

Dynamical Characterization & Analysis of the Optimization Algorithms: Linearized Bregman and Iterative Shrinkage Thresholding Algorithm

Yihua Xu and Ariba Khan Mentor: Dr. Rachel Kuske

Gradient Descent (GD)

<u>What is GD?</u> One of the most common optimization algorithm used to calculate the minimum value of a differentiable function.

$$x_{k+1} = x_k - t_k A^T (Ax_k - b)$$

Basic Algorithm Process:

- 1. x is a minimizer of ||Ax-b||
- 2. Set random initial condition for x
- 3. Iterate over the time step
- 4. Find the local minimum



 $\begin{array}{ll} \underline{2 \text{ Ways of Measuring Algorithm Behavior:}}\\ \text{Normalized Residual} & \mathcal{R}(x_k) = \frac{\|Ax_k - b\|_2}{\|b\|_2}.\\ \text{Model Error} & M(x_k) = \frac{\|x_{True} - x_k\|_2}{\|x_{True}\|_2} \end{array}$

**Other optimization algorithms have been developed based on GD that deal with different data types

Methods Introduction

LB & ISTA are both based on Gradient Descent but deal with sparse solutions

	Linearized Bregman (LB)	Iterative Shrinkage Thresholding Algorithm (ISTA)
Iterative Formula	$z_{k+1} = z_k - t_k A^T (Ax_k - b) x_{k+1} = S_\lambda (z_{k+1})$	$z_{k+1} = x_k - t_k A^T (Ax_k - b) x_{k+1} = S_\lambda (z_{k+1})$
Time Step	$t_k = \frac{\ Ax_k - b\ _2^2}{\ A^\top (Ax_k - b)\ _2^2}.$	$t_k = \frac{1}{\ A\ _2^2}$
Shrinkage Process	$[S_{\lambda}(z_{k+1})]_i = \begin{cases} z_i - \lambda \\ 0 \\ z_i + \lambda \end{cases}$	$\begin{array}{ll} \lambda & \text{if } z_i > \lambda \\ & \text{if } -\lambda \leq z_i \leq \lambda \\ \lambda & \text{if } z_i < -\lambda \end{array}$

Variable Definitions & Data Setup

Variable	Definition	Data Setup
k	Number of iterations	Chose varying k, we will specify in the following slides
A	A matrix; normally a tall matrix because in the real world there are often many variables to consider	800*70 matrix of normally distributed random variables
b	Created by the multiplication of A and xTrue, with added noise	A * xTrue + (0.1 * noise)
xTrue	The correct value, the algorithm's final output should be around this vector value	A size 70 vector created with a random exponential decay, including large, small, and zero entries; this is done to mimic sparse solutions in the real world
Noise	Source of chatter	random normal distribution with mean=0 and variance=1
t _k	Time step	Dynamic time step for LB and constant time step for ISTA
x ₀	Initial guess for x	Vector of zeros

Preview of Our Findings & Applications

- Linearized Bregman performs best with large entries and a large $\boldsymbol{\lambda}$
- ISTA performs best with small entries and small λ
- Linearized Bregman outperforms ISTA when subsampling
- Hybrid Method is a combination of both methods and outperforms both LB and ISTA when dealing with a mix of entry sizes





Comparing λ

0.01

Choices concerning hyperparameter λ

1.0

Comparison of different values of lambda with respect to LB

LB

Though converging slightly quicker early on when λ is small, it would lead to overfitting/ fluctuation later on. Overall, a moderate or comparatively large λ is most suitable for LB, a reasonable range could be from 0.5 - 1.

$$z_{k+1} = z_k - t_k A^T (Ax_k - b) x_{k+1} = S_\lambda (z_{k+1})$$



Comparison of different values of lambda with respect to ISTA

Obviously, if we put a large
$$\lambda$$
 in ISTA, the residual would decrease to a very limited extent. (Basically predicting most entries as zero). ISTA needs a very small λ to reach a low residual. A good value for λ_{ISTA} is around 0.01.

$$z_{k+1} = x_k - t_k A^T (Ax_k - b) x_{k+1} = S_\lambda (z_{k+1})$$

Comparing λ



For large entries, The difference between different λ is very small early on in the iterations and overall convergence rate is quick, while comparatively larger λ could result in smaller chatter later on.



For small entries and small λ , ISTA tends to predict well and converge quickly without fluctuation in later iterations.

Introduction to Hybrid Method & Initial Results

The Hybrid Method

<u>General Idea:</u> x_{hybrid} follows x_{lb} for the first K iterations, then smaller entries follow ISTA and larger entries follow LB



Subsampling .vs. No Subsampling

<u>What is subsampling?</u> Taking a different random subset of the data by choosing k rows out of entire set for each iteration.

<u>Why subsampling?</u> In real world setting, iterating over the entire dataset usually mean a very time-consuming process. By subsampling, we strive to solve problems in shorter periods of time with only few full datapasses. Why is LB more suited for subsampling? ISTA leaves out certain information, because it uses x_k , which undergoes a shrinkage process every iteration, also since A_k is mostly different from iteration to iteration when subsampling, when throwing out information, the impact would be more severe than using the complete dataset.

$$egin{aligned} &z_{k+1} = \overbrace{X_k}^{-} - t_k \; A_k^T \; (A_k x_k - b_k) \ &x_{k+1} = S_\lambda \; (z_{k+1}) \end{aligned}$$

while LB retains all information which naturally makes LB more fitted for subsampling

$$z_{k+1} = z_k - t_k A_k^T (A_k x_k - b_k) x_{k+1} = S_\lambda (z_{k+1})$$

After Subsampling



.

LB ISTA Hybrid

- K= 16(top) / 50 (bottom)
- K_{max}= 48(top) / 150 (bottom)
- Rows per iteration = 50 (top)/ 16 (bottom)
- Noise = 0.1 * random normal variable

Observations from both setups:

- Consistent with the analytical result, Linearized Bregman converges much quicker than ISTA when subsampling
- Hybrid method has a lower residual than ISTA and LB since it is able to perform well with both large and small entries

Comparing Entries



Submatrix size for subsampling



Flat subsampling mean residual: 0.0194 Tall subsampling mean residual: 0.0255 Square subsampling mean residual: 0.0218

Flat subsampling mean residual: 0.05381 Tall subsampling mean residual: 0.02483 Square subsampling mean residual: 0.0311

Conclusion & References

Through this dynamical study, we observed that Linearized Bregman functions best with large entries and large λ while ISTA operates best with small entries and small λ . This prompted us to create the Hybrid algorithm which benefits from both quick convergence as well as less chatter. Furthermore, we discovered that LB outperforms ISTA when subsampled. To extend on this research, more studies regarding the parameter choices (eg. noise, initial guess x_0 , data setup, specific rows when subsampling) could be conducted.

Emmanoil Daskalakis, Felix J. Herrman, and Rachel Kuske, Accelerating Sparse Recovery by Reducing Chatter (2020)

W. Yin, S. Osher, D. Goldfarb, and J. Darbon, Bregman iterative algorithms for l1- minimization with applications to compressed sensing, SIAM Journal on Imaging Sciences, 1 (2008), pp. 143–168, https://doi.org/10.1137/070703983, http://dx.doi.org/10. 1137/070703983, https://arxiv.org/abs/http://dx.doi.org/10.1137/070703983.

D. A. Lorenz, F. Schöpfer, and S. Wenger, The linearized Bregman method via split feasibility problems: Analysis and generalizations, SIAM Journal on Imaging Sciences, 7 (2014), pp. 1237–1262, https://doi.org/10.1137/130936269, http://dx.doi.org/10.1137/130936269, https://arxiv.org/abs/http://dx.doi.org/10.1137/130936269.

D. A. Lorenz, S. Wenger, F. Schöpfer, and M. Magnor, A sparse Kaczmarz solver and a linearized Bregman method for online compressed sensing, in 2014 IEEE International Conference on Image Processing (ICIP), Oct 2014, pp. 1347–1351, https://doi.org/10.1109/ ICIP.2014.7025269.