

AI505
Optimization

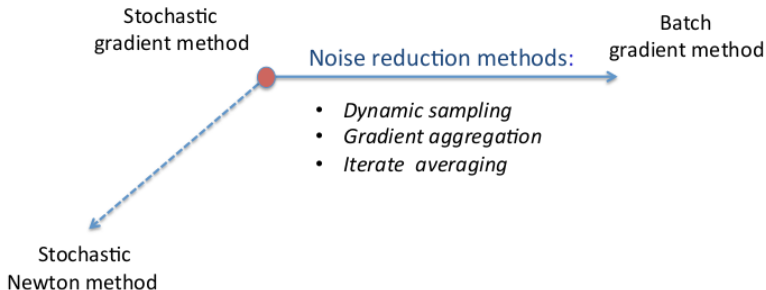
Optimization in Machine Learning

Marco Chiarandini

Department of Mathematics & Computer Science
University of Southern Denmark

Noise Reduction Methods

- SG as the ideal optimization approach for large-scale applications.
- SG suffers from the adverse effect of **noisy gradient estimates**.
when fixed stepsizes are used it prevents SG from converging to the solution
when a diminishing stepsize sequence $\{\alpha_k\}$ is employed it leads to a slow, sublinear rate of convergence.
- Remedies:



Achieve linear rate of convergence to the optimal value using a fixed stepsize.

- **Dynamic sampling methods** achieve noise reduction by gradually increasing the mini-batch size used in the gradient computation, thus employing increasingly more accurate gradient estimates as the optimization process proceeds.
- **Gradient aggregation methods** improve the quality of the search directions by storing gradient estimates corresponding to samples employed in previous iterations, updating one (or some) of these estimates in each iteration, and defining the search direction as a weighted average of these estimates.

Rate of convergence remains sublinear but reduces variance of iterates

- **iterate averaging methods** maintain an average of iterates computed during the optimization process and employ a more aggressive stepsize sequence—of order $O(1/\sqrt{k})$ rather than $O(1/k)$.

Reducing Noise at a Geometric Rate

rate of decrease in noise that allows a stochastic-gradient-type method to converge at a linear rate.

Consequence of Lipschitz assumption with ℓ constant:

$$\mathbb{E}_{\xi_k}[F(\mathbf{w}_{k+1})] - F(\mathbf{w}_k) \leq -\alpha_k \nabla F(\mathbf{w}_k)^T \mathbb{E}_{\xi_k}[g(\mathbf{w}_k, \xi_k)] + \frac{1}{2} \alpha_k^2 \ell \mathbb{E}_{\xi_k}[\|g(\mathbf{w}_k, \xi_k)\|_2^2]$$

We want to make the left hand side small (sequence of expected optimality gaps).

Theorem 5.1 (Strongly Convex Objective, Noise Reduction)

The SG method with a fixed stepsize $\bar{\alpha}$ and previous assumptions plus a variance of the stochastic vectors that decreases geometrically

$$\text{Var}_{\xi_k}[g(\mathbf{w}_k, \xi_k)] \leq M \zeta^{k-1}$$

has a sequence of expected optimality gaps that vanishes at a linear rate:

$$\mathbb{E}[F(\mathbf{w}_k) - F^*] \leq \omega \rho^{k-1}$$

Dynamic Sample Size Methods

Can we design efficient optimization methods attaining the critical bound on the variance?

Mini-batch stochastic gradient:

$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \bar{\alpha} g(\mathbf{w}_k, \xi_k)$$

where the stochastic directions are computed for some $\tau > 1$ as

$$g(\mathbf{w}_k, \xi_k) \stackrel{\text{def}}{=} \frac{1}{n_k} \sum_{i \in \mathcal{S}_k} \nabla f(\mathbf{w}_k; \xi_{k,i}) \quad \text{with } n_k \stackrel{\text{def}}{=} |\mathcal{S}_k| = \lceil \tau^{k-1} \rceil.$$

the mini-batch size increases geometrically as a function of the iteration counter k

Corollary 5.2. Let $\{\mathbf{w}_k\}$ be the iterates generated with unbiased gradient estimates, i.e., $\mathbb{E}_{\xi_{k,i}}[\nabla f(\mathbf{w}_k; \xi_{k,i})] = \nabla F(\mathbf{w}_k)$ for all $k \in \mathbb{N}$ and $i \in \mathcal{S}_k$. Then, the variance condition is satisfied, and if all other assumptions of Theorem 5.1 hold, then the expected optimality gap vanishes linearly.

Dynamic Sample Size Methods

Note: we described a method as linearly convergent but the per-iteration cost increases without bound.

Recall that SG method needs $\mathcal{T}(n, \epsilon) \leq 1/\epsilon$ evaluations to guarantee $\mathbb{E}[F(\mathbf{w}_k) - F^*] \leq \epsilon$

Theorem 5.3 Suppose that the dynamic sampling SG method is run with a stepsize $\bar{\alpha}$ satisfying (5.2) and some τ . In addition, suppose that all previous Assumptions hold. Then, the total number of evaluations of a stochastic gradient of the form $\nabla f(\mathbf{w}_k; \xi_{k,i})$ required to obtain $\mathbb{E}[F(\mathbf{w}_k) - F^*] \leq \epsilon$ is $O(\epsilon^{-1})$.

given the rate of convergence of a batch optimization algorithm on strongly convex functions (i.e., linear, superlinear, etc.), what should be the sampling rate so that the overall algorithm is **efficient** in the sense that it results in the lowest computational complexity?

- if the optimization method has a sublinear rate of convergence, then there is no sampling rate that makes the algorithm “efficient”;
- if the optimization algorithm is linearly convergent, then the sampling rate must be geometric (with restrictions on the constant in the rate) for the algorithm to be “efficient”;
- for superlinearly convergent methods, increasing the sample size at a rate that is slightly faster than geometric will yield an “efficient” method.

Design in Practice

- presetting the sampling rate, ie, $\tau > 1$ before running the optimization algorithm, requires some experimentation. Care must be put in preventing the full sample set from being employed too soon

- adaptive mechanisms to produce descent directions sufficiently often

- any direction $g(\mathbf{w}_k, \xi_k)$ is a descent direction for F at \mathbf{w}_k if, for some $\chi \in [0, 1)$, one has

$$\delta(\mathbf{w}_k, \xi_k) \stackrel{\text{def}}{=} \|g(\mathbf{w}_k, \xi_k) - \nabla F(\mathbf{w}_k)\|_2 \leq \chi \|g(\mathbf{w}_k, \xi_k)\|_2$$

verifying the inequality may be costly because involves the evaluation of $\nabla F(\mathbf{w}_k)$, one can estimate the left-hand side $\delta(\mathbf{w}_k, \xi_k)$, and then choose n_k so it holds sufficiently often.

- The sample variance obtained by sampling without replacements is bounded above by $\chi^2 \|g(\mathbf{w}_k, \xi_k)\|_2^2$
 - If this condition is not satisfied, then increase the sample size—either immediately in iteration k to a size that one might predict would satisfy such a condition.
 - no guarantee that the size n_k increases at a geometric rate. Remedy: if the adaptive increases the sampling rate more slowly than a preset geometric sequence, then a growth in the sample size is imposed.

Gradient Aggregation

- Rather than compute increasingly more **new** stochastic gradient information in each iteration, achieve a lower variance by **reusing and/or revising** previously computed information
- achieve a linear rate of convergence on strongly convex problems.
- improved rate is achieved primarily by either an increase in computation or an increase in storage.
- works on finite sums like R_n

SVRG

Procedure SVRG ; # Methods for Minimizing an Empirical Risk R_n

Choose an initial iterate $\mathbf{w}_1 \in \mathbb{R}^d$, stepsize $\alpha > 0$ and a positive integer m ;

for $k = 1, 2, \dots$ **do**

Compute the batch gradient $\nabla R_n(\mathbf{w}_k)$;

Initialize $\tilde{\mathbf{w}}_1 \leftarrow \mathbf{w}_k$;

for $j = 1, \dots, m$ **do**

$\tilde{\mathbf{g}}_j \leftarrow \nabla f_{i_j}(\tilde{\mathbf{w}}_j) - (\nabla f_{i_j}(\mathbf{w}_k) - \nabla R_n(\mathbf{w}_k))$; # $\nabla R_n(\mathbf{w}_k)$ from batch gradient

$\tilde{\mathbf{w}}_{j+1} \leftarrow \tilde{\mathbf{w}}_j - \alpha \tilde{\mathbf{g}}_j$;

Option (a): Set $\mathbf{w}_{k+1} = \tilde{\mathbf{w}}_{m+1}$;

Option (b): Set $\mathbf{w}_{k+1} = \frac{1}{m} \sum_{j=1}^m \tilde{\mathbf{w}}_{j+1}$;

Option (c): Choose j uniformly from $\{1, \dots, m\}$ and set $\mathbf{w}_{k+1} = \tilde{\mathbf{w}}_{j+1}$;

- since $E_{i_j}[\nabla f_{i_j}(\mathbf{w}_k)] = \nabla R_n(\mathbf{w}_k)$, one can view $\nabla f_{i_j}(\mathbf{w}_k) - \nabla R_n(\mathbf{w}_k)$ as the bias in the gradient estimate $\nabla f_{i_j}(\mathbf{w}_k)$.
- sampled gradient $\nabla f_{i_j}(\tilde{\mathbf{w}}_j)$ is corrected based on a perceived bias. Overall, $\tilde{\mathbf{g}}_j$ represents an unbiased estimator of $\nabla R_n(\tilde{\mathbf{w}}_j)$, but with a variance that one can expect to be smaller than as in simple SG

SAGA

in each iteration, it computes a stochastic vector \mathbf{g}_k as the average of stochastic gradients evaluated at previous iterates.

Procedure SAGA ; # Method for Minimizing an Empirical Risk R_n
Choose an initial iterate $\mathbf{w}_1 \in \mathbb{R}^d$ and stepsize $\alpha > 0$;
for $i = 1, \dots, n$ **do**
 Compute $\nabla f_i(\mathbf{w}_1)$;
 Store $\nabla f_i(\mathbf{w}_{[i]}) \leftarrow \nabla f_i(\mathbf{w}_1)$; # $\mathbf{w}_{[i]}$ represents the latest iterate at which ∇f_i
for $k = 1, 2, \dots$ **do**
 Choose j uniformly in $\{1, \dots, n\}$;
 Compute $\nabla f_j(\mathbf{w}_k)$;
 Set $\mathbf{g}_k \leftarrow \nabla f_j(\mathbf{w}_k) - \nabla f_j(\mathbf{w}_{[j]}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}_{[i]})$;
 Store $\nabla f_j(\mathbf{w}_{[j]}) \leftarrow \nabla f_j(\mathbf{w}_k)$;
 Set $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \alpha \mathbf{g}_k$;

As in SVRG, the method employs unbiased gradient estimates, but with variances that are expected to be less than the stochastic gradients that would be employed in a basic SG routine

SAGA

- Same per-iteration costs as basic SG
- on strongly convex R_n can achieve a linear rate of convergence but needs knowledge of at least ℓ .
- More effective initialization instead of evaluating all the gradients $\{\nabla f_i\}_{i=1}^n$ at the initial point. For example, one could perform one epoch of simple SG, or one can assimilate iterates one-by-one and compute \mathbf{g}_k only using the gradients available up to that point.
- SAGA needs to store n stochastic gradient vectors
- for very large n , gradient aggregation methods are comparable to batch algorithms and therefore cannot beat SG in this regime

Iterated Averaging Methods