradient(x,A,lambda):
    '''

    Parameters
    -----
    x : numpy.ndarray
        shape: (row*col,1)
        DESCRIPTION: The independent variable of gradient function.

    Returns
    -----
    type: numpy.ndarray
        shape: (row*col,1)
        DESCRIPTION: The solution of gradient function.

    '''

    v=0.5
    row=A.shape[0]
    col=A.shape[1]
    solution1=np.asarray(np.concatenate(list(map(lambda i:x[i*row:(i+1)*row]-
    solution2=np.zeros([row*col,1])
    for i in range(col-1):
        for j in index_list[i]:
            n=np.linalg.norm(x[i*row:(i+1)*row]-x[(j+i)*row:(j+i+1)*row])
            w=np.exp((-v)*A_dist_list[i][j])
            solution2[i*row:(i+1)*row]+=w*(x[i*row:(i+1)*row]-x[(j+i)*row:(j+i+1)*row])
```

```
#solution2[i*row:(i+1)*row] += w*(x[i*row:(i+1)*row]-x[(j+i)*row:(j+1)*row])
return solution1+lamd*solution2
```

In []:

```
def find_k_near(A,k):
    row=A.shape[0]
    col=A.shape[1]
    A_dist_list = list()
    index_list = list()
    for i in range(col-1):
        if k > col-i-1:
            k = col-i-1
        A_diff = sparse.hstack([A[:,i] for _ in range(col-i-1)]) - A[:, i+1:]
        A_dist = np.array(A_diff.power(2).sum(axis=0)).squeeze(axis=0)
        A_dist_list.append(A_dist)
        index = np.argpartition(A_dist,k)[k:]
        index_list.append(index)
    return A_dist_list, index_list
```

In []:

```
from numpy.linalg import norm
def W_AGM(obj,grad,x,options):
    global A_dist_list
    global index_list

    alpha = options['alpha']
    tol = options['tol']
    isprint = options['isprint']
    #####
    A_dist_list, index_list = find_k_near(A,7)
    #####
    k = 0
    x_pre = x
    t = 1
    t_pre = 1

    gradient = grad(x,A,lambda1,delta)
    x_lst = [x]
    gradient_lst = [norm(gradient)]
    while norm(gradient) >= tol:
        # accelerate
        beta = (t_pre -1)/t
        y = x + beta * (x - x_pre)
        # update
        x_pre = x
        x = y - alpha * grad(y,A,lambda1,delta)
        t_pre = t
        t = 1/2 * (1 + np.sqrt(1+4*t**2))
        gradient = grad(x,A,lambda1,delta)
        x_lst.append(x)
        gradient_lst.append(norm(gradient))
        k = k + 1
        if isprint:
            print("Iteration:",k,"obj:",obj(x,A,lambda1,delta),"norm_gradient",norm(gradient))
    return x,x_lst,gradient_lst
```

In []:

```
def invoke_W_AGM(f,g,x_init,options):
    global lambda1
    global delta
    lambda1 = options['lambda']
    delta = options['delta']
```

```

L = 1+x_init.shape[0]*lambda1/delta
alpha = 1/L
options["alpha"] = alpha
solution, x_lst, gradient_lst = W_AGM(f,g,x_init,options)

print("Stepsize:",alpha,"NumIteration:",len(x_lst)-1)
return solution,x_lst,gradient_lst

```

In []:

```

#test wine
T1 = time.time()
options = {
    "lambda":0.05,
    "delta":1e-4,
    "tol":1e-1,
    "isprint":True
}
A=wine_data
x0=np.concatenate(list(map(lambda i:A.getcol(i).toarray(),range(A.shape[1]))))
solution,x_lst,gradient_lst = invoke_W_AGM(function,gradient,x0,options)
T2 = time.time()
run_time = (T2-T1)/60

```

In []:

```
print(run_time)
```

test find_k_near

In []:

```

A=csc_matrix((np.array([1,1,1]),(np.array([0,1,1]),np.array([0,0,1])),shape=
x=np.concatenate(list(map(lambda i:A.getcol(i).toarray(),range(A.shape[1]))),
np.array(A.todense()))

```

In []:

```
A_dist_list, index_list = find_k_near(A,2)
```

Other

In []:

```

import numpy as np
from numpy.linalg import norm

class UnionFindSet:
    def __init__(self, n):
        self.parent = list(range(n))

    #合并index1和index2所属集合
    def union(self, index1: int, index2: int):
        self.parent[self.find(index2)] = self.find(index1)

    #查找index结点的父结点 (含路径压缩)
    def find(self, index) -> int:
        if self.parent[index] != index:
            self.parent[index] = self.find(self.parent[index])
        return self.parent[index]

X = solution4
#print(X)

```

```

uf = UnionFindSet(X.shape[0])
# clustering
tolerance = 3 # TBD

for i in range(X.shape[0]):
    for j in range(i+1, X.shape[0]):
        # 两两比较距离
        if norm(X[i] - X[j]) < tolerance:
            uf.union(i, j)
# 聚类结果
label_result4 = list()
for i in range(X.shape[0]):
    label_result4.append(uf.find(i)) # 第i条数据的聚类类别

# Visualization
print(label_result4)

```

In []: print(label_generation4)

In []:

```

#original data
x=data_generation4[:,0]
y=data_generation4[:,1]
c=label_generation4

fig, ax = plt.subplots()
scatter = ax.scatter(x, y, c=c, cmap=plt.cm.RdYlBu)
plt.show()

```

In []:

```

#predict data
x=data_generation4[:,0]
y=data_generation4[:,1]
c=label_result4

fig, ax = plt.subplots()
scatter = ax.scatter(x, y, c=c, cmap=plt.cm.RdYlBu)
plt.show()

```

In []:

```

#test data_generation1
options = {
    "lambda":0.05,
    "delta":1e-4,
    "tol":1e-2,
    "isprint":True
}
solution_1,x1_lst,gradient1_lst = invoke_AGM(data_generation1,options)

```

In []:

```

X = solution_1
#print(X)
uf = UnionFindSet(X.shape[0])
# clustering
tolerance = 4 # TBD

for i in range(X.shape[0]):
    for j in range(i+1, X.shape[0]):
        # 两两比较距离
        if norm(X[i] - X[j]) < tolerance:
            uf.union(i, j)

```

```
# 聚类结果
label_result1 = list()
for i in range(X.shape[0]):
    label_result1.append(uf.find(i)) # 第i条数据的聚类类别

# Visualization
print(label_result1)
print(set(label_result1))
```

In []:

```
#orginal data
x=data_generation1[:,0]
y=data_generation1[:,1]
c=label_generation1

fig, ax = plt.subplots()
scatter = ax.scatter(x, y, c=c, cmap=plt.cm.RdYlBu)
plt.show()
```

In []:

```
#predict data
x=data_generation1[:,0]
y=data_generation1[:,1]
c=label_result1

fig, ax = plt.subplots()
scatter = ax.scatter(x, y, c=c, cmap=plt.cm.RdYlBu)
plt.show()
```

```

1 # -*- coding: utf-8 -*-
2 """
3 Created on Sun Dec 12 14:29:25 2021
4
5 @author: Evelyn
6 """
7 #如果要运行程序, 请把本文件放到解压的数据文件下
8
9 from scipy.io import loadmat
10 import numpy as np
11 from scipy.sparse import csc_matrix
12 import time
13
14 # -----my version(slow version)-----
15 def function(x):
16     """
17     Parameters
18     -----
19     x : numpy.ndarray
20         shape: (row*col,1)
21         DESCRIPTION: The independent variable of function 'f_clust'.
22
23     Returns
24     -----
25     type: float
26         DESCRIPTION: The solution of f_clust.
27
28     PS: A, row, col, lamd, delta all store as global variables.
29     """
30     f1=[np.power(np.linalg.norm(x[i*row:(i+1)*row]-A[:, i].reshape(-1, 1)),2) for i
31     in range(col)]
32     f2=0
33     for i in range(col):
34         for j in range(i+1,col):
35             n=np.linalg.norm(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
36             f2+=0.5/delta*np.power(n,2) if n<=delta else n-delta*0.5
37     return 0.5*sum(f1)+lamd*f2
38
39 def gradient(x):
40     """
41     Parameters
42     -----
43     x : numpy.ndarray
44         shape: (row*col,1)
45         DESCRIPTION: The independent variable of gradient function.
46     Returns
47     -----
48     type: numpy.ndarray
49         shape: (row*col,1)
50         DESCRIPTION: The solution of gradient function.
51
52     """
53
54     solution1=np.asarray(np.concatenate(list(map(lambda i:x[i*row:(i+1)*row]-A[:, i].reshape(-1, 1),range(col)))),axis=0)
55     solution2=np.zeros([row*col,1])
56     for i in range(col):
57         for j in range(col):
58             n=np.linalg.norm(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
59             solution2[i*row:(i+1)*row]+=(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])/delta

```

```

 60     if n<=delta else (x[i*row:(i+1)*row]-x[j*row:(j+1)*row])/n
 61     return solution1+lamd*solution2
 62
 63 def hessian(x):
 64     """
 65     Parameters
 66     -----
 67     x : numpy.ndarray
 68         shape: (row*col,1)
 69         DESCRIPTION: The independent variable of hessian funciton
 70     Returns
 71     -----
 72     type: numpy.ndarray
 73         shape: (row*col,row*col)
 74         DESCRIPTION: The solution of hessian function.
 75
 76     """
 77     mat=[ ]
 78     hub={}
 79     for i in range(col):
 80         hubij=np.zeros([row,row])
 81         for j in range(col):
 82             n=np.linalg.norm(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
 83             hubij+=np.identity(row)/delta if n<=delta else np.identity(row)/n-
 84 (x[i*row:(i+1)*row]-x[j*row:(j+1)*row])@(x[i*row:(i+1)*row]-x[j*row:
 85 (j+1)*row]).T/np.power(n,3)
 86         hub[i]=hubij
 87
 88     for i in range(col):
 89         r=[]
 90         for j in range(col):
 91             if i==j:
 92                 r.append(np.identity(row)+lamd*hub[i])
 93             else:
 94                 n=np.linalg.norm(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
 95                 hubij=np.identity(row)/delta if n<=delta else np.identity(row)/n-
 96 (x[i*row:(i+1)*row]-x[j*row:(j+1)*row])@(x[i*row:(i+1)*row]-x[j*row:
 97 (j+1)*row]).T/np.power(n,3)
 98                 r.append(-lamd*hubij)
 99     mat.append(np.concatenate(r,axis=1))
100
101     return np.concatenate(mat,axis=0)

# -----new version(quick)-----
102
103 def function1(x):
104     """
105     Parameters
106     -----
107     x : numpy.ndarray
108         shape: (row*col,1)
109         DESCRIPTION: The independent variable of function 'f_clust'.
110
111     Returns
112     -----
113     type: float
114         DESCRIPTION: The solution of f_clust.
115
116     PS: A, row, col, lamd, delta all store as global variables.
117     """
118
119     f1=[np.power(np.linalg.norm(x[i*row:(i+1)*row]-A[:, i].reshape(-1, 1)),2) for i
120 in range(col)]
121     f2=0

```

```

116     for i in range(col-1):
117         xi_xj_norm=np.asarray(np.linalg.norm([x[i*row:(i+1)*row]-x[j*row:(j+1)*row]
118 for j in range(i+1,col)],axis=1,ord=2))
119         mask=xi_xj_norm<=delta
120         f2+=sum(xi_xj_norm[mask]**2*0.5/delta)
121         f2+=sum(xi_xj_norm[~mask]-delta*0.5)
122     return 0.5*sum(f1)+lamd*f2
123
124
125
126 def gradient1(x):
127     """
128     Parameters
129     -----
130     x : numpy.ndarray
131         shape: (row*col,1)
132         DESCRIPTION: The independent variable of gradient function.
133     Returns
134     -----
135     type: numpy.ndarray
136         shape: (row*col,1)
137         DESCRIPTION: The solution of gradient function.
138
139     solution1=np.asarray(np.concatenate(list(map(lambda i:x[i*row:(i+1)*row]-A[:,i].reshape(-1, 1),range(col))),axis=0))
140     solution2=np.zeros([row*col,1])
141     for i in range(col):
142         xi_xj=np.asarray([x[i*row:(i+1)*row]-x[j*row:(j+1)*row] for j in
143 range(col)]).reshape((col,row),order='C')
144         m=np.all(np.equal(xi_xj, 0), axis=1)
145         xi_xj=xi_xj[~m] #去掉0项
146         xi_xj_norm=np.linalg.norm(xi_xj,axis=1,ord=2)
147         mask=xi_xj_norm<=delta
148         solution2[i*row:(i+1)*row]+=(np.sum(xi_xj[mask]/delta, axis=0)).reshape(row,1)
149     solution2[i*row:(i+1)*row]+=
150     (np.sum(xi_xj[~mask]/np.expand_dims(xi_xj_norm[~mask],axis=1),axis=0)).reshape(row,1)
151     return solution1+lamd*solution2
152
153
154
155 def hessian1(x):
156     """
157     Parameters
158     -----
159     x : numpy.ndarray
160         shape: (row*col,1)
161         DESCRIPTION: The independent variable of hessian funciton
162     Returns
163     -----
164     type: numpy.ndarray
165         shape: (row*col, row*col)
166         DESCRIPTION: The solution of hessian function.
167
168     mat=[ ]
169     hub={}
170     for i in range(col):
171         hubij=np.zeros([row, row])
172         xi_xj=np.asarray([x[i*row:(i+1)*row]-x[j*row:(j+1)*row] for j in
173 range(col)]).reshape((col, row),order='C')
174         m=np.all(np.equal(xi_xj, 0), axis=1)

```

```

173     xi_xj=xi_xj[~m] #去掉0项
174     xi_xj_norm=np.linalg.norm(xi_xj, axis=1, ord=2)
175     mask=xi_xj_norm<=delta
176     hubij+=np.sum([(np.identity(row)/delta if judge==True else
177 np.identity(row)/xi_xj_norm[j]-xi_xj[j].reshape(row,1)@xi_xj[j].reshape(1,row)/(xi_xj_norm[j]**3)) for j,judge in
178 enumerate(mask)],axis=0)
179     hub[i]=hubij
180
181     for i in range(col):
182         r=[]
183         for j in range(col):
184             if i==j:
185                 r.append(np.identity(row)+lamd*hub[i])
186             else:
187                 n=np.linalg.norm(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
188                 hubij=np.identity(row)/delta if n<=delta else np.identity(row)/n-
189 (x[i*row:(i+1)*row]-x[j*row:(j+1)*row])@(x[i*row:(i+1)*row]-x[j*row:
190 (j+1)*row]).T/np.power(n,3)
191                 r.append(-lamd*hubij)
192             mat.append(np.concatenate(r, axis=1))
193     return np.concatenate(mat, axis=0)

194 # 测试以上函数
195 # print(gradient(x))
196 # print(hessian(x))
197 # print(function(x))
198 # xk=gradient(x)
199 # B=hessian(x)
200 # print((xk.T@B@xk)[0][0])

201 class Newton_CG:
202     def __init__(self, func, grad, hess, x0, tol, maxiter):
203         self.func=func
204         self.grad=grad
205         self.hess=hess
206         self.x0=x0      #initial point
207         self.x_sequence=[x0] #sequence of {x_k}
208         self.norm_grad_x_sequence=[np.linalg.norm(self.grad(x0))] #sequence of
209         {||grad_x_k||}
210         self.cg_max=10 #the max iterations of CG
211         self.epsilon=tol #the tolerance of newton-CG
212         self.maxiter=maxiter #the max iteration of newton-CG

213     def backtracking(self, x, d):
214         alpha=s
215         x_1=x+alpha*d
216         while self.func(x_1)>self.func(x)+gamma*alpha*np.dot(self.grad(x).T,d)[0,0]:
217             alpha=sigma*alpha
218             x_1=x+alpha*d
219         return alpha

220     def cal_dk(self, xk):
221         B=self.hess(xk)
222         vj=np.zeros([row*col,1])
223         rj=self.grad(xk)
224         pj=-rj
225         cg_tol=min(1,np.power(np.linalg.norm(rj),0.1))*np.linalg.norm(rj)
226         dk=-rj
227         for j in range(0, self.cg_max):
228             judge=(pj.T@B@pj)[0][0]
229             if judge<=0:

```

```

229         dk=vj if j>0 else dk
230         break
231     norm_rj=np.linalg.norm(rj) #the norm of old rj
232     sigmaj=np.power(norm_rj,2)/judge
233     vj+=sigmaj*pj
234     rj+=sigmaj*B@pj #new rj
235     norm_rj_1=np.linalg.norm(rj) #the norm of new rj
236     if norm_rj_1<cg_tol:
237         dk=vj
238         break
239     betaj_1=np.power(norm_rj_1,2)/np.power(norm_rj,2)
240     pj=-rj+betaj_1*pj
241     return dk
242
243
244 def Newton_CG(self):
245     xk=self.x0
246     print("INITIAL OBJ:",self.func(x0))
247     for k in range(self.maxiter):
248         dk=self.cal_dk(xk)
249         alphak=self.backtracking(xk,dk)
250         xk=xk+alphak*dk
251         self.x_sequence.append(xk)
252         norm_grad_xk=np.linalg.norm(self.grad(xk))
253         self.norm_grad_x_sequence.append(norm_grad_xk)
254         print("ITER:{:0>6d}".format(len(self.x_sequence)-1),end='      ')
255         # print("OBJ:",self.func(xk),end='      ')
256         print("NORM:",norm_grad_xk,end='\n')
257         if norm_grad_xk<=self.epsilon:
258             return xk
259
260 T1 = time.time() #计时
261
262 # 测试Newton_CG
263
264 #设置一些全局变量, 可自行修改
265 lamd=0.5
266 delta=0.001
267 s=1
268 sigma=0.5
269 gamma=0.1
270 maxiter=10000
271
272 #稀疏矩阵
273 file='.\wine\wine_data.mat'
274 dataA=loadmat(file,mat_dtype=True)
275 file='.\wine\wine_label.mat'
276 datab=loadmat(file,mat_dtype=True)
277 A=dataA['A']
278 b=datab['b']
279 row=A.shape[0]
280 col=A.shape[1]
281 x0=np.concatenate(list(map(lambda i:A.getcol(i).toarray(), range(col))),axis=0)
282 tol=0.001
283 maxiter=10000
284
285 # -----测试wine
286 lamda=0.05
287 delta=1e-4
288 temp=Newton_CG(function, gradient, hessian, x0, tol, maxiter)

```

```
289 temp.Newton_CG()
290 T2 = time.time()
291 print('程序运行时间:{:.2f}分' .format((T2 - T1)/60))
292 np.savez('.\\compare\\wine_no_weight',temp.x_sequence, temp.norm_grad_x_sequence)
293
294 # -----测试一个自己建的很小的数据集
295 # A=np.array([[0.,0.5,1.,0.5,1.,0.5]]*4)
296 # col=A.shape[1]
297 # row=A.shape[0]
298 # x0=A.reshape([-1,1],order='F')
299 # temp=Newton_CG(function, gradient, hessian, x0, tol, maxiter)
300 # temp.Newton_CG()
301
302 # -----测试人工生成数据集
303 # data=np.loadtxt('.\\op_data\\200_circle_2.txt')
304 # A=data[0:2]
305 # x0=A.reshape([-1,1],order='F')
306 # row=A.shape[0]
307 # col=A.shape[1]
308 # cir_200=Newton_CG(function, gradient, hessian, x0, tol, maxiter)
309 # x_star=cir_200.Newton_CG()
310 # T2 = time.time()
311 # print('程序运行时间:{:.2f}分' .format((T2 - T1)/60))
312 # np.savez('cir_200',cir_200.x_sequence,cir_200.norm_grad_x_sequence)
313
314 # data=np.loadtxt('.\\op_data\\200_dot_3.txt')
315 # A=data[0:2]
316 # x0=A.reshape([-1,1],order='F')
317 # row=A.shape[0]
318 # col=A.shape[1]
319 # dot_200=Newton_CG(function, gradient, hessian, x0, tol, maxiter)
320 # dot_200.Newton_CG()
321 # T2 = time.time()
322 # print('程序运行时间:{:.2f}分' .format((T2 - T1)/60))
323 # np.savez('dot_200',dot_200.x_sequence,dot_200.norm_grad_x_sequence)
324
325 # T1 = time.time()
326 # data=np.loadtxt('.\\op_data\\200_moon_2.txt')
327 # A=data[0:2]
328 # x0=A.reshape([-1,1],order='F')
329 # row=A.shape[0]
330 # col=A.shape[1]
331 # moon_200=Newton_CG(function, gradient, hessian, x0, tol, maxiter)
332 # moon_200.Newton_CG()
333 # T2 = time.time()
334 # print('程序运行时间:{:.2f}分' .format((T2 - T1)/60))
335 # np.savez('moon_200',moon_200.x_sequence,moon_200.norm_grad_x_sequence)
336
337 # T1 = time.time()
338 # data=np.loadtxt('.\\op_data\\200_random_4.txt')
339 # A=data[0:2]
340 # x0=A.reshape([-1,1],order='F')
341 # row=A.shape[0]
342 # col=A.shape[1]
343 # random_200=Newton_CG(function, gradient, hessian, x0, tol, maxiter)
344 # random_200.Newton_CG()
345 # T2 = time.time()
346 # print('程序运行时间:{:.2f}分' .format((T2 - T1)/60))
347 # np.savez('random_200',random_200.x_sequence,random_200.norm_grad_x_sequence)
348
```

```
349
350 # T1 = time.time()
351 # data=np.loadtxt('.\\op_data\\200_dot_6.txt')
352 # A=data[0:2]
353 # x0=A.reshape([-1,1],order='F')
354 # row=A.shape[0]
355 # col=A.shape[1]
356 # dot_200_6=Newton_CG(function, gradient, hessian, x0, tol, maxiter)
357 # dot_200_6.Newton_CG()
358 # T2 = time.time()
359 # print('程序运行时间:{:.2f}分' .format((T2 - T1)/60))
360 # np.savez('dot_200_6',dot_200_6.x_sequence, dot_200_6.norm_grad_x_sequence)
361
362 # T1 = time.time()
363 # lamda_list=[0.05,0.1,0.2,0.5,0.8,1,3,5,7.5,10.]
364 # data=np.loadtxt('.\\op_data\\200_dot_6.txt')
365 # A=data[0:2]
366 # x0=A.reshape([-1,1],order='F')
367 # row=A.shape[0]
368 # col=A.shape[1]
369 # tol=0.1
370 # for i in lamda_list:
371 #     lamd=i
372 #     dot_2000=Newton_CG(function, gradient, hessian, x0, tol, maxiter)
373 #     dot_2000.Newton_CG()
374 #     np.savez('dot_200_6_'+str(i),dot_2000.x_sequence,dot_2000.norm_grad_x_sequence)
375 # T2 = time.time()
376 # print('程序运行时间:{:.2f}分' .format((T2 - T1)/60))
377 # np.savez('dot_2000',dot_2000.x_sequence,dot_2000.norm_grad_x_sequence)
378
379 # -----测试真实数据
380 # file='.\vowel\vowel_data.mat'
381 # dataA=loadmat(file,mat_dtype=True)
382 # A=dataA['A']
383 # row=A.shape[0]
384 # col=A.shape[1]
385 # x0=np.concatenate(list(map(lambda i:A.getcol(i).toarray(), range(col))),axis=0)
386 # tol=0.001
387 # vowel=Newton_CG(function, gradient, hessian, x0, tol, maxiter)
388 # vowel.Newton_CG()
389 # T2 = time.time()
390 # print('程序运行时间:{:.2f}分' .format((T2 - T1)/60))
391 # np.savez('vowel',vowel.x_sequence, vowel.norm_grad_x_sequence)
```

```

1 # -*- coding: utf-8 -*-
2 """
3 Created on Sun Dec 26 22:49:35 2021
4
5 @author: Evelyn
6 """
7 from scipy.io import loadmat
8 import numpy as np
9 from scipy.sparse import csc_matrix
10 import time
11
12 def get_weighted_sparse(k):
13     for i in range(col-1):
14         w=np.array([np.exp(-v*np.power(np.linalg.norm((A[:, i].reshape(-1,1)-A[:, j]).reshape(-1,1)),ord=2),2)) for j in range(i+1,col)])
15         if k>=col-i-1:
16             k=col-i-1
17         index=np.argpartition(w,-k)[-k:]
18         for j in range(k):
19             Weight[(i,index[j]+i+1)]=w[index[j]]
20             Weight[(index[j]+i+1,i)]=w[index[j]]
21     return Weight
22
23 def get_weighted(k):
24     for i in range(col-1):
25         w=np.array([np.exp(-v*np.power(np.linalg.norm((A[:, i].reshape(-1,1)-A[:, j]).reshape(-1,1))),2)) for j in range(i+1,col)])
26         if k>=col-i-1:
27             k=col-i-1
28         index=np.argpartition(w,-k)[-k:]
29         for j in range(k):
30             Weight[(i,index[j]+i+1)]=w[index[j]]
31             Weight[(index[j]+i+1,i)]=w[index[j]]
32     return Weight
33
34 # -----my version(slow version)-----
35 def function(x):
36     """
37     Parameters
38     -----
39     x : numpy.ndarray
40         shape: (row*col,1)
41         DESCRIPTION: The independent variable of function 'f_clust'.
42
43     Returns
44     -----
45     type: float
46         DESCRIPTION: The solution of f_clust.
47
48     PS: A, row, col, lamd, delta all store as global variables.
49     """
50     f1=[np.power(np.linalg.norm(x[i*row:(i+1)*row]-A[:, i].reshape(-1, 1)),2) for i in range(col)]
51     f2=0
52     for i in range(col):
53         for j in range(i+1,col):
54             n=np.linalg.norm(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
55             f2+=Weight.get((i,j),0)*0.5/delta*np.power(n,2) if n<=delta else
56             Weight.get((i,j),0)*(n-delta*0.5)
57     return 0.5*sum(f1)+lamd*f2
58 def gradient(x):

```

```

59 """
60
61     Parameters
62     -----
63     x : numpy.ndarray
64         shape: (row*col,1)
65         DESCRIPTION: The independent variable of gradient function.
66     Returns
67     -----
68     type: numpy.ndarray
69         shape: (row*col,1)
70         DESCRIPTION: The solution of gradient function.
71
72 """
73
74     solution1=np.asarray(np.concatenate(list(map(lambda i:x[i*row:(i+1)*row]-A[:,i].reshape(-1, 1),range(col)))),axis=0))
75     solution2=np.zeros([row*col,1])
76     for i in range(col):
77         for j in range(col):
78             n=np.linalg.norm(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
79             solution2[i*row:(i+1)*row]+=Weight.get((i,j),0)*(x[i*row:(i+1)*row]-
x[j*row:(j+1)*row])/delta if n<=delta else Weight.get((i,j),0)*(x[i*row:(i+1)*row]-
x[j*row:(j+1)*row])/n
80     return solution1+lamd*solution2
81
82 def hessian(x):
83 """
84
85     Parameters
86     -----
87     x : numpy.ndarray
88         shape: (row*col,1)
89         DESCRIPTION: The independent variable of hessian funciton
90     Returns
91     -----
92     type: numpy.ndarray
93         shape: (row*col,row*col)
94         DESCRIPTION: The solution of hessian function.
95
96 """
97     mat=[]
98     hub={}
99     for i in range(col):
100         hubij=np.zeros([row,row])
101         for j in range(col):
102             n=np.linalg.norm(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
103             hubij+=Weight.get((i,j),0)*np.identity(row)/delta if n<=delta else
Weight.get((i,j),0)*np.identity(row)/n-(x[i*row:(i+1)*row]-x[j*row:
(j+1)*row])@(x[i*row:(i+1)*row]-x[j*row:(j+1)*row]).T/np.power(n,3)
104         hub[i]=hubij
105
106     for i in range(col):
107         r=[]
108         for j in range(col):
109             if i==j:
110                 r.append(np.identity(row)+lamd*hub[i])
111             else:
112                 n=np.linalg.norm(x[i*row:(i+1)*row]-x[j*row:(j+1)*row])
113                 hubij=Weight.get((i,j),0)*np.identity(row)/delta if n<=delta else
Weight.get((i,j),0)*np.identity(row)/n-(x[i*row:(i+1)*row]-x[j*row:
(j+1)*row])@(x[i*row:(i+1)*row]-x[j*row:(j+1)*row]).T/np.power(n,3)

```

```

114         r.append(-lamd*hubij)
115     mat.append(np.concatenate(r, axis=1))
116 return np.concatenate(mat, axis=0)
117
118 class Newton(CG:
119     def __init__(self, func, grad, hess, x0, tol, maxiter):
120         self.func=func
121         self.grad=grad
122         self.hess=hess
123         self.x0=x0      #initial point
124         self.x_sequence=[x0] #sequence of {x_k}
125         self.norm_grad_x_sequence=[np.linalg.norm(self.grad(x0))] #sequence of
{||grad_x_k||}
126         self.cg_max=10 #the max iterations of CG
127         self.epsilon=tol #the tolerance of newton-CG
128         self.maxiter=maxiter #the max iteration of newton-CG
129
130     def backtracking(self, x, d):
131         alpha=s
132         x_1=x+alpha*d
133         while self.func(x_1)>self.func(x)+gamma*alpha*np.dot(self.grad(x).T,d)[0,0]:
134             alpha=sigma*alpha
135             x_1=x+alpha*d
136         return alpha
137
138     def cal_dk(self, xk):
139         B=self.hess(xk)
140         vj=np.zeros([row*col,1])
141         rj=self.grad(xk)
142         pj=-rj
143         cg_tol=min(1,np.power(np.linalg.norm(rj),0.1))*np.linalg.norm(rj)
144         dk=-rj
145         for j in range(0,self.cg_max):
146             judge=(pj.T@B@pj)[0][0]
147             if judge<=0:
148                 dk=vj if j>0 else dk
149                 break
150             norm_rj=np.linalg.norm(rj) #the norm of old rj
151             sigmaj=np.power(norm_rj,2)/judge
152             vj+=sigmaj*pj
153             rj+=sigmaj*B@pj #new rj
154             norm_rj_1=np.linalg.norm(rj) #the norm of new rj
155             if norm_rj_1<cg_tol:
156                 dk=vj
157                 break
158             betaj_1=np.power(norm_rj_1,2)/np.power(norm_rj,2)
159             pj=-rj+betaj_1*pj
160         return dk
161
162
163     def Newton(CG(self):
164         xk=self.x0
165         print("INIT OBJ:",self.func(x0))
166         for k in range(self.maxiter):
167             dk=self.cal_dk(xk)
168             alphak=self.backtracking(xk,dk)
169             xk=xk+alphak*dk
170             self.x_sequence.append(xk)
171             norm_grad_xk=np.linalg.norm(self.grad(xk))
172             self.norm_grad_x_sequence.append(norm_grad_xk)
173             print("ITER:{:0>6d}{}".format(len(self.x_sequence)-1),end=' ')

```

```
174     print("OBJ:",self.func(xk),end='    ')
175     print("NORM:",norm_grad_xk,end='\n')
176     if norm_grad_xk<=self.epsilon:
177         return xk
178
179 #读取数据, 只用wine数据进行了测试, 其他可自行修改
180 file='.\wine\wine_data.mat'
181 dataA=loadmat(file,mat_dtype=True)
182
183 file='.\wine\wine_label.mat'
184 datab=loadmat(file,mat_dtype=True)
185
186 A=dataA['A']
187 b=datab['b']
188
189 #设置一些全局变量, 可自行修改
190 lamd=0.05
191 delta=0.001
192 row=A.shape[0]
193 col=A.shape[1]
194 Weight={}
195 s=1
196 sigma=0.5
197 gamma=0.1
198 v=0.5
199
200 T1 = time.time() #计时
201
202 # 测试Newton_CG(Weighted 版本)
203 # -----测试wine
204 x0=np.concatenate(list(map(lambda i:A.getcol(i).toarray(), range(col))),axis=0)
205 tol=0.001
206 delta=1e-4
207 maxiter=10000
208
209 Weight=get_weighted_sparse(5)
210 temp=Newton_CG(function, gradient, hessian, x0, tol, maxiter)
211 temp.Newton_CG()
212 T2 = time.time()
213 print('程序运行时间:{:.2f}分' .format((T2 - T1)/60))
214 np.savez('..\compare\wine_weight',temp.x_sequence, temp.norm_grad_x_sequence)
215
216 # -----测试一个自己建的很小的数据集
217 # A=np.array([[0.,0.5,1.,0.5,1.,0.5]]*4)
218 # col=A.shape[1]
219 # row=A.shape[0]
220 # Weighted=get_weighted(3)
221 # x0=A.reshape([-1,1])
222 # temp=Newton_CG(function, gradient, hessian, x0, tol, maxiter)
223 # temp.Newton_CG()
```

```
1 import numpy as np
2 from numpy.linalg import norm
3 import matplotlib.pyplot as plt
4 import os
5 from scipy.io import loadmat
6 from scipy.sparse import csc_matrix
7 import time
8
9 # load data
10 data_200_dot_3 = np.loadtxt('op_data/200_dot_3.txt')
11 data_200_dot_6 = np.loadtxt('op_data/200_dot_6.txt')
12 data_2000_dot_3 = np.loadtxt('op_data/2000_dot_6.txt')
13 wine_data_mat = loadmat("wine_data.mat",mat_dtype=True)
14 wine_label_mat = loadmat("wine_label.mat",mat_dtype=True)
15 wine_data=wine_data_mat['A']
16 wine_label=wine_label_mat['b']
17
18
19 def function(x):
20     """
21     Parameters
22     -----
23     x : numpy.ndarray
24         shape: (row*col,1)
25         DESCRIPTION: The independent variable of function 'f_clust'.
26
27     Returns
28     -----
29     type: float
30         DESCRIPTION: The solution of f_clust.
31
32     PS: A, row, col, lamd, delta all store as global variables.
33     """
34     f1=[np.power(np.linalg.norm(x[i*row:(i+1)*row]-A[:, i].reshape(-1, 1)),2) for i in range(col)]
35     f2=0
36     for i in range(col-1):
37         xi_xj_norm=np.asarray(np.linalg.norm([x[i*row:(i+1)*row]-x[j*row:(j+1)*row]
38 for j in range(i+1,col)],axis=1,ord=2))
39         mask=xi_xj_norm<=delta
40         f2+=sum(xi_xj_norm[mask]**2*0.5/delta)
41         f2+=sum(xi_xj_norm[~mask]-delta*0.5)
42     return 0.5*sum(f1)+lamd*f2
43
44 def gradient(x):
45     """
46
47     Parameters
48     -----
49     x : numpy.ndarray
50         shape: (row*col,1)
51         DESCRIPTION: The independent variable of gradient function.
52     Returns
53     -----
54     type: numpy.ndarray
55         shape: (row*col,1)
56         DESCRIPTION: The solution of gradient function.
57
58     """
59
```

```

60     solution1=np.asarray(np.concatenate(list(map(lambda i:x[i*row:(i+1)*row]-A[:, i].reshape(-1, 1),range(col))),axis=0))
61     solution2=np.zeros([row*col,1])
62     for i in range(col):
63         xi_xj=np.asarray([x[i*row:(i+1)*row]-x[j*row:(j+1)*row] for j in
range(col)]).reshape((col,row),order='C')
64         m=np.all(np.equal(xi_xj, 0), axis=1)
65         xi_xj=xi_xj[~m] #去掉0项
66         xi_xj_norm=np.linalg.norm(xi_xj, axis=1, ord=2)
67         mask=xi_xj_norm<=delta
68         solution2[i*row:(i+1)*row]+=(np.sum(xi_xj[mask]/delta, axis=0)).reshape(row,1)
69         solution2[i*row:(i+1)*row]+=
    (np.sum(xi_xj[~mask]/np.expand_dims(xi_xj_norm[~mask],axis=1),axis=0)).reshape(row,1)
70     return solution1+lamd*solution2
71
72
73 def plot_convergence_figure(obj, x_lst, gradient_lst, x_init):
74
75     assert len(x_lst) == len(gradient_lst)
76     iteration_num = len(gradient_lst)
77
78     #The norm of the gradient VS #iterations
79     plt.figure(dpi = 300, figsize=(6, 6.5))
80     ax1 = plt.subplot(2, 1, 1)
81     scatter = ax1.scatter(range(iteration_num), gradient_lst, s = 5)
82     plt.title("The norm of the gradient VS #Iterations")
83     ax1.set_xlabel('#iterations')
84     ax1.set_ylabel('gradient norm')
85
86     #Relative error VS #Iterations
87     ax2 = plt.subplot(2, 1, 2)
88     best_function_value = obj(x_lst[-1])
89     scatter = ax2.scatter(range(iteration_num), [abs(obj(each) - best_function_value)
/ (max(1, abs(best_function_value))) for each in x_lst], s = 5)
90     plt.title("Relative error VS #Iterations")
91     ax2.set_xlabel('#iterations')
92     ax2.set_ylabel('relative error')
93
94     plt.tight_layout()
95     plt.show()
96
97
98
99 class UnionFindSet:
100     def __init__(self, n):
101         self.parent = list(range(n))
102
103     #合并index1和index2所属集合
104     def union(self, index1: int, index2: int):
105         self.parent[self.find(index2)] = self.find(index1)
106
107     #查找index结点的父结点 (含路径压缩)
108     def find(self, index: int) -> int:
109         if self.parent[index] != index:
110             self.parent[index] = self.find(self.parent[index])
111         return self.parent[index]
112
113
114 def BFGS(func, grad, x0, options):
115     epsilon = options['tol']
116     maxiter = options['maxiter']

```

```

117     k = 0 #number of iterations
118     n = np.shape(x0)[0]
119     Bk = np.eye(n) #initial Hessian matrix is I
120     alpha = s
121     x_lst = [x0]
122     norm_grad_lst = [norm(grad(x0))]
123
124
125     while k < maxiter:
126         gk1 = grad(x0)
127         if norm(gk1) < epsilon:
128             break
129         dk = -1.0*np.linalg.solve(Bk,gk1)
130         while func(x0 + alpha * dk) > func(x0) + gamma * alpha * np.dot(gk1.T, dk)
131 [0,0]:
132             alpha = sigma * alpha
133
134             x1 = x0 + alpha * dk
135             x_lst.append(x1)
136             norm_gk1 = norm(grad(x1))
137             norm_grad_lst.append(norm_gk1)
138             print("ITER:{:>6d}".format(len(x_lst)-1),end='   ')
139             print("OBJ:", func(x1),end='   ')
140             print("NORM:", norm_gk1,end='\n')
141
142             #BFGS
143             sk = x1 - x0
144             yk = gk1 - grad(x0)
145
146             if np.dot(sk.T,yk) > 0:
147                 Bs = np.dot(Bk,sk)
148                 ys = np.dot(yk.T,sk)
149                 sBs = np.dot(np.dot(sk.T,Bk),sk)
150                 Bk = Bk - Bs.reshape((n,1))*Bs/sBs + yk.reshape((n,1))*yk/ys
151
152                 k += 1
153                 x0 = x1
154             x = x0
155             return x_lst ,norm_grad_lst, x, k
156
157 def LBFGS(fun, grad, x0, options):
158     epsilon = options['tol']
159     maxiter = options['maxiter']
160     rho = 0.55
161     sigma = 0.4
162
163     H0 = np.eye(np.shape(x0)[0])
164
165     x_lst = [x0]
166     norm_grad_lst = [norm(grad(x0))]
167
168     #s和y用于保存近期m个, 这里m取6
169     s = []
170     y = []
171     m = 6
172
173     k = 1
174     gk = np.mat(grad(x0))#计算梯度
175     dk = -H0 * gk
176     while (k < maxiter):

```

```

177     gk1 = grad(x0)
178     if norm(gk1) < epsilon:
179         break
180     n = 0
181     mk = 0
182     gk = np.mat(grad(x0))#计算梯度
183     while (n < 20):
184         newf = fun(x0 + rho ** n * dk)
185         oldf = fun(x0)
186         if (newf < oldf + sigma * (rho ** n) * (gk.T * dk)[0, 0]):
187             mk = n
188             break
189         n = n + 1
190
191     #LBFGS校正
192     x = x0 + rho ** mk * dk
193     #print x
194
195     #保留m个
196     if k > m:
197         s.pop(0)
198         y.pop(0)
199
200     #计算最新的
201     sk = x - x0
202     yk = grad(x) - gk
203
204     s.append(sk)
205     y.append(yk)
206
207     #two-loop的过程
208     t = len(s)
209     qk = grad(x)
210     a = []
211     for i in range(t):
212         alpha = (s[t - i - 1].T * qk) / (y[t - i - 1].T * s[t - i - 1])
213         qk = qk - alpha[0, 0] * y[t - i - 1]
214         a.append(alpha[0, 0])
215     r = H0 * qk
216
217     for i in range(t):
218         beta = (y[i].T * r) / (y[i].T * s[i])
219         r = r + s[i] * (a[t - i - 1] - beta[0, 0])
220
221
222     if (yk.T * sk > 0):
223         dk = -r
224
225     k = k + 1
226     x_lst.append(x)
227     norm_gk1 = norm(grad(x))
228     norm_grad_lst.append(norm_gk1)
229     print("ITER:{:>6d}".format(len(x_lst)-1),end='    ')
230     print("OBJ:", fun(x),end='    ')
231     print("NORM:", norm_gk1,end='\n')
232     x0 = x
233
234     x = x0
235     return x_lst, norm_grad_lst, x, k
236

```

```
237
238 #global variables
239 lamd = 0.5
240 delta = 0.001
241
242 #global variables for backtracking
243 s=1
244 sigma=0.5
245 gamma=0.1
246
247
248 options = {
249     "tol": 1e-3,
250     "maxiter": 1e5
251 }
252
253 # wine BFGS
254 A = wine_data
255 row = A.shape[0]
256 col = A.shape[1]
257 x0 = np.concatenate(list(map(lambda i:A.getcol(i).toarray(), range(col))),axis=0)
258
259
260 T1 = time.time()
261 temp_wine_1 = BFGS(function, gradient, x0, options)
262 T2 = time.time()
263 run_time_wine_1 = (T2-T1)/60
264
265 run_time_wine_1
266
267 plot_convergence_figure(function, temp_wine_1[0], temp_wine_1[1], x0)
268
269 # wine L-BFGS
270 T1 = time.time()
271 temp_wine_2 = LBFGS(function, gradient, x0, options)
272 T2 = time.time()
273 run_time_wine_2 = (T2-T1)/60
274
275 run_time_wine_2
276
277 plot_convergence_figure(function, temp_wine_2[0], temp_wine_2[1], x0)
278
279 # data_200_dot_3 BFGS
280 A = data_200_dot_3[0:2]
281 x0 = A.reshape([-1,1],order='F')
282 row = A.shape[0]
283 col = A.shape[1]
284
285 T1 = time.time()
286 x0=A.reshape([-1,1],order='F')
287 temp_200_3 = BFGS(function, gradient, x0, options)
288 T2 = time.time()
289 run_time_200_3 = (T2-T1)/60
290
291 run_time_200_3
292
293 plot_convergence_figure(function, temp_200_3[0], temp_200_3[1], x0)
294
295 X = temp_200_3[2].reshape([-1,2], order='F')
296 #print(X)
297 uf = UnionFindSet(X.shape[0])
```

```
298 # clustering
299 tolerance = 3 # TBD
300
301 for i in range(X.shape[0]):
302     for j in range(i+1, X.shape[0]):
303         # 两两比较距离
304         if norm(X[i] - X[j]) < tolerance:
305             uf.union(i, j)
306
307 # 聚类结果
308 label_result_200_3 = list()
309 for i in range(X.shape[0]):
310     label_result_200_3.append(uf.find(i)) # 第i条数据的聚类类别
311
312 #orginal data
313 x=data_200_dot_3[0, :]
314 y=data_200_dot_3[1, :]
315 c=label_result_200_3
316
317 fig, ax = plt.subplots()
318 scatter = ax.scatter(x, y, c=c, cmap=plt.cm.RdYlBu)
319 plt.show()
320
321 #data_200_dot_6 BFGS
322 A = data_200_dot_6[0:2]
323 x0 = A.reshape([-1,1],order='F')
324 x_dot_200_6 = data_200_dot_6[0]
325 y_dot_200_6 = data_200_dot_6[1]
326 row = A.shape[0]
327 col = A.shape[1]
328
329 T1 = time.time()
330 x0 = A.reshape([-1,1],order='F')
331 temp_200_6 = BFGS(function, gradient, x0, options)
332 T2 = time.time()
333 run_time_200_6 = (T2-T1)/60
334
335 run_time_200_6
336
337 plot_convergence_figure(function, temp_200_6[0], temp_200_6[1], x0)
338
339
340 X = temp_200_6[2].reshape([-1,2], order='F')
341 #print(X)
342 uf = UnionFindSet(X.shape[0])
343
344 # clustering
345 tolerance = 6 # TBD
346
347 for i in range(X.shape[0]):
348     for j in range(i+1, X.shape[0]):
349         # 两两比较距离
350         if norm(X[i] - X[j]) < tolerance:
351             uf.union(i, j)
352
353 # 聚类结果
354 label_result_200_6 = list()
355 for i in range(X.shape[0]):
356     label_result_200_6.append(uf.find(i)) # 第i条数据的聚类类别
357
358 fig, ax = plt.subplots()
359 scatter = ax.scatter(x_dot_200_6, y_dot_200_6, c=label_result_200_6,
```

2021/12/28 下午10:26

/Users/ronajiang/Desktop/Optimization/final project/final project/codes/BFGS&LBFGS.py

```
|cmap=plt.cm.RdYlBu)
358 plt.show()
```

```
1 #!/usr/bin/env python
2 # coding: utf-8
3
4 # In[11]:
5
6
7 #Load libraries
8 import pandas as pd
9 import numpy as np
10 from sklearn.datasets import load_wine
11 from sklearn.model_selection import train_test_split
12 from sklearn.neighbors import KNeighborsClassifier
13 from scipy.io import loadmat
14 from scipy.sparse import csc_matrix
15 import tensorflow as tf
16 import matplotlib.pyplot as plt
17 from numpy.linalg import norm
18 import random
19 import math
20 import time
21
22
23 # In[2]:
24
25
26 #load data 200 two-dimensional points, 6 cluster centers
27 data_dot_200_6 = np.loadtxt(r"C:\Users\HAOHAN\Desktop\优化\data\op_data\200_dot_6.txt")
28 x_dot_200_6 = data_dot_200_6[0]
29 y_dot_200_6 = data_dot_200_6[1]
30 label_dot_200_6 = data_dot_200_6[2]
31
32 fig, ax = plt.subplots(dpi=200)
33 plt.title("Artificial Data Set(200 two-dimensional points, 6 cluster centers)")
34 scatter = ax.scatter(x_dot_200_6, y_dot_200_6, c=label_dot_200_6, cmap=plt.cm.RdYlBu)
35 plt.show()
36
37
38 # In[61]:
39
40
41 #load wine data
42 file=r"C:\Users\HAOHAN\Desktop\优化\data\datasets\datasets\wine\wine_data.mat"
43 data_wine = loadmat(file,mat_dtype=True)[ "A"].todense().T
44
45
46 # In[6]:
47
48
49 class UnionFindSet:
50     def __init__(self, n):
51         self.parent = list(range(n))
52
53     #合并index1和index2所属集合
54     def union(self, index1: int, index2: int):
55         self.parent[self.find(index2)] = self.find(index1)
56
57     #查找index结点的父结点 (含路径压缩)
58     def find(self, index: int) -> int:
59         if self.parent[index] != index:
```

```

60         self.parent[index] = self.find(self.parent[index])
61     return self.parent[index]
62
63
64 # In[136]:
65
66
67 def f(x,a,lambda1,delta):
68     f_1 = np.sum((norm(x[i]-a[i]))**2 for i in range(len(x)))
69     x_diff = list()
70     f_2 = 0
71     for i in range(len(a)):
72         # for j in range(i+1,len(a)):
73         #     x_diff.append(x[i]-x[j])
74         xi_diff = x[i] - x[i+1:]
75         xi_diff_norm = norm(xi_diff, ord=2, axis=1)
76         mask = xi_diff_norm <= delta
77         f_2 += np.sum(xi_diff_norm[mask]**2 / (2*delta))
78         f_2 += np.sum(xi_diff_norm[~mask] - delta/2)
79
80     #f_2 = phi(x_diff,delta)
81     return 1/2*(f_1) + lambda1 * f_2
82
83 def g(x,a,lambda1,delta):
84     grad_1 = x - a
85     grad_2 = np.zeros(x.shape)
86     for i in range(x.shape[0]):
87         xi_diff = x[i] - x[i+1:]
88         xi_diff_norm = norm(xi_diff, ord=2, axis=1)
89         mask = xi_diff_norm <= delta
90         grad_2[i] += lambda1 * np.sum(xi_diff[mask] / delta, axis=0)
91         grad_2[i] += lambda1 * np.sum(xi_diff[~mask] /
92             np.expand_dims(xi_diff_norm[~mask], axis=1), axis=0)
93
94     return grad_1 + grad_2
95
96 def g_MiniBatch(x,a,lambda1,delta,batch_size):
97     data_index_array = np.random.randint(0, x.shape[0], batch_size)
98     grad_1 = x - a
99     grad_2 = np.zeros(x.shape)
100    for i in data_index_array:
101        xi_diff = x[i] - x[i+1:]
102        xi_diff_norm = norm(xi_diff, ord=2, axis=1)
103        mask = xi_diff_norm <= delta
104        grad_2[i] += lambda1 * np.squeeze(np.asarray(np.sum(xi_diff[mask] / delta,
axis=0)))
105        grad_2[i] += lambda1 * np.squeeze(np.asarray(np.sum(xi_diff[~mask] /
np.expand_dims(xi_diff_norm[~mask], axis=1), axis=0)))
106
107    return grad_1 + grad_2
108
109 # In[279]:
110
111
112 def GradientMethod(obj,grad,x,options,a,MiniBatch=False):
113     T1 = time.time() #计时
114     global lambda1
115     global delta
116     lambda1 = options['lambda']
117     delta = options['delta']
118     alpha = options['alpha']

```

```
119     tol = options['tol']
120     isprint = options['isprint']
121     max_iter = options["max_iter"]
122     if MiniBatch:
123         batch_size = options["batch_size"]
124         gradient = grad(x,a,lambda1,delta,batch_size)
125     else:
126         gradient = grad(x,a,lambda1,delta)
127     x_lst = [x]
128     gradient_lst = [norm(gradient)]
129     k = 0
130     while norm(gradient) >= tol and k < max_iter:
131         k += 1
132         x = x - alpha * gradient / k
133         if MiniBatch: gradient = grad(x,a,lambda1,delta,batch_size)
134         else: gradient = grad(x,a,lambda1,delta)
135         if isprint and (k % 10000 == 1):
136             print("Iteration:",k-
1,"obj:",obj(x,a,lambda1,delta),"norm_gradient:",norm(gradient), "time:
137 {::2f}s".format(time.time() - T1))
138             x_lst.append(x)
139             gradient_lst.append(norm(gradient))
140             print("Iteration:",k-
1,"obj:",obj(x,a,lambda1,delta),"norm_gradient:",norm(gradient))
141             T2 = time.time() #计时
142             print('程序运行时间:{:.4f}分'.format((T2 - T1)/60))
143             return x,x_lst,gradient_lst
144
145 # In[14]:
146
147
148 def plot_convergence_figure(obj, x_lst, gradient_lst, options, x_init):
149     lambda1 = options['lambda']
150     alpha = options['alpha']
151
152     assert len(x_lst) == len(gradient_lst)
153     iteration_num = len(gradient_lst)
154
155     #The norm of the gradient VS #iterations
156     plt.figure(dpi = 300, figsize=(6, 6.5))
157     ax1 = plt.subplot(2, 1, 1)
158     scatter = ax1.scatter(range(iteration_num), gradient_lst, s = 5)
159     plt.title("The norm of the gradient VS #Iterations")
160     ax1.set_xlabel('#iterations')
161     ax1.set_ylabel('gradient norm')
162
163     #Relative error VS #Iterations
164     ax2 = plt.subplot(2, 1, 2)
165     best_function_value = obj(x_lst[-1],x_init,lambda1,delta)
166     scatter = ax2.scatter(range(iteration_num), [abs(obj(each,x_init,lambda1,delta) -
167     best_function_value) / (max(1, abs(best_function_value))) for each in x_lst], s = 5)
168     plt.title("Relative error VS #Iterations")
169     ax2.set_xlabel('#iterations')
170     ax2.set_ylabel('relative error')
171
172     plt.tight_layout()
173     plt.show()
174
175 # In[ ]:
176
```

```

177
178 data = data_dot_200_6.T[:, :2]
179 delta_value = 1e-4
180 solution_list_all, x_list_all, gradient_list_all, function_value_list_all = [], [], []
181 for batch_size in [1, 25, 50, 100, 200]:
182     options = {
183         "lambda": lambda_val,
184         "delta": delta_value,
185         "tol": 1e-1,
186         "isprint": True,
187         "alpha": 1,
188         "max_iter": 30000,
189         "batch_size": batch_size
190     }
191     print(batch_size)
192     solution_list, x_list, gradient_list =
193     GradientMethod(f, g_MiniBatch, data, options, data, MiniBatch=True)
194     tmp = []
195     for each in x_list:
196         tmp.append(f(each, data, lambda_val, delta_value))
197     function_value_list_all.append(tmp)
198     solution_list_all.append(solution_list)
199     x_list_all.append(x_list)
200     gradient_list_all.append(gradient_list)
201
202 # In[ ]:
203
204
205 plt.figure(dpi = 300)
206 k_list = [1, 25, 50, 100, 200]
207 for i in range(len(function_value_list_all)):
208     plt.plot([each * 50 for each in range(len(function_value_list_all[i]) - [10:5000:50])], [np.sqrt(norm(each - min(function_value_list_all[i][10:5000:50])) / max(1, abs(min(function_value_list_all[i][10:5000:50])))) for each in
209     function_value_list_all[i][10:5000:50]], label="k=" + str(k_list[i]))
210 plt.legend()
211 plt.xlabel("#Iterations")
212 plt.ylabel("Relative Error")
213 plt.title("Relative Error VS #Iterations in artifitial dataset")
214 plt.show()
215
216 # In[224]:
217
218
219 X = solution_data_dot_200_6_SGD
220 #print(X)
221 uf = UnionFindSet(X.shape[0])
222 # clustering
223 tolerance = 0.7 # TBD
224
225 for i in range(X.shape[0]):
226     for j in range(i+1, X.shape[0]):
227         # 两两比较距离
228         if norm(X[i] - X[j]) < tolerance:
229             uf.union(i, j)
230 # 聚类结果
231 label_data_dot_200_6 = list()
232 for i in range(X.shape[0]):
233

```

```

    label_data_dot_200_6.append(uf.find(i)) # 第i条数据的聚类类别
234
235 # Visualization
236 print(len(set(label_data_dot_200_6)))
237
238 #predict data
239 fig, ax = plt.subplots()
240 scatter = ax.scatter(x_dot_200_6, y_dot_200_6, c=label_data_dot_200_6,
cmap=plt.cm.RdYlBu)
241 plt.show()
242
243
244 # In[280]:
245
246
247 data = data_wine
248 delta = delta_value
249 solution_list_all_wine, x_list_all_wine, gradient_list_all_wine = [], [], []
250 for batch_size in [1, 50, 100, 178]:
251     options = {
252         "lambda":lambda_val,
253         "delta":delta_value,
254         "tol":1e-1,
255         "isprint":True,
256         "alpha": 0.1,
257         "max_iter": 20000,
258         "batch_size": batch_size
259     }
260     print(batch_size)
261     solution_list, x_list, gradient_list =
GradientMethod(f,g_MiniBatch,data,options,data,MiniBatch=True)
262     solution_list_all_wine.append(solution_list)
263     x_list_all_wine.append(x_list)
264     gradient_list_all_wine.append(gradient_list)
265
266
267 # In[281]:
268
269
270 function_value_list_all_wine = []
271 for i in range(len(x_list_all_wine)):
272     tmp = []
273     for each in x_list_all_wine[i]:
274         tmp.append(f(each,data,lambda_val,delta_value))
275     function_value_list_all_wine.append(tmp)
276
277
278 # In[285]:
279
280
281 plt.figure(dpi = 300)
282 k_list = [1, 50, 100, 178]
283 for i in range(len(function_value_list_all_wine)):
284     plt.plot([each * 50 for each in range(len(function_value_list_all_wine[i])[:20000:50])], [np.sqrt(norm(each - min(function_value_list_all_wine[i][:20000:50])) / max(1, abs(min(function_value_list_all_wine[i][:20000:50])))) for each in
function_value_list_all_wine[i][:20000:50]], label="k="+str(k_list[i]))
285 plt.legend()
286 plt.xlabel("#Iterations")
287 plt.ylabel("Relative Error")
288 plt.title("Relative Error VS #Iterations in wine dataset")
289 plt.show()

```

```
290
291
292 # In[289]:
293
294
295 plt.figure(dpi = 300)
296 k_list = [1, 50, 100, 178]
297 for i in range(len(function_value_list_all_wine)):
298     plt.plot([each * 50 for each in range(len(gradient_list_all_wine[i])
299 [:20000:50])], gradient_list_all_wine[i][:20000:50], label="k="+str(k_list[i]))
300 plt.legend()
301 plt.xlabel("#Iterations")
302 plt.ylabel("Gradient Norm")
303 plt.title("Gradient Norm VS #Iterations in wine dataset")
304 plt.show()
305
306 # In[ ]:
```