

STOCHASTIC GRADIENT ALGORITHMS FROM ODE SPLITTING PERSPECTIVE

Daniil Merkulov & Ivan Oseledets

Center for Computational and Data-Intensive Science and Engineering

Skolkovo Institute of Science and Technology

Bolshoy Boulevard 30, bld. 1, Moscow, Russia, 121205

daniil.merkulov@skolkovotech.ru, i.oseledets@skoltech.ru

ABSTRACT

We present a different view on stochastic optimization, which goes back to the splitting schemes for approximate solutions of ODE. In this work, we provide a connection between stochastic gradient descent approach and first-order splitting scheme for ODE. We consider the special case of splitting, which is inspired by machine learning applications and derive a new upper bound on the global splitting error for it. We present, that the Kaczmarz method is the limit case of the splitting scheme for the unit batch SGD for linear least squares problem. We support our findings with systematic empirical studies, which demonstrates, that a more accurate solution of local problems leads to the stepsize robustness and provides better convergence in time and iterations on the softmax regression problem.

1 INTRODUCTION

A lot of practical problems arising in machine learning require minimization of a finite sample average which can be written in the form

$$f(\theta) = \frac{1}{n} \sum_{i=1}^n f_i(\theta) \rightarrow \min_{\theta \in \mathbb{R}^p}, \quad (1)$$

where the sum goes over the *minibatches* of the original dataset. Vanilla stochastic gradient descent (SGD) method Robbins & Monro (1951) consists sequential steps in the direction of the gradient of $f_i(\theta)$, where i is to be chosen randomly from 1 to n without replacement.

$$\theta_{k+1} = \theta_k - h_k \nabla f_i. \quad (2)$$

Gradient descent method Cauchy (1847) can be considered as an Euler discretization of the ordinary differential equation (ODE) of the form of the gradient flow

$$\frac{d\theta}{dt} = -\nabla f(\theta). \quad (3)$$

In continuous time, SGD is often analyzed by introducing noise into the right-hand side of (3). However, for a real dataset, the distribution of the noise obtained by replacing the full gradient by its minibatch variant is not known and can be different for different problems. Instead, we propose a new view on the SGD as a *first-order splitting scheme* for (3), thus shedding a new light on SGD-type algorithms. This representation allows using more efficient local problem solvers for the approximation of the full gradient flow.

Contributions

- We show, that vanilla SGD could be considered as a splitting scheme for a full gradient flow and highlight connection between learning rate, batch size and size of the approximation step of SGD in continuous time.
- We propose new optimization scheme, which uses numerical integration of simple ODE at each step instead of stochastic gradient calculation and show empirically, that such approach can be considered as a stepsize-robust alternative to SGD for some practical ML problems.
- We present, that the Kaczmarz method is the limit case of the splitting scheme for the unit batch SGD for linear least squares problem.

2 SGD AS A SPLITTING SCHEME

We firstly consider simple ODE, where we can apply splitting idea and corresponding minimization problem. The best example to start from is simple ODE with right-hand-side, consisting of two summands:

$$\frac{d\theta}{dt} = -\frac{1}{2} (g_1(\theta) + g_2(\theta)) \quad (4)$$

Suppose, we want to find the solution $\theta(h)$ of (4) via integrating it on the small timestep h . The first order splitting scheme defined by solving first $\frac{d\theta}{dt} = -\frac{1}{2}g_1(\theta)$, $\theta(0) = \theta_0$ with exact solution $\theta_1(h)$ at the moment h , followed by $\frac{d\theta}{dt} = -\frac{1}{2}g_2(\theta)$, $\theta(0) = \theta_1(h)$ with exact solution $\theta_2(h)$ at the moment h . Thus, the first order approximation could be written as a combinations of both solutions $\theta^I(h) = \theta_2(h) \circ \theta_1(h) \circ \theta_0$.

It is interesting to study how the pure splitting scheme Marchuk (1968); Strang (1968) corresponds to the SGD approach. For this purpose, we consider an illustrative example of Gradient Flow equation 5, where the right-hand side of ODE is just the sum of operators acting on θ , which allows us to apply splitting scheme approximation directly.

$$\frac{d\theta}{dt} = -\frac{1}{2} \sum_{i=1}^2 \nabla f_i(\theta) = -\frac{1}{2} \nabla f_1(\theta) - \frac{1}{2} \nabla f_2(\theta) \quad (5)$$

Table 1: The table describes the correspondence between splitting scheme for discretized Gradient Flow ODE and epoch of SGD

Splitting step	Euler discretization	SGD Epoch	First-order splitting
$\frac{d\theta}{dt} = -\frac{1}{2} \nabla f_1(\theta)$	$\tilde{\theta}_I = \theta_0 - \frac{h}{2} \nabla f_1(\theta_0)$	$\tilde{\theta}_{SGD} = \theta_0 - h \nabla f_1(\theta_0)$	$\tilde{\theta}_I = \theta_0 - \frac{h}{2} \nabla f_1(\theta_0)$
$\frac{d\theta}{dt} = -\frac{1}{2} \nabla f_2(\theta)$	$\theta_I = \tilde{\theta}_I - \frac{h}{2} \nabla f_2(\tilde{\theta}_I)$	$\theta_{SGD} = \tilde{\theta}_{SGD} - h \nabla f_2(\tilde{\theta}_{SGD})$	$\theta_I = \tilde{\theta}_I - \frac{h}{2} \nabla f_2(\tilde{\theta}_I)$

Thus, we can conclude, that *one epoch of SGD is just the splitting scheme for the discretized Gradient Flow ODE with $2 \cdot h$ step size ($m \cdot h$ in case of m batches)*

Indeed, in SGD we go in the direction of the batch gradient, which stands for the Euler discretization of batch gradient flow ODE or *local ODE*. This idea gives additional intuition on the method. Given information about the Euler scheme limitation (first-order accuracy, stability issues), we propose to solve each local problem more precisely.

3 OPTIMIZATION STEP WITH ODE SOLVER

We propose to integrate local problem more precisely instead of Euler step in SGD. Solution of the local ODE problem involves replacing gradient in the right-hand side of gradient flow ODE 4 with batch gradient version. In our experiments the explicit Runge-Kutta method Dormand & Prince (1980); Shampine (1986) was used via scipy Virtanen et al. (2020) function odeint.

Table 2: The table presents ODE, which we need to solve at each step of the algorithm. The last column shows the ODE, which is needed to be solved at each iteration of the algorithm for each given problem.

Problem	Loss function	Batch gradient	Initial local ODE
Linear Least Squares	$f(\theta) = \frac{1}{n} \sum_{i=1}^m \ X_i \theta - y_i\ _2^2$	$\frac{1}{b} X_i^\top (X_i \theta - y_i)$	$\frac{d\theta}{dt} = -\frac{1}{n} X_i^\top (X_i \theta - y_i)$
Binary logistic regression	$f(\theta) = -\frac{1}{n} \sum_{i=1}^n \left(y_i \ln \sigma(\theta^\top x_i) + (1 - y_i) \ln (1 - \sigma(\theta^\top x_i)) \right)$	$\frac{1}{b} X_i^\top (\sigma(X_i \theta) - y_i)$	$\frac{d\theta}{dt} = -\frac{1}{n} X_i^\top (\sigma(X_i \theta) - y_i)$
One FC Layer + softmax	$f(\Theta) = -\frac{1}{n} \sum_{i=1}^n \log \left(\frac{y_i^\top e^{\Theta^\top x_i}}{\mathbf{1}^\top e^{\Theta^\top x_i}} \right)$	$\frac{1}{b} X_i^\top (s(\Theta^\top X_i^\top) - Y_i)^\top$	$\frac{d\Theta}{dt} = -\frac{1}{n} X_i^\top (s(\Theta^\top X_i^\top) - Y_i)^\top$

Algorithm 1: Splitting optimization

 θ_0 - initial parameter; b - batch size; α - learning rate; m - total number of batches
 $h := \alpha m$ $t := 0$ **for** $k = 0, 1, \dots$ **do** **for** $i = 1, 2, \dots, m$ **do** Formulate local ODE problem \mathcal{P}_i^k $\theta_{t+1} = \text{integrate } \mathcal{P}_i^k \text{ given an initial value } \theta(0) = \theta_t \text{ to the step } h$ $t := t + 1$ **end****end**

Typical machine learning problems involves dealing with mini-batch of size b , which is often less, than the number of trainable parameters p , which allows us to reduce dimensionality of the dynamic system via QR decomposition of each batch data matrix $X_i^\top = Q_i R_i$ (see details in the Appendix) and substitution $\eta_i = Q_i^\top \theta$. Note, that QR decomposition is only needed to be performed once before the training.

Table 3: The table shows initial local ODE and paired \mathcal{P}_i^k . Note, that $\eta_i \in \mathbb{R}^b$, while $\theta \in \mathbb{R}^p$

Initial local ODE	\mathcal{P}_i^k	Integration
$\frac{d\theta}{dt} = -\frac{1}{n} X_i^\top (X_i \theta - \mathbf{y}_i)$	$\frac{d\eta_i}{dt} = -\frac{1}{n} R_i (R_i^\top \eta_i - \mathbf{y}_i), \eta_i = Q_i^\top \theta$	analytical
$\frac{d\theta}{dt} = -\frac{1}{n} X_i^\top (\sigma(X_i \theta) - \mathbf{y}_i)$	$\frac{d\eta_i}{dt} = -\frac{1}{n} R_i (\sigma(R_i^\top \eta_i) - \mathbf{y}_i), \eta_i = Q_i^\top \theta$	odeint
$\frac{d\Theta}{dt} = -\frac{1}{n} X_i^\top (s(\Theta^\top X_i^\top) - Y_i)^\top$	$\frac{dH_i}{dt} = -\frac{1}{n} R_i (s(H_i^\top R) - Y_i)^\top, H_i = Q_i^\top \Theta$	odeint

There is an analytical solution for each local ODE in linear least squares case:

Theorem 1. For any matrix $\mathbf{x}_i \in \mathbb{R}^{b \times p}$, $b \leq p$, $\text{rank} X_i = b$, any vector of right-hand side $\mathbf{y}_i \in \mathbb{R}^b$ and initial vector of parameters θ_0 , there is a solution of the $\frac{d\theta}{dt} = -\frac{1}{n} X_i^\top (X_i \theta - \mathbf{y}_i)$, given by formula:

$$\theta(h) = Q_i e^{-\frac{1}{n} R_i R_i^\top h} (Q_i^\top \theta_0 - R_i^{-\top} \mathbf{y}_i) + Q_i R_i^{-\top} \mathbf{y}_i + (I - Q_i Q_i^\top) \theta_0, \quad (6)$$

where $Q_i \in \mathbb{R}^{p \times b}$ and $R_i \in \mathbb{R}^{b \times b}$ stands for the QR decomposition of the matrix \mathbf{X}_i^\top , $\mathbf{X}_i^\top = Q_i R_i$.

It is interesting to mention, that the splitting approach immediately leads to the Kaczmarz Kaczmarz. (1937); Strohmer & Vershynin (2009); Gower & Richtárik (2015) method for solving linear system in the same setting with unit batch size.

$$\lim_{h \rightarrow \infty} \theta(h) = \frac{(y_i - \mathbf{x}_i^\top \theta_0)}{\|\mathbf{x}_i\|^2} \mathbf{x}_i + \theta_0, \quad (7)$$

which is exact formula for Kaczmarz method for solving linear system. This result correlates with the statements of Needell et al. (2014), but provides us with a new sense of similarity between SGD and Kaczmarz method.

4 RESULTS

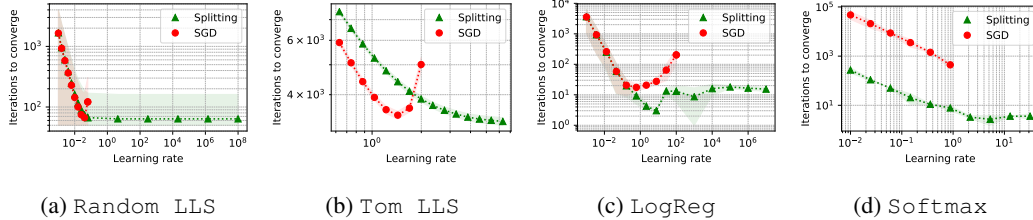
In this section, we describe the experimental setting. The majority of computations were performed on the NVIDIA DGX-2 cluster with 80 CPUs and 512 Gb RAM. We restricted the number of CPU usage per each experiment with an upper limit of 5 CPUs per experiment. All time measurements were done with the time library for Python. All experiments were done with the fixed random seed for reproducibility. For each experiment we performed 30 runs with random initialization and plotted trend line with the standard deviation.

Linear Least Squares Both random and the real linear systems were tested. For random linear system (`random lls`) we generated 10000×500 matrix with additive Gaussian noise of magnitude 0.01. Presented figures correspond to the batch size equals to 20. The real linear system (`tom lls`) is the standard tomography data from AIRTools II Hansen & Jørgensen (2018). Solution of the linear system is the 50×50 image reconstructed from solving 12780×2500 linear system. Presented figures correspond to the batch size equals to 60. Relative error 10^{-3} was used as the stopping criterion.

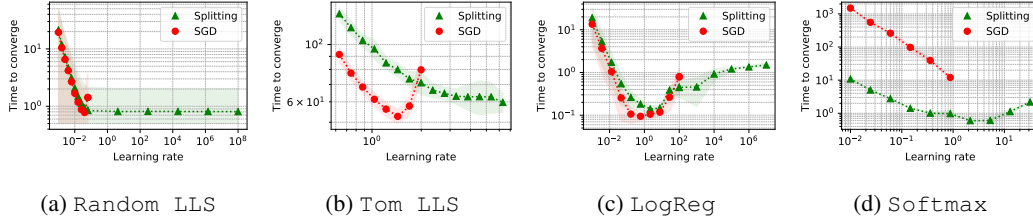
Binary Logistic Regression (`logreg`) In our experiments we used two classes from MNIST LeCun et al. (1998) dataset, which corresponds to the 0 and 1 digits. The size of the batch for presented figure is 50. Test error 0.001 was used as the stopping criterion.

Softmax Logistic Regression (`softmax`) We took Fashion MNIST Xiao et al. (2017) dataset with 60000 grayscale pictures from 10 classes. Each example is 28×28 image. The size of the batch for presented figure is 64. Test error 0.25 was used as the stopping criterion.

On the figures below we have two labels: SGD and `Splitting`, which stands for batch stochastic gradient descent and proposed algorithm. We use different constant learning rates to perform our experiments. All the learning rates tested for both algorithms. Lack of point of one algorithm on the graph means reaching the limit of iterations without achieving the termination rule.



As it is expected, SGD diverges starting from some value of learning rate, which is specific for each problem. While we can see comparative robustness of the proposed splitting optimization approach.



5 RELATED WORK

In this work, we presented another point of view on the nature of stochasticity in the stochastic gradient algorithms. From this perspective, different splitting schemes yield different stochastic gradient algorithms. We focused on the first-order splitting scheme for ODE, which corresponds to the SGD with the constant learning rate. Given this tractable setting, we performed a systematic empirical study of the local problem integration influence on the quality of the approximation scheme in machine learning problems. While the question of using these ideas to make general-purpose optimizer remains open, splitting optimization approach showed itself quite robust to the hyperparameter tuning for particular practical problems. Appendix to the paper contains proofs of the theorems and a new global error upper bounds for the first-order splitting for the special case. In Su et al. (2014) authors introduced second order ODE, which is equivalent (in the limit sense) to the gradient descent with Nesterov momentum Nesterov (1983). Generalization of these ideas was presented in Wibisono et al. (2016) with an arbitrary polynomial acceleration using the same parameter in ODE. General overview of the interplay between continuous-time and discrete-time points of view on dynamical systems and iterative optimization methods is covered in Helmke & Moore (2012), Evtushenko & Zhadan (1994)

REFERENCES

- Cauchy, A. Méthode générale pour la résolution des systèmes d'équations simultanées. *Comp. Rend. Sci. Paris*, 25(1847):536–538, 1847.
- Dormand, J. R. and Prince, P. J. A family of embedded runge-kutta formulae. *Journal of computational and applied mathematics*, 6(1):19–26, 1980.
- Evtushenko, Y. G. and Zhadan, V. G. Stable barrier-projection and barrier-newton methods in linear programming. *Computational Optimization and Applications*, 3(4):289–303, 1994.
- Gower, R. M. and Richtárik, P. Randomized iterative methods for linear systems. *SIAM Journal on Matrix Analysis and Applications*, 36(4):1660–1690, 2015.
- Hansen, P. C. and Jørgensen, J. S. Air tools ii: algebraic iterative reconstruction methods, improved implementation. *Numerical Algorithms*, 79(1):107–137, 2018.
- Helmke, U. and Moore, J. B. *Optimization and dynamical systems*. Springer Science & Business Media, 2012.
- Kaczmarz, S. Angenäherte auflösung von systemen linearer gleichungen. *Bull. Internat. Acad. Polon.Sci. Lettres A*, pp. 335–357, 1937.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Marchuk, G. I. Some application of splitting-up methods to the solution of mathematical physics problems. *Aplikace matematiky*, 13(2):103–132, 1968.
- Needell, D., Ward, R., and Srebro, N. Stochastic gradient descent, weighted sampling, and the randomized kaczmarz algorithm. In *Advances in neural information processing systems*, pp. 1017–1025, 2014.
- Nesterov, Y. E. A method of solving a convex programming problem with convergence rate $o(k^2)$. In *Doklady Akademii Nauk*, volume 269, pp. 543–547. Russian Academy of Sciences, 1983.
- Osher, S., Ruan, F., Xiong, J., Yao, Y., and Yin, W. Sparse recovery via differential inclusions. *Applied and Computational Harmonic Analysis*, 41(2):436–469, 2016.
- Robbins, H. and Monro, S. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
- Shampine, L. F. Some practical runge-kutta formulas. *Mathematics of computation*, 46(173):135–150, 1986.
- Sheng, Q. Global error estimates for exponential splitting. *IMA Journal of Numerical Analysis*, 14(1):27–56, 1994.
- Strang, G. On the construction and comparison of difference schemes. *SIAM Journal on Numerical Analysis*, 5(3):506–517, 1968.
- Strohmer, T. and Vershynin, R. A randomized kaczmarz algorithm with exponential convergence. *Journal of Fourier Analysis and Applications*, 15(2):262, 2009.
- Su, W., Boyd, S., and Candes, E. A differential equation for modeling nesterov’s accelerated gradient method: Theory and insights. In *Advances in Neural Information Processing Systems*, pp. 2510–2518, 2014.
- Virtanen, P., Gommers, R., Oliphant, T. E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S. J., Brett, M., Wilson, J., Jarrod Millman, K., Mayorov, N., Nelson, A. R. J., Jones, E., Kern, R., Larson, E., Carey, C., Polat, İ., Feng, Y., Moore, E. W., Van der Plas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E. A., Harris, C. R., Archibald, A. M., Ribeiro, A. H., Pedregosa, F., van Mulbregt, P., and Contributors, S. . . SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: <https://doi.org/10.1038/s41592-019-0686-2>.
- Wibisono, A., Wilson, A. C., and Jordan, M. I. A variational perspective on accelerated methods in optimization. *proceedings of the National Academy of Sciences*, 113(47):E7351–E7358, 2016.
- Xiao, H., Rasul, K., and Vollgraf, R. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

A UPPER BOUND ON THE GLOBAL SPLITTING ERROR

Suppose, that we have only two batches, and the problem (17) is consistent, i.e. there exists an exact solution θ_* such as $X\theta_* = \mathbf{y}$. The GD flow has the form

$$\begin{aligned}\frac{d\theta}{dt} &= -X^\top(X\theta - \mathbf{y}) = -X^\top X(\theta - \theta_*) = \\ &= -(X_1^\top X_1 + X_2^\top X_2)(\theta - \theta_*),\end{aligned}\tag{8}$$

i.e. the splitting scheme corresponds to a linear operator splitting

$$A = A_1 + A_2, \quad A = -X^\top X, \quad A_i = -X_i^\top X_i, \quad i = 1, 2.$$

Both A_1 and A_2 are symmetric non-negative definite matrices. Without loss of generality, we can assume that $\theta_* = 0$,

Suppose that the rank of A is r_1 and the rank of A_2 is r_2 . Then, we can write them as

$$A_i = Q_i B_i Q_i^*,$$

where Q_i is an $N \times r_i$ matrix with orthonormal columns. The following Lemma gives the representation of the matrix exponents of such matrices.

Lemma 1. *Let $A = QBQ^*$, where Q is an $N \times r$ matrix with orthonormal columns, and B is an $r \times r$ matrix. Then,*

$$e^{tA} = (I - QQ^*) + Qe^{tB}Q^*.\tag{9}$$

To prove (15) we note that

$$\begin{aligned}e^{tA} &= \sum_{k=0}^{\infty} \frac{t^k A^k}{k!} = \sum_{k=0}^{\infty} \frac{t^k QB^k Q^*}{k!} = \\ &= I - QQ^* + QQ^* + Q \sum_{k=1}^{\infty} \frac{t^k B^k}{k!} Q^* = \\ &= (I - QQ^*) + Qe^{tB}Q^*.\end{aligned}$$

Lemma 2. *Let $A_1, A_2 \in \mathbb{S}_+^p$ be the square negative semidefinite matrices, that don't have full rank, i.e. $\text{rank } A_1 \leq p$ and $\text{rank } A_2 \leq p$. While the sum of those matrices has full rank, i.e. $A = A_1 + A_2, \text{rank } A = p$. Then, the global upper bound error will be written as follows:*

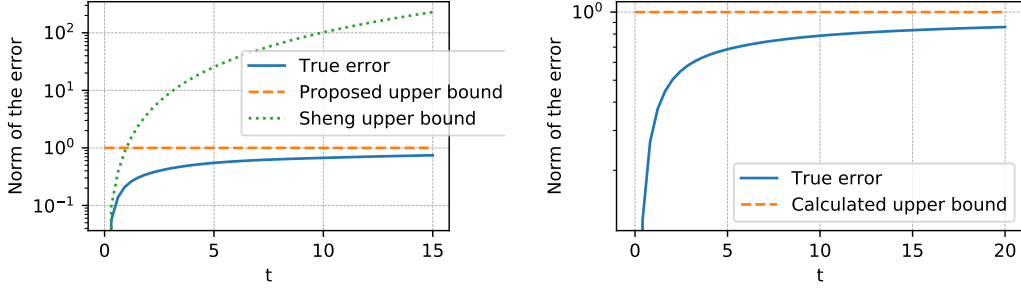
$$\lim_{t \rightarrow \infty} \|e^{A_2 t} e^{A_1 t} - e^{At}\| = \|(I - Q_2 Q_2^*)(I - Q_1 Q_1^*)\|\tag{10}$$

Proof. The proof is straightforward. We will use the low rank matrix exponential decomposition from the Lemma 3

$$e^{A_i t} = \Pi_i + Q_i e^{B_i t} Q_i^*, \text{ where } \Pi_i = I - Q_i Q_i^*; i = 1, 2$$

$$\begin{aligned}\lim_{t \rightarrow \infty} \|e^{A_2 t} e^{A_1 t} - e^{At}\| &= \\ &= \lim_{t \rightarrow \infty} \|(\Pi_2 + Q_2 e^{B_2 t} Q_2^*)(\Pi_1 + Q_1 e^{B_1 t} Q_1^*) - e^{At}\| = \\ &= \lim_{t \rightarrow \infty} \|\Pi_2 \Pi_1 + Q_1 e^{B_1 t} Q_1^* \Pi_2 + \Pi_1 Q_2 e^{B_2 t} Q_2^* + \\ &\quad + Q_1 e^{B_1 t} Q_1^* Q_2 e^{B_2 t} Q_2^* - e^{At}\| = \\ &= \Pi_2 \Pi_1\end{aligned}$$

Since all matrices B_1, B_2, A are negative all the matrix exponentials are decaying: $\|e^{At}\| \leq e^{t\mu(A)} \forall t \geq 0$, where $\mu(A) = \lambda_{\max}\left(\frac{A+A^\top}{2}\right)$ - the logarithmic norm. \square



(a) Global error of the splitting scheme. Initial random full rank matrix $X \in \mathbb{R}^{100 \times 100}$ was splitted by rows. $X_1, X_2 \in \mathbb{R}^{50 \times 100}$. Target matrices were obtained the following way: $A_1 = -X_1^* X_1, A_2 = -X_2^* X_2, A = -X^* X$. So A_1, A_2 are negative and lacking full rank, while $A = A_1 + A_2$ has full rank.

(b) Global upper bound on the splitting scheme in case of 40 summands in the right-hand side.

The graph presented on the Figure 3a describes . One can easily see significant difference between existing global upper bounds for that case Sheng (1994) and derived upper bound.

Theorem 2. Let $A_1, A_2, \dots, A_b \in \mathbb{S}_+^p$ be the square negative semidefinite matrices, that don't have full rank, i.e. $\text{rank } A_i \leq p, \forall i = 1, \dots, b$. While the sum of those matrices has full rank, i.e. $A = \sum_{i=1}^b A_i, \text{rank } A = p$. Then, the global upper bound error will be written as follows:

$$\lim_{t \rightarrow \infty} \|e^{A_b t} \dots e^{A_1 t} - e^{A t}\| = \left\| \prod_{i=1}^b \Pi_{b-i+1} \right\|, \quad (11)$$

where $\Pi_i = I - Q_i Q_i^*$ and $A_i = Q_i B_i Q_i^*$ and Q_i is a matrix with orthonormal columns.

The graph on the Figure 3b shows empirical validity of the presented upper bound.

B PROOFS

Theorem 1. For any matrix $\mathbf{x}_i \in \mathbb{R}^{b \times p}, b \leq p, \text{rank } X_i = b$, any vector of right-hand side $\mathbf{y}_i \in \mathbb{R}^b$ and initial vector of parameters θ_0 , there is a solution of the $\frac{d\theta}{dt} = -\frac{1}{n} X_i^\top (X_i \theta - \mathbf{y}_i)$, given by formula:

$$\theta(h) = Q_i e^{-\frac{1}{n} R_i R_i^\top h} (Q_i^\top \theta_0 - R_i^{-\top} \mathbf{y}_i) + Q_i R_i^{-\top} \mathbf{y}_i + (I - Q_i Q_i^\top) \theta_0, \quad (6)$$

where $Q_i \in \mathbb{R}^{p \times b}$ and $R_i \in \mathbb{R}^{b \times b}$ stands for the QR decomposition of the matrix $\mathbf{X}_i^\top, \mathbf{X}_i^\top = Q_i R_i$.

Proof. Given $X_i^\top = Q_i R_i$, we have $(I - Q_i Q_i^\top) X_i^\top = 0$. Note, that Q_i is left unitary matrix, i.e. $Q_i^\top Q_i = I$.

$$\begin{aligned} \frac{d\theta}{dt} &= -\frac{1}{n} X_i^\top (X_i \theta - \mathbf{y}_i) \\ (I - Q_i Q_i^\top) \frac{d\theta}{dt} &= 0 \\ \frac{d\theta}{dt} &= Q_i \frac{d(Q_i^\top \theta)}{dt} \quad Q_i^\top \theta = \eta_i \\ \frac{d\theta}{dt} &= Q_i \frac{d\eta_i}{dt} \quad \text{integrate from 0 to } h \\ \theta(h) &= Q_i (\eta_i(h) - \eta_i(0)) + \theta_0 \end{aligned} \quad (12)$$

On the other hand:

$$\begin{aligned}
\frac{d\boldsymbol{\eta}_i}{dt} &= Q_i^\top \frac{d\boldsymbol{\theta}}{dt} = -\frac{1}{n} Q_i^\top X_i^\top (X_i \boldsymbol{\theta} - \mathbf{y}_i) = \\
&= -\frac{1}{n} Q_i^\top Q_i R_i (R_i^\top Q_i^\top \boldsymbol{\theta} - \mathbf{y}_i) = \\
&= -\frac{1}{n} (R_i R_i^\top \boldsymbol{\eta}_i - R_i \mathbf{y}_i)
\end{aligned} \tag{13}$$

Consider the moment of time $t = \infty$. $\frac{d\boldsymbol{\eta}_i}{dt} = 0$, since $\exists \boldsymbol{\theta}^*, Q_i^\top \boldsymbol{\theta}^* = \boldsymbol{\eta}_i^*$. Also consider (13):

$$\begin{aligned}
\frac{d\boldsymbol{\eta}_i}{dt} = 0 &= -\frac{1}{n} (R_i R_i^\top \boldsymbol{\eta}_i^* - R_i \mathbf{y}_i) \\
R_i \mathbf{y}_i &= R_i R_i^\top \boldsymbol{\eta}_i^*
\end{aligned} \tag{14}$$

Now we look at the (13) with the replacement, given in (14):

$$\begin{aligned}
\frac{d\boldsymbol{\eta}_i}{dt} &= -\frac{1}{n} (R_i R_i^\top \boldsymbol{\eta}_i - R_i R_i^\top \boldsymbol{\eta}_i^*) \\
\frac{d\boldsymbol{\eta}_i}{dt} &= -\frac{1}{n} R_i R_i^\top (\boldsymbol{\eta}_i - \boldsymbol{\eta}_i^*) \quad \text{integrate from 0 to } h \\
\boldsymbol{\eta}_i(h) - \boldsymbol{\eta}_i^* &= e^{-\frac{1}{n} R_i R_i^\top h} (\boldsymbol{\eta}_i(0) - \boldsymbol{\eta}_i^*) \\
&\quad \text{while } \boldsymbol{\eta}_i^* = R_i^{-\top} \mathbf{y}_i, \boldsymbol{\eta}_i(0) = Q_i^\top \boldsymbol{\theta}_0 \\
\boldsymbol{\eta}_i(h) &= e^{-\frac{1}{n} R_i R_i^\top h} (Q_i^\top \boldsymbol{\theta}_0 - R_i^{-\top} \mathbf{y}_i) + R_i^{-\top} \mathbf{y}_i
\end{aligned}$$

Using (12) we obtain the target formula

$$\begin{aligned}
\boldsymbol{\theta}(h) &= Q_i e^{-\frac{1}{n} R_i R_i^\top h} (Q_i^\top \boldsymbol{\theta}_0 - R_i^{-\top} \mathbf{y}_i) + \\
&\quad + Q_i R_i^{-\top} \mathbf{y}_i + (I - Q_i Q_i^\top) \boldsymbol{\theta}_0,
\end{aligned}$$

□

Lemma 3. Let $A = QBQ^*$, where Q is an $N \times r$ matrix with orthonormal columns, and B is an $r \times r$ matrix. Then,

$$e^{tA} = (I - QQ^*) + Qe^{tB}Q^*. \tag{15}$$

To prove (15) we note that

$$\begin{aligned}
e^{tA} &= \sum_{k=0}^{\infty} \frac{t^k A^k}{k!} = \sum_{k=0}^{\infty} \frac{t^k QB^kQ^*}{k!} = \\
&= I - QQ^* + QQ^* + Q \sum_{k=1}^{\infty} \frac{t^k B^k}{k!} Q^* = \\
&= (I - QQ^*) + Qe^{tB}Q^*.
\end{aligned}$$

Lemma 4. Let $A_1, A_2 \in \mathbb{S}_+^p$ be the square negative semidefinite matrices, that don't have full rank, i.e. $\text{rank } A_1 \leq p$ and $\text{rank } A_2 \leq p$. While the sum of those matrices has full rank, i.e. $A = A_1 + A_2, \text{rank } A = p$. Then, the global upper bound error will be written as follows:

$$\lim_{t \rightarrow \infty} \|e^{A_2 t} e^{A_1 t} - e^{At}\| = \|(I - Q_2 Q_2^*)(I - Q_1 Q_1^*)\| \tag{16}$$

Proof. The proof is straightforward. We will use the low rank matrix exponential decomposition from the Lemma 3

$$e^{A_i t} = \Pi_i + Q_i e^{B_i t} Q_i^*, \text{ where } \Pi_i = I - Q_i Q_i^*; i = 1, 2$$

$$\begin{aligned}
& \lim_{t \rightarrow \infty} \|e^{A_2 t} e^{A_1 t} - e^{A t}\| = \\
& = \lim_{t \rightarrow \infty} \|(\Pi_2 + Q_2 e^{B_2 t} Q_2^*)(\Pi_1 + Q_1 e^{B_1 t} Q_1^*) - e^{A t}\| = \\
& = \lim_{t \rightarrow \infty} \|\Pi_2 \Pi_1 + Q_1 e^{B_1 t} Q_1^* \Pi_2 + \Pi_1 Q_2 e^{B_2 t} Q_2^* + \\
& + Q_1 e^{B_1 t} Q_1^* Q_2 e^{B_2 t} Q_2^* - e^{A t}\| = \\
& = \Pi_2 \Pi_1
\end{aligned}$$

Since all matrices B_1, B_2, A are negative all the matrix exponentials are decaying: $\|e^{A t}\| \leq e^{t\mu(A)} \forall t \geq 0$, where $\mu(A) = \lambda_{\max}\left(\frac{A+A^\top}{2}\right)$ - the logarithmic norm. \square

C APPLICATIONS

C.1 LINEAR LEAST SQUARES

C.1.1 PROBLEM

Let $f_i(\boldsymbol{\theta}) = \|\mathbf{x}_i^\top \boldsymbol{\theta} - y_i\|^2$, then problem (1) is the linear least squares problem, which can be written as

$$f(\boldsymbol{\theta}) = \frac{1}{n} \|X\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \frac{1}{n} \sum_{i=1}^s \|X_i \boldsymbol{\theta} - \mathbf{y}_i\|_2^2 \rightarrow \min_{\boldsymbol{\theta} \in \mathbb{R}^p}, \quad (17)$$

where $X \in \mathbb{R}^{n \times p}$ and $\mathbf{y} \in \mathbb{R}^p$ and the second part of the equation stands for s mini-batches with size b regrouping ($b \cdot s = n$): $X_i \in \mathbb{R}^{b \times p}, \mathbf{y}_i \in \mathbb{R}^b$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \nabla f(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^s X_i^\top (X_i \boldsymbol{\theta} - \mathbf{y}_i) \quad (18)$$

The gradient flow equation will be written as follows:

$$\frac{d\boldsymbol{\theta}}{dt} = -\frac{1}{n} \sum_{i=1}^s X_i^\top (X_i \boldsymbol{\theta} - \mathbf{y}_i) \quad (19)$$

C.1.2 EXACT SOLUTION OF THE LOCAL PROBLEM

Theorem 1 gives us explicit formula for the local solution:

$$\boldsymbol{\theta}(h) = Q_i e^{-\frac{1}{n} R_i R_i^\top h} (Q_i^\top \boldsymbol{\theta}_0 - R_i^{-\top} \mathbf{y}_i) + Q_i R_i^{-\top} \mathbf{y}_i + (I - Q_i Q_i^\top) \boldsymbol{\theta}_0$$

C.1.3 KACZMARZ AS THE LIMIT CASE OF SPLITTING

Kaczmarz method Kaczmarz. (1937), Strohmer & Vershynin (2009), Gower & Richtárik (2015) is a well-known iterative algorithm for solving linear systems. It is interesting to mention, that splitting approach immediately leads to the Kaczmarz method for solving linear system in the same setting with unit batch size.

When the batch size is equal to one, we need to do n QR decompositions for each transposed batch matrix, which is just column vector \mathbf{x}_i in our case:

$$\mathbf{x}_i = \mathbf{q}_i \mathbf{r}_i = \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|} \|\mathbf{x}_i\| \quad (20)$$

Now, we need to use (6) to derive analytic local solution in that case:

$$\begin{aligned}
\boldsymbol{\theta}(h) &= \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|} e^{-\frac{\|\mathbf{x}_i\|^2 h}{n}} \left(\frac{\mathbf{x}_i^\top}{\|\mathbf{x}_i\|} \boldsymbol{\theta}_0 - \frac{y_i}{\|\mathbf{x}_i\|} \right) + \\
&+ \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|^2} y_i + \left(I - \frac{\mathbf{x}_i \mathbf{x}_i^\top}{\|\mathbf{x}_i\|^2} \right) \boldsymbol{\theta}_0 = \\
&= \frac{(y_i - \mathbf{x}_i^\top \boldsymbol{\theta}_0)}{\|\mathbf{x}_i\|^2} \left(1 - e^{-\frac{\|\mathbf{x}_i\|^2 h}{n}} \right) \mathbf{x}_i + \boldsymbol{\theta}_0
\end{aligned}$$

It can be easily seen, that:

$$\lim_{h \rightarrow \infty} \boldsymbol{\theta}(h) = \frac{(y_i - \mathbf{x}_i^\top \boldsymbol{\theta}_0)}{\|\mathbf{x}_i\|^2} \mathbf{x}_i + \boldsymbol{\theta}_0, \quad (21)$$

which is exact formula for Kaczmarz method for solving linear system. This result correlates with the statements of Needell et al. (2014), but provides us with a new sense of similarity between SGD and Kaczmarz method.

C.2 BINARY LOGISTIC REGRESSION

C.2.1 PROBLEM

In this classification task then problem (1) takes the following form:

$$-\frac{1}{n} \sum_{i=1}^n \left(y_i \ln \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) + (1 - y_i) \ln(1 - \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i)) \right) \rightarrow \min_{\boldsymbol{\theta} \in \mathbb{R}^p}, \quad (22)$$

where $\sigma(x) = \frac{1}{1+e^{-x}}$ is the sigmoid function, while $y_i \in \{0, 1\}$ stands for the label of the object class.

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \nabla f(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i (\sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) - y_i) \quad (23)$$

The gradient flow equation will be written as follows:

$$\frac{d\boldsymbol{\theta}}{dt} = -\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i (\sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) - y_i) \quad (24)$$

Our particular interest lies in mini-batch reformulation of the given problem. We consider s mini-batches with size b regrouping ($b \cdot s = n$): $X_i \in \mathbb{R}^{b \times p}$, $\mathbf{y}_i \in \mathbb{R}^b$ and $\sigma(\mathbf{x})$ stands for the element-wise sigmoid function.

$$\frac{d\boldsymbol{\theta}}{dt} = -\frac{1}{n} \sum_{i=1}^s X_i^\top (\sigma(X_i \boldsymbol{\theta}) - \mathbf{y}_i) \quad (25)$$

C.2.2 SPLITTING SCHEME AND LOCAL PROBLEM

Since we are applying splitting scheme to find the approximate solution of the (25), each local problem should be written as follows:

$$\frac{d\boldsymbol{\theta}}{dt} = -\frac{1}{n} X_i^\top (\sigma(X_i \boldsymbol{\theta}) - \mathbf{y}_i) \quad (26)$$

Note, that this is not linear equation and cannot be solved as easy as in Theorem 1. However, we can apply the same technique to reduce the dimension of ODE, which is needed to be solved numerically.

Suppose, we have QR decomposition of each batch data matrix $X_i^\top = Q_i R_i$, then we can multiply both sides of (26) on the $(I - Q_i Q_i^\top)$ on the left.

$$\begin{aligned}
(I - Q_i Q_i^\top) \frac{d\boldsymbol{\theta}}{dt} &= (I - Q_i Q_i^\top) \frac{1}{n} X_i^\top (\mathbf{y}_i - \sigma(X_i \boldsymbol{\theta})) \\
\frac{d\boldsymbol{\theta}}{dt} &= Q_i \frac{d(Q_i^\top \boldsymbol{\theta})}{dt} \quad Q_i^\top \boldsymbol{\theta} = \boldsymbol{\eta}_i \\
\frac{d\boldsymbol{\theta}}{dt} &= Q_i \frac{d\boldsymbol{\eta}_i}{dt} \quad \text{integrate from 0 to } h \\
\boldsymbol{\theta}(h) &= Q_i (\boldsymbol{\eta}_i(h) - \boldsymbol{\eta}_i(0)) + \boldsymbol{\theta}_0
\end{aligned} \tag{27}$$

On the other hand:

$$\begin{aligned}
\frac{d\boldsymbol{\eta}_i}{dt} &= Q_i^\top \frac{d\boldsymbol{\theta}}{dt} = -\frac{1}{n} Q_i^\top X_i^\top (\sigma(X_i \boldsymbol{\theta}) - \mathbf{y}_i) = \\
&= -\frac{1}{n} Q_i^\top Q_i R_i (\sigma(X_i \boldsymbol{\theta}) - \mathbf{y}_i) = \\
&= -\frac{1}{n} R_i (\sigma(X_i \boldsymbol{\theta}) - \mathbf{y}_i)
\end{aligned}$$

Recall, that each hypothesis function depends on linear function $\mathbf{x}_i^\top \boldsymbol{\theta}$, which means, that in batch reformulation it is just entries of the vector $X_i \boldsymbol{\theta}$. Since we have QR decomposition of X_i^\top , we can write: $X_i \boldsymbol{\theta} = R_i^\top Q_i^\top \boldsymbol{\theta} = R_i^\top \boldsymbol{\eta}_i$. In other words:

$$\frac{d\boldsymbol{\eta}_i}{dt} = -\frac{1}{n} R_i (\sigma(R_i^\top \boldsymbol{\eta}_i) - \mathbf{y}_i), \tag{28}$$

To sum it up, we need to solve (28) (which is much simpler, than original differential equation (26)), than substitute it to the (27) with $\boldsymbol{\eta}_i(0) = Q_i^\top \boldsymbol{\theta}_0$. Note, that matrices Q_i and R_i can be computed only once before the training.

C.3 SOFTMAX REGRESSION

C.3.1 PROBLEM

In this classification task then problem (1) takes the following form:

$$-\frac{1}{n} \sum_{i=1}^n \log \left(\frac{\mathbf{y}_i^\top e^{\Theta^\top \mathbf{x}_i}}{\mathbf{1}^\top e^{\Theta^\top \mathbf{x}_i}} \right) \rightarrow \min_{\Theta \in \mathbb{R}^{p \times K}}, \tag{29}$$

where $e^{\mathbf{x}}$ is element-wise exponential function, while $\mathbf{y}_i \in \mathbb{R}^K$ stands for the one-hot encoding of the i -th object label.

$$\nabla_{\Theta} f(\Theta) = -\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \left(\mathbf{y}_i - \frac{e^{\Theta^\top \mathbf{x}_i}}{\mathbf{1}^\top e^{\Theta^\top \mathbf{x}_i}} \right)^\top \tag{30}$$

$$\nabla_{\Theta} f(\Theta) = -\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i (\mathbf{y}_i - s(\Theta^\top \mathbf{x}_i))^\top \tag{31}$$

Here we use $s(\mathbf{x})$ as a softmax function of a vector \mathbf{x} , i.e. $s(\mathbf{x}) = \frac{e^{\mathbf{x}}}{\mathbf{1}^\top e^{\mathbf{x}}}$. While mini-batch reformulation will take the following form:

$$\nabla_{\Theta} f(\Theta) = -\frac{1}{n} \sum_{i=1}^s X_i^\top (Y_i - s(\Theta^\top X_i^\top))^\top, \tag{32}$$

where $s(X) = \begin{bmatrix} s(\mathbf{x}_{(1)}) & s(\mathbf{x}_{(2)}) & \cdots & s(\mathbf{x}_{(b)}) \end{bmatrix}$ is a column-wise softmax function. Indeed, in a very similar manner to the binary logistic regression we can write down gradientflow ODE for softmax regression in a mini-batch form:

$$\frac{d\Theta}{dt} = -\frac{1}{n} \sum_{i=1}^s X_i^\top (s(\Theta^\top X_i^\top) - Y_i)^\top \quad (33)$$

Splitting method requires the local problem, which is focused on a single minibatch:

$$\frac{d\Theta}{dt} = -\frac{1}{n} X_i^\top (s(\Theta^\top X_i^\top) - Y_i)^\top \quad (34)$$

$$\begin{aligned} (I - Q_i Q_i^\top) \frac{d\Theta}{dt} &= (I - Q_i Q_i^\top) \frac{1}{n} X_i^\top (Y_i - s(\Theta^\top X_i^\top))^\top \\ \frac{d\Theta}{dt} &= Q_i \frac{d(Q_i^\top \Theta)}{dt} \quad Q_i^\top \Theta = H_i \\ \frac{d\Theta}{dt} &= Q_i \frac{dH_i}{dt} \quad \text{integrate from 0 to } h \\ \Theta(h) &= Q_i (H_i(h) - H_i(0)) + \Theta_0 \end{aligned} \quad (35)$$

On the other hand:

$$\begin{aligned} \frac{dH_i}{dt} &= Q_i^\top \frac{d\Theta}{dt} = -\frac{1}{n} Q_i^\top X_i^\top (s(\Theta^\top X_i^\top) - Y_i)^\top = \\ &= -\frac{1}{n} Q_i^\top Q_i R_i (s(\Theta^\top X_i^\top) - Y_i)^\top = \\ &= -\frac{1}{n} R_i (s(\Theta^\top X_i^\top) - Y_i)^\top = \\ &= -\frac{1}{n} R_i (s(H_i^\top R) - Y_i)^\top \end{aligned}$$

Now we need to solve ODE of variable of the size $b \times k$, rather, than $p \times k$.