
Contents

1	Probability and training algorithms	3
1.1	Stochastic gradient descent method and convergence theory	3
1.1.1	Convergence of SGD	3
1.1.2	SGD with mini-batch	5

1

Probability and training algorithms

1.1 Stochastic gradient descent method and convergence theory

The next optimization problem is the most common case in machine learning.

Problem 1.

$$(1.1) \quad \min_{x \in \mathbb{R}^n} f(x),$$

where

$$(1.2) \quad f(x) = \frac{1}{N} \sum_{i=1}^N f_i(x).$$

One version of stochastic gradient descent (SGD) algorithm is:

Algorithm 1 SGD

Input: initialization x_0 , learning rate η_t .

For: $t = 0, 1, 2, \dots$

Randomly pick $i_t \in \{1, 2, \dots, N\}$ independently with probability $\frac{1}{N}$

$$(1.3) \quad x_{t+1} = x_t - \eta_t \nabla f_{i_t}(x_t).$$

1.1.1 Convergence of SGD

Theorem 1. Assume that each $f_i(x)$ is λ -strongly convex and $\|\nabla f_i(x)\| \leq M$ for some $M > 0$. If we take $\eta_t = \frac{a}{\lambda(t+1)}$ with sufficiently large a such that

$$(1.4) \quad \|x_0 - x^*\|^2 \leq \frac{a^2 M^2}{(a-1)\lambda^2}$$

then

$$(1.5) \quad \mathbb{E}e_t^2 \leq \frac{a^2 M^2}{(a-1)\lambda^2(t+1)}, \quad t \geq 1,$$

where $e_t = \|x_t - x^*\|$.

Proof. The L^2 error of SGD can be written as

$$(1.6) \quad \begin{aligned} \mathbb{E}\|x_{t+1} - x^*\|^2 &\leq \mathbb{E}\|x_t - \eta_t \nabla f_{i_t}(x_t) - x^*\|^2 \\ &\leq \mathbb{E}\|x_t - x^*\|^2 - 2\eta_t \mathbb{E}(\nabla f_{i_t}(x_t) \cdot (x_t - x^*)) + \eta_t^2 \mathbb{E}\|\nabla f_{i_t}(x_t)\|^2 \\ &\leq \mathbb{E}\|x_t - x^*\|^2 - 2\eta_t \mathbb{E}(\nabla f(x_t) \cdot (x_t - x^*)) + \eta_t^2 M^2 \\ &\leq \mathbb{E}\|x_t - x^*\|^2 - \eta_t \lambda \mathbb{E}\|x_t - x^*\|^2 + \eta_t^2 M^2 \\ &= (1 - \eta_t \lambda) \mathbb{E}\|x_t - x^*\|^2 + \eta_t^2 M^2 \end{aligned}$$

The third line comes from the fact that

$$(1.7) \quad \begin{aligned} \mathbb{E}(\nabla f_{i_t}(x_t) \cdot (x_t - x^*)) &= \mathbb{E}_{i_1 i_2 \dots i_t}(\nabla f_{i_t}(x_t) \cdot (x_t - x^*)) \\ &= \mathbb{E}_{i_1 i_2 \dots i_{t-1}} \frac{1}{N} \sum_{i=1}^N \nabla f_i(x_t) \cdot (x_t - x^*) \\ &= \mathbb{E}_{i_1 i_2 \dots i_{t-1}} \nabla f(x_t) \cdot (x_t - x^*) \\ &= \mathbb{E} \nabla f(x_t) \cdot (x_t - x^*), \end{aligned}$$

and

$$(1.8) \quad \mathbb{E}\|\nabla f_{i_t}(x_t)\|^2 \leq \mathbb{E}M^2 = M^2.$$

Note when $t = 0$, we have

$$(1.9) \quad \mathbb{E}e_0^2 = \|x_0 - x^*\|^2 \leq \frac{a^2 M^2}{(a-1)\lambda},$$

based on the assumption.

In the case of SDG, by the inductive hypothesis,

$$(1.10) \quad \begin{aligned} \mathbb{E}e_{t+1}^2 &\leq (1 - \eta_t \lambda) \mathbb{E}e_t^2 + \eta_t^2 M^2 \\ &\leq (1 - \frac{a}{t+1}) \frac{a^2 M^2}{(a-1)\lambda^2(t+1)} + \frac{a^2 M^2}{\lambda^2(t+1)^2} \\ &\leq \frac{a^2 M^2}{(a-1)\lambda^2} \frac{1}{(t+1)^2} (t+1 - a + a - 1) \\ &= \frac{a^2 M^2}{(a-1)\lambda^2} \frac{t}{(t+1)^2} \\ &\leq \frac{a^2 M^2}{(a-1)\lambda^2(t+2)} \cdot \left(\frac{t}{(t+1)^2} \leq \frac{1}{t+2} \right) \end{aligned}$$

□

1.1.2 SGD with mini-batch

Firstly, we will introduce a natural extended version of the SGD discussed above with introducing mini-batch.

Algorithm 2 SGD with mini-batch

Input: initialization x_0 , learning rate η_t .

For: $t = 0, 1, 2, \dots$

Randomly pick $B_t \subset \{1, 2, \dots, N\}$ independently
with probability $\frac{1}{\binom{N}{m}}$ and $\#B_t = m$.

$$(1.11) \quad x_{t+1} = x_t - \eta_t g_t(x_t).$$

where

$$g_t(x_t) = \frac{1}{m} \sum_{i \in B_t} \nabla f_i(x_t)$$

Now we introduce the SGD algorithm with mini-batch without replacement which is the most commonly used version of SGD in machine learning.

Algorithm 3 Shuffle SGD with mini-batch

Input: learning rate η_k , mini-batch size m , parameter initialization x_0 and denote $M = \lceil \frac{N}{m} \rceil$.

For Epoch $k = 1, 2, \dots$

Randomly pick $B_i \subset \{1, 2, \dots, N\}$ without replacement
with $\#B_i = m$ for $i = 1, 2, \dots, t$.

For mini-batch $t = 1 : M$

Compute the gradient on B_t :

Update x :

$$x \leftarrow x - \eta_k g_t(x),$$

where

$$g_t(x) = \frac{1}{m} \sum_{i \in B_t} \nabla f_i(x)$$

EndFor

EndFor

To “randomly pick $B_i \subset \{1, 2, \dots, N\}$ without replacement with $\#B_i = m$ for $i = 1, 2, \dots, t$ ”, we usually just randomly shuffle the index set first and then consec-

utively pick every m elements in the shuffled index set. That is the reason why we would like to call the algorithm as shuffled SGD while this is the mostly used version of SGD in machine learning.

Remark 1. Let recall a general machine learning loss function

$$(1.12) \quad L(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(h(X_i; \theta), Y_i),$$

where $\{(X_i, Y_i)\}_{i=1}^N$ correspond to these data pairs. For example, $\ell(\cdot, \cdot)$ takes cross-entropy and $h(x; \theta) = p(x; \theta)$ as we discussed in Logistic regression section.

Thus, we have the following corresponding relation

$$\begin{aligned} f(x) &\leftrightarrow L(\theta) \\ f_i(x) &\leftrightarrow \ell(h(X_i; \theta), Y_i). \end{aligned}$$