

# CONVERGENCE RATE FOR REGULARISED LINEAR LEAST SQUARES USING DISTRIBUTED FORWARD-BACKWARD METHOD

Beilun Zhang, supervised by Matthew Tam

University of Melbourne



## Introduction

Consider the problem where we want to find the most sparse vector  $x$  that satisfies  $Ax = b$  (i.e. maximise  $\|x\|_0$ , the number of zero entries in  $x$ ). This problem can be approximated as an optimisation problem as follow:  
Let  $\lambda > 0$ ,  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ , find

$$\arg \min_{x \in \mathbb{R}^n} \lambda \|x\|_1 + \frac{1}{2} \|Ax - b\|_2^2 \quad (1)$$

This problem is called the  $\ell_1$ -regularised linear least squares. There are already well-established algorithms to solve this problem (such as the forward-backward method), which works well in a centralised setting, however, here we will consider an algorithm for this problem that is optimised for a distributed setting of ring networks

## Model

In [1], Arago-Artacho, Malitsky, Tam and Torregrosa-Belen proposed an algorithms to solve monotone inclusion problems in a real Hilbert space  $\mathcal{H}$  in the form of:

$$\text{find } x \in \mathcal{H} \text{ such that } 0 \in \left( \sum_{i=1}^n A_i + \sum_{i=1}^m B_i \right)(x) \quad (2)$$

where  $A_1, \dots, A_n : \mathcal{H} \rightrightarrows \mathcal{H}$  are maximally monotone operator and  $B_1, \dots, B_m : \mathcal{H} \rightarrow \mathcal{H}$  are  $\frac{1}{L}$ -cocoercive. By using zero operators if necessary, the algorithm can always assume that  $m = n - 1$ . The algorithm can be expressed as fixed point iteration  $z^{k+1} = T(z^k)$ ,  $T : \mathcal{H}^{n-1} \rightarrow \mathcal{H}^{n-1}$  is given by

$$T(z) := z + \gamma \begin{pmatrix} x_2 - x_1 \\ x_3 - x_2 \\ \vdots \\ x_n - x_{n-1} \end{pmatrix} \quad (3)$$

where  $x = (x_1, \dots, x_n) \in \mathcal{H}^n$  depends on  $z = (z_1, \dots, z_{n-1}) \in \mathcal{H}^{n-1}$  and is given by

$$\begin{cases} x_1 = J_{\alpha A_1}(z_1) \\ x_i = J_{\alpha A_i}(z_i + x_{i-1} - z_{i-1} - \alpha B_{i-1}(x_{i-1})) \quad \forall i \in [2, n-1] \\ x_n = J_{\alpha A_n}(x_1 + x_{n-1} - z_{n-1} - \alpha B_{n-1}(x_{n-1})) \end{cases}$$

Where  $J_f$  is the resolvent for operator  $f$ . After finding  $z^* \in \text{Fix } T$  through fixed point iteration, solution to (2) is given by  $x^* = J_{\alpha A_1}(z_1^*)$ .

In a distributed system with  $n$  agents, agent  $i$  is responsible for updating  $x_i$  and sending  $x_i$  to agent  $i - 1, i + 1 \pmod n$

We can apply this to our problems of  $\ell_1$ -regularised linear least squares. The equation in (1) is equivalent to

$$\sum_{i=1}^n \lambda |x_i| + \sum_{i=1}^m \frac{1}{2} |A_i x - b_i|^2 = \sum_{i=1}^n g_i(x) + \sum_{i=1}^m f_i(x) \quad (4)$$

Since the problem is convex, finding  $\arg \min_{x \in \mathbb{R}^n}$  for (4) is the same as finding  $x \in \mathbb{R}^n$  such that  $0 \in (\sum_{i=1}^n \partial g_i + \sum_{i=1}^m \nabla f_i)(x)$ , where  $\partial g_i$  is the subdifferential of  $g_i$ , thus we have  $A_i = \partial g_i$  and  $B_i = \nabla f_i$ , with

$$J_{\alpha \partial g_i} = \begin{bmatrix} x_1 \\ \vdots \\ S(x_i; \alpha \lambda) \\ \vdots \\ x_n \end{bmatrix} \text{ and } \nabla f_i(x) = A_i^T (A_i x - b_i)$$

Where  $A_i$  is the  $i$ th row of matrix  $A$  and

$$S(y; k) = \begin{cases} y + k & y < -k \\ 0 & -k \leq y \leq k \\ y - k & y > k \end{cases}$$

## Results

We explored 3 methods/conditions that affects the convergence rate

Values for  $\alpha$  and  $\gamma$

There are 3 free parameters in the model:  $\lambda$  (regularization parameter),  $\alpha$  (step size) and  $\gamma$  (relaxation parameter). We will focus only on the parameter  $\alpha$  and  $\gamma$ , as  $\lambda$  controls how good the approximation is and has little effect on the convergence rate. From [1], the algorithm would converge when  $\alpha \in (0, \frac{2}{L})$  and  $\gamma \in (0, 1 - \frac{\alpha L}{2})$ . Fig. 1 shows the normalized number of iterations the algorithm takes when we set  $\alpha$  and  $\gamma$  a certain percentage of their upper bound. The figure shows that the number of normalized iterations is minimum when  $\alpha$  is approximately 50% of  $\frac{2}{L}$  and  $\gamma$  is close to 100% of  $1 - \frac{\alpha L}{2}$ .

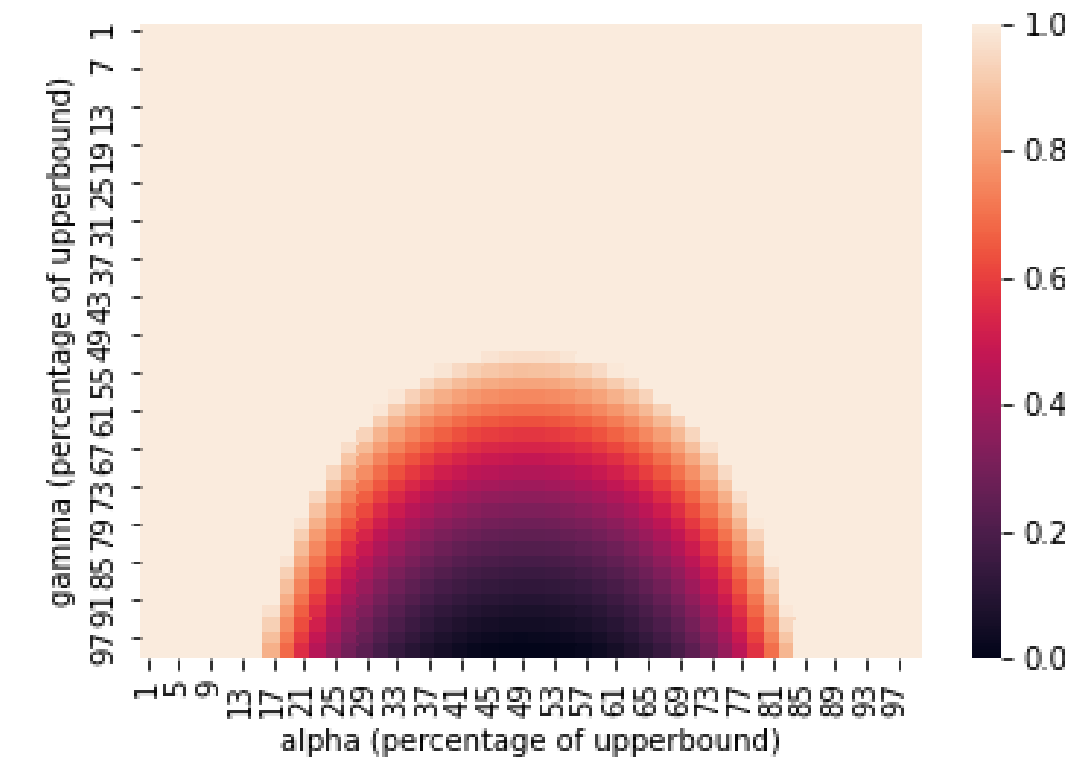


Fig. 1: Axis represent the % of theoretical upper bound  $\alpha$ ,  $\gamma$  is at

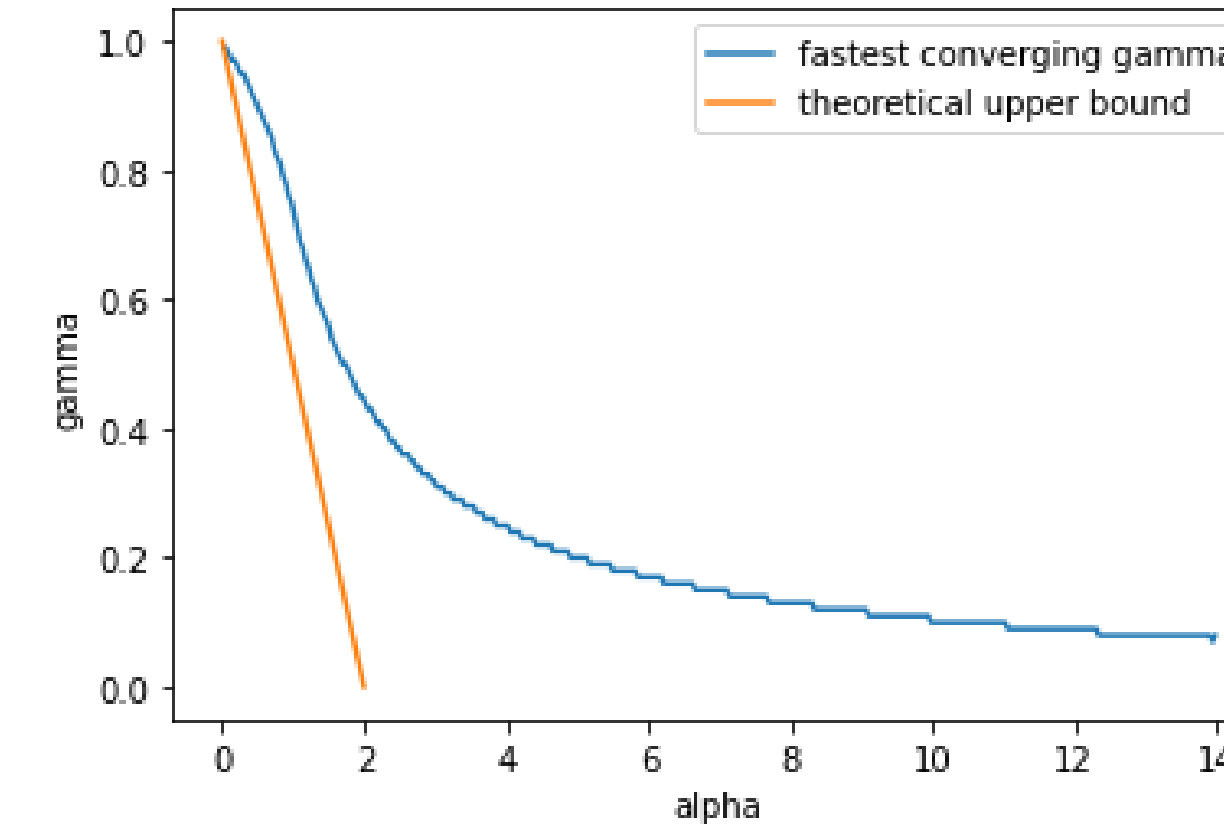


Fig. 2: Experimental upper bound for  $\gamma$  vs theoretical upper bound

However, it turns out that the upper bound for  $\gamma$  and  $\alpha$  could be improved. Fig. 2 shows the  $\gamma$  values that gives the lowest number of iterations for each  $\alpha$ , depicts that the algorithm still converges when the parameters are not within their theoretical upper bound. If we go beyond the theoretical upper bound, the optimal  $\alpha$  value is still somewhere around 50% of  $\frac{2}{L}$  ( $L = 1$  in this example), as shown in Fig. 4, with the optimal  $\gamma$  value being around 0.7

### Preconditioning

Another way to improve the convergence rate is the scale the Lipschitz constant  $L$  so that each  $\nabla f_i$  is at most  $L$ -Lipschitz (i.e. they are not  $(L - \epsilon)$ -Lipschitz). Note due to Baillon-Haddad theorem,  $\nabla f_i$  is  $\frac{1}{L}$ -cocoercive iff it is  $L$ -Lipschitz.  $\nabla f_i$  is  $L$ -Lipschitz if:

$$\|\nabla f_i(x) - \nabla f_i(y)\| = \|A_i^T A_i(x - y)\| \leq L \|x - y\| \quad \forall x, y \in \mathbb{R}^n \quad (5)$$

As  $\|A_i^T A_i(x - y)\| \leq \sigma_{\max}(A_i^T A_i) \|x - y\|$ , where  $\sigma_{\max}(A_i^T A_i)$  is the maximum singular value of  $A_i^T A_i$ ,  $\nabla f_i$  is thus at most  $\sigma_{\max}(A_i^T A_i)$ -Lipschitz. We can multiply each row  $A_i$  with  $\sqrt{\frac{L}{\sigma_{\max}(A_i^T A_i)}}$  to have each  $\nabla f_i$  at most  $L$ -Lipschitz. This is equivalent to multiplying  $A$  with a diagonal matrix  $D$ , and our new problem has  $A_{\text{new}} = DA$ , and  $b_{\text{new}} = Db$ , note that this won't change the solution of our algorithm.

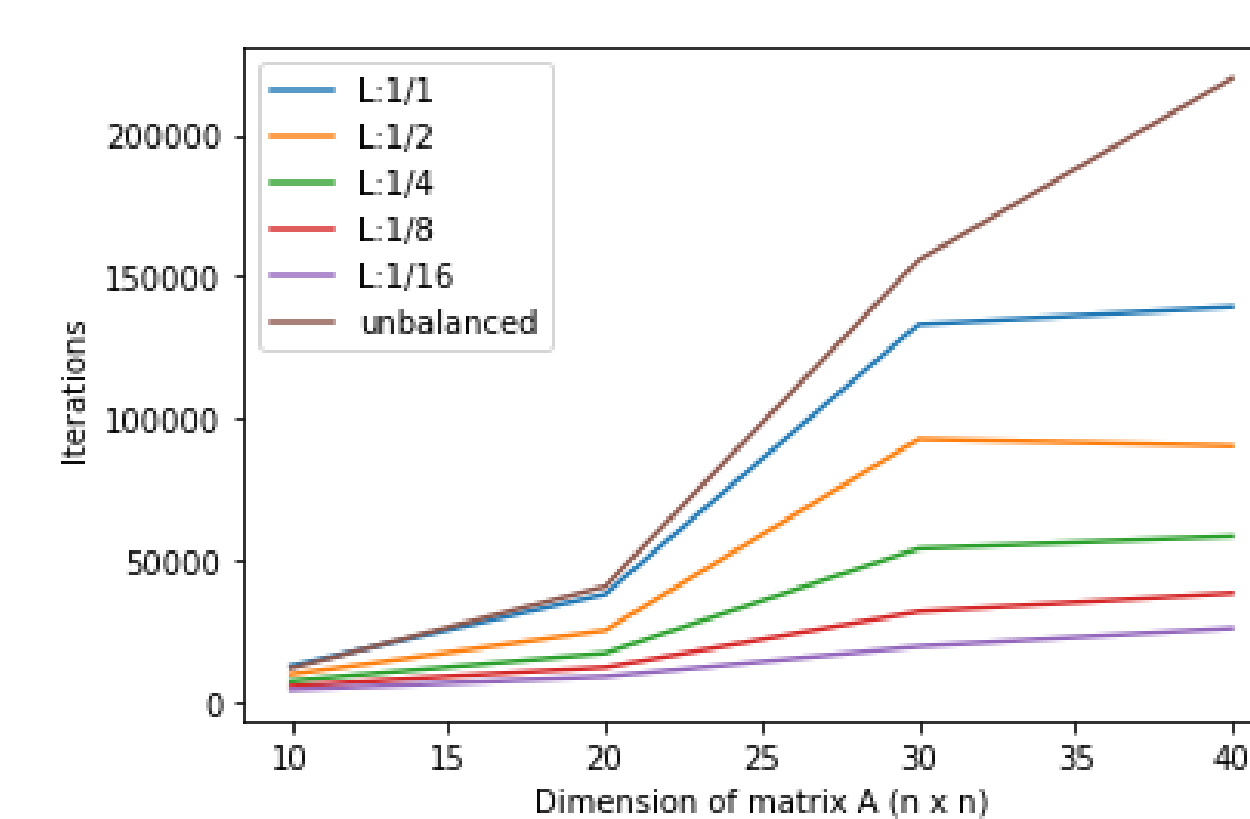


Fig. 3: Convergence rate for preconditioning with different  $L$  "unbalanced" is when we have no preconditioning

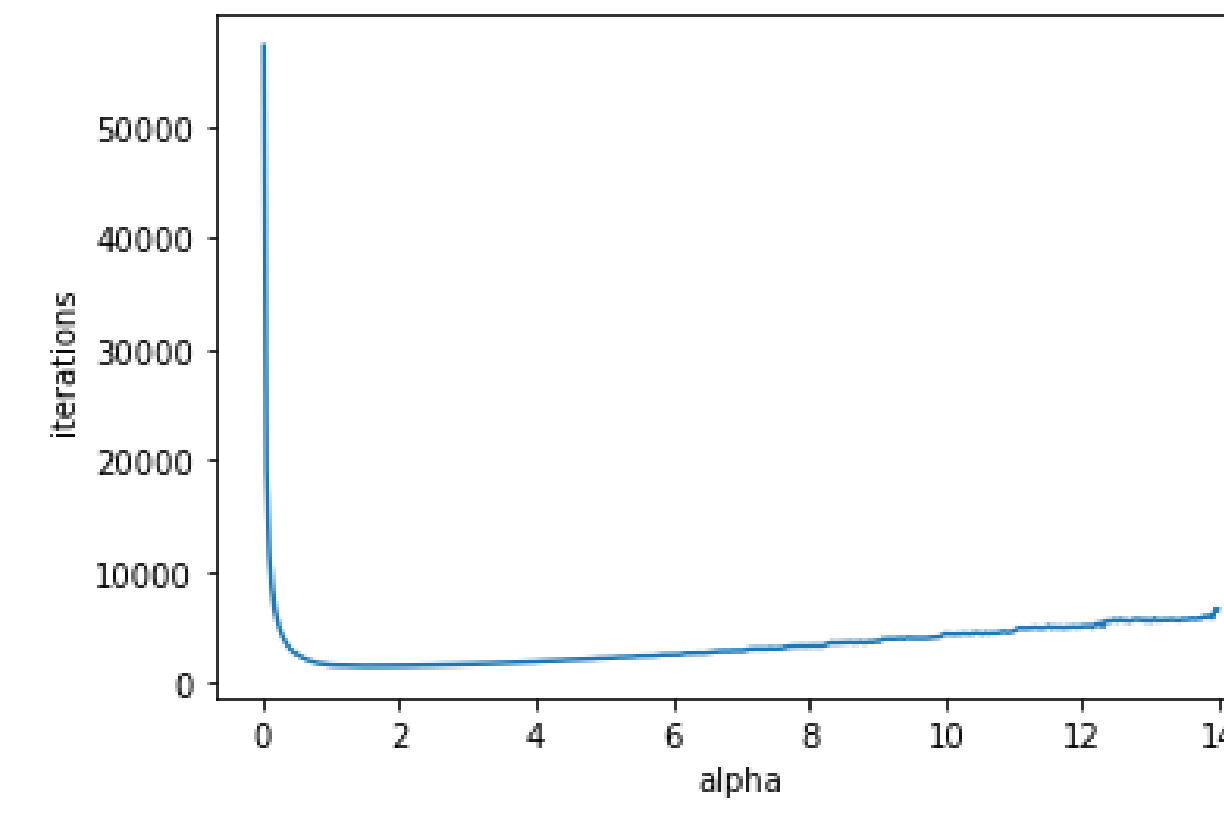


Fig. 4: iterations for the algorithm when the optimal  $\gamma$  value is used for that  $\alpha$

Preconditioning this way drastically improves the convergence rate, as shown in Fig. 3. In general, the lower  $L$  is, the faster the algorithm converges, which is to be expected as the

upper bound for  $\alpha$  is larger so we can take a larger step for each iteration. However, having a small  $L$  comes at a cost - the  $\lambda$  parameter tends to be more sensitive when we have a small  $L$  so it would require more work to calibrate it in order for the approximation to be decent.

### Different splitting methods

In equation (4), we set  $g_i(x) = \lambda |x_i|$ , which uses one coordinate of  $x$ , and  $f_i(x) = \frac{1}{2} |A_i x - b_i|^2$ , which uses one row of  $A$ . However, we could also split  $\lambda \|x\|_1$  and  $\frac{1}{2} \|Ax - b\|_2^2$  into  $g_i(x), f_i(x)$  differently, so that  $g_i(x)$  uses  $k$  coordinates from  $x$  and  $f_i(x)$  uses  $k$  rows from  $A$  (e.g. when  $k = 2$ ,  $g_1(x) = \lambda(|x_1| + |x_2|)$ ,  $f_1(x) = \frac{1}{2}(|A_1 x - b_1|^2 + |A_2 x - b_2|^2)$ )

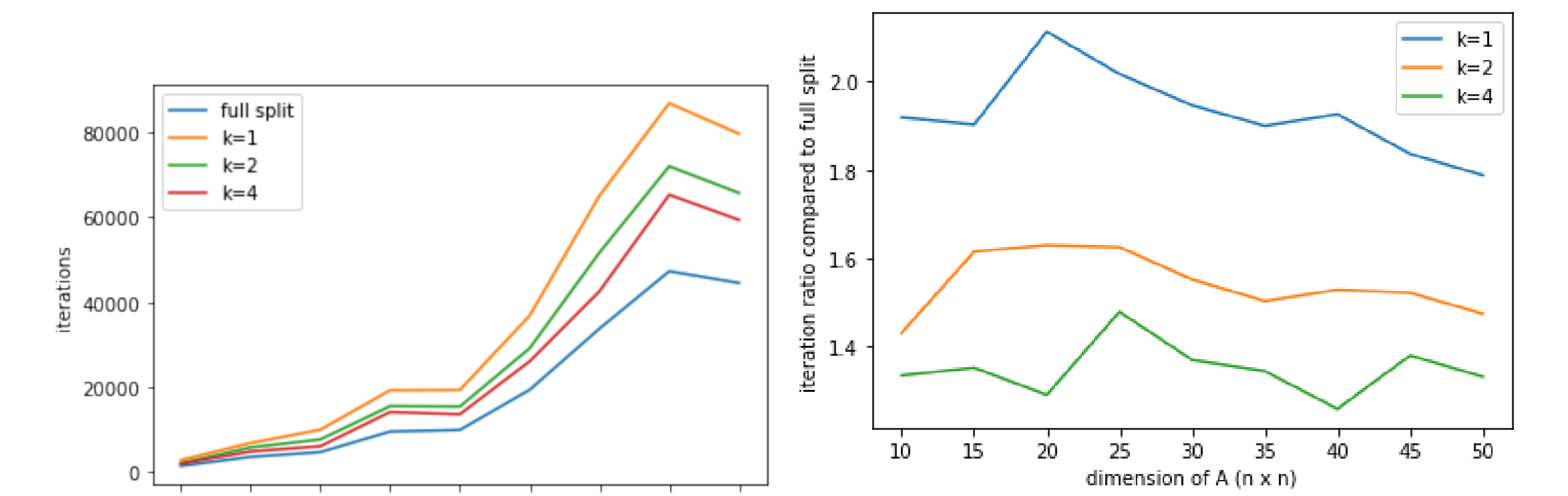


Fig. 6: Ratio of iteration between full split and other split methods

Fig. 5: Convergence rate of different splitting methods

Some notable observations are that:

- When  $g_1(x) = \lambda \|x\|_1, g_2(x) = 0, f_1(x) = \frac{1}{2} \|Ax - b\|_2^2$  (which we call the "full split"), the algorithm is the same as the Davis-Yin splitting
- The less number of  $g_i, f_i$  we split, the better the convergence rate
- Compared to full split (which can be seen as the a centralized algorithm), splitting with  $k = 1$  (which can be seen as distributed algorithm) needs to be approximately 2 times faster to outperform it, which is rather manageable in a distributed setting as each agents do calculations in parallel

## Remark

- We make the algorithm terminate when  $\|z^{k+1} - z^k\| \leq 1e - 3$  or when  $k = 10,000$ , whichever comes first
- Trials were done by generating a random sparse  $x$  vector with dimension  $n$ , and then randomly generates a  $n \times n$  matrix  $A$ .  $b$  is generated using  $b = Ax$

## Future Work

- Find the formula for the actual upper bound of  $\gamma$  and  $\alpha$
- Implement the distributed version of the algorithm

## References

- [1] Francisco J. Aragon-Artacho et al. *Distributed forward-backward methods for Ring Networks*. July 2022. URL: <https://doi.org/10.48550/arXiv.2112.00274>.
- [2] Yu. Malitsky. *Projected reflected gradient methods for monotone variational inequalities*. Feb. 2015. URL: <https://arxiv.org/abs/1502.04968>.
- [3] Ernest K. Ryu and Wotao Yin. *Large-scale convex optimization: Algorithms and analyses via monotone operators*. Cambridge University Press, 2023.