

MDS5210: Machine Learning Final project—Robust Matrix Factorization

Peng Deng

School of Data Science
220041042

Shihan Wang

School of Data Science
220041083

Han Yang

School of Data Science
220041046

Jing Sun

School of Data Science
116010199

Sicheng Liu

School of Data Science
220041011

Abstract—This paper* is about the application of subgradient method and A-IRLS algorithms to solve the optimization problem of matrix factorization/PCA and apply them into different situations. We first evaluate our algorithms on artificially generated dataset to validate their effectiveness. Then we move to the real dataset and apply the model on the problem of face shadow removal and video frames static background extraction. And we found that both algorithms provide satisfactory performance while A-IRLS algorithm reaches smaller error and converges in much faster speed while subgradient method has much shorter running time.†

Index Terms—Machine Learning, PCA, Subgradient, A-IRLS

I. INTRODUCTION

PCA/matrix factorization is a very useful and popular method in the field of computer vision. In our project, we will adopt this method and build a machine learning system to solve the problem of removing the shadow of face images and extracting backgrounds from video frames.

PCA/matrix factorization aims to factorize the data matrix $\mathbf{X} \in \mathbb{R}^{d \times n}$ into \mathbf{UV}^\top . The underlying idea for applying low-rank matrix factorization is that the data matrix \mathbf{X} possesses intrinsically low-rank structures. Under the sparsity assumption of the residual error, one possible learning problem formulation is as formula (1).

$$\underset{\mathbf{U} \in \mathbb{R}^{d \times r}, \mathbf{V} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad \|\mathbf{X} - \mathbf{UV}^\top\|_0 \quad (1)$$

Assume that $\mathbf{X} = \mathbf{L}^* + \mathbf{S}^*$, \mathbf{L}^* is a low-rank matrix with rank $r \ll \min(d, n)$ and $\mathbf{S}^* \in \mathbb{R}^{d \times n}$ is a sparse matrix (the residual error). \mathbf{L}^* can be factorized as $\mathbf{L}^* = \mathbf{U}^* \mathbf{V}^{*\top}$, where $\mathbf{U}^* \in \mathbb{R}^{d \times r}$, $\mathbf{V}^* \in \mathbb{R}^{n \times r}$. Thus, our goal is equivalent to find the underlying \mathbf{U}^* and \mathbf{V}^* given \mathbf{X} . Plugging \mathbf{U}^* and \mathbf{V}^* into (1), we get $\|\mathbf{S}^*\|_0$. The fact that \mathbf{S}^* is sparse, together with the sparsity promoting nature of l_0 -quasinorm, suggests that $(\mathbf{U}^*, \mathbf{V}^*)$ is likely to be one optimal solution of (1). Thus, solving (1) to the globally optimal solution is likely to provide us the underlying $(\mathbf{U}^*, \mathbf{V}^*)$. However, the l_0 -quasinorm makes the learning problem (1) intractable, so that we use the l_1 -norm

to replace the l_0 -quasinorm and the final tractable learning problem is as formula (2).

$$\underset{\mathbf{U} \in \mathbb{R}^{d \times r}, \mathbf{V} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad \|\mathbf{X} - \mathbf{UV}^\top\|_1 \quad (2)$$

In our project, we will implement **subgradient method** and **alternating iteratively reweighted least squares (A-IRLS)** algorithms to solve the learning problem.

II. ALGORITHMS DESCRIPTION

A. Subgradient Method

The loss function $\mathcal{L} = \|\mathbf{X} - \mathbf{UV}^\top\|_1$ is not differentiable and we cannot employ gradient descent method directly. Thus, we apply subgradient method to solve the problem. We first calculate the subgradient at \mathbf{U} and \mathbf{V} :

$$\begin{aligned} \partial \mathcal{L}(\mathbf{U}) &= -\text{sign}(\mathbf{X} - \mathbf{UV}^\top) \mathbf{V} \\ \partial \mathcal{L}(\mathbf{V}) &= -\text{sign}(\mathbf{X} - \mathbf{UV}^\top)^\top \mathbf{U} \end{aligned} \quad (3)$$

Then set an initial α and use diminishing stepsize, so in each iteration, we have:

$$\begin{aligned} \mathbf{U}_{k+1} &= \mathbf{U}_k - \frac{\alpha}{\sqrt{n}} \cdot \partial \mathcal{L}(\mathbf{U}_k) \\ \mathbf{V}_{k+1} &= \mathbf{V}_k - \frac{\alpha}{\sqrt{n}} \cdot \partial \mathcal{L}(\mathbf{V}_k) \end{aligned} \quad (4)$$

After certain iteration, we can obtain the final estimation of $(\mathbf{U}^*, \mathbf{V}^*)$ as $(\hat{\mathbf{U}}, \hat{\mathbf{V}})$.

The subgradient method to solve this learning problem can be summarized as Algorithm (1).

Algorithm 1: The Subgradient Method

- 1 Initialization: Choose $\mathbf{U}_0 \in \mathbb{R}^{d \times r}$, $\mathbf{V}_0 \in \mathbb{R}^{n \times r}$ with $\mathbf{U}_0^{(i,j)} \sim N(0, 1)$, $\mathbf{V}_0^{(i,j)} \sim N(0, 1)$. Choose α .
 - for** $k = 0, 1, 2, \dots, N$ **do**
 - 2 Update \mathbf{U} and \mathbf{V} :

$$\begin{aligned} \mathbf{U}_{k+1} &= \mathbf{U}_k - \frac{\alpha}{\sqrt{n}} \cdot \partial \mathcal{L}(\mathbf{U}_k) \\ \mathbf{V}_{k+1} &= \mathbf{V}_k - \frac{\alpha}{\sqrt{n}} \cdot \partial \mathcal{L}(\mathbf{V}_k) \end{aligned}$$
 - return** $(\mathbf{U}_N, \mathbf{V}_N)$
-

*Group Member and Responsibilities:

Peng Deng: A-IRLS Algorithm, Data Processing, Result Generation

Shihan Wang: Subgradient method, Data Generation

Han Yang: A-IRLS Algorithm, Report Writing

Jing Sun: Subgradient method, Slides Making

Sicheng Liu: A-IRLS Algorithm, Report Writing

†https://github.com/Dpaul891/Machine_Learning_Project

B. A-IRLS Algorithm

The method of Alternating iteratively reweighed least squares (A-IRLS) is used to solve certain optimization problems in the form of p -norm iteratively in which every iteration involves solving a subproblem of weighted least squares [1].

Specifically, in our project $p = 1$ and in each step, we should fix U and V respectively and solve the weighted least squares problem of another variable. Thus, in each iteration:

$$\begin{aligned} V_{k+1}(i, :) &= (U_k^\top W_k^{(v,i)} U_k)^{-1} U_k^\top W_k^{(v,i)} X(:, i) \\ U_{k+1}(i, :) &= (V_{k+1}^\top W_k^{(u,i)} V_{k+1})^{-1} V_{k+1}^\top W_k^{(u,i)} X^\top(:, i) \end{aligned} \quad (5)$$

where $W_k^{(u,i)}$ and $W_k^{(v,i)}$ are the diagonal matrix of weights at the k^{th} iteration for $U(i, :)$ and $V(i, :)$ respectively. They are initialized with $W_0^{(u,i)} = I_n$; $W_0^{(v,i)} = I_d$. Let δ be some small value like 0.0001 then the update formula is:

$$\begin{aligned} W_k^{(u,i)} &= \text{diag} \left(\frac{1}{\max\{\delta, |X^\top(:, i) - V_k U_k(i, :)|\}} \right) \\ W_k^{(v,i)} &= \text{diag} \left(\frac{1}{\max\{\delta, |X(:, i) - U_k V_k(i, :)|\}} \right) \end{aligned} \quad (6)$$

The A-IRLS Algorithm to solve this learning problem can be summarized as Algorithm (2).

Algorithm 2: The A-IRLS Algorithm

- 1 Initialization: Choose $U_0 \in \mathbb{R}^{d \times r}$, $V_0 \in \mathbb{R}^{n \times r}$ with $U_0^{(i,j)} \sim N(0, 1)$, $V_0^{(i,j)} \sim N(0, 1)$. Choose $W_0^{(u,i)} = I_n$; $W_0^{(v,i)} = I_d$. Choose $\delta = 0.0001$.
 - for $k = 0, 1, 2, \dots, N$ do
 - 2 Fix U_k and update V_{k+1} :
 - for $i = 0, 1, 2, \dots, n$ do

$$\begin{aligned} V_{k+1}(i, :) &= (U_k^\top W_k^{(v,i)} U_k)^{-1} U_k^\top W_k^{(v,i)} X(:, i) \\ W_k^{(v,i)} &= \text{diag} \left(\frac{1}{\max\{\delta, |X(:, i) - U_k V_k(i, :)|\}} \right) \end{aligned}$$
 - 3 Fix V_{k+1} and update U_{k+1} :
 - for $i = 0, 1, 2, \dots, d$ do

$$\begin{aligned} U_{k+1}(i, :) &= (V_{k+1}^\top W_k^{(u,i)} V_{k+1})^{-1} V_{k+1}^\top W_k^{(u,i)} X^\top(:, i) \\ W_k^{(u,i)} &= \text{diag} \left(\frac{1}{\max\{\delta, |X^\top(:, i) - V_k U_k(i, :)|\}} \right) \end{aligned}$$
 - return (U_N, V_N)
-

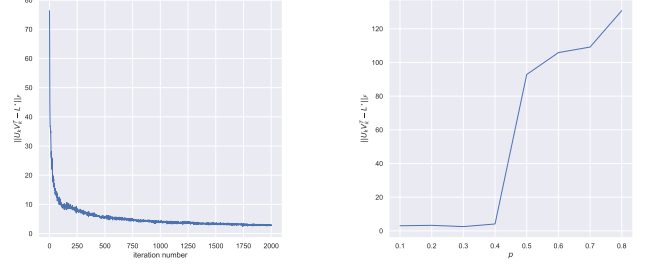
III. EVALUATION ON ARTIFICIAL DATASET

In this part, we evaluate our algorithms on artificially generated dataset. We first generate the underlying low-rank matrix $L^* = U^* V^{*\top}$ by generating $U^* \in \mathbb{R}^{d \times r}$, $V^* \in \mathbb{R}^{n \times r}$ with *i.i.d.* standard Gaussian entries. To generate $S^* \in \mathbb{R}^{d \times n}$ we first randomly select pdn locations. Then, we fill each of the selected locations with an *i.i.d.* mean 0 and variance 100 Gaussian entry, while the remaining locations are set to 0. Here

p is the ratio of the nonzero elements in S^* that represent the sparsity of the residual error.

We set $d = n = 50$, $r = 5$ and first set $p = 0.3$ to run the two algorithms on data X to plot the error $\|U_k V_k^\top - L^*\|_F$ in the process of iteration. After that, we vary $p \in \{0.1, 0.2, 0.3, \dots, 0.8\}$ and run algorithms on X with different p and plot the final error $\|\hat{U} \hat{V}^\top - L^*\|_F$ for each algorithm to compare the performance under different sparsity of residual errors S^* . The results are shown below:

A. Subgradient Method



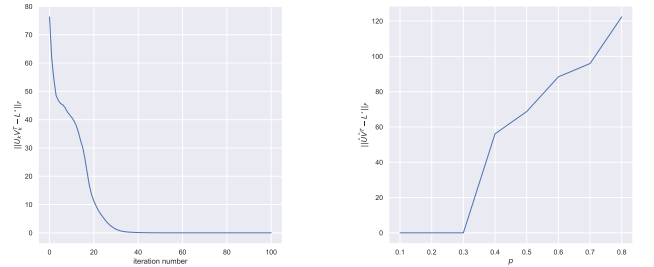
(a) Iterated errors of Subgradient with $p = 0.3$ (b) Final errors of Subgradient with different p

Fig. 1. The subgradient method

From Figure 1(a), we found that the error of sub-gradient method will decrease rapidly during first 250 iterations and the rate of descent will be slower in the following iterations. Finally, the error will fluctuate around 3.

Based on Figure 1(b), the final error will generally increase when p increases. When p ranges from 0.1 to 0.4, the achieved error is near 0 which is satisfactory. When p is larger than 0.4, the final error will increase a lot. This is because when p is larger, the matrix S^* is not sparse, so the algorithm is not suitable at this moment.

B. A-IRLS Algorithm



(a) Iterated errors of A-IRLS with $p = 0.3$ (b) Final errors of A-IRLS with different p

Fig. 2. The A-IRLS method

As shown in Figure 2(a), the error of A-IRLS algorithm reaches 0 in the 60th iteration which is much faster than sub-gradient method.

From Figure 2(b), we found that when p ranges from 0.1 to 0.3, the final error will achieve 0 in the first 100 iterations and the error will increase a lot when p is larger than 0.3, which is very similar to the sub-gradient method.

According to the iterated errors plot, we can see that both algorithms have successfully made the errors decrease while the descent speed of A-IRLS is much better. And from the final errors plot, we can conclude that both algorithms perform well when p ranges from 0.1 to 0.3.

IV. SHADOW AND ILLUMINATION REMOVING

Next, we will apply our two algorithms on the Yale-B02 face image dataset to remove the shadow and illumination. The main idea is that the clean face images of the same person are similar to each other such that they can form a low-rank matrix (regard each face image as one column), while the illumination and shadow is sparse. The results are shown below:

A. Subgradient Method

We set the iteration number equals to 2000 to get $L^* = U^*V^{*\top}$ and we can see that the loss decreases quickly in the first 250 iterations and then fluctuates around a value in later iterations as in Figure 3.

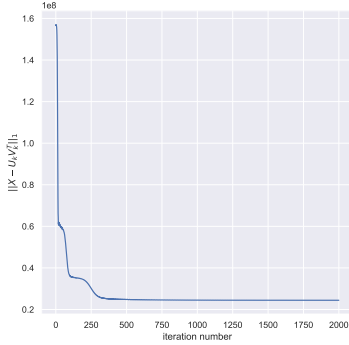


Fig. 3. Iteration Process of Subgradient in Face Image Data

The results of face shadow removal are shown as Figure 4. X is the original images, \hat{L} is the processed images and \hat{S} is the removed shadow.

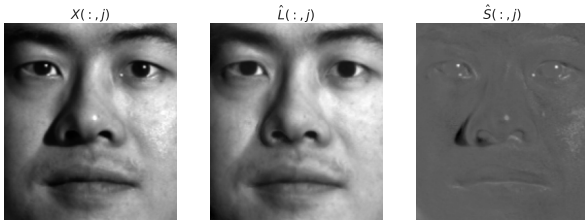


Fig. 4. Face Shadow Removal of Subgradient

B. A-IRLS Algorithm

We set the iteration number equals to 100 to get $L^* = U^*V^{*\top}$ and we can see the loss appears to be decreasing with iteration as Figure 5.

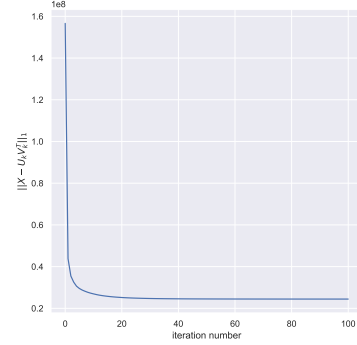


Fig. 5. Iteration Process of A-IRLS in Face Image Data

The results of face shadow removal are shown as Figure 6. X is the original images, \hat{L} is the processed images and \hat{S} is the removed shadow. We can see the shadow are removed from the original image.

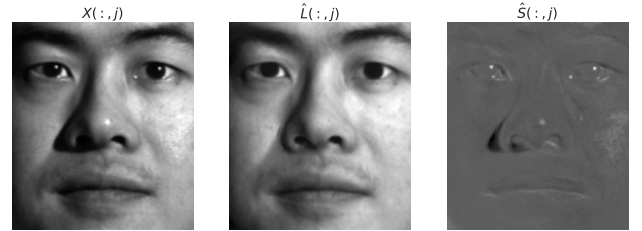


Fig. 6. Face Shadow Removal of A-IRLS

According to the figures shown above, the performance of both algorithms for face shadow removal are well while sub-gradient method needs more iterations to achieve similar result.

V. STATIC BACKGROUND EXTRACTION

In this part, we apply our algorithms to the escalator video dataset of a moving escalator and some pedestrians. The csv-format video is composed of 200 frames each with 130×160 pixels. Our goal is to split the video into non-static foreground and static background. The background has similar components as the original image, thus can be expressed by low rank matrices. Thus, the process can also be transformed into the minimization problem we have clarified above.

A. Subgradient Method

We set the iteration number equals to 2000 and the iteration process is similar as previous display of subgradient method. The loss still decreases rapidly in first 250 iterations and then fluctuates around a value in future iterations as in Figure 7.

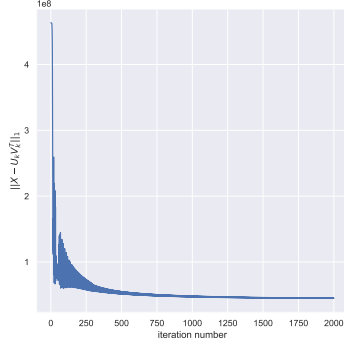


Fig. 7. Iteration Process of Subgradient in Escalator Data

After implementation of algorithm, we can obtain \hat{L} and \hat{S} . Each column of X , \hat{L} and \hat{S} represents an image of original picture, estimated background and foreground respectively. Here is an example as Figure 8.



Fig. 8. Example of Static Background Extraction by Subgradient

And by combining the 200 images together, we can obtain the whole video of background and foreground which is attached in the submitted folder.

B. A-IRLS Algorithm

We set the iteration number equals to 50 and the result comes as expected. The loss appears to be decreasing with iteration and we can see the process has the rise fluctuation in about 6th iteration. The results are shown as Figure 9.

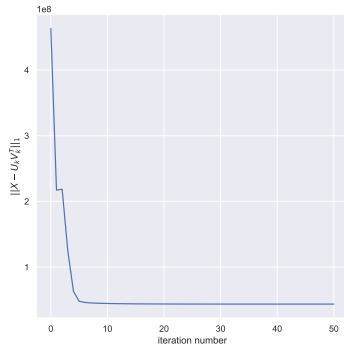


Fig. 9. Iteration Process of A-IRLS in Escalator Data

After implementation of algorithm, we can obtain \hat{L} and \hat{S} . Each column of X , \hat{L} and \hat{S} represents an image of original picture, estimated background and foreground respectively. Here is an example as Figure 10.



Fig. 10. Example of Static Background Extraction by A-IRLS

And by combining the 200 images together, we can obtain the whole video of background and foreground which is attached in the submitted folder.

VI. CONCLUSION

In this report, we apply subgradient method and A-IRLS algorithm on three matrix factorization/PCA problems: artificial datasets, face shadow removal and Static Background Extraction. In the artificial dataset evaluation, we find out good sparsity is useful for our model; In the process of face shadow removal, we have successfully removed a proportion of shadow of original image. However, due to limited iteration number the results are not so perfect; In the problem of escalator video, we have divided background and foreground of each image of the video and then combined them as a new video.

Comparatively, both algorithms provide satisfactory performance while A-IRLS algorithm reaches a litter bit smaller error and converges in a much faster speed.

In the future model tuning and development, we may try different algorithms and increase the iteration numbers in the process of algorithms to get better performance and results.

REFERENCES

- [1] Wikipedia, "Iteratively reweighted least squares," https://en.wikipedia.org/wiki/Iteratively_reweighted_least_squares.

APPENDIX

All the codes and supporting materials can be found in our github repository.