



Optimization for Data Science: final project
Frank-Wolfe methods for Recommender Systems

Kokot Maksim (ID: 2072065)

Date: 17/09/2023

1 Introduction

In this project we explore two important algorithms designed for constraint optimization tasks: Frank-Wolfe and Projected Gradient.

The project features theoretical analysis of the algorithms referencing to the existing papers dedicated to their computational cost and application. Moreover, the project is aimed to demonstrate the algorithms performance for project real world Recommender System task. In order to do that, we have selected Netflix Prize as such a task, formalized it as constrained and convex optimization problem, applied the algorithms to that problem.

2 Formalization

In this section we formulate convex and constrained optimization problem which solution is aimed to build Recommender System. Further analysis is mainly based on [1] and [2]. [1] describes algorithms' computational costs per iteration with respect to different convex sets commonly used in optimization. Meanwhile, [2] gives stronger convergence results Frank-Wolfe algorithm.

2.1 Low rank matrix completion

The majority of Recommender System problems provide user-item rating table, which contains large amount of missing values. In this case, the problem could be formalized as matrix completion task.

First of all, we will consider the constrained optimization problem

$$\min_{x \in C} f(x), \tag{1}$$

where:

x – solution of the problem

$C \subset R^n$ – a compact convex set

n – dimension of vector x

f – function to minimize

For given sparse rating matrix M , we are required to obtain matrix X whose entries suitably match to corresponding entries of matrix M in case if those entries are not missed. Moreover, we expect the entries of matrix X , which correspond to missing values in matrix M , to be plausible predictions of unknown ratings.

For this purpose, we impose low-rank constraint to matrix X . This type of constraint is aimed to preserve interdependency between the entries.

Thus, we may transform (1) into

$$\min_{\text{rank}(x) \leq r} f(x), \quad (2)$$

where:

r – threshold for rank of x

However, the low-rank constraint is hard to handle and it makes the problem non-convex. In order to overcome this problem, we can replace it with the nuclear norm constraint as follows

$$\min_{\|x\|_{\text{nuc}} \leq \beta} f(x), \quad (3)$$

where:

β – threshold for nuclear norm of x

By imposing this constraint, we obtain a convex optimization problem. The set \mathcal{C} satisfying this constrain could be described as a convex hull of points described in expression bellow

$$\mathcal{C} = \{X \in R^{n \times m} \mid \|X\|_{\text{nuc}} \leq \beta\} = \{\beta uv^T \mid u \in R^n, v \in R^m, \|u\|_2 = \|v\|_2 = 1\} \quad (4)$$

2.2 Loss function

As a function to minimize f , we can use squared Frobenius norm of difference between matrices M and X . Note, that matrix $M - X$ is also sparse and its entry is defined if and only if corresponding entry of matrix M is defined.

$$f(X) = \|M - X\|_F^2, \quad (5)$$

Thus, expression (3) could be rewritten as

$$\min_{\|x\|_{\text{nuc}} \leq \beta} \|M - X\|_F^2, \quad (6)$$

2.3 Gradient

The gradient matrix for current solution is calculated as follows:

$$\nabla f(X) = -2(M - X) \quad (7)$$

2.4 Frank-Wolfe algorithm

In this project we apply Frank-Wolfe algorithm to the convex minimization problem described in expression (6). This is first-order algorithm which is purposed for solving constrained optimization problems. It is described in Algorithm 1.

Algorithm 1 Frank-Wolfe

Input: Start point $\mathbf{x}_0 \in \mathcal{C}$, step-size strategy $(\gamma_t)_{t \in N} \subset [0, 1]$

for $t = 0$ to $T - 1$ **do**

$$\mathbf{w}_t := \arg \min_{\mathbf{v} \in \mathcal{C}} \langle \mathbf{v}, \nabla f(\mathbf{x}_t) \rangle$$

$$\mathbf{x}_{t+1} := \mathbf{x}_t + \gamma_t(\mathbf{w}_t - \mathbf{x}_t)$$

End for

For matrix completion with the nuclear norm constraint, we can clarify the computation of \mathbf{w}_t for each iteration. First, we need to obtain \mathbf{u}_t and \mathbf{v}_t which are the top left and right singular vectors of $-\nabla f(\mathbf{x}_t)$. This could be done by using truncated singular value decomposition (truncated SVD) with ARPACK solver. Second, we need to compute \mathbf{w}_t which equals $\beta \mathbf{u}_t \mathbf{v}_t^T$. The whole algorithm is described in Algorithm 2.

Algorithm 2 Frank-Wolfe for matrix completion with the nuclear norm constraint

Input: Start point $\mathbf{x}_0 \in \mathcal{C}$, step-size strategy $(\gamma_t)_{t \in N} \subset [0, 1]$

for $t = 0$ to $T - 1$ **do**

 Compute $(\mathbf{u}_t, \mathbf{v}_t)$, the top singular vectors of $-\nabla f(\mathbf{x}_t)$

$$\mathbf{w}_t := \beta \mathbf{u}_t \mathbf{v}_t^T$$

$$\mathbf{x}_{t+1} := \mathbf{x}_t + \gamma_t(\mathbf{w}_t - \mathbf{x}_t)$$

End for

The huge advantage of Frank-Wolfe algorithm is that it uses linear minimizations over given convex set and it is projection-free algorithm, hence its usage does not suppose solving quadratic programming problem [1]. Moreover, it doesn't require to use full SVD and relies on using truncated SVD, which is suitable for sparse data.

During the classes it has been proven that if f is a convex function with Lipschitz continuous gradient having constant L and for each iteration t we use $\gamma_t = \frac{2}{t+2}$, we have the following convergence result:

$$f(x_{t+1}) - f(x^*) \leq \frac{2LD^2}{t+1}, \quad (8)$$

where:

x^* – optimal solution of the problem

L – Lipschitz constant

D – diameter of the convex set \mathcal{C}

Diameter of the convex set \mathcal{C} can be computed by following formula:

$$D = \max_{x,y \in \mathcal{C}} \|x - y\|_2 \quad (9)$$

According to paper [1], linear minimization over nuclear norm-ball performed each iteration has a following complexity:

$$O(v \ln(m+n) \sqrt{\sigma_1} / \sqrt{\varepsilon}), \quad (10)$$

where:

v – the number of nonzero entries of $-\nabla f(x_t)$

σ_1 – the top singular value of $-\nabla f(x_t)$

ε – the additive error which satisfies $\sigma_1 - u^T(-\nabla f(x_t))v \leq \varepsilon$

Note that $v \leq mn$. As the rule, in Recommender System tasks matrix $-\nabla f(x_t)$ is highly sparse, so $v \ll mn$. The expression (10) shows that sparsity of $-\nabla f(x_t)$ reduces the complexity of linear minimization. This implies effectiveness of Frank-Wolfe for Recommender System tasks.

2.5 Projected Gradient algorithm

In addition, we are going to apply Projected Gradient algorithm to the same problem (6). Unlike Frank-Wolfe algorithm, this algorithm is projection-based and it requires to compute projection onto \mathcal{C} every iteration. Projected Gradient algorithm is described in Algorithm 3.

Algorithm 3 Projected Gradient

Input: Start point $\mathbf{x}_0 \in \mathcal{C}$, step-size strategies $(s_t)_{t \in N} > \mathbf{0}$, $(\alpha_t)_{t \in N} \subset [0, 1]$
for $t = 0$ **to** $T - 1$ **do**

$$\hat{\mathbf{x}}_t := \text{proj}(\mathbf{x}_t - s_t \nabla f(\mathbf{x}_t))$$

$$\mathbf{x}_{t+1} := \mathbf{x}_t + \alpha_t (\hat{\mathbf{x}}_t - \mathbf{x}_t)$$

End for

A projection onto the nuclear norm-ball task can be transformed into projection onto the simplex task. For this purpose, we need first to obtain $(U, \sigma, V) \in R^{n \times k} \times R^k \times R^{m \times k}$, $k = \min\{n, m\}$ as the results of SVD of $\mathbf{x}_t - s_t \nabla f(\mathbf{x}_t)$. Then we need to compute $\hat{\sigma}$ as projection of σ onto the simplex $\Delta_k = \{x \in R^k \mid x_l \geq 0 \text{ for } 1 \leq l \leq k, \sum_{l=1}^k x_l = \beta\}$. Finally, a projection onto the nuclear norm-ball could be calculated as $\hat{\mathbf{x}}_t := U \text{diag}(\hat{\sigma}) V^T$. The whole algorithm is described in Algorithm 4.

Algorithm 4 Projected Gradient for matrix completion with the nuclear norm constraint

Input: Start point $\mathbf{x}_0 \in \mathcal{C}$, step-size strategies $(s_t)_{t \in N} > \mathbf{0}$, $(\alpha_t)_{t \in N} \subset [0, 1]$
for $t = 0$ **to** $T - 1$ **do**

$$\mathbf{w}_t := \mathbf{x}_t - s_t \nabla f(\mathbf{x}_t)$$

Compute (U, σ, V) , the results of SVD of \mathbf{w}_t

Compute $\hat{\sigma} = \text{proj}(\sigma, \Delta_k)$

$$\hat{\mathbf{x}}_t := U \text{diag}(\hat{\sigma}) V^T$$

$$\mathbf{x}_{t+1} := \mathbf{x}_t + \alpha_t (\hat{\mathbf{x}}_t - \mathbf{x}_t)$$

End for

The projection of σ onto the simplex Δ_k can be obtained by using an algorithm proposed here [3]. In practice, this algorithm has complexity $O(n)$.

The course materials provides us a proof for a statement that if f is a convex function with Lipschitz continuous gradient having constant L and for each iteration t we use $\alpha_t = \frac{1}{L}$, we have the following convergence result:

$$f(x_{t+1}) - f(x^*) \leq \frac{2L}{t} \|x_1 - x^*\|^2 \quad (11)$$

Finding projection on nuclear norm-ball performed each iteration has a following complexity [1]:

$$O(mn \min\{m, n\} + \min\{m^3, n^3\}) \quad (12)$$

3 Data

In order to demonstrate developed solutions, we chose MovieLens dataset for Netflix Prize competition. This dataset has several available variants. We used the smallest variant [4] due to limited memory. Thus, selected dataset contains 100 000 ratings provided by 600 users to 9000 movies.

The set of unique available ratings is $\{0.5, 1, 1.5, 2, 2.5, 3, 3.5, 4, 4.5, 5\}$.

The project is mainly aimed to demonstrate the way selected algorithms converge to optimal solution. However, we believe that showing predictive ability of obtained solution is also necessary. For this reason, we split given rating into train and test part in such a way that test size equals 25% of all dataset.

4 Conducting experiment

4.1 Experiment's configuration

Before we start conducting the experiment, we have to properly select several algorithms' input parameters such as number of iteration and step-size strategies.

For each algorithm step-size step-size strategy has been carefully chosen.

Multiple papers including [2] describes Frank-Wolfe variant with step-size strategy $\gamma_t = \frac{2}{t+2}$ where t is number of iteration. In our case, γ_t values obtained this way turned out to be large and validation curves looked noisy. In order to make validation curves look smooth and obtain better convergence results, we have selected step-size strategy $\gamma_t = \frac{0.4}{t+2}$.

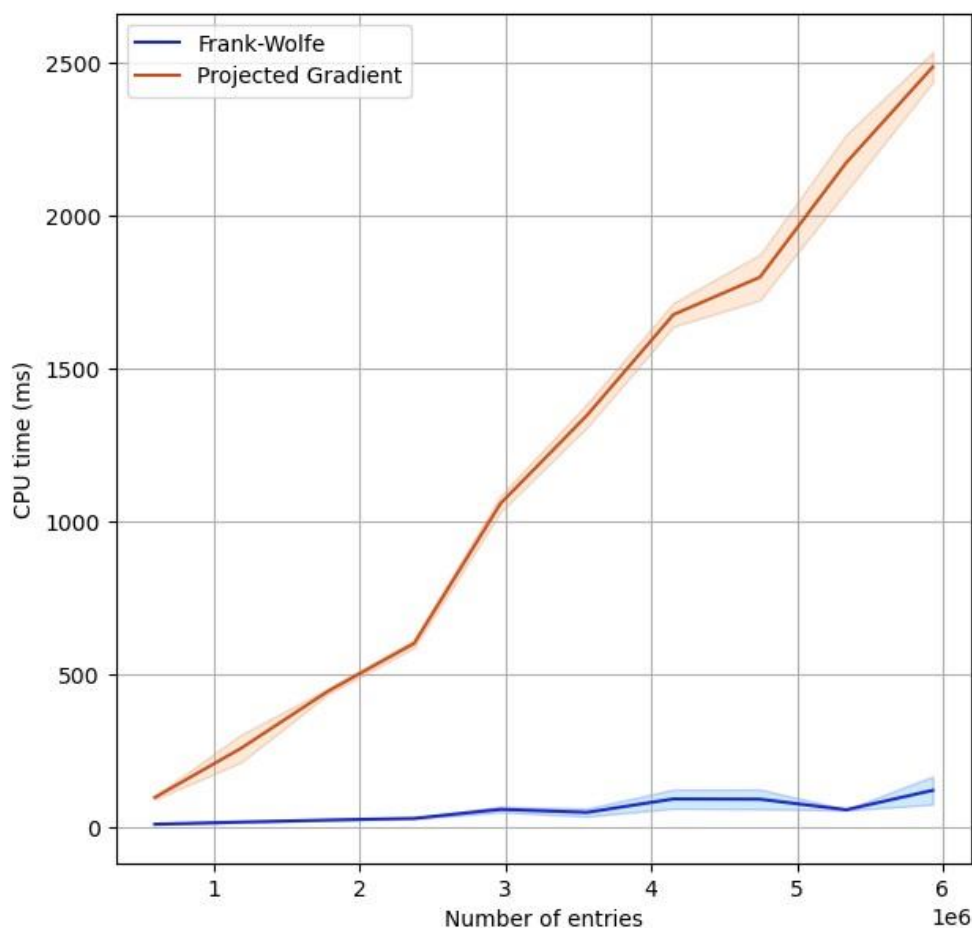
Minimized loss function f is quadratic convex function with Lipschitz constant $L = 2$. During the classes we proved faster convergence rate of Projected Gradient algorithm for the case when we minimize a convex function with Lipschitz continuous gradient and fix $s_t = \alpha_t = 1/L$. For this reason, we set $s_t = \alpha_t = 1/2$.

It's possible to show that at least 2000 iterations are sufficient to reach a solution closed to optimal one for both algorithms given introduced step-size strategies.

Finally, we have to set parameter β for with the nuclear norm constraint. Empirically we have selected $\beta = 10000$. This value is enough to make the predictions match corresponding ratings if they exist. Moreover, this value is small enough for preserving interdependency between entries, which leads to plausibility of predictions corresponding to unknown ratings.

4.2 Cost per iteration

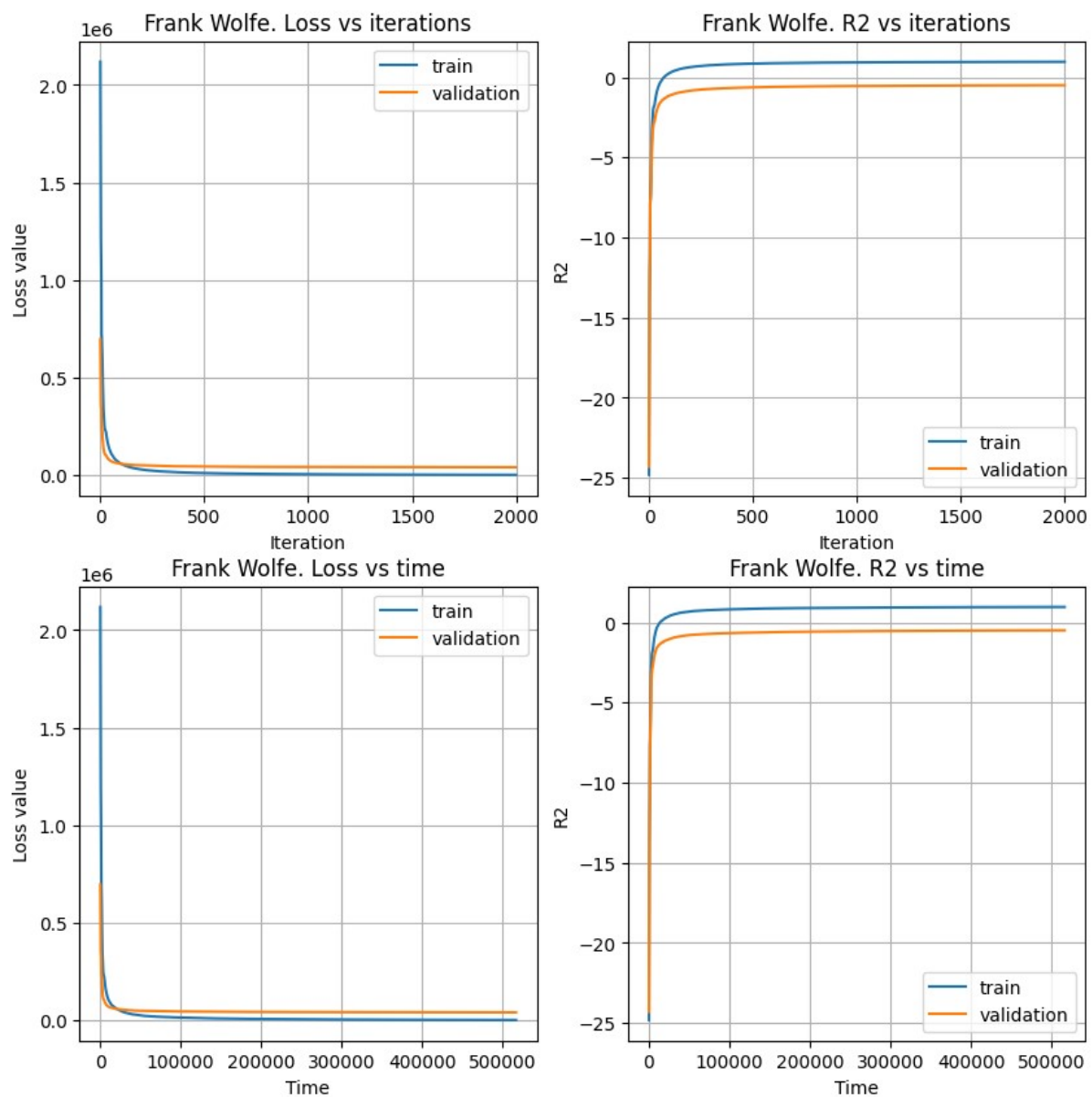
First we are going to compare computational costs per iteration for both algorithms. The comparison will be held with respect to number of entries in matrix X . After conducted analysis we get the following plot:

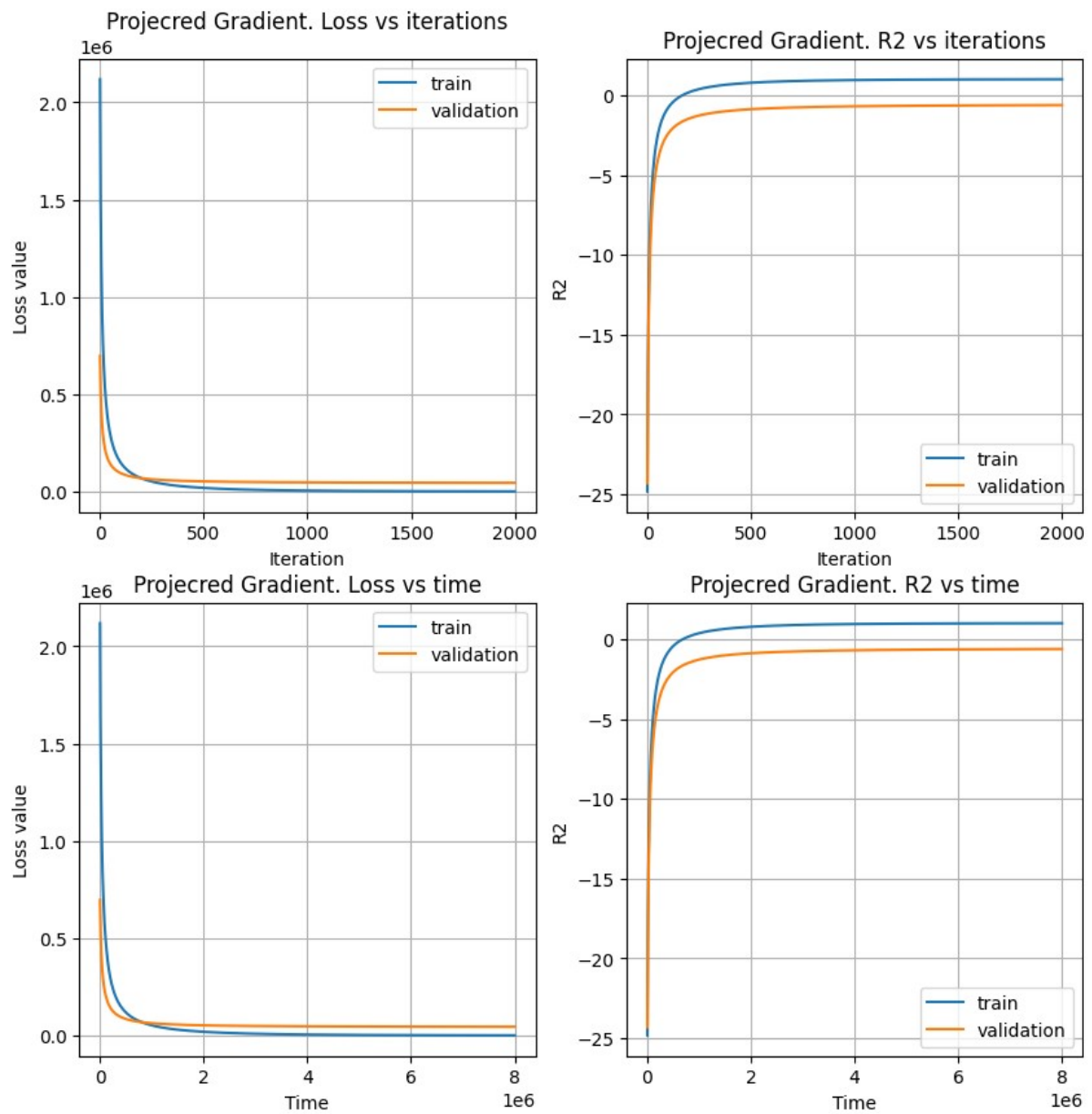


Here we can clearly observe that Projected Gradient iteration costs much more than Frank-Wolfe iteration, especially if we deal with large scale problem. Those results are consisted with theoretical analysis conducted bellow.

4.3 Loss function and R2 values

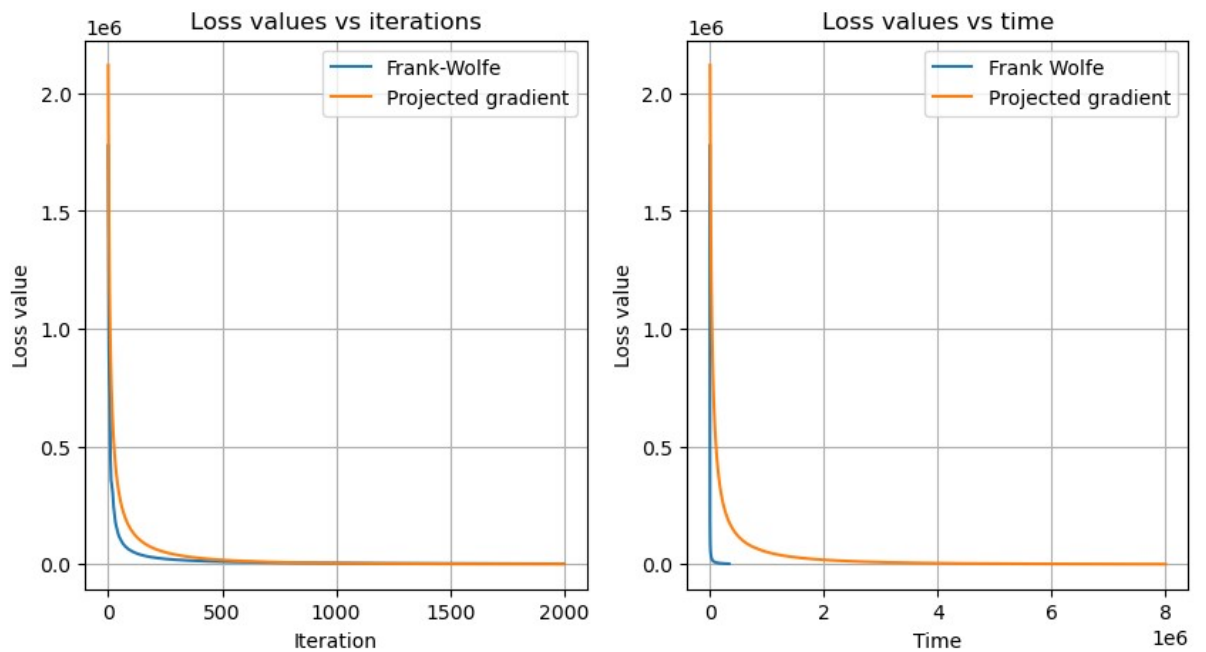
Now we are going to see the performance of both algorithms related to chosen Recommender System problem. For this purpose, we plot loss function and R2 values against number of iterations and CPU time. In addition, we make those plots for train and test datasets separately.





During 2000 iterations, both algorithms lead to sufficient decrease loss function value.

Now we compare convergence results of the algorithms.



Frank-Wolfe algorithm seems to have better convergence results in terms of iterations. However, it's impossible to notice, that final loss function value after using Projected Gradient algorithm (480) is smaller than one obtained after using Frank-Wolfe algorithm (2073).

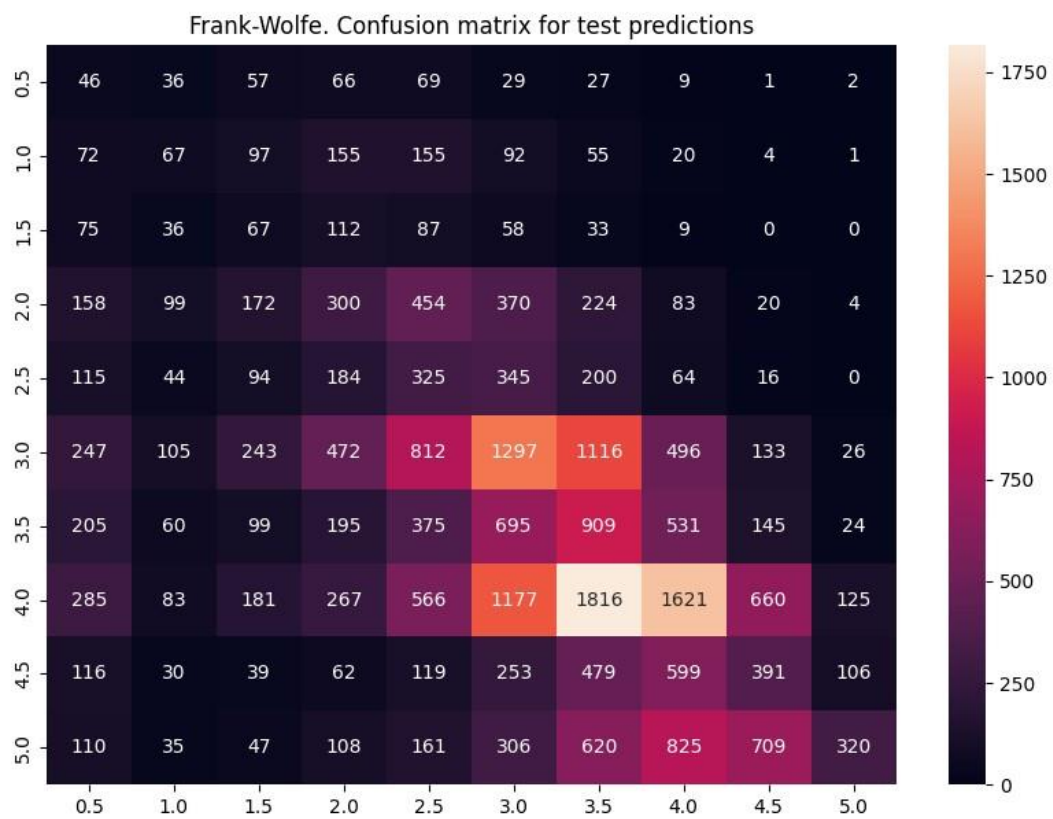
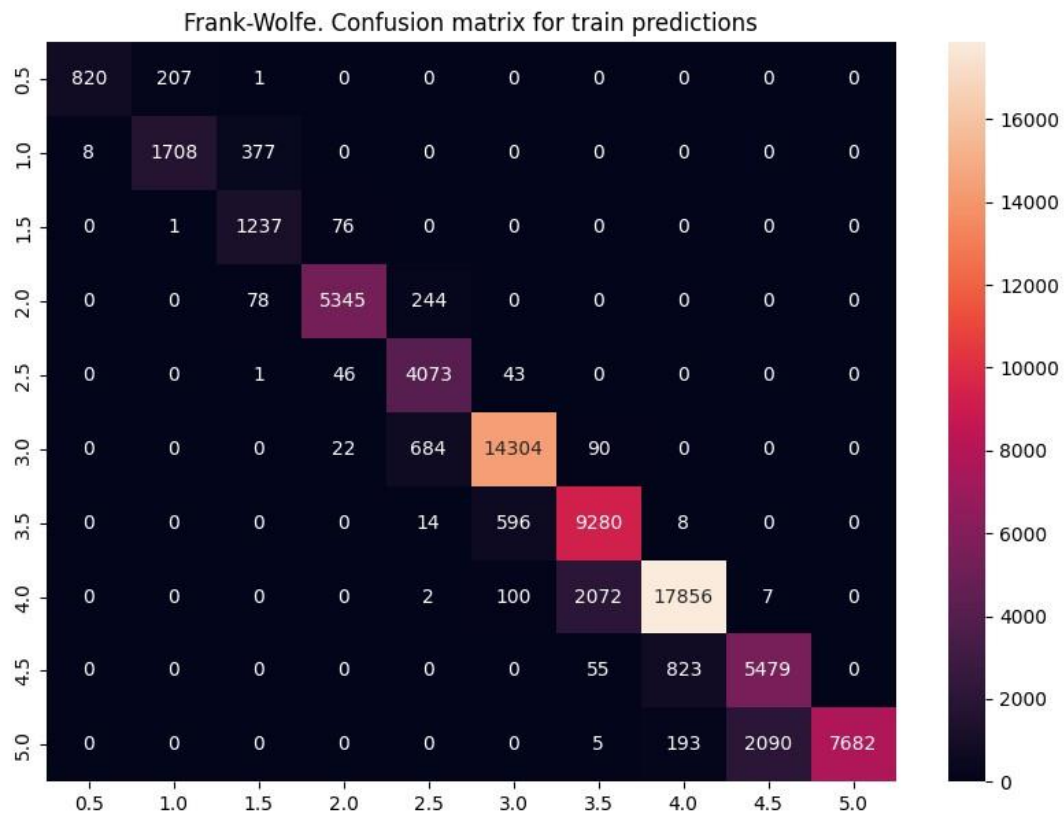
The plot above suggests that Projected Gradient requires less iterations to reach good solution for the problem. However, due to the fact that Projected Gradient iteration is costlier than Frank-Wolfe iteration, Frank-Wolfe algorithm consumes much less CPU time to find the solution.

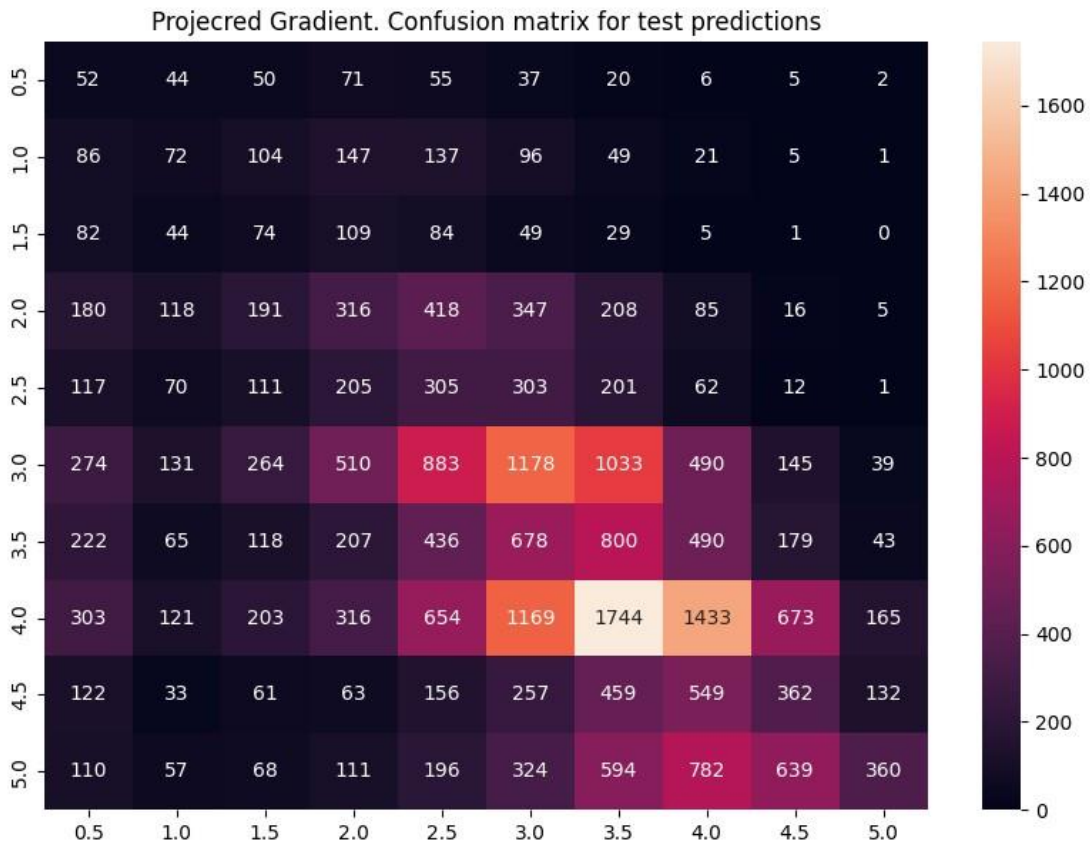
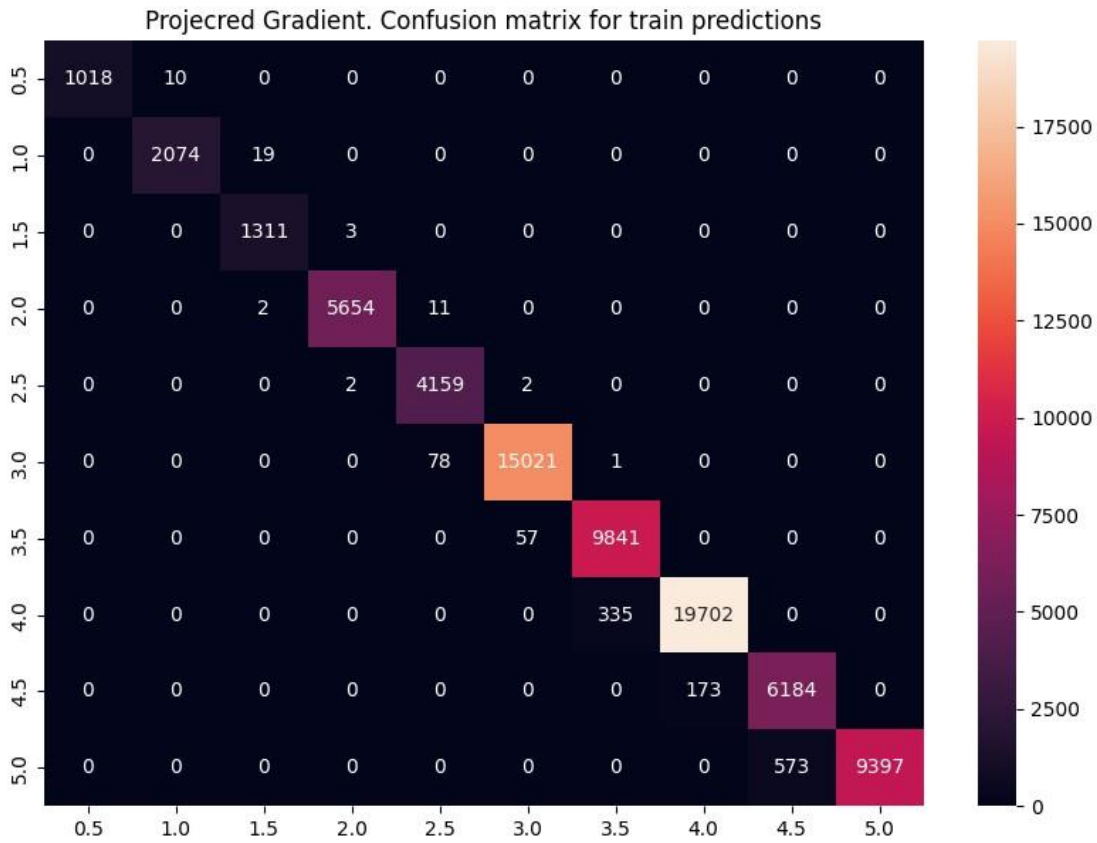
Final R2 scores are listed in the table below:

	Dataset	
Algorithm	train	test
Frank-Wolfe	0.97	-0.37
Projected Gradient	0.99	-0.49

4.4 Confusion matrices

In order to explicitly estimate the plausibility of obtained solutions in terms of providing recommendations, we plot confusion matrices for train and test predictions.





By looking at main diagonals of confusion matrices built for train prediction, we can conclude, that values of matrix X corresponding to train ratings

of matrix M have sufficiently become closer to given ratings. However, this could not be said about test predictions, even though it's possible to see moderate values close to part of main diagonals related to good and neutral ratings. Generally, such results are frequently obtained collaborative by using collaborative filtering methods for real world Recommender System problems.

5 Conclusion

In this project we explored performance of Frank-Wolfe and Projected Gradient Algorithms in constrained and convex optimization problem related to Recommender Systems. By analyzing papers describing algorithms, their complexity, convergence results and conducting computational experiments, we may conclude, that Projected Gradient algorithm is computationally expensive with respect to one single iteration compared with Frank-Wolfe algorithm due to several reasons. First, Projected Gradient algorithm requires to find projection which is quadratic programming problem, whereas Frank-Wolfe is projection-free algorithm. Second, Projected Frank-Wolfe exploits truncated SVD on sparse matrix, and this is computationally less expensive than using full SVD for Projected Gradient algorithm. Those aspects make Projected Gradient algorithm converge slower to optimal solution in term of CPU time. Meanwhile, one Projected Gradient iteration gives a new solution X which is closer to optimal solution than the one given by Frank-Wolfe iteration. However, this advantage doesn't make Projected Gradient algorithm provide the fastest CPU time solution.

References

- [1] Cyrille W. Combettes, Sebastian Pokutta. «Complexity of Linear Minimization and Projection on Some Sets». In: (2021).
- [2] Martin Jaggi. «Revisiting Frank-Wolfe: Projection-Free Sparse Convex Optimization».
- [3] L. Condat. «Fast projection onto the simplex and the l_1 ball». In: (2016).
- [4] MovieLens dataset. URL: <https://files.grouplens.org/datasets/movielens/ml-latest-small.zip>.
- [5] Lijun Ding, Madeleine Udell. «Frank-Wolfe Style Algorithms for Large Scale Optimization».