

Optimization for Machine Learning

November 4, 2022

Program: Advanced stochastic gradient methods

→ Batch variants

→ Variance reduction

→ SG methods in practice

① Batch stochastic gradient methods

Pb minimize $w \in \mathbb{R}^d$ $f(w) = \frac{1}{n} \sum_{i=1}^m f_i(w)$, $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$
 C^1

f_i depends on the i th data point of a dataset

Stochastic Gradient:

$$w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$$

i_k drawn at random in $\{1, \dots, m\}$

Gradient Descent:

$$w_{k+1} = w_k - \alpha_k \nabla f(w_k)$$
$$= w_k - \frac{\alpha_k}{n} \sum_{i=1}^m \nabla f_i(w_k)$$

Batch Stochastic Gradient:

$$w_{k+1} = w_k - \frac{\alpha_k}{|S_k|} \sum_{i \in S_k} \nabla f_i(w_k)$$

stochastic gradient

(approximation of $\nabla f(w_k)$)

where S_k is a set of indices

randomly drawn in $\{1, \dots, m\}$

→ choosing $S_k = \{1, \dots, m\}$ (drawing n indices without replacement) gives GD

→ choosing $|S_k| = 1$ recovers SG

For $1 < |S_k| < n$, can identify two regimes of batch methods: depending on the batch size $|S_k|$

1) Large batch regime: $1 \ll |S_k| \ll n$

⊕ Converges to a tight neighborhood of a "solution" (global minimum / first-order stationary point)

⊖ High per-iteration cost \Rightarrow convergence curves are similar to that of GD

⊕ Still cheaper than GD but significantly more expensive than SG

2) Mini-batch regime: $1 < |S_k| \ll n$

⊕ Still much cheaper than GD, not so expensive compared to a SG iteration

⊕ Unlike SG, possibility to evaluate $\{f_i(x_k)\}_{i \in S_k}$ in parallel (with the right computing architecture/hardware)

⊖ Not necessarily better than SG, more sensitive to redundancies

$$\text{Ex) } \mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \quad \text{and} \quad \mathcal{D}' = \underbrace{\{\mathcal{D}, \dots, \mathcal{D}\}}_{\times 10}$$

SG with uniform sampling on \mathcal{D} and \mathcal{D}'
 \Rightarrow same probability of drawing a point corresponding to (x_1, y_1)

For batch SG with batch size > 1 , this is not true

NB: Batch methods also reduce the variability between several runs of the same method (compared to SG)

Connections with CV rates

For SG, under the assumptions that f was $C_{2,1}^{1,1}$ and μ -strongly convex, and assuming that the $\{i_k\}$ were independent with

$$\mathbb{E}_{i_k} [\nabla f_{i_k}(w_k)] = \nabla f(w_k) \quad \text{and} \quad \mathbb{E}_{i_k} [\|\nabla f_{i_k}(w_k)\|^2] \leq \|\nabla f(w_k)\|^2 + \sigma^2,$$

we proved that SG with stepsize $\alpha_k = \alpha$ converges to a neighborhood of the optimal value $\mathbb{E}[f(w_k)] \rightarrow [f^*, f^* + \frac{L\alpha\sigma^2}{2\mu}]$

$$(f^* = \min_{w \in \mathcal{D}} f(w))$$

↳ For batch SG: Suppose that every batch S_k consists of indices drawn in an iid fashion (e.g. uniformly at random in $\{1, \dots, n\}$ with replacement). Then, the

$$\mathbb{E}_{S_k} \left[\frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(w_k) \right] = \nabla f(w_k) \quad \text{and}$$

$$\mathbb{E}_{S_k} \left[\left\| \frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(w_k) \right\|^2 \right] \leq \|\nabla f(w_k)\|^2 + \frac{\sigma^2}{m_b} \quad \text{where } m_b = |S_k|$$

Lower "variance" for batch SG estimates

$$\frac{\sigma^2}{m_b} < \sigma^2 \quad \text{when } m_b > 1$$

$$\Rightarrow \text{CV to } \left[f^*, f^* + \frac{L\alpha\sigma^2}{2\mu m_b} \right] < \left[f^*, f^* + \frac{L\alpha\sigma^2}{2\mu} \right]$$

⚠ By using a batch size $m_b > 1$, the number of accesses to data points per iteration increases. For a given epoch budget, a batch method with batch size $m_b > 1$ will perform less iterations than SG.

1 epoch \equiv 1 iteration of GD
 \equiv n iterations of SG
 \equiv $\frac{n}{m_b}$ iterations of batch SG with batch size m_b

\Rightarrow Tradeoff: \rightarrow Increasing the batch size reduces the variance of the stochastic gradient estimates
 \rightarrow But increasing the batch size increases the number of data points accessed during 1 iteration

(2) Variance reduction techniques

Goal: Reduce the variance of stochastic gradients to avoid steps that are too different from gradient steps.
 \Rightarrow Reduce the overall variance within the algorithm

2.1. Increase the batch size

\rightarrow More expensive stochastic gradient estimates
 \rightarrow Reduces the variance of those estimates and the variability between algorithmic runs

2.2 Averaging

SG iteration: $w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$

Averaging = maintaining an averaged iterate $\bar{w}_k = \frac{1}{k+1} \sum_{i=0}^k w_i$

$\{\bar{w}_k\}$ has better guarantees than $\{w_k\}$ (statistical theory)
 in the convex setting

Computing $\{\bar{w}_k\}$ requires an additional vector

$$\bar{w}_{k+1} = \frac{1}{k+2} \left((k+1)\bar{w}_k + w_{k+1} \right)$$

Many iterations of SG $\Rightarrow \frac{1}{k+1}$ gets small very quickly
(numerical issues with computing \bar{w}_k)

2.3. Gradient aggregation methods

\hookrightarrow Idea: Take full gradient steps every once in a while to reduce the variance of the method

SVRG (Stochastic Variance Reduced Gradient) 2013

\hookrightarrow Two loops within the algorithm

* Outer loop (k): Compute a full gradient $\nabla f(w_k)$

* Inner loop (j): Perform m iterations

$$\tilde{w}_0 = w_k$$

$$\tilde{w}_{j+1} = \tilde{w}_j - \alpha_j \tilde{g}_j$$

Full gradient
computed once per
outer loop

$$\tilde{g}_j = \nabla f_{i_j}(\tilde{w}_j) - \nabla f_{i_j}(w_k) + \nabla f(w_k)$$

Stochastic
gradient
from iteration j
of the inner loop

where i_j drawn at random in $\{1, \dots, m\}$

already computed upon
calculating $\nabla f(w_k) = \frac{1}{m} \sum_{i=1}^m \nabla f_i(w_k)$

$E_{i_j}[\tilde{g}_j] = \nabla f(\tilde{w}_j)$ and \tilde{g}_j has a lower variance
than $\nabla f_{i_j}(\tilde{w}_j)$

Cost of an outer iteration = 1 iteration of GD + m iterations of SG

⊕ 1 iteration of SURG $\Rightarrow m$ updates of the iterate

⊕ Reduced variance

⊖ Need to store $\nabla f(w_k)$ and $\{\nabla f_i(w_k)\}_i$

⊕ Better w theory than SG

(in particular, $\mathbb{E}[f(w_k)] \rightarrow f^*$
at a linear rate for f
strongly convex \approx GD rate)

Beneficial with a (very) large budget
and efficient storage of $\nabla f_i(w_k)$ (not always possible)

SAGA (Stochastic Average Gradient++)

Approach: Improve over SURG by putting the emphasis on storing information rather than computing full gradients.

Algorithm: Iteration 0: Compute $\{\nabla f_i(w_0)\}_{i=1}^m$ and define $\nabla f_i(w_{[i]}) = \nabla f_i(w_0) \forall i$

Iteration k: Draw i_k at random in $\{1, \dots, m\}$

Compute $g_k = \nabla f_{i_k}(w_k) - \nabla f_{i_k}(w_{[i_k]}) + \frac{1}{m} \sum_{i=1}^m \nabla f_i(w_{[i]})$

and $w_{k+1} = w_k - \alpha_k g_k$

Set $\nabla f_{i_k}(w_{[i_k]}) = \nabla f_{i_k}(w_k)$

- ⊕ Same theoretical guarantees than SVRG
- ⊕ Same iteration cost (in terms of accesses to data points) than SG except for iteration 0
- ⊖ Need to store $\{\nabla f_i(w_{[i]})\}_{i=1}^m$ (and need to compute $\frac{1}{m} \sum_{i=1}^m \nabla f_i(w_{[i]})$)

Takeaways:

- Gradient aggregation methods have nice theoretical guarantees and converge faster than SG in terms of epochs
- But their storage cost and the need for full gradient evaluations make them prohibitive in certain applications like deep learning
- SAGA implemented in `skit-learn`, and has had some recent success in reinforcement learning.

③ Stochastic gradient methods for deep learning

General framework

minimize $f(w) = \frac{1}{m} \sum_{i=1}^m f_i(w)$
 $w \in \mathbb{R}^d$

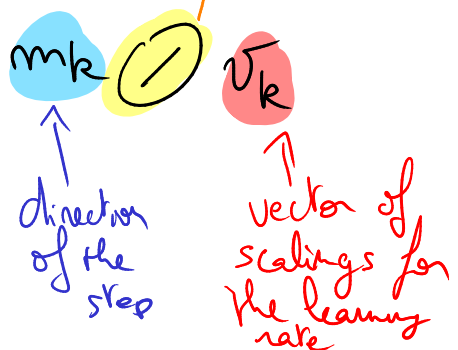
At every iteration k , i_k is drawn at random in $\{1, \dots, m\}$ and $\nabla f_{i_k}(w_k)$ is used to compute w_{k+1}

Iteration k : $w_{k+1} = w_k - \alpha$

$\forall i=1..d$

$$[w_{k+1}]_i := [w_k]_i - \frac{\alpha}{[w_k]_i} [m_k]_i$$

$\alpha > 0$
 (constant for simplicity)



↳ We recover SG by setting $m_k = \nabla f_{i_k}(w_k)$

$$\text{and } v_k = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^d$$

↳ The variants used in practice use more elaborate formulas for m_k and/or v_k that depend on the previous iterations

3.1. Stochastic gradient with momentum ("SGD with momentum")

$$\begin{aligned} \text{↳ } v_k &= \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^d, & m_k &= \beta_1 m_{k-1} + (1-\beta_1) \nabla f_{i_k}(w_k) \\ & & (m_{-1} &= 0) & \beta_1 &\in (0, 1) \\ & & & & (\beta_1 = 0 &\Rightarrow \text{SG}) \end{aligned}$$

↳ Related to accelerated gradient techniques, and in particular the heavy-ball method

$$(w_{k+1} = w_k - \alpha \nabla f(w_k) + \beta_k (w_k - w_{k-1}))$$

⇒ This method was used in the breakthrough results on training NN for image processing (~2012)

⇒ Very difficult to get theoretical guarantees for this method, especially in a (nonconvex) deep learning

3.2. AdaGrad (~2011)

↳ Issue: Scaling of learning rate per coordinate

$$\text{↳ } m_k = \nabla f_{i_k}(w_k)$$

$$\forall i=1..d, \quad [R_k^-]_i = [R_{k-1}^-]_i + [\nabla f_{i_k}(w_k)]_i^2 \quad \left. \vphantom{[R_k^-]_i} \right\} \rightarrow \text{sum of squares of all previous values of } [\nabla f_{i_k}(w_k)]_i$$

$$[v_k]_i = \sqrt{[R_k]_i + \epsilon} \Rightarrow \text{estimate of the magnitude of } [\nabla f(w)]_i$$

For numerical stability (e.g. 10^{-8})

$\Rightarrow \forall k, \forall i$, the learning rate at iteration k for the i th coordinate is

$$\frac{\alpha}{\sqrt{\sum_{j=0}^k [\nabla f(w_j)]_i^2 + \epsilon}}$$

- ⊕ Theory, very efficient for problems with sparse gradients (a lot of zero coordinates) \Rightarrow AdaGrad is efficient for recommender systems
- ⊖ If the problem is not sparse, AdaGrad might drive the stepsize too quickly towards 0

3.3. RMSProp

\hookrightarrow For training very deep networks

$\hookrightarrow m_k = \nabla f(w_k)$,

$$[R_k]_i = \beta_2 [R_k]_i + (1 - \beta_2) [\nabla f(w_k)]_i^2$$

$\beta_2 \in (0, 1)$
weighted average of

$$[v_k]_i = \sqrt{[R_k]_i + \epsilon}$$

current previous and information

3.4. Adam (2015)

\hookrightarrow Idea: Combine momentum and learning rate scaling through geometric averages

$$m_k = \frac{(1-\beta_1) \sum_{j=0}^k \beta_1^{k-j} \nabla f_j(w_j)}{1-\beta_1^{k+1}}$$

$\forall i=1 \dots d$

$$[v_k]_i = \sqrt{\frac{(1-\beta_2) \sum_{j=0}^k \beta_2^{k-j} [\nabla f_j(w_j)]_i^2}{1-\beta_2^{k+1}}} + \epsilon$$

→ The method of choice today for training deep models, especially in NLP (Adam & AdamW)

→ CV proof in 2015 was found incorrect in 2018