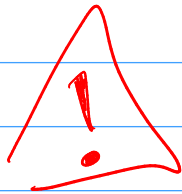


Optimization for Machine Learning

Stochastic Gradient pt 3

Today: Advanced methods (via notebook)



IMPORTANT: COURSE PROJECT

- 5 possible projects online
- Each student should indicate their 3 preferred projects
- Each student will then get 1 project assigned to them according to:
 - preferences
 - response time
 - balance number of students / project

↳ Regardless of the project

- You have to pick a dataset

- Deliverables: Python code / Notebook that can be run by the instructor + Short PDF report

- Projects are individual

↳ Deadline: TBA (last year January 15, likely around that time)

↳ We'll start allocating projects next week (~Oct 30) according to expressed preferences

Back to the exercise from last session

Recall: $f \in C_L^{1,1}$, μ -strongly convex, $f^* = \min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$

SG: $x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k)$ i_k random index in $\{1, \dots, n\}$

batch SG: $x_{k+1} = x_k - \frac{\alpha_k}{|S_k|} \sum_{i \in S_k} \nabla f_i(x_k)$ S_k set of m_b random indices drawn iid as in SG (e.g. uniformly) $i \in \{1, \dots, n\}$

CV rate for SG with $\alpha_k = \frac{1}{L}$

HKEM,

$$\mathbb{E}[f(x_k) - f^*] \leq \frac{\sigma^2}{2\mu} + \left(1 - \frac{\mu}{L}\right)^k \left(f(x_0) - f^* - \frac{\sigma^2}{2\mu}\right)$$

CV rate for batch SG with $\alpha_k = \frac{1}{L}$

$$\mathbb{E}[f(x_k) - f^*] \leq \frac{\sigma^2}{2\mu m_b} + \left(1 - \frac{\mu}{L}\right)^k \left(f(x_0) - f^* - \frac{\sigma^2}{2\mu m_b}\right)$$

1) Is the second guarantee better than for SG?

→ Rate is identical $\left(1 - \frac{\mu}{L}\right)^k$

→ Batch SG converges (in function value) in

$$\left[f^*, f^* + \frac{\sigma^2}{2\mu m_b}\right] \subseteq \left[f^*, f^* + \frac{\sigma^2}{2\mu}\right]$$

⇒ Better (smaller) neighborhood

→ Every iteration of batch SG is more expensive when $m_b > 1$

$$\frac{m_b}{m} \text{ epoch} \quad \text{VS} \quad \frac{1}{m} \text{ epoch}$$

For a fixed number of epochs N_E , the rate will be $(1 - \frac{\mu}{L})^{N_E \times m}$ for SG

$$(1 - \frac{\mu}{L})^{N_E \times \frac{m}{m_b}} \text{ for batch SG}$$

\Rightarrow Better (faster) for SG!

CONCLUSION : IT DEPENDS!

2) What about running SG with $\alpha_k = \frac{1}{m_b L}$?

Generic formula with constant α

$$\mathbb{E}[f(x_k) - f^*] \leq \frac{L\sigma^2\alpha}{2\mu} + (1 - \alpha\mu)^k \left[f(x_0) - f^* - \frac{L\sigma^2\alpha}{2\mu} \right]$$

$$\alpha = \frac{1}{m_b L}$$

$$\mathbb{E}[f(x_k) - f^*] \leq \frac{\sigma^2}{2\mu m_b} + \left(1 - \frac{\mu}{L m_b}\right)^k \left[f(x_0) - f^* - \frac{\sigma^2}{2\mu m_b} \right]$$

$\frac{\sigma^2}{2\mu m_b}$: same neighborhood

$\left(1 - \frac{\mu}{L m_b}\right)^k$: slower ∇ rate!
(smaller stepsize)

Comparison between SG $\alpha = \frac{1}{L}$ and $\alpha = \frac{1}{L m_b}$: IT depends!

ADVANCED STOCHASTIC GRADIENT METHODS

① Variance reduction methods

Basic SG:

- Randomized algorithm
⇒ How much do we deviate from the average performance in practice?
- Based on stochastic gradient approximations
⇒ How bad can those approximations be?

In the analysis, these variance considerations are represented by σ^2

For SG:
$$\mathbb{E}_{i_h} [\|\nabla f_{i_h}(x_h)\|^2] - \underbrace{\|\mathbb{E}_{i_h} [\nabla f_{i_h}(x_h)]\|^2}_{=\|\nabla f(x_h)\|^2 \text{ under additional assumption}} \leq \sigma^2$$

Q) What can we change in the algorithm to reduce that variance?

a) Use a batch

↳ If $\mathbb{E}_{i_h} [\|\nabla f_{i_h}(x_h)\|^2] - \|\nabla f(x_h)\|^2 \leq \sigma^2$ for SG
and a batch method uses m_b indices drawn iid as in SG,

then $\mathbb{E}_{S_k} [\|\uparrow\|^2] - \|\nabla f(x_k)\|^2 \leq \frac{\sigma^2}{n_b}$

$\frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(x_k)$

\uparrow
 $< \sigma^2$
if $n_b > 1$

b) Iterate averaging

$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ "average of the f_i functions"

\Rightarrow Under some assumptions, we proved convergence in expected value ($\mathbb{E}[f(x_k) - f^*]$)

on average!

Averaging:

$\bullet x_{k+1} = x_k - \alpha_k \nabla f_i(x_k)$

$\bullet \hat{x}_{k+1} = \frac{1}{k+1} \sum_{l=0}^k x_l$

$= \frac{k}{k+1} \hat{x}_k + \frac{1}{k+1} x_k$

c) Gradient aggregation

Idea: Combine SG steps with a full gradient estimate

First important method: SVRG (2013)

Iteration k of SVRG (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 $i_j \sim \mathcal{U}\{1, \dots, m\}$

Set $\tilde{x}_{j+1} = \tilde{x}_j - \alpha_j \tilde{g}_j$, where $\alpha_j > 0$

and

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

- Draw $j_k \sim \mathcal{U}\{0, \dots, m-1\}$ and
set $x_{k+1} = \tilde{x}_{j_k+1}$

Motivation: In the inner loop, the stochastic gradient $\nabla f_{i_j}(\tilde{x}_j)$ is corrected using gradient information from the outer loop

^{outer}
1 iteration of SVRG = 1 full gradient $\nabla f(x_k)$
+ m stochastic gradients $\nabla f_{i_j}(\tilde{x}_j)$

\Rightarrow A lot more expensive than SG and GD
(per iteration)

\Rightarrow But SVRG has much better w/guarantees

Ex) On $C_{L,1}^{1,1}$, μ -strongly convex f ,
SG with stepsize $\frac{1}{L}$

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

SVRG stepsize $1/L$ $\mathbb{E}[f(x_k)] \rightarrow f^*$

Remark:

$\forall k, \forall j,$

$$\begin{aligned}\mathbb{E}_{ij}[\tilde{g}_j] &= \mathbb{E}_{ij}[\nabla f_{ij}(\tilde{x}_j) - \nabla f_{ij}(x_k) + \nabla f(x_k)] \\ &= \nabla f(\tilde{x}_j)\end{aligned}$$

Main drawback of SVRG: per-iteration cost

\hookrightarrow Multiple attempts to improve SVRG

\hookrightarrow One important method: SAGA

(Bach, Le Roux,
Schmidt ~2015)

SAGA

≡ full
gradient
calculation

$$x_0 \in \mathbb{R}^d$$

For $i=1 \dots m$

$$\text{Set } \nabla f_{[i]}(x_{[i]}) = \nabla f_i(x)$$

↪ store m component
gradients

Main loop → For $k=0, 1, \dots$

Compute new
 ∇f_i

Draw $i_k \sim \mathcal{U}(\{1, \dots, m\})$

Compute $\nabla f_{i_k}(x_k) \rightarrow$ computed at iteration k

perform step

$$\text{Set } g_k = \nabla f_{i_k}(x_k) - \nabla f_{i_k}(x_{[i_k]}) + \frac{1}{m} \sum_{i=1}^m \nabla f_{[i]}(x_{[i]})$$

$$\text{and } x_{k+1} = x_k - \alpha_k g_k$$

Update
 $\nabla f_{[i_k]}()$

$$\text{Set } \nabla f_{[i_k]}(x_{[i_k]}) = \nabla f_{i_k}(x_k)$$

↓
every $\nabla f_{[i]}(x_{[i]})$

contains the most recently computed

$$\nabla f_i(-)$$

$\nabla f_{[i_k]}(\text{orange}) = \text{last gradient } \nabla f_{i_k}() \text{ that}$
was computed

↳ In terms of access to data points, SAGA and SG have the same per-iteration cost (1) except for iteration 0 ($n+1$ for SAGA vs 1 for SG)

↳ In terms of storage, SAGA requires to store n vectors of dimension d ($\nabla f_{[1]}(x_{[1]}), \dots, \nabla f_{[n]}(x_{[n]})$)

(Notebook: $n=1000$, $d=50$
50000)

- Difficult to implement in a large-scale regime

- But efficient with moderate n and d

⇒ Good scikit-learn implementation

+ Ad hoc implementations for certain problems

Ex) Linear regression problem

$$f(x) = \frac{1}{2n} \|Ax - y\|^2 = \frac{1}{n} \sum_{i=1}^n \frac{(\overbrace{a_i^T x - y_i}^{f_i(x)})^2}{2}$$

$\forall i=1..n, \quad \nabla f_i(x) = (\overbrace{a_i^T x - y_i}^{\in \mathbb{R}}) \overbrace{a_i}^{\substack{\uparrow \\ i^{\text{th}} \text{ data point}}}$

⇒ If we store $a_i^T x$, we can recompute the gradient when needed using an access to (a_i, y_i)

⇒ During a SAGA iteration, you will access

(x_{in}, y_{in}) and can use that to compute
 $\nabla f_{in}(x_{in})$ and $\nabla f_{in}(x_{in})$

\Rightarrow with this approach, can go from $n \times d$
storage to $n + d$

③ Stochastic gradient methods for deep learning

\hookrightarrow Used to train deep learning models

\hookrightarrow Implemented in PyTorch / JAX / TensorFlow

General formula (for SG, can be generalized to batch)

$$\forall k \in \mathbb{N}, \quad x_{k+1} = x_k - \overset{\alpha_k > 0}{\alpha_k} \odot \frac{m_k}{\|v_k\|}$$

$$\Leftrightarrow [x_{k+1}]_j = [x_k]_j - \alpha_k \frac{[m_k]_j}{[v_k]_j}$$

for every $j = 1 \dots d$

m_k : direction (computed using a stochastic gradient $\nabla f_{in}(x_k)$)

v_k : vector of normalization for the learning rate

Ex: SG fits that framework with
 $m_k = \nabla f_{ih}(x_k)$ and $v_k = 1_{\mathbb{R}^d} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

① SG with momentum / SGD with momentum

$$v_k = 1_{\mathbb{R}^d}$$

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

↑
momentum
step

↑
stochastic
gradient step

$$m_{k-1} = \frac{1}{\alpha_k} (x_k - x_{k-1})$$

$\beta_1 = 0$
 \rightarrow SG

↓

$$\beta_1 \in [0, 1)$$

→ Used in 2012 for the breakthrough paper
 on NN for ImageNet

⇒ Typically $\beta_1 = 0.9$ and $\alpha_k = 0.01$

② Adagrad (2011 ~ 2014)

$$m_k = \nabla f_{ih}(x_k) \quad (\text{as in SG})$$

$$\forall j=1 \dots d, \quad [v_k]_j = \sqrt{\sum_{l=0}^k [m_l]_j^2}$$

$$= \sqrt{\sum_{l=0}^k [\nabla f_{ih}(x_l)]_j^2}$$

\Rightarrow For coordinate j , the stepsize is set according to the j^{th} coordinate of all previous stochastic gradients

\Rightarrow Useful for problems with sparse gradients
= lots of zero coordinates

common in training recommender systems

③ RMSProp (2012)

$$m_k = \nabla f_k(x_k)$$

$$j=1..d \quad [v_k]_j = \sqrt{\beta_2 [v_{k-1}]_j^2 + (1-\beta_2) [m_k]_j^2}$$

$$\beta_2 \in [0, 1) \quad (\beta_2 = 0 \Rightarrow \text{Adagrad})$$

\Rightarrow Good performance on very deep networks

\Rightarrow Gives better stepsizes than Adagrad

4

Adam (2015) → Kingma & Ba
(most cited optimization paper!)

$$0 < \beta_1 < 1$$

$$0 < \beta_2 < 1$$

$$m_k = \frac{\underbrace{\nabla f_{i_2}(x_k) + \beta_1 \nabla f_{i_{k-1}}(x_{k-1}) + \dots + \beta_1^{k-1} \nabla f_{i_0}(x_0)}_{\sum_{l=0}^{k-1} \beta_1^{k-l} \nabla f_{i_l}(x_l)}}{1 - \beta_1^{k+1}}$$

$$\forall j=1 \dots d, \quad [v_k]_j = \sqrt{\frac{(1 - \beta_2) \sum_{l=0}^{k-1} \beta_2^{k-l} [\nabla f_{i_l}(x_l)]_j^2}{1 - \beta_2^{k+1}}}$$

⇒ Gives more weight to the most recent stochastic gradient through a geometric average

⇒ Empirical estimates of the mean of the stochastic gradients and the variance of all stochastic gradient coordinates

↳ THE method of choice today to train large networks, especially those in NLP (LLMs, Transformers, Diffusion models, etc)

- Maybe modern architectures are tuned to give good results when trained with Adam
- Default training algorithm in many platforms (PyTorch, JAX)
- CV proof from 2015 was found to be wrong in 2018