

Stochastic Trust Region Inexact Newton Method for Large-scale Machine Learning

Vinod Kumar Chauhan * · Anuj Sharma ·
Kalpana Dahiya

Received: date / Accepted: date

Abstract Nowadays stochastic approximation methods are one of the major research direction to deal with the large-scale machine learning problems. From stochastic first order methods, now the focus is shifting to stochastic second order methods due to their faster convergence and availability of computing resources. In this paper, we have proposed a novel Stochastic Trust RegiOn Inexact Newton method, called as STRON, to solve large-scale learning problems which uses conjugate gradient (CG) to inexactly solve trust region subproblem. The method uses progressive subsampling in the calculation of gradient and Hessian values to take the advantage of both, stochastic and full-batch regimes. We have extended STRON using existing variance reduction techniques to deal with the noisy gradients and using pre-conditioned conjugate gradient (PCG) as subproblem solver, and empirically proved that they do not work as expected, for the large-scale learning problems. Finally, our empirical results prove efficacy of the proposed method against existing methods with bench marked datasets.

Keywords stochastic optimisation · subsampling · second order methods · inexact newton · large-scale learning

1 Introduction

Machine learning involves data intensive optimization problems which have large number of component functions corresponding to large amount of available data. Traditional/classical methods, like Newton method, fail to perform well on such large-scale learning problems (i.e., problems with large number of data points) due to large per-iteration complexity. So nowadays one of the major challenge in machine learning is to develop scalable and efficient algorithms to deal with these large-scale learning problems [40, 11, 12].

To solve the machine learning problems, gradient descent (GD) [10] is the classical method of choice

*Most of this work was done when author was doing his PhD at Panjab University Chandigarh, India.

*This is a pre-print of an article published in International Journal of Machine Learning and Cybernetics. The final authenticated version is available online at: <https://doi.org/10.1007/s13042-019-01055-9>

Vinod Kumar Chauhan
Department of Engineering
University of Cambridge, UK
E-mail: vk359@cam.ac.uk
Homepage: <http://www.eng.cam.ac.uk/profiles/vk359>

Anuj Sharma
Computer Science & Applications
Panjab University Chandigarh, INDIA
E-mail: anujs@pu.ac.in
Homepage: <https://sites.google.com/view/anujsharma>

Kalpana Dahiya
University Institute of Engineering and Technology
Panjab University Chandigarh, INDIA
E-mail: kalpanas@pu.ac.in

but it trains slowly while dealing with large-scale learning problems due to high per-iteration cost. Stochastic approximation based methods [33] can be quite effective in such situations but they converge slowly due to the noisy approximations of gradient. So a variety of stochastic variance reduction techniques came to existence, e.g., [25, 17, 34, 2, 19]. But the major limitation of these methods is that they can converge up to linear rate only.

Newton method is another classical method to solve optimization problems, which can give up to quadratic convergence rate [7]. But again (pure) Newton method is not feasible with large-scale learning problems due to huge per-iteration computational complexity and need to store a huge Hessian matrix. So nowadays one of the most significant open question in optimization for machine learning is: “Can we develop stochastic second order methods with quadratic convergence, like Newton method but has low per-iteration complexity, like stochastic approximation methods?”. After success in the stochastic first order methods, the research is shifting its focus towards the stochastic second order methods to leverage the faster convergence of second order methods and the available computing power.

Inexact Newton (also called truncated Newton or Hessian free) methods and quasi-Newton methods are among the major research directions for developing second order methods [5]. Inexact Newton methods try to solve the Newton equation approximately without calculating and storing the Hessian matrix. On the other hand, quasi-Newton methods try to approximate the Hessian inverse and avoid the need to store the Hessian matrix. Thus both the methods try to resolve the issues with Newton method for solving large-scale learning problems. The stochastic variants of inexact Newton and quasi-Newton, further reduce the complexity of these methods by using subsampled gradient and Hessian calculations. In this paper, we have proposed a novel **stochastic trust region inexact Newton (STRON)** method to solve the large-scale learning problems, which introduces subsampling to gradient and Hessian calculations. It uses progressive subsampling to enjoy the benefits of both the regimes, stochastic approximation and full-batch learning. We further extend the method using existing variance reduction techniques to deal with noise produced by subsampling of gradient values, and by proposing PCG for solving the trust region subproblem.

1.1 Optimization Problem

We consider unconstrained convex optimization problem of expected risk, as given below:

$$\min_w R(w) = \mathbb{E}[f(w; \xi)], \quad (1)$$

where $f(w; \xi) = f(w; x_i, y_i) = f(h(w; x_i), y_i)$ is a smooth composition of linear prediction model h and loss function f over randomly selected data point (x_i, y_i) from the unknown distribution $P(x_i, y_i)$, parameterized by the model parameter $w \in \mathbb{R}^n$. Since it is not feasible to solve (1) as P is unknown so the model is approximated by taking a set $N = \{(x_1, y_1), \dots, (x_l, y_l)\}$ of l data points from the unknown distribution P and then solving the empirical risk minimization problem, as given below:

$$\min_w F(w) = \frac{1}{l} \sum_{i=1}^l f(w; x_i, y_i). \quad (2)$$

For simplicity, we take $f(w; x_i, y_i) = f_i(w)$. Finite sum optimization problems of type (2) exists across different fields, like signal processing, statistics, operation research, data science and machine learning, e.g., logistic regression and SVM in machine learning.

1.2 Solution Techniques

Simple iterative classical method to solve (2) is gradient descent (GD) [10], as given below:

$$w^{k+1} = w^k - \alpha_k \nabla F(w_k), \quad (3)$$

where $(k + 1)$ is iteration number and α_k is called learning rate or step size. The complexity of this iteration is $O(nl)$ which is large for large-scale learning problems due to large number of data points.

Stochastic gradient descent (SGD) [33] is very effective to deal with such problems due to its low per-iteration complexity, as given below:

$$w^{k+1} = w^k - \alpha_k \nabla f_{i_k}(w_k), \quad (4)$$

where i_k is randomly selected data point. But convergence can't be guaranteed in SGD because of noisy gradient values.

Another classical second order method to solve (2) is Newton method, as given below:

$$w^{k+1} = w^k - \alpha_k \nabla^2 F(w_k)^{-1} \nabla F(w_k). \quad (5)$$

The complexity of this iteration is $O(n^2l + n^3)$ and it needs to store and invert the Hessian matrix $\nabla^2 F(w_k)$, which is computationally very expensive and needs large memory for large-scale learning problems, respectively. That's why first order methods and their stochastic variants have been studied very extensively, during the last decade, to solve the large-scale learning problems but not second order methods. As stochastic first order methods have hit their limits and due to good availability of computing power, the main focus is shifting towards stochastic second order methods, and nowadays one important open question is to find stochastic second order methods with quadratic convergence rates.

There are two major research directions for second order methods: quasi-Newton methods and inexact Newton methods, both of which try to resolve the issues associated with the Newton method. Quasi-Newton methods try to approximate the Hessian matrix during each iteration, as given below:

$$w^{k+1} = w^k - \alpha_k B_k \nabla F(w_k), \quad (6)$$

where B_k is an approximate of $\nabla^2 F(w_k)^{-1}$, e.g., Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm is one such method [20]. On the other hand, inexact Newton methods try to solve the Newton system approximately, e.g., Newton-CG [38]. Both the methods try to resolve the issues related with Newton method but still their complexities are large for large-scale problems. So a number of stochastic variants of these methods have been proposed, e.g., [8, 9, 6, 3] which introduce subsampling to gradient and Hessian calculations.

1.3 Contributions

The contributions of the manuscript are listed below:

- The manuscript highlights the recent shift of stochastic methods from first order to second order methods and raises an open question that we need stochastic second order methods with quadratic convergence rate to solve the large-scale machine learning problems.
- We have proposed a novel subsampled variant of trust region inexact Newton method for solving large-scale learning problems, which is called STRON and is the first stochastic variant of trust region inexact Newton methods. STRON uses progressive subsampling scheme for gradient and Hessian calculations to enjoy the benefits of both stochastic and full batch regimes. STRON can converge up to quadratic rate and answers the raised open question.
- STRON has been extended using existing variance reduction techniques to deal with the noisy approximations of the gradient calculations. The extended method uses stochastic variance reduced gradient (SVRG) for variance reduction with static batching for gradient calculations and progressive batching for the Hessian calculations. The empirical results prove that this does not work as expected for large-scale learning.
- We further extend STRON and use PCG, instead of CG method, to solve the Newton system inexactly. We have used weighted average of identity matrix and diagonal matrix as the preconditioner. But even this fails to work as expected for large-scale learning.
- Finally, our empirical experiments prove the efficacy of STRON against existing techniques with bench marked datasets.

2 Literature Review

Stochastic approximation methods are very effective to deal with the large-scale learning problems due to their low per-iteration cost, e.g., SGD [33], but they lead to slow convergence rates due to the noisy approximation. To deal with the noise issue, a number of techniques have been proposed and some of the important techniques (as discussed in [16]) are: (a) decreasing learning rates [37], (b) using mini-batching [13], (c) importance sampling [16], and (d) variance reduction [23]. The variance reduction methods can further be classified into three categories: primal methods [34, 15], dual methods [36] and primal-dual methods [39]. The variance reduction techniques are effective to deal with the large-scale learning problems because of low per-iteration complexity, like SGD, and have fast linear convergence, like GD. These techniques exploit the best of SGD and GD but for these stochastic variants of first order methods the convergence is limited to linear rate only, unlike the second order methods which can give up to quadratic rate.

Second order methods utilize the curvature information to guide the step direction towards the solution and exhibit faster convergence than the first order methods. But huge per-iteration cost due to the huge Hessian matrix and its inversion make the training of models slow for large-scale problems. So certain techniques have been developed to deal with the issues related to Hessian matrix, e.g., quasi-Newton methods and inexact Newton methods are two major directions to deal with the huge computational cost of Newton method. Quasi-Newton methods approximate the Hessian matrix during each iteration, e.g., BFGS [20] and its limited memory variant, called L-BFGS [28], are examples of the quasi-Newton class which use gradient and parameter values from the previous iterations to approximate the Hessian inverse. L-BFGS uses only recent information from previous M -iterations. On the other hand, inexact Newton methods try to solve the Newton system approximately, e.g., Newton-CG [38].

Recently, several stochastic variants of BFGS and L-BFGS have been proposed to deal with large-scale problems. Schraudolph et al. [35] proposed stochastic variants of BFGS and L-BFGS for the online setting, known as oBFGS. Mokhtari and Ribeiro [30] extended oBFGS by adding regularization which enforces upper bound on the eigen values of the approximate Hessian, known as RES (Regularized Stochastic BFGS). Stochastic quasi-Newton (SQN) [9] is another stochastic variant of L-BFGS which collects curvature information at regular intervals, instead of at each iteration. Variance-reduced Stochastic Newton (VITE) [29] extended RES and proposed to use variance reduction for the subsampled gradient values for solving smoothly strongly convex problems. Kolte et al. [24] provided another stochastic L-BFGS method with variance reduction using SVRG (referred as SVRG-LBFGS). Moritz et al. [31] proposed Stochastic L-BFGS (SLBFGS) using SVRG for variance reduction and using Hessian-vector product to approximate the gradient differences for calculating the Hessian approximations, also referred as SVRG-SQN. SVRG-LBFGS and SVRG-SQN differ in how the curvature is approximated (i.e. how Hessian is approximated), former collects the curvature information once during each epoch (outer iterations) using gradient differences but later collects the curvature information after regular intervals inside the inner-iterations using Hessian-vector products. Berahas et al. [4] proposed multi-batch scheme into stochastic L-BFGS where batch sample changes with some overlaps with previous iteration. Bollapragada et al. [6] proposed progressive batching, stochastic line search and stable Newton updates for L-BFGS. Bollapragada et al. [5] studies the conditions on the subsample sizes to get the different convergence rates.

Stochastic inexact Newton methods are also explored extensively. Byrd et al. [8] proposed stochastic variants of Newton-CG along with L-BFGS method. Bollapragada et al. [5] studies subsampled Newton methods and find conditions on subsample sizes and forcing term (constant used with the residual condition), for linear convergence of Newton-CG method. Bellavia et al. [3] studies the effect of forcing term and line search to find linear and super-linear convergence of Newton-CG method. Newton-SGI (stochastic gradient iteration) is another way of solving the linear system approximately and is studied in Agarwal et al. [1].

Trust Region Newton (TRON) method is one of the most efficient solver for solving large-scale linear classification problems [27]. This is trust region inexact Newton method which does not use any subsampling and is present in LIBLINEAR library [18]. Hsia et al. [21] extends TRON by improving the trust region radius value. Hsia et al. [22], further extends TRON and uses preconditioned conjugate gradient (PCG) which uses weighted average of identity matrix and diagonal matrix as a preconditioner, to solve the trust region subproblem. Since subsampling is an effective way to deal with the large-scale problems so in this paper, we have proposed a stochastic variant of trust region inexact Newton method, which have not been studied so far to the best of our knowledge.

3 Trust Region Inexact Newton Method

Inexact Newton methods, also called as Truncated Newton or Hessian free methods, solve the Newton equation (linear system) approximately. CG method is a commonly used technique to solve the trust region subproblem approximately. In this section, we discuss inexact Newton method and its trust region variation.

3.1 Inexact Newton Method

The quadratic model $m_k(p)$ obtained using Taylor's theorem is given below:

$$F(w_k + p) - F(w_k) \approx m_k(p) \equiv \nabla F(w_k)^T p + \frac{1}{2} p^T \nabla^2 F(w_k) p. \quad (7)$$

Taking derivative of $m_k(p)$ w.r.t. p and equating to zero, we get,

$$\nabla^2 F(w_k) p = -\nabla F(w_k), \quad (8)$$

which is Newton system and its solution gives Newton method, as given below:

$$w^{k+1} = w^k + p_k = w^k - \nabla^2 F(w_k)^{-1} \nabla F(w_k). \quad (9)$$

The computational complexity of this iteration is $O(n^2l + n^3)$ which is very expensive. This iteration involves the calculation and inversion of a large Hessian matrix which is not only very expensive to calculate but expensive to store also. CG method approximately solves the subproblem (8) without forming the Hessian matrix, which solves the issues related to large computational complexity and need to store the large Hessian matrix. Each iteration runs for a given number of CG iterations or until the residual condition is satisfied, as given below:

$$\|r_k\| \leq \eta'_k \|\nabla F(w_k)\|, \quad (10)$$

where $r_k = \nabla^2 F(w_k) p + \nabla F(w_k)$ and η'_k is a small positive value, known as forcing term [32].

3.2 Trust Region Inexact Newton Method

Trust region is a region in which the approximate quadratic model of the given function gives correct approximation for that function. In trust region methods, we don't need to calculate the step size (also called learning rate) directly but they indirectly adjust the step size as per the trust region radius. Trust region method solves the following subproblem to get the step direction p_k :

$$\min_p m_k(p) \quad \text{s.t.} \quad \|p\| \leq \Delta_k, \quad (11)$$

where $m_k(p)$ is a quadratic model of $F(w_k + p) - F(w_k)$, as given in (7) and Δ_k is the trust region radius. This subproblem can be solved similar to Newton-CG, except that now we need to take care of the extra constraint of p . TRON (trust region Newton method) [27] is one of the most famous and widely used such method, which is used in LIBLINEAR [18] to solve l_2 -regularized logistic regression and l_2 -SVM problems. Hsia et al. [21] extends TRON by proposing better trust region radius. Hsia et al. [22] further extends TRON using PCG subproblem solver which uses average of identity matrix and diagonal matrix as preconditioner, to solve the trust region subproblem and shows that PCG could be effective to solve ill-conditioned problems.

Then the ratio of actual and predicted reductions of the model is calculated, as given below:

$$\rho_k = \frac{F(w_k + p_k) - F(w_k)}{m_k(p_k)}. \quad (12)$$

The parameters are updated for the $(k + 1)$ th iteration as given below:

$$w_{k+1} = \begin{cases} w_k + p_k, & \text{if } \rho_k > \eta_0, \\ w_k, & \text{if } \rho_k \leq \eta_0, \end{cases} \quad (13)$$

where $\eta_0 > 0$ is a given constant. Then the trust region radius Δ_k is updated as per the ratio of actual reduction and predicted reduction, and a framework for updating Δ_k as given in [26], is given below:

$$\Delta_{k+1} \in \begin{cases} [\gamma_1 \min\{\|p_k\|, \Delta_k\}, \gamma_2 \Delta_k], & \text{if } \rho_k \leq \eta_1, \\ [\gamma_1 \Delta_k, \gamma_3 \Delta_k], & \text{if } \rho_k \in (\eta_1, \eta_2), \\ [\Delta_k, \gamma_3 \Delta_k], & \text{if } \rho_k \geq \eta_2, \end{cases} \quad (14)$$

where $0 < \eta_1 < \eta_2 \leq 1$ and $0 < \gamma_1 < \gamma_2 < 1 < \gamma_3$. If $\rho_k \leq \eta_1$ then the Newton step is considered unsuccessful and the trust region radius is shrunk. On the other hand if $\rho_k \geq \eta_2$ then the step is successful and the trust region radius is enlarged. We have implemented this framework as given in the LIBLINEAR library [18] and chose the following pre-defined values for the above constants: $\eta_0 = 1e - 4$, $\eta_1 = 0.25$, $\eta_2 = 0.75$, $\gamma_1 = 0.25$, $\gamma_2 = 0.5$ and $\gamma_3 = 4$.

4 STRON

STRON exploits the best of both, stochastic and batch regimes, using progressive subsampling to solve the large-scale learning problems. As stochastic gradient descent (SGD) trains faster for large-scale learning problems than gradient descent (GD) due to low computations per iteration but GD is more accurate than SGD due to batch calculations, similarly STRON takes benefit of low computation during initial iterations and as it reaches the solution region it uses batch calculations to find accurate solution like TRON. The major challenge with STRON is to decide when to switch from stochastic to full-batch regime, i.e., to tune the subsampling rate.

STRON introduces stochasticity into the trust region inexact Newton method and calculates subsampled function, gradient and Hessian values to solve the trust region subproblem, as given below:

$$\min_p m_k(p) = \nabla F_{X_k}(w_k)^T p + \frac{1}{2} p^T \nabla^2 F_{S_k}(w_k) p, \quad \text{s.t. } \|p\| \leq \Delta_k, \quad (15)$$

where $\nabla^2 F_{S_k}(w_k)$ and $\nabla F_{X_k}(w_k)$ are subsampled Hessian and gradient values over the subsamples S_k and X_k , respectively, as defined below:

$$\begin{aligned} \nabla^2 F_{S_k}(w_k) &= \frac{1}{|S_k|} \sum_{i \in S_k} \nabla^2 f_i(w_k), \\ \nabla F_{X_k}(w_k) &= \frac{1}{|X_k|} \sum_{i \in X_k} \nabla f_i(w_k), \\ F_{X_k}(w_k) &= \frac{1}{|X_k|} \sum_{i \in X_k} f_i(w_k), \end{aligned} \quad (16)$$

where subsamples are increasing, i.e., $|X_k| < |X_{k+1}|$, $|S_k| < |S_{k+1}|$ and F_{X_k} is subsampled function value used for calculating ρ_k . STRON solves (15) approximately for given number of CG iterations or until the following residual condition is satisfied:

$$\|r_k\| \leq \eta'_k \|\nabla F_{X_k}(w_k)\|, \quad (17)$$

where $r_k = \nabla^2 F_{S_k}(w_k) p + \nabla F_{X_k}(w_k)$.

STRON is presented by Algorithm 1. It randomly selects subsamples S_k and X_k for the k th iteration (outer iterations). X_k and S_k are used for calculating the gradient and Hessian values, respectively. Then it solves the trust region subproblem using CG solver (inner iterations) which uses subsampled Hessian in calculating Hessian-vector products. CG stops when residual condition, same as (17), satisfies, it reaches maximum #CG iterations or it reaches the trust region boundary. The ratio of reduction in actual and predicted reduction is calculated similar to (12) but using subsampled function, and is used for updating the parameters as given in (13). Then trust region radius Δ_k is updated as per ρ_k as given in (14) and these steps are repeated for a given number of iterations or until convergence.

STRON uses progressive subsampling, i.e., dynamic subsampling to calculate function, gradient and Hessian values, and solves the Newton system approximately. It is effective to deal with large-scale problems since it uses subsampling and solves the subproblem approximately, without forming the Hessian matrix but using only Hessian-vector products. So it handles the complexity issues related with the Newton method.

Algorithm 1 STRON

```
1: Input:  $w_0$ 
2: Result:  $w = w_k$ 
3: for  $k = 0, 1, \dots$  do
4:   Randomly select subsamples  $S_k$  and  $X_k$ 
5:   Calculate subsampled gradient  $\nabla F_{X_k}(w_k)$ 
6:   Solve the trust region subproblem using Algorithm 2, to get the step direction  $p_k$ 
7:   Calculate the ratio  $\rho_k = (F_{X_k}(w_k + p_k) - F_{X_k}(w_k)) / m_k(p_k)$ 
8:   Update the parameters using (13)
9:   Update the trust region radius  $\Delta_k$  using (14)
10: end for
```

Algorithm 2 CG Subproblem Solver

```
1: Inputs:  $\Delta_k > 0, \eta'_k \in (0, 1)$ 
2: Result:  $p_k = p_j$ 
3: Initialize  $p_0 = 0, r_0 = d_0 = -\nabla F_{X_k}(w_k)$ 
4: for  $j = 1, 2, \dots$  do
5:   if  $\|r_{j-1}\| < \eta'_k \|\nabla F_{X_k}(w_k)\|$  then
6:     return  $p_k = p_{j-1}$ 
7:   end if
8:   Calculate subsampled Hessian-vector product  $v_j = \nabla^2 F_{S_k}(w_k) d_{j-1}$ 
9:    $\alpha_j = \|r_{j-1}\|^2 / (d_{j-1}^T v_j)$ 
10:   $p_j = p_{j-1} + \alpha_j d_{j-1}$ 
11:  if  $\|p_j\| \geq \Delta_k$  then
12:    Calculate  $\tau_j$  such that  $\|p_{j-1} + \tau_j d_{j-1}\| = \Delta_k$ 
13:    return  $p_k = p_{j-1} + \tau_j d_{j-1}$ 
14:  end if
15:   $r_j = r_{j-1} - \alpha_j v_j,$ 
16:   $\beta_j = \|r_j\|^2 / \|r_{j-1}\|^2, d_j = r_j + \beta_j d_{j-1}$ 
17: end for
```

4.1 Complexity

The complexity of trust region inexact Newton (TRON) method depends on function, gradient and CG subproblem solver. This is dominated by CG subproblem solver and is given by $O(nl) \times \#CG$ iterations, and for sparse data, $O(\#nnz) \times \#CG$ iterations, where $\#nnz$ is number of non-zeros values in the dataset. For subsampled trust region inexact Newton (STRON) method, the complexity per-iteration is given by $O(n|S_k|) \times \#CG$ iterations, and for sparse data, $O(\#nnz_{S_k}) \times \#CG$ iterations, where $\#nnz_{S_k}$ is number of non-zeros values in the subsample S_k . Since $\#CG$ iterations taken by TRON and STRON do not differ much so the per-iteration complexity of STRON is smaller than TRON in the initial iterations and later becomes equal to TRON due to progressive subsampling, i.e., when $|S_k| = N$.

4.2 Analysis

STRON method uses progressive subsampling with the assumption that eventually mini-batch size becomes equal to the whole dataset, i.e., for some value of $k \geq \bar{k} > 0$, STRON becomes TRON. So we can follow the theoretical results given in TRON, which itself refers the results from [26]. For $k \geq \bar{k}$, we get different convergence depending upon the value of η'_k : If $\eta'_k < 1$ then STRON converges Q-linearly, if $\eta'_k \rightarrow 0$ as $k \rightarrow \infty$ then STRON has Q-superlinear convergence and when $\eta'_k \leq \kappa_0 \|\nabla F(w_k)\|$ for $\kappa_0 > 0$ then STRON converges at quadratic rate.

5 Experimental Results

In this section, we discuss experimental settings and results. The experiments have been conducted with the bench marked binary datasets as given in the Table 1, which are available for download from LibSVM website¹. We use following methods in experimentation:

¹ <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

Table 1: Datasets used in experimentation

Dataset	#features	#datapoints
gisette	5000	6,000
rcv1.binary	47,236	20,242
webspam (unigram)	254	350,000
covtype.binary	54	581,012
ijcnn1	22	49,990
news20.binary	1,355,191	19,996
real-sim	20,958	72,309
Adult	123	32,561
mushroom	112	8124

TRON [22]: This is a trust region inexact Newton method without any subsampling. It uses preconditioned CG method to solve the trust region subproblem and it is present in the current version of LIBLINEAR library [18].

STRON: This is the proposed stochastic trust region inexact Newton method with progressive subsampling technique for gradient and Hessian calculations. It uses CG method to solve the trust region subproblem.

STRON-PCG: This is an extension of STRON using PCG for solving the trust region subproblem, as discussed in the Subsection A.1.

STRON-SVRG: This is another extension of STRON using variance reduction for subsampled gradient calculations, as discussed in Subsection A.2.

Newton-CG [8]: This is stochastic inexact Newton method which uses CG method to solve the subproblem. It uses progressive subsampling similar to STRON.

SVRG-SQN [31]: This is stochastic L-BFGS method with variance reduction for gradient calculations.

SVRG-LBFGS [24]: This is another stochastic L-BFGS method with variance reduction. It differs from SVRG-SQN method in approach by which Hessian information is sampled.

5.1 Experimental Setup

The datasets have been divided into 80% and 20%, for training and testing datasets, respectively to plot the convergence, and 5-fold cross-validation has been used to present results in Table 2. We have used $\lambda = 1/l$ for all methods because generally it gives good results and at the same time help to reduce number of tunable hyper-parameters. #CG iterations has been set to a sufficiently large value of 25, as all the inexact Newton methods use 5-10 iterations and hardly go beyond 20 iterations. Progressive batching scheme uses initial batch size of 1% for all datasets except ijcn which uses 10% of dataset and the subsample size is increased linearly for all datasets, except covtype where exponential rate is used. Moreover we set the rate such that subsample size equals dataset size in 5 epochs because, generally second order methods reach the solution region in 4-5 epochs and converge in 8-10 epochs. For the sake of simplicity and avoid extra sampling, we take same subsample for Hessian and gradient calculations, i.e., $S_k = X_k$. Quasi-Newton methods (SVRG-SQN and SVRG-LBFGS) use mini-batch size of 10% with stochastic backtracking line search to find the step size and same size mini-batches are taken for gradient and Hessian subsampling. Memory of $M = 5$ is used in quasi-Newton methods with $L = 5$ as update frequency of Hessian inverse approximation for SVRG-SQN method. All the algorithms use similar exit criterion of decrease in gradient value where we calculate gradient ($\|g_0\|$) at w_0 and run the algorithms until $\|g_k\| \leq \epsilon \|g_0\|$, where $\|g_k\|$ is gradient value at k^{th} -iteration and ϵ is the given tolerance level. Moreover all the methods are implemented in C++² with MATLAB interface and experiments have been executed on MacBook Air (8 GB 1600 MHz DDR3, 1.6 GHz Intel Core i5 and 256GB SSD).

5.2 Comparative Study

The experiments have been performed with strongly convex and smooth l_2 -regularized logistic regression problem as given below:

² experimental results can be reproduced using the library LIBS2ML [14].

$$\min_w F(w) = \frac{1}{l} \sum_{i=1}^l \log(1 + \exp(-y_i w^T x_i)) + \frac{\lambda}{2} \|w\|^2. \quad (18)$$

The results have been plotted as optimality ($F(w) - F(w^*)$) versus training time (in seconds) and accuracy versus training time for high ($\epsilon=10^{-10}$)-accuracy solutions, as given in the Figs. 1, 2 and 3. As it is clear from the results, STRON converges faster than all other methods and shows improvement against TRON on accuracy vs. time plots. Moreover, quasi-Newton methods converges slower than inexact Newton methods as already established in the literature [27]. As per the intuitions, STRON takes initial advantage over TRON due to subsampling and as it reaches the solution region the progressive batching scheme reaches the full batching scheme and converges with same rate as TRON. That's why, in most of the figures, we can observe STRON and TRON converging in parallel lines. Moreover, we observe a horizontal line for accuracy vs. time plot with covtype dataset because all methods give 100% accuracy. Generally, the models are trained for low ($\epsilon=10^{-02}$)-accuracy solutions. So we present the results for

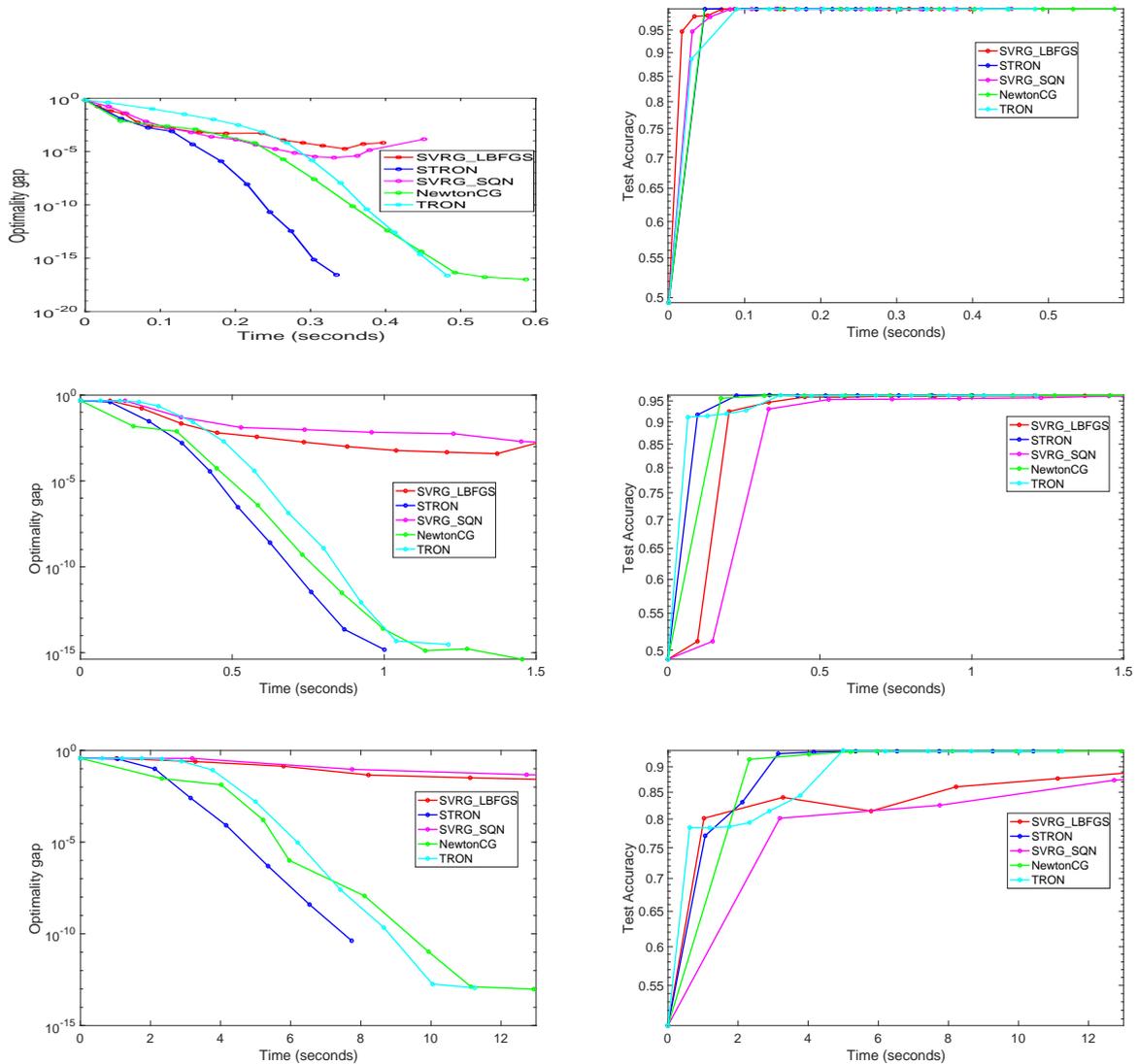


Fig. 1: First column presents optimality versus training time (in seconds) and second column presents accuracy versus training time, on mushroom (first row), rcv1 (second row) and news20 (third row) datasets.

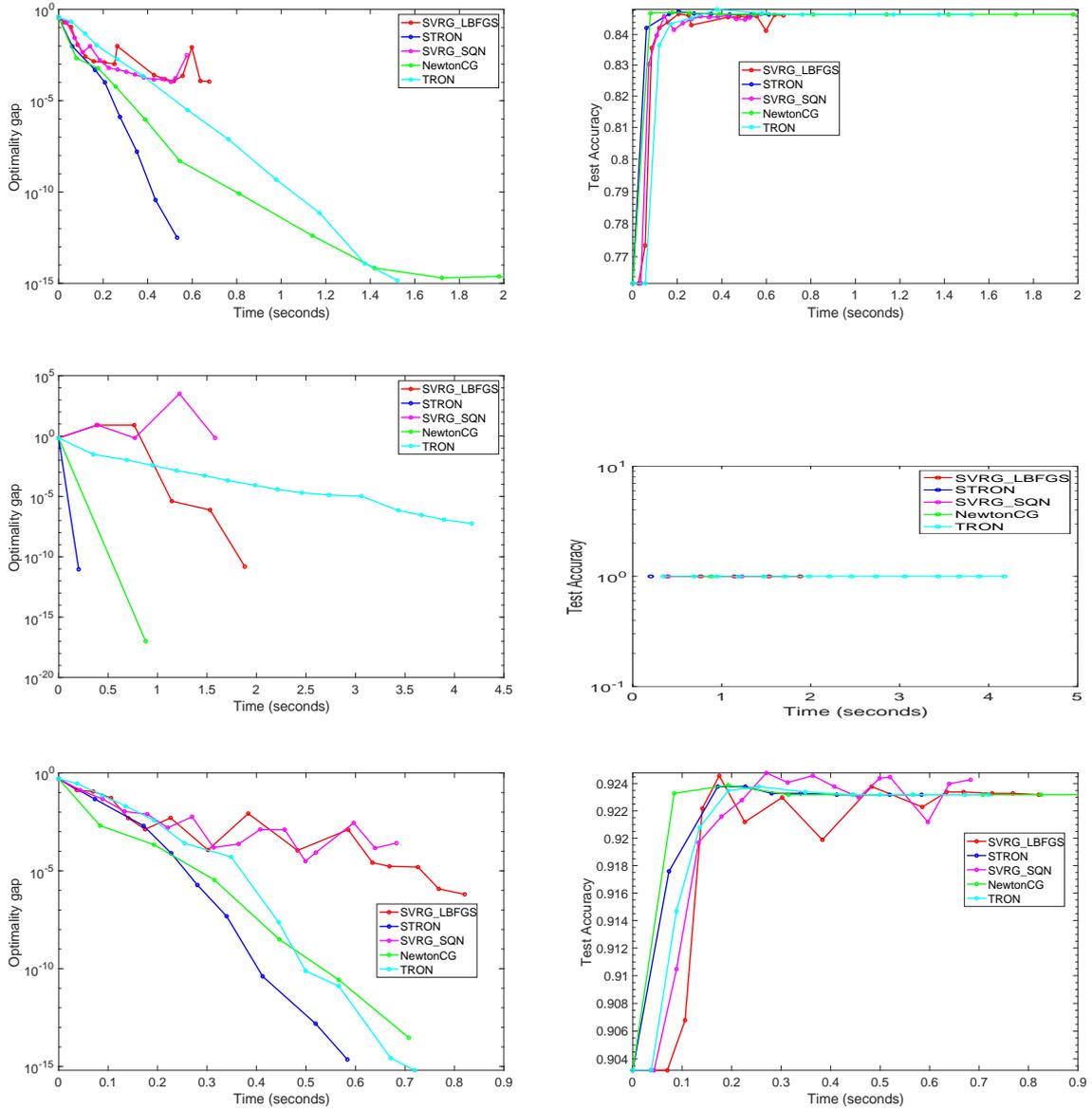


Fig. 2: First column presents optimality versus training time (in seconds) and second column presents accuracy versus training time, on Adult (first row), covtype (second row) and ijcnn1 (third row).

such a case using Table 2, which reports results using 5-fold cross validation. As it is clear from the table, STRON either outperforms other solvers or shows results pretty close to the best method.

5.3 Results with SVM

We extend STRON to solve l_2 -SVM problem which is a non smooth problem, as given below:

$$\min_w F(w) = \frac{1}{l} \sum_{i=1}^l \max(0, 1 - y_i w^T x_i)^2 + \frac{\lambda}{2} \|w\|^2. \quad (19)$$

The results are reported in the Fig. 4 with news20 and real-sim datasets. As it is clear from the figure, STRON shows results similar to logistic regression problem and outperforms all other methods.

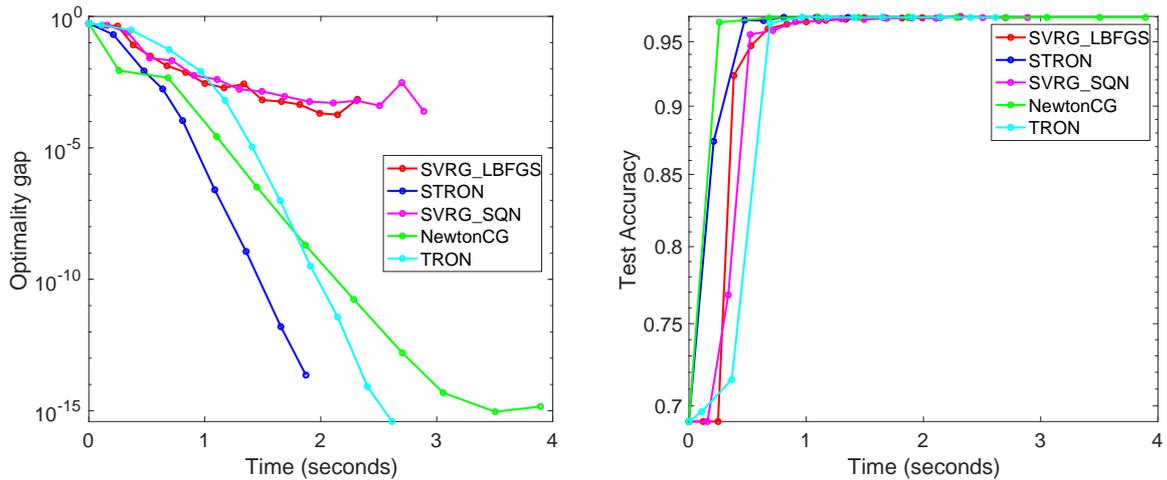


Fig. 3: First column presents optimality versus training time (in seconds) and second column presents accuracy versus training time, on real-sim dataset.

Table 2: Comparison of Training Time (seconds) for low accuracy ($\epsilon=0.01$) solution

Datasets↓	Methods→	SVRG_LBFGS	STRON	SVRG_SQN	Newton_CG	TRON
covtype	Time	0.8440±0.1770	0.0620±0.0045	0.4140±0.0397	0.9740±0.1907	0.7294±0.0162
	Accuracy	1.00±0.0	1.00±0.0	1.00±0.0	1.00±0.0	1.00±0.0
real-sim	Time	1.5920±0.5867	0.9340±0.1258	2.3480±0.5909	0.68±0.01	0.88±0.10
	Accuracy	0.9670±0.0055	0.9697±0.0017	0.9670 ±0.0029	0.9688±0.0022	0.9698±0.0015
rcv1	Time	2.1780±0.1813	0.5700±0.0200	3.2840±0.0270	0.5860±0.0167	0.6589±0.0866
	Accuracy	0.9641±0.0023	0.9640±0.0025	0.9637±0.0024	0.9640±0.0024	0.9639±0.0027
news20	Time	42.8960±2.5993	6.9620±1.3534	84.0520±2.2915	6.764±0.0074	7.2470±0.3567
	Accuracy	0.9333±0.0079	0.9337±0.0076	0.9339±0.0072	0.9338±0.0074	0.9338±0.0077
mushroom	Time	0.1540±0.01340	0.0820±0.0268	0.1820 ±0.01920	0.072±0.0045	0.1316±0.0169
	Accuracy	0.9995±0.0007	0.9992±0.0005	0.9986±0.0026	0.9996±0.0005	0.9995±0.0005
ijcnn1	Time	0.3460±0.1390	0.2560±0.03780	0.3600±0.1158	0.3120±0.0795	0.2564±0.0185
	Accuracy	0.9235±0.0022	0.9232±0.0021	0.9237±0.0018	0.9238±0.0025	0.9233±0.0185
Adult	Time	0.3620±0.0981	0.2000±0.0187	0.3100±0.0644	0.2380±0.0084	0.2139±0.0046
	Accuracy	0.8479±0.0028	0.8470±0.0028	0.8469±0.0031	0.8473±0.0030	0.8477±0.0028
gisette	Time	9.5620±0.9891	8.1800±0.4260	10.7020±1.1314	41.6480±36.1190	10.4277±0.1709
	Accuracy	0.9725±0.0036	0.9730±0.0045	0.9690±0.0026	0.9701±0.0057	0.9755±0.0028
webspam	Time	6.0480±2.2814	6.3720±0.9942	7.6080±2.1462	3.272±0.2141	9.0548±1.2572
	Accuracy	0.9139±0.0033	0.9214±0.0019	0.9158±0.0036	0.9219±0.0017	0.9227±0.0009

6 Conclusion

We proposed a novel stochastic trust region inexact Newton method, called as STRON, to solve the large-scale learning problems. The proposed method used progressive batching scheme to deal with noisy approximations of gradient and Hessian, and enjoyed the benefits of both stochastic and full batch regimes. STRON has been extended to use preconditioned CG as trust region subproblem solver and to use variance reduction for noisy gradient calculations. Our empirical results proved the efficacy of STRON against the state-of-art techniques with bench marked datasets.

Acknowledgements First author is thankful to Ministry of Human Resource Development, Government of INDIA, to provide fellowship (University Grants Commission - Senior Research Fellowship) to pursue his PhD. We are also thankful to the anonymous reviewers for their constructive comments to improve the quality of our manuscript.

Conflict of interest

The authors declare that they have no conflict of interest.

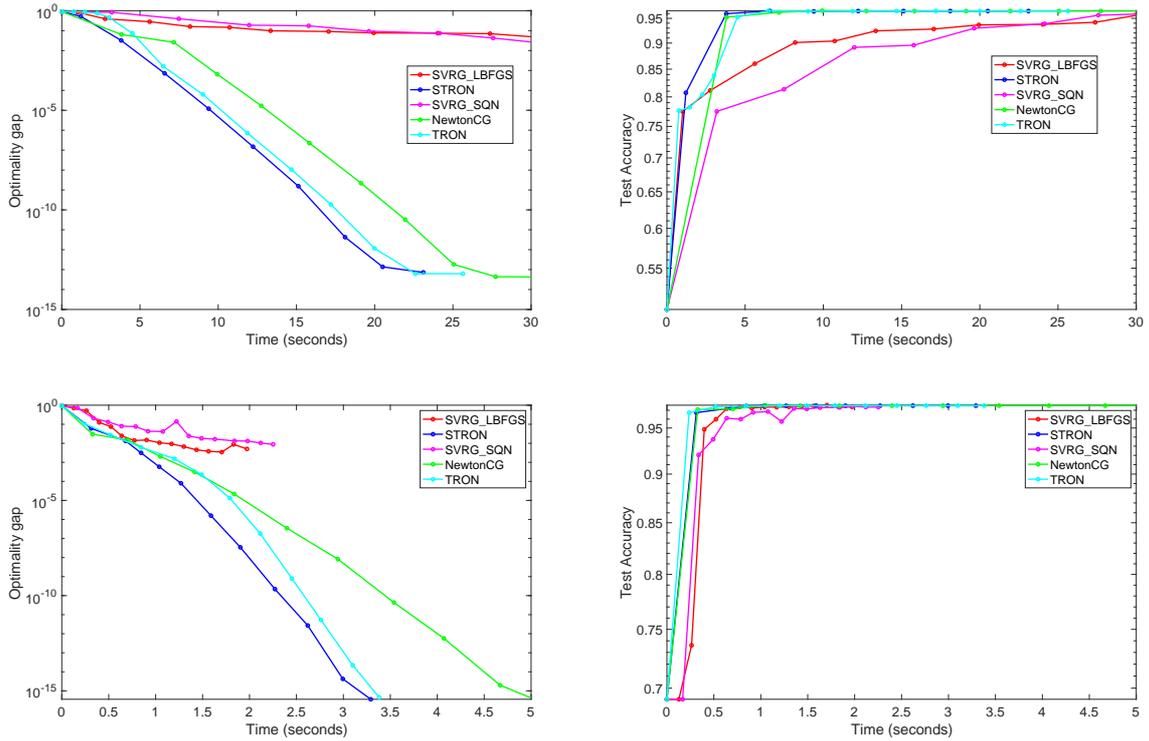


Fig. 4: Experiments with l_2 -SVM on news20 (first row) and real-sim (second row) datasets.

A Extensions

In this section, we discuss extensions of the proposed method with PCG for solving the trust region subproblem, and with variance reduction technique.

A.1 PCG Subproblem Solver

Number of iterations required by CG method to solve the subproblem depend on the condition number of the Hessian matrix. So for ill-conditioned problems CG method converges slowly. To avoid such situations, generally a non-singular matrix M , called preconditioner, is used as follow. For the linear system $\nabla^2 F(w)p = -\nabla F(w)$, we solve following system:

$$M^{-1}\nabla^2 F(w)p = -M^{-1}\nabla F(w). \quad (20)$$

Generally, $M^{-1} = LL^T$ is taken to ensure the symmetry and positive definiteness of $M^{-1}\nabla^2 F(w)$. PCG can be useful for solving the ill-conditioned problems but it involves extra computational overhead. We follow [22] to use PCG as a weighted average of identity matrix and diagonal matrix of Hessian, as given below:

$$M = \alpha \times \text{diag}(H) + (1 - \alpha) \times I, \quad (21)$$

where H is a Hessian matrix and $0 \leq \alpha \leq 1$. For $\alpha = 0$, there is no preconditioning and for $\alpha = 1$ it is a diagonal preconditioner. In the experiments, we have taken $\alpha = 0.01$ for TRON and STRON-PCG [22]. To apply PCG to trust region subproblem, we can use Algorithm 2 without any modifications, after changing the trust region subproblem (11), as given below [38]:

$$\min_{\hat{p}} (L^{-1}\nabla F_{X_k}(w_k))^T \hat{p} + \frac{1}{2} \hat{p}^T (L^{-1}\nabla^2 F_{S_k}(w_k) L^{-T}) \hat{p}, \quad \text{s.t. } \|\hat{p}\| \leq \Delta_k, \quad (22)$$

where $\hat{p} = L^T p$. STRON using PCG as a trust region subproblem solver is denoted by STRON-PCG and the results are reported in Fig. 5. It compares TRON, STRON and STRON-PCG on news20 and rcv1 datasets. As it is clear from the figure, both STRON and STRON-PCG outperform TRON.

PCG trust region subproblem solver involves extra cost for calculating the preconditioner, and for TRON the overhead due to preconditioner is given by

$$O(n) \times \#\text{CG iterations} + O(nl). \quad (23)$$

And for STRON-PCG, preconditioner involves extra cost as given below:

$$O(n) \times \#\text{CG iterations} + O(n|S_k|). \quad (24)$$

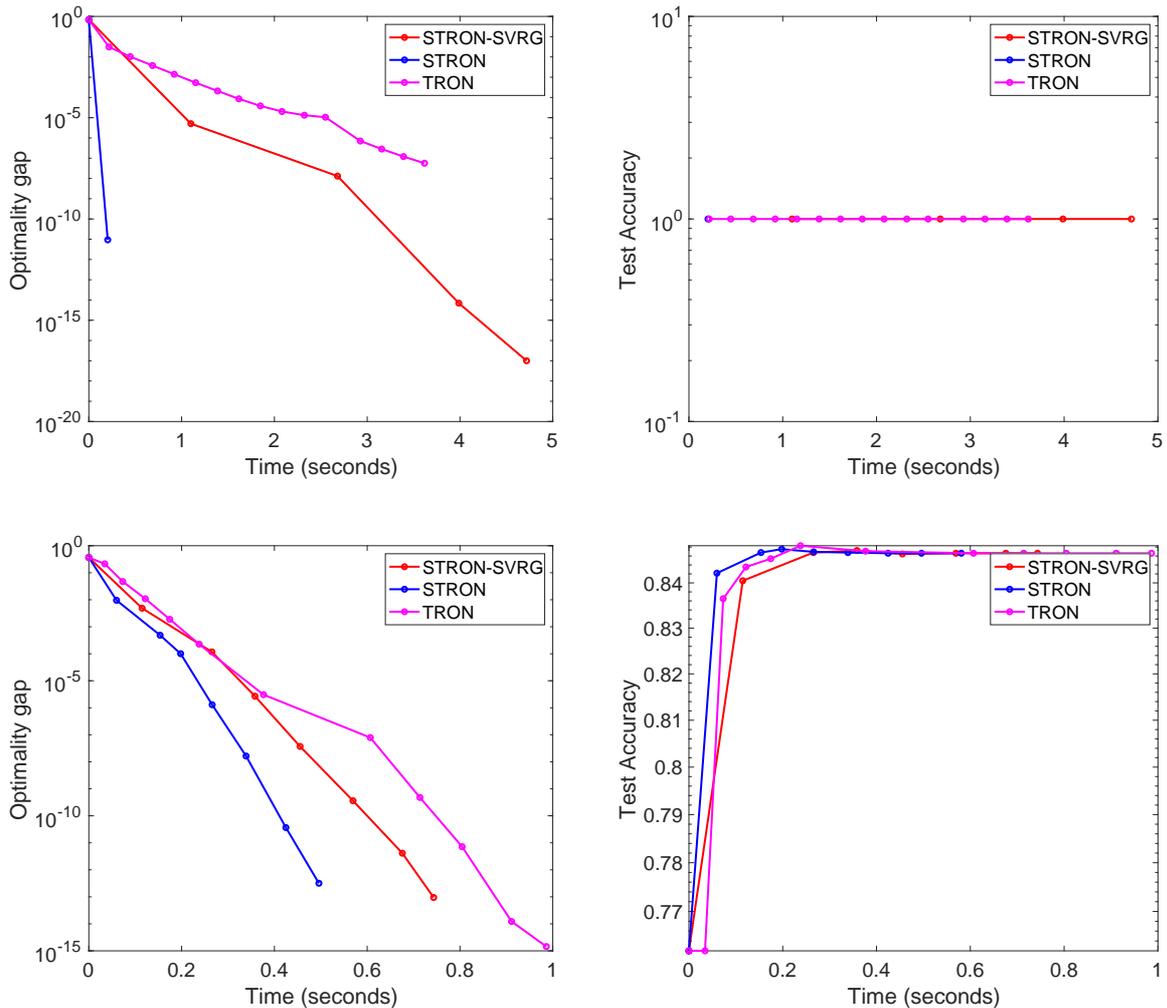


Fig. 5: Comparative study of STRON_PCG, STRON and TRON on covtype (first row) and Adult (second row) datasets.

A.2 Stochastic Variance Reduced Trust Region Inexact Newton Method

Recently, researchers have proposed stochastic variants of second order methods with variance reduction. So it is an interesting question to know that how will variance reduction work with stochastic trust region inexact Newton methods, as this is not studied yet. Our empirical results prove that variance reduction does not work in this case, even after using progressive subsampling for Hessian calculation.

To improve the quality of search direction, we have used SVRG as variance reduction technique for gradient calculations, as given below:

$$g_k = \nabla F_{S_k}(w_k) - \nabla F_{S_k}(\bar{w}) + \nabla F(\bar{w}), \quad (25)$$

where \bar{w} is parameter value at the start of outer iteration. STRON-SVRG uses variance reduction for gradient calculations and progressive batching for Hessian calculations, as given in the Algorithm 3. The experimental results are presented in Fig. 6 with news20 and rcv1 datasets. As it is clear from the figures, STRON-SVRG lags behind STRON and TRON, i.e., variance reduction in STRON-SVRG is not sufficient to beat the progressive batching in gradient calculations of STRON. This is because both, STRON-SVRG and STRON, are stochastic/subsampled variants of TRON and to compensate for the noisy gradient calculations, former uses well-known variance reduction strategy but later uses progressive subsampling strategy. STRON is able to beat TRON but STRON-SVRG fails and lags behind both.

B More Results

In this appendix, we provide more results and study the effect of regularization coefficient on the methods.

Algorithm 3 STRON with Variance Reduction

```

1: Inputs:  $w_0, m$ 
2: Result:  $w = w_k$ 
3: for  $i = 0, 1, 2, \dots$  do
4:   Calculate  $\nabla F(w_i)$  and set  $\bar{w} = w_i$ 
5:   for  $k = 0, 1, \dots, (m - 1)$  do
6:     Randomly select subsamples  $S_k$  and  $X_k$ 
7:     Calculate subsampled gradient  $\nabla F_{X_k}(w_k)$ 
8:     Calculate variance reduced gradient using (25)
9:     Solve the trust region subproblem using Algorithm 2 with variance reduced gradient, instead of subsampled
        gradient, to get the step direction  $p_k$ 
10:    Calculate the ratio  $\rho_k = (F_{X_k}(w_k + p_k) - F_{X_k}(w_k)) / m_k(p_k)$ 
11:    Update the parameters using (13)
12:    Update the trust region  $\Delta_k$  using (14)
13:  end for
14: end for
  
```

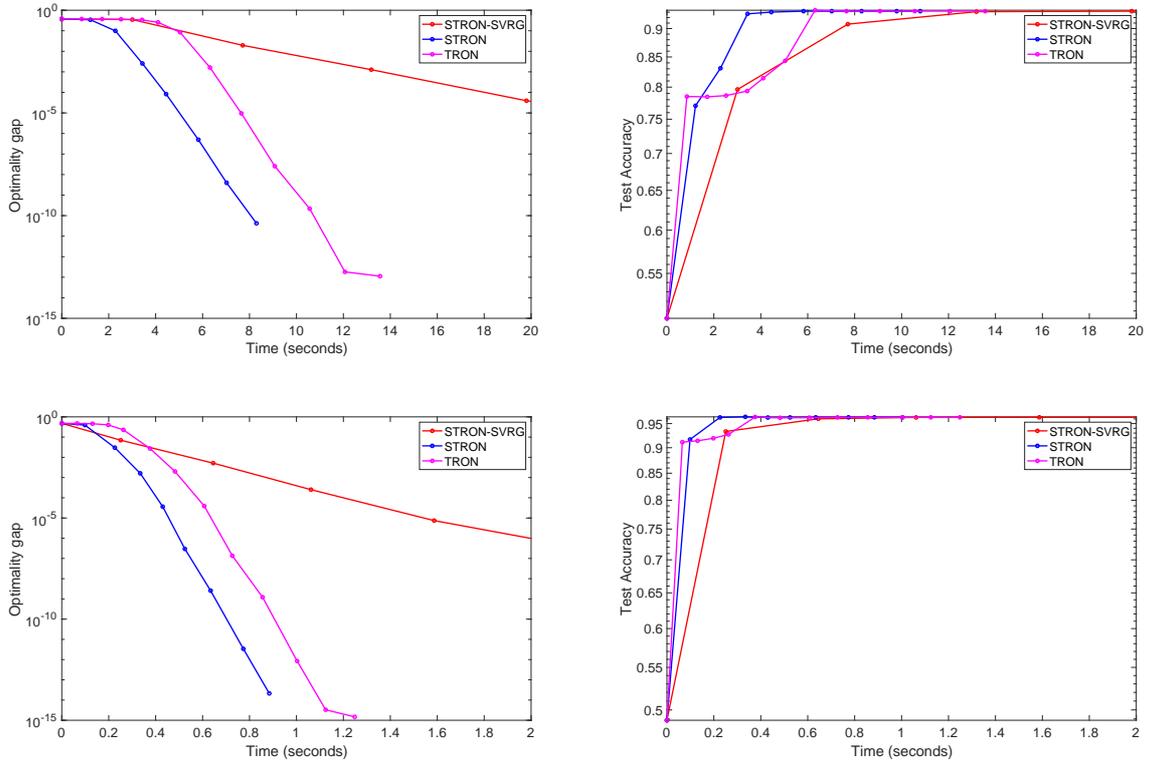


Fig. 6: Comparative study of STRON-SVRG, STRON and TRON on news20 (first row) and rcv1 (second row) datasets.

B.1 More Experiments

Fig. 7 presents more results on gisette and webspam datasets on l_2 -SVM and l_2 -regularized logistic regression, respectively. We observe results similar to Figs. 1–4, which show that STRON outperforms all other techniques.

B.2 Effect of Regularization Coefficient

Here we study the effect of the value of regularization coefficient (λ) on the convergence and accuracy of STRON and TRON methods. Fig. 8 presents the results with a series of values for $\lambda = \{1/l, 1e-1, 1e-3, 1e-5, 1e-7\}$ using news20 dataset with SVM. From the figure, it is clear that both the methods are affected by the choice of value λ . For larger values of $\lambda = \{1e-1, 1e-3\}$, both the methods converge to the less accurate solution, as depicted in terms of optimality gap and accuracy plots. On the other hand, for values of $\lambda < 1e-3$, there do not seem to be much difference in the accuracy of both

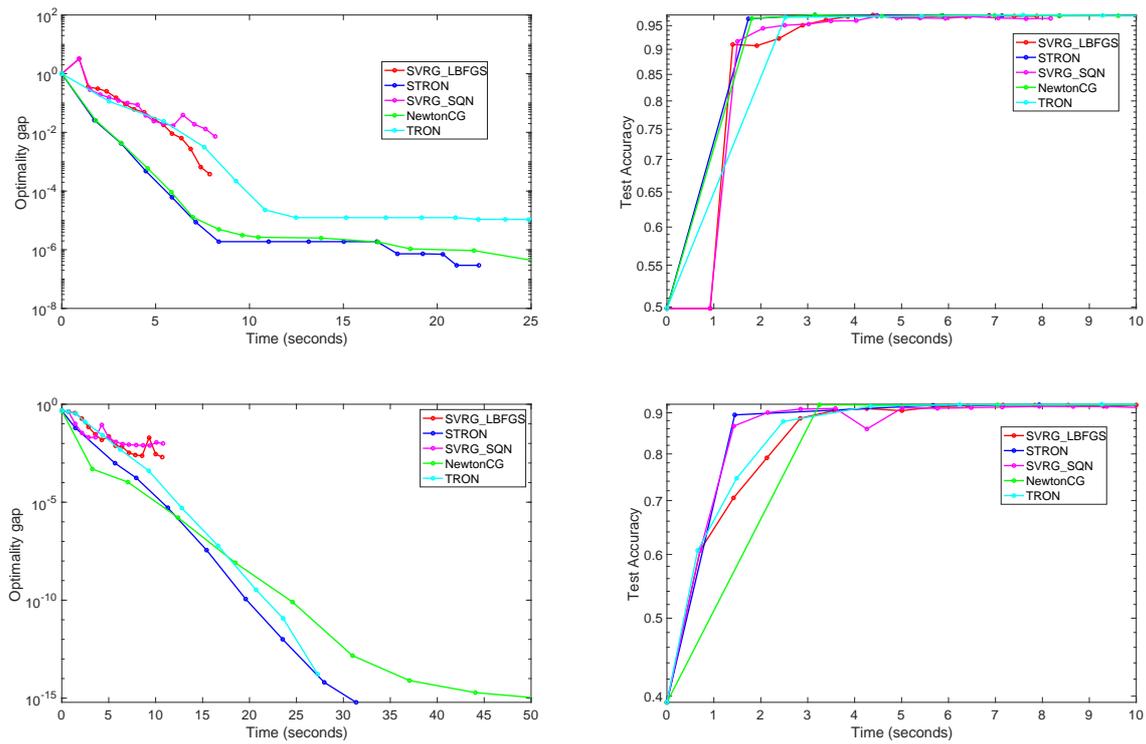


Fig. 7: First row presents results with gisette on SVM and second row presents results with webspam on logistic regression.

the methods. But, in terms of optimality gap, it is clear that smaller the value of λ , better is the solution, although there seem to be no difference on the corresponding accuracy plot. Moreover, for all the values of λ , STRON clearly outperforms TRON method.

Generally, machine learning problems involve tuning a lot of hyper-parameters which are quite difficult to tune. So to reduce number of parameters to tune, we set $\lambda = 1/l$, which works well in practice, as it sets the value relative to the size of the dataset and gives good results, as is clear from the figure.

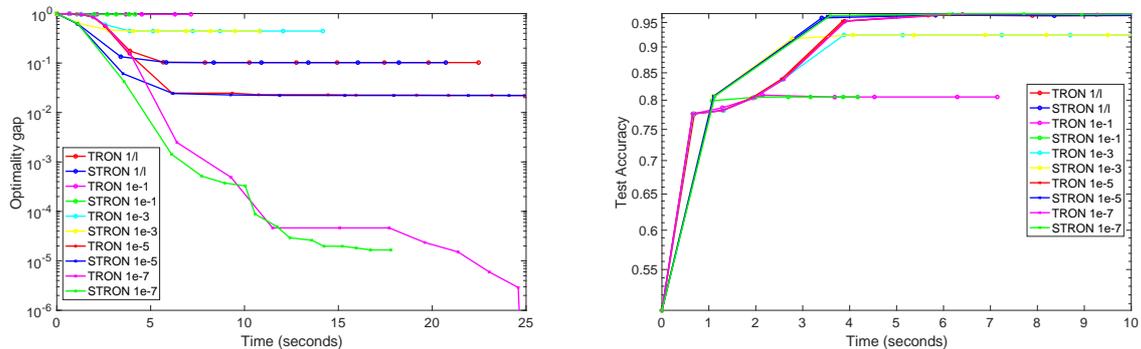


Fig. 8: Effect of regularization coefficient on STRON and TRON methods.

References

1. Agarwal N, Bullins B, Hazan E (2017) Second-order stochastic optimization for machine learning in linear time. *Journal of Machine Learning Research* 18(116):1–40
2. Allen-Zhu Z (2017) Katyusha: The First Direct Acceleration of Stochastic Gradient Methods. *Journal of Machine Learning Research* (to appear) Full version available at <http://arxiv.org/abs/1603.05953>
3. Bellavia S, Krejic N, Jerinkic NK (2018) Subsampled inexact newton methods for minimizing large sums of convex functions. *Optimization Online* URL http://www.optimization-online.org/DB_HTML/2018/01/6432.html
4. Berahas AS, Nocedal J, Takac M (2016) A multi-batch l-bfgs method for machine learning. In: *Advances in Neural Information Processing Systems* 29, pp 1055–1063
5. Bollapragada R, Byrd R, Nocedal J (2016) Exact and Inexact Subsampled Newton Methods for Optimization. arXiv URL <https://arxiv.org/abs/1609.08502>
6. Bollapragada R, Nocedal J, Mudigere D, Shi HJ, Tang PTP (2018) A progressive batching l-BFGS method for machine learning. In: *Proceedings of the 35th International Conference on Machine Learning*, PMLR, *Proceedings of Machine Learning Research*, vol 80, pp 620–629
7. Boyd S, Vandenberghe L (2004) *Convex Optimization*. Cambridge University Press, New York, NY, USA
8. Byrd R, Chin G, Neveitt W, Nocedal J (2011) On the use of stochastic hessian information in optimization methods for machine learning. *SIAM Journal on Optimization* 21(3):977–995, DOI 10.1137/10079923X
9. Byrd RH, Hansen SL, Nocedal J, Singer Y (2016) A stochastic quasi-newton method for large-scale optimization. *SIAM Journal on Optimization* 26(2):1008–1031
10. Cauchy AL (1847) Méthode générale pour la résolution des systèmes d’équations simultanées. *Compte Rendu des S’éances de L’Acad’emie des Sciences XXV S’erie A*(25):536–538
11. Chauhan VK, Dahiya K, Sharma A (2017) Mini-batch block-coordinate based stochastic average adjusted gradient methods to solve big data problems. In: *Proceedings of the Ninth Asian Conference on Machine Learning*, PMLR, vol 77, pp 49–64, URL <http://proceedings.mlr.press/v77/chauhan17a.html>
12. Chauhan VK, Dahiya K, Sharma A (2018) Problem formulations and solvers in linear svm: a review. *Artificial Intelligence Review* DOI 10.1007/s10462-018-9614-6, URL <https://doi.org/10.1007/s10462-018-9614-6>
13. Chauhan VK, Sharma A, Dahiya K (2018) Faster learning by reduction of data access time. *Applied Intelligence* 48(12):4715–4729, DOI 10.1007/s10489-018-1235-x
14. Chauhan VK, Sharma A, Dahiya K (2019) LIBS2ML: A Library for Scalable Second Order Machine Learning Algorithms. arXiv URL <https://arxiv.org/abs/1904.09448>, 1904.09448
15. Chauhan VK, Sharma A, Dahiya K (2019) Saags: Biased stochastic variance reduction methods for large-scale learning. *Applied Intelligence* DOI 10.1007/s10489-019-01450-3
16. Csiba D, Richt P (2016) Importance Sampling for Minibatches pp 1–19, arXiv:1602.02283v1
17. Defazio A, Bach F, Lacoste-Julien S (2014) Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In: *Proceedings of the 27th International Conference on Neural Information Processing Systems*, MIT Press, Cambridge, MA, USA, NIPS’14, pp 1646–1654
18. Fan R, Chang K, Hsieh C, Wang X, Lin C (2008) Liblinear: A library for large linear classification. *JMLR* 9:1871–1874
19. Fanhua S, Zhou K, Cheng J, Tsang IW, Zhang L, Tao D (2018) Vr-sgd: A simple stochastic variance reduction method for machine learning. arXiv URL <https://arxiv.org/abs/1802.09932>
20. Fletcher R (1980) *Practical methods of optimization*, vol. 1, unconstrained optimization
21. Hsia CY, Zhu Y, Lin CJ (2017) A study on trust region update rules in newton methods for large-scale linear classification. In: *Proceedings of the Ninth Asian Conference on Machine Learning*, PMLR, *Proceedings of Machine Learning Research*, vol 77, pp 33–48
22. Hsia CY, Chiang WL, Lin CJ (2018) Preconditioned conjugate gradient methods in truncated newton frameworks for large-scale linear classification. In: *Proceedings of the Tenth Asian Conference on Machine Learning*, PMLR, *Proceedings of Machine Learning Research*
23. Johnson R, Zhang T (2013) Accelerating stochastic gradient descent using predictive variance reduction. In: *Burges CJC, Bottou L, Welling M, Ghahramani Z, Weinberger KQ (eds) Advances in Neural Information Processing Systems 26*, Curran Associates, Inc., pp 315–323
24. Kolte R, Erdogdu M, Ozgur A (2015) Accelerating svrg via second-order information. In: *NIPS Workshop on Optimization for Machine Learning*
25. Le Roux N, Schmidt M, Bach F (2012) A Stochastic Gradient Method with an Exponential Convergence Rate for Strongly-Convex Optimization with Finite Training Sets. Tech. rep., INRIA
26. Lin C, Mor J (1999) Newton’s method for large bound-constrained optimization problems. *SIAM Journal on Optimization* 9(4):1100–1127, DOI 10.1137/S1052623498345075
27. Lin CJ, Weng RC, Keerthi SS (2008) Trust region newton method for logistic regression. *JMLR* 9:627–650
28. Liu DC, Nocedal J (1989) On the limited memory bfgs method for large scale optimization. *Mathematical Programming* 45(1):503–528
29. Lucchi A, McWilliams B, Hofmann T (2015) A variance reduced stochastic newton method. arXiv URL <http://arxiv.org/abs/1503.08316>
30. Mokhtari A, Ribeiro A (2014) Res: Regularized stochastic bfgs algorithm. *IEEE Transactions on Signal Processing* 62(23):6089–6104
31. Moritz P, Nishihara R, Jordan MI (2016) A linearly-convergent stochastic l-bfgs algorithm. In: *AISTATS*
32. Nocedal, Wright S (1999) *Numerical Optimization*. Springer, New York
33. Robbins H, Monro S (1951) A stochastic approximation method vol-22:pp. 400–407
34. Schmidt M, Le Roux N, Bach F (2016) Minimizing finite sums with the stochastic average gradient. *Math Program* pp 1–30
35. Schraudolph NN, Yu J, Gnter S (2007) A stochastic quasi-newton method for online convex optimization. In: *Meila M, Shen X (eds) Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics*, PMLR,

Proceedings of Machine Learning Research, vol 2, pp 436–443

36. Shalev-Shwartz S, Zhang T (2013) Stochastic dual coordinate ascent methods for regularized loss. *J Mach Learn Res* 14(1):567–599
37. Shalev-Shwartz S, Singer Y, Srebro N (2007) Pegasos: Primal estimated sub-gradient solver for svm. In: Proceedings of the 24th International Conference on Machine Learning, ACM, New York, NY, USA, ICML '07, pp 807–814
38. Steihaug T (1983) The conjugate gradient method and trust regions in large scale optimization. *SIAM Journal on Numerical Analysis* 20(3):626–637
39. Zhang Y, Xiao L (2015) Stochastic primal-dual coordinate method for regularized empirical risk minimization. In: Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37, ICML'15, pp 353–361
40. Zhou L, Pan S, Wang J, Vasilakos AV (2017) Machine learning on big data: Opportunities and challenges. *Neurocomputing* 237:350 – 361, DOI <https://doi.org/10.1016/j.neucom.2017.01.026>