

OPTIMIZATION FOR MACHINE LEARNING

November 7, 2024

Today: Advanced stochastic gradient (with numerical illustration)

Next week (Nov. 14): Lab session on Stochastic gradient

Back to the exercise

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimize}} \quad f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{x}) \quad f_i \in C^{1,1}_{L_i}$$

Variant of SG: $x_{k+1} = x_k - \frac{\alpha_k}{C_{ik}} \nabla f_{i_k}(x_k)$

where i_k is drawn randomly in $\{1, \dots, m\}$

$$\forall i \in \{1, \dots, m\}, \quad \Pr(i_k = i) = \frac{c_i}{\sum_{j=1}^m c_j}$$

with $c_i = \frac{m L_i}{\sum_{j=1}^m L_j}$

① Show that $\mathbb{E}_{i_k} \left[\frac{1}{c_{i_k}} \nabla f_{i_k}(x_k) \right] = \nabla f(x_k)$

$$\mathbb{E}_{i_k} \left[\frac{1}{c_{i_k}} \nabla f_{i_k}(x_k) \right] = \sum_{i=1}^m \Pr(i_k = i) \times \frac{1}{c_i} \nabla f_i(x_k)$$

does not depend on i $\xrightarrow{\sum_{i=1}^m \frac{c_i}{\sum_{j=1}^m c_j} \times \frac{1}{c_i} \nabla f_i(x_k)}$

$$= \frac{1}{\sum_{j=1}^m c_j} \sum_{i=1}^m \nabla f_i(x_k)$$

$$\sum_{j=1}^m c_j = \sum_{j=1}^m \frac{m L_j}{\sum_{l=1}^m L_l} = \frac{1}{\sum_{l=1}^m L_l} \sum_{j=1}^m m L_j = m \frac{\sum_{j=1}^m L_j}{\sum_{l=1}^m L_l} = m$$

Hence $\mathbb{E}_{i_k} \left[\frac{1}{c_{i_k}} \nabla f_{i_k}(x_k) \right] = \frac{1}{m} \sum_{i=1}^m \nabla f_i(x_k) = \nabla f(x_k)$

$$\textcircled{2} \quad f \in C_L^{1,1} \quad L = \frac{\sum_{i=1}^m L_i}{m}$$

a) If $\alpha_n = \frac{1}{L}$, what is the value of $\frac{\alpha_k}{c_{ik}}$?

$$\frac{\alpha_k}{c_{ik}} = \frac{1}{L} \times \frac{\sum_{j=1}^m L_j}{m L_{ik}} = \frac{m}{\sum_{j=1}^m L_j} \times \frac{\sum_{j=1}^m L_j}{m} \times \frac{1}{L_{ik}} = \frac{1}{L_{ik}}$$

b) Interest of importance sampling?

Stepsize is tuned to the stochastic gradient / data point that is used at every iteration

$$L_{ik} \gg \frac{1}{m} \sum_{j=1}^m L_j = L \Rightarrow \frac{\alpha_n}{c_{ik}} \ll \alpha_n$$

ADVANCED STOCHASTIC GRADIENT METHODS

① Variance reduction methods

↳ Theory (and practical relevance) of SG relies on the stochastic gradients not varying too much from the true gradient

⇒ In the analysis, represented by the condition

$$\mathbb{E}_{ik} \left[\|\nabla f_{ik}(x_k)\|^2 \right] - \|\nabla f(x_k)\|^2 \leq \sigma^2$$

↑
 Variance
 parameter

Q) what can we change in the SG method to reduce the variance of the stochastic gradient estimates?
 (without reverting to GD!)

a) Use a batch

$$\nabla f_{i_k}(x_k) \rightarrow \frac{1}{m_b} \sum_{i \in S_k} \nabla f_i(x_k)$$

where S_k is a set of m_b indices drawn with or without replacement in $\{1, -1\}$

If the indices in S_k are drawn iid from a distribution

that satisfies $E_{i_k} \left[\| \nabla f_{i_k}(x_k) \|^2 \right] - \| \nabla f(x_k) \|^2 \leq \sigma^2$

when used in SG, then

$$E_{S_k} \left[\left\| \frac{1}{m_b} \sum_{i \in S_k} \nabla f_i(x_k) \right\|^2 \right] - \| \nabla f(x_k) \|^2 \leq \frac{\sigma^2}{m_b}$$

Theoretical consequences

- With 1 sample at every iteration and stepsize $\alpha > 0$, we showed that SG converges in function value (and in expectation) to a neighborhood of the optimum

$f(x)$
 μ-strongly
 convex

$$E \left[f(x_k) - f^* \right] \rightarrow \left[0, \frac{\sigma^2}{2\mu} \right]$$

at a rate $(1 - \alpha_u)^K$

$$\frac{\sigma^2}{m_b} \leq \sigma^2$$

• with a batch size m_b , can show

$$\mathbb{E}[f(x_k) - f^*] \rightarrow [0, \frac{Lx}{\alpha u} \frac{\sigma^2}{m_b}]$$

at a slower rate than SG $(1 - \frac{\alpha}{m_b} \mu)^K$

Practical consequences: Runs of batch SG tend to vary less than runs of Vanilla SG

b) Iterate averaging

$$\hookrightarrow \text{SG: } x_{k+1} = x_k - \alpha_u \nabla f_{i_k}(x_k)$$

$$\mathbb{E}[f(x_k) - f^*] \rightarrow \text{neighborhood of } 0$$

If we analyze the behavior of the average iterate, can

$$\text{show } \mathbb{E}\left[f\left(\frac{1}{K+1} \sum_{k=0}^K x_k\right) - f^*\right] \rightarrow 0$$

⊕ Smoother behavior (like for subgradient methods)

⊖ Memory cost

Naive implementation: store $K+1$ iterates

Better implementation: store 2 vectors (may still be too much for large models)

$$x_k, x_{k+1}$$

$$\hat{x}_k = \frac{1}{k+1} \sum_{l=0}^k x_l \quad \Rightarrow \quad \hat{x}_{k+1} = \frac{k+1}{k+2} \hat{x}_k + \frac{1}{k+2} x_{k+1}$$

(Running average of the iterates)

c) Gradient aggregation methods

→ Idea: Combine SG steps with full gradient calculations

. First major algorithm of that form: SVRG (2013)
Stochastic Variance-Reduced Gradient

Outer loop
Iteration k (x_k)

- Compute $\nabla f(x_k) = \frac{1}{m} \sum_{i=1}^m \nabla f_i(x_k)$
- Set $\tilde{x}_0 = x_k$
- Inner loop
 - For $j=0, \dots, m-1$ ($m \in \mathbb{N}, m \geq 1$)
 - Draw i_j uniformly in $\{1, \dots, m\}$
 - Set $\tilde{x}_{j+1} = \tilde{x}_j - \alpha_j \tilde{g}_j$ with $\alpha_j > 0$
 - And $\tilde{g}_j = \nabla f_{i_j}(\tilde{x}_j) - \nabla f_{i_j}(x_k) + \nabla f(x_k)$

1 iteration of
the inner loop
=
1 stochastic
gradient calculation

Picking an
iterate
at
random
↓

Better behavior
than the last iterate

$$\mathbb{E}_{j_k} [x_{k+1}] = \frac{1}{m} \sum_{j=0}^{m-1} \tilde{x}_{j+1}$$

$$\text{Key property : } \mathbb{E}_{ij} [\tilde{g}_j] = \nabla f(x_j)$$

→ Can also show that this method has better convergence

properties than SG

$$\mathbb{E}[f(x_k) - f^*] \rightarrow 0$$

vs $\rightarrow (0, \dots)$ for SG

Downside of SVRG : → Need to be able to compute $\nabla f(x_k)$!

→ 1 iteration costs $m + m$ sample gradients \Rightarrow higher than 1 iteration of GD!

An important alternative to SVRG : SAGA (2015)
Bach, Le Roux, Schmidt

↳ Do ① full gradient calculation at iteration 0

$$x_0 \rightarrow \{\nabla f_1(x_0), \dots, \nabla f_m(x_0)\}$$

At iteration k , you know x_k as well as an "old" gradient

for every component $\nabla f_{[1]}(x_{[1]}), \dots, \nabla f_{[m]}(x_{[m]})$

- (new) stochastic gradient calculation
- Sample i_h uniformly in $\{1, \dots, m\}$
 - Compute $\nabla f_{i_h}(x_h)$
 - Set

$$x_{h+1} = x_h - \alpha_h g_h$$

where $g_h = \nabla f_{i_h}(x_h) - \nabla f_{[i_h]}(x_{[i_h]}) + \frac{1}{m} \sum_{i=1}^m \nabla f_{[i]}(x_{[i]})$

- Set $\nabla f_{[i_h]}(x_{[i_h]}) = \nabla f_{i_h}(x_h)$

The last gradient of the form $\nabla f_{i_h}(\cdot)$ that was computed

- Properties:
- Except at iteration 0, the cost of 1 iteration is the same than an iteration of SG
 - $E_{\text{th}}[g_a] = \nabla f(x_a)$
 - Smaller convergence guarantees than SVRG
(better than SG)

Caveat

- Memory cost (store n gradients)
- Can be slow effectively for certain problems

Ex) Linear regression $f_i(x) = \frac{1}{2}(\hat{a}_i^T x - y_i)^2$

$$\nabla f_i(x) = (\hat{a}_i^T x - y_i) \hat{a}_i$$

can store $\hat{a}_i^T x \in \mathbb{R}$ and
recompute $\nabla f_i(x)$ when needed

\Rightarrow Implementation of SAGA in scikit-learn

Refs:

R.M. Gower, M. Schmidt, F. Bach and P. Richtárik
Variance-reduced methods for machine learning
(2020)

(2) Stochastic gradient methods for deep learning

→ Most popular methods used to train deep learning models
that you can find in PyTorch, JAX, ...

Generic iteration

$$x_{k+1} = x_k - \alpha_k m_k v_k$$

$\alpha_k > 0$
 $m_k \in \mathbb{R}^d$
 $v_k \in \mathbb{R}^d$

$$\Leftrightarrow [x_{k+1}]_j = [x_k]_j - \alpha_k \frac{[m_k]_j}{[v_k]_j} \quad \forall j=1..d$$

SG special case : $v_k = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$, $m_k = \nabla f_{i_k}(x_k)$ *i_k random*

Batch SG — $v_k = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$, $m_k = \frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(x_k)$ *S_k random*

⇒ For every method below (presented using SG), there is a batch counterpart

(1) SG with momentum

$$v_k = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$m_k = \beta_1 m_{k-1} + (1-\beta_1) \nabla f_{i_k}(x_k)$$

$\beta_1 \in [0, 1]$

$$\beta_1 = 0 \Rightarrow SG$$

→ Inspired by heavy ball

→ Default settings : PyTorch : $\beta_1=0$ JAX : $\beta_1=0.9$

→ Was the method of choice for training neural networks in 2012

② Adagrad (2011~2014)

$$m_k = \nabla f_k(x_k) \text{ (like SG)} + \text{diagonal scaling}$$

$$\forall j=1..d, [v_k]_j = \sqrt{\sum_{l=0}^k [m_l]_j^2} = \sqrt{\sum_{l=0}^k [\nabla f_l(x_l)]_j^2}$$

⇒ At every iteration, the stepsize along the j^{th} coordinate is normalized according to all past values of stochastic gradients

$$\Rightarrow \text{Equivalent to } x_{k+1} = x_k - \alpha_k H_k^{-1/2} \nabla f_k(x_k)$$

$$H_k = \text{diag}\left(\sum_{l=0}^k m_l m_l^T\right)$$

$$\text{diag}(A) = \begin{bmatrix} A_{11} & & \\ 0 & \ddots & \\ & & A_{dd} \end{bmatrix}$$

→ Efficient on problems in which the (stochastic) gradients have coordinates that differ in magnitude or when the gradients are sparse (lots of zero coefficients)

→ Has been applied successfully to problems from recommendation systems

⚠ Tends to produce exceedingly small stepsizes very quickly

→ Theoretical guarantees!

③ RMSProp

(Root Mean Square Propagation,
contemporary to Adagrad)

$$m_k = \nabla f_{ik}(x_k)$$

$$\forall j=1..d, [v_k]_j = \sqrt{\beta_2 [v_{k-1}]_j^2 + (1-\beta_2) [\nabla f(x)]_j^2}$$

for some $\beta_2 \in [0,1)$

$$\text{Adagrad} \approx \beta_2 = \frac{1}{2}$$

→ PyTorch: $\beta_2 = 0.95$, JAX $\beta_2 = 0.9$

→ RMSProp was found effective for training very deep neural networks (early 2010s)

→ Some theory, including recent results (2022-2024)

④ Adam (2015) → Most cited optimization paper (Kingma & Ba)

$$0 < \beta_1 < 1 \quad 0 < \beta_2 < 1$$

$$m_k = \frac{(1-\beta_1)}{1-\beta_1^{k+1}} \sum_{l=0}^k \beta_1^{k-l} \nabla f_{il}(x_l)$$

$\nabla f_{il}(x_l)$
stochastic gradient at iteration l



Acts as an estimate of the average stochastic gradient

↑
Geometric weighted average:

More weight given to the most recent estimates

$$\forall j=1..d, \quad [\nabla_k]_{:j} = \sqrt{\frac{1-\beta_2}{1-\beta_2^{k+1}} \sum_{l=0}^k \beta_2^{k-l} [\nabla f_{\text{st}}(x_l)]_{:j}^2}$$

estimate
 of the
 variance of
 stochastic gradients in the j^{th} coordinate

Geometric average

→ THE method used today (Transformers, LMs, etc)
 + Batch, stepsize tuning

→ PyTorch / JAX / original paper: $\beta_1 = 0.9$, $\beta_2 = 0.999$

→ * Original paper had a proof of convergence (2015)
 * Flaw was found in 2018 → People kept using it!

M. Hardt and B. Recht (2022)
 Patterns, predictions and actions

Recommendations for SG. (+ momentum)

- ① Pick as large a batch size as possible given your computer's RAM
 \Rightarrow Memory / parallelism concerns
- ② Set $\beta_1 = 0$ or $\beta_1 = 0.9$
 \Rightarrow Do not tune β_1 too much
 (same holds for β_2)

(3) Find the largest constant stepsize such that SG doesn't diverge

(4) Run SG with that stepsize until the loss plateaus



(5) Reduce the stepsize by a constant factor (e.g. 10)

$$0.1 \rightarrow 0.01 \rightarrow 0.001$$

(6) Repeat (4) and (5) until convergence/budget exhausted.

Exercise

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

$$\text{Batch SG: } x_{k+1} = x_k - \frac{\alpha_n}{|S_k|} \sum_{i \in S_k} \nabla f_i(x_k)$$

$$\text{where } |S_k| = m_b \quad \forall k \quad 1 \leq m_b \leq n$$

$$\forall S \subseteq \{1, \dots, m\}, P(S_k = S) = \begin{cases} 0 & \text{if } |S| \neq m_b \\ \frac{1}{\binom{m}{m_b}} = \frac{m_b! (n-m_b)!}{n!} & \text{if } |S| = m_b \end{cases}$$

a) Show that $E_{S_k} \left[\frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(x_k) \right] = \nabla f(x_k)$

b) Describe two modifications to the algorithm that guarantee $\mathbb{E}[f(x_n) - f^*] \xrightarrow{n \rightarrow \infty} 0$ for $f \in \mathcal{L}_{\text{Lip}}, \mu\text{-smooth}$ convex

Very final note :

SG = Stochastic Gradient

SGD = Stochastic Gradient Descent

SG is not a descent method like GD
 \Rightarrow only guarantee descent on average