# OPTIMIZATION FOR MACHINE LEARNING

This week: Stochastic gradient!

Today : Basics          (+ early stopping?)

Course project: Coming up!

# STOCHASTIC GRADIENT METHODS
AKA SGD

<u>Motivation:</u> • Most problems in ML involve data, and have the typical form

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad f(x) = \frac{1}{m} \sum_{i=1}^{m} \ell\left( h(a_i ; x), y_i \right)$$

Loss function: measures agreement between model & data

Model parameterized by $x$

Data ($a_i$ to be matched with $y_i$)

↳ If all functions in the sum are $C^1$, then $f$ is $C^1$
and $\nabla f(x) = \frac{1}{m} \sum_{i=1}^{m} \nabla f_i(x)$ with $f_i(x) = \ell\left( h(a_i ; x), y_i \right)$

→ Computing $\nabla f(x)$ (for, e.g., performing an iteration of GD) requires to compute all gradients $\nabla f_i(x)$ (sample gradients)

⇒ Can be expensive if $m \gg 1$

⇒ If $\{(a_i, y_i)\}_{i=1..m}$ are sampled from some data distribution, then a good approximation for $\nabla f(x)$ may be found using less than $m$ sample gradients

→ More generally, this observation applies to stochastic optimization problems of the form
$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad f(x) = \mathbb{E}_{(a,y)}\left[ \ell\left( h(a ; x), y \right) \right]$$

# ① Basics of stochastic gradient

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

$m \geq 1$

Implicitly, $f_i$ depends on the $i^{th}$ data point in a dataset of size $m$.

## Gradient Descent

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

Iteration $k$:    $x_{k+1} = x_k - \alpha_k \nabla f(x_k) = x_k - \frac{\alpha_k}{m} \sum_{i=1}^{m} \nabla f_i(x_k)$

$\alpha_k > 0$

choosing all
$\nabla f_i$ s deterministically

## Stochastic Gradient (basic form)  1952
### aka Vanilla SG

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

Iteration $k$:    $x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k) = x_k - \frac{\alpha_k}{1} \nabla f_{i_k}(x_k)$

$\alpha_k > 0$

index drawn randomly
within $\{1, -, m\}$

## Example: The perceptron algorithm

$a_i \in \mathbb{R}^d, \ y_i \in \{-1, 1\}$

Problem    minimize    $\frac{1}{m} \sum_{i=1}^{m} \max(1 - y_i \, a_i^T x, 0)$
$x \in \mathbb{R}^d$

Perceptron algorithm (Rosenblatt)    • Start with $x_0$, $\alpha > 0$
  • Iteration $k$:    → Pick a data point $(a_{i_k}, y_{i_k})$ at random
          → If $\max(1 - y_{i_k} a_{i_k}^T x_k, 0) = 0$, set $x_{k+1} = x_k$.

$\longrightarrow$ If $\max(1 - y_{i_n} a_{i_n}^T x_n, 0) = 1 - y_{i_n} a_{i_n}^T x_n > 0$

then set $x_{n+1} = x_n + \alpha \, y_{i_n} a_{i_n}$

$\Rightarrow$ This is a stochastic (sub) gradient method!

NB: Stochastic gradient techniques can be generalized to nonsmooth functions by using subgradients instead of gradients

# The stochastic gradient family

- Choosing the random index at every iteration

$\longrightarrow$ Choose them in cyclic order $\{1, 2, -, m, 1, -, m\}$

$\Rightarrow$ Deterministic

$\Rightarrow$ Guarantee that every sequence of

"epoch" $\Big[$ $m$ iterations corresponds to looking at the entire dataset

Random Re-shuffling $\Bigg\{$ $\longrightarrow$ Draw a random permutation of $\{1, -, m\}$ $\{\sigma(1), -, \sigma(m)\}$, then use $i_0 = \sigma(1)$
$i_1 = \sigma(2)$
$\vdots$
$i_{m-1} = \sigma(m)$

Then repeat the process every $m$ iterations.

$\Rightarrow$ Every block of $m$ iterations corresponds to a pass over the entire dataset ("epoch")

$\longrightarrow$ Draw $\{i_k\}_k$ iid uniformly at random

$$\forall k, \quad \mathbb{P}\left(i_k = 1\right) = \ldots = \mathbb{P}\left(i_k = m\right) = \frac{1}{m}$$

$\Rightarrow$ Sampling with replacement: no guarantee that the entire dataset has been seen after $m$ iterations

$\Rightarrow$ On average, however, likely that all data points will be seen in comparable proportions

$\Rightarrow$ In that setting, an epoch $= m$ iterations

$\longrightarrow$ Why not sample more than 1 index?

$\Rightarrow$ Batch stochastic gradient methods

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

Iteration $k$

"batch"

• Sample a set $\widehat{S_k}$ of indices within $\{1, \ldots, m\}$ with or without replacement.

• $x_{k+1} = x_k - \dfrac{\alpha_k}{|S_k|} \displaystyle\sum_{i \in S_k} \nabla f_i(x_k)$

$\alpha_k > 0$

• $S_k = \{1, \ldots, m\} \; \forall k$

($m$ samples without replacement)

$\Rightarrow$ GD / "batch (S)G"

$$x_{k+1} = x_k - \frac{\alpha_k}{m} \sum_{i=1}^{m} \nabla f_i(x_k)$$

• $S_k = \{i_k\}$ (1 index per iteration)

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

$\Rightarrow$ Vanilla SG!

$\longrightarrow$ Batch methods include SG and GD as special cases

$\longrightarrow$ Other cases of interest:

$|S_k|$: batch size

• $|S_k| \ggg 1$ but not exactly $m$ ("large batch regime")

• $1 < |S_k| \lll m$ ("mini batch regime")

$\longrightarrow$ Often used when evaluations of $\nabla f_i$ can be parallelized
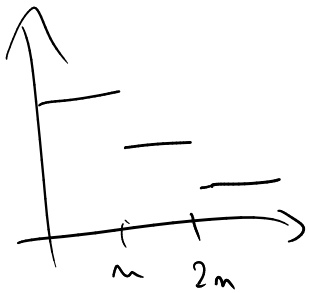
**. Choose the stepsize (aka learning rate)**

→ A constant stepsize can be surprisingly effective in practice!

→ Decreasing stepsizes typically core with the best theoretical guarantees ( proposed for the first stochastic gradient method: Pick $\{\alpha_n\}_n$ such that $\sum_{n=0}^{\infty} \alpha_n = \infty$

$$(\text{ex}) = \frac{1}{n+1}$$

$$\sum_{n=0}^{\infty} \alpha_n^2 < \infty \quad )$$

→ Hybrid strategies:



. Fix $\alpha_n$ to a constant value for $m$ iterations ("an epoch")

. Decrease that value before starting the next epoch

<u>Variant</u>: Use $\alpha$ for $T$ epochs

then $\alpha/2$ for $2T$ epochs

then $\alpha/4$ for $4T$ epochs

⋮

→ line search?

<u>Recall</u>: For GD, line search consists in testing values $\{\alpha, \theta\alpha, \theta^2\alpha, \dots\}$ for some $\theta \in (0,1)$ until a value $\theta^i\alpha$ is found such that

$$f(x - \theta^i\alpha \nabla f(x)) < f(x) - c\,\theta^i\alpha \|\nabla f(x)\|^2$$

→ For a stochastic gradient method, line search can be

defined in two ways

- Looking for $\alpha > 0$ such that $f(x_k - \alpha \nabla f_{i_k}(x_k)) < f(x_k)$

  $\longrightarrow$ Requires to evaluate $f = \frac{1}{m} \sum_{i=1}^{m} f_i$, hence to look at all data points, which defeats the purpose of SG (might have to evaluate $\nabla f$ as well)

  $\longrightarrow$ Typically not considered in SG methods

- Looking for $\alpha > 0$ such that $f_{i_k}(x_k - \alpha \nabla f_{i_k}(x_k)) < f_{i_k}(x_k)$

  $\longrightarrow$ No guarantee that $f(x_k - \alpha \nabla f_{i_k}(x_k)) < f(x_k)$

  $\longrightarrow$ Some success in overparameterized models (when the solution satisfies $\nabla f_i(x) = 0$ $\forall i = 1 \dots m$)

  $\longrightarrow$ Used in the Pytorch implementation of L-BFGS

  - For every sample/batch, there is an inner loop of iterations using the same batch and line search

- <span style="color:red">Advanced aspects (see Thursday's lectures)</span>

  $\longrightarrow$ Momentum (like in heavy ball)

  $\longrightarrow$ Variance reduction (Combine gradient and stochastic gradient steps)

  $\longrightarrow$ Output statistics of the iterates

  $\hookrightarrow$ Instead of $x_K$, output $\frac{1}{K} \sum_{k=0}^{K-1} x_k$

  (as in subgradient methods!)

$\hookrightarrow$ Instead of $x_K$, return $x_k$ $(k \leq K)$

with probability $\dfrac{\alpha_k}{\sum\limits_{h=0}^{K-1} \alpha_h}$.

---

# ② Theory of stochastic gradient

Setup:

$\qquad$ minimize $f(x) = \dfrac{1}{n} \sum\limits_{i=1}^{n} f_i(x)$ , $\quad f_i \; C^1$
$\quad x \in \mathbb{R}^d$

$\qquad\qquad f \; C_L^{1,1}$

$\qquad\qquad$ (Df $L$-Lipschitz continuous)

$\qquad \{x^*\} := \underset{x \in \mathbb{R}^d}{\text{argmin}} \; f(x)$

$\qquad\qquad f \; \mu$-strongly convex

$\qquad f^* := f(x^*) = \underset{x \in \mathbb{R}^d}{\min} f(x)$

$\qquad\qquad (\mu \leq L)$

$f \; C_L^{1,1} \implies \forall (x, w) \in (\mathbb{R}^d)^2,$

$\qquad f(w) \leq f(x) + \nabla f(x)^T (w - x) + \dfrac{L}{2} \|w - x\|^2$

$\qquad$ Instrumental in proving convergence for GD

$\qquad x = x_k \quad , \quad w = x_k - \alpha_k \nabla f(x_k) = x_{k+1}$

$\qquad f(x_{k+1}) \leq f(x_k) - \alpha_k \|\nabla f(x_k)\|^2 + \dfrac{L \alpha_k^2}{2} \|\nabla f(x_k)\|^2$

Q) what happens when $x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k)$ ?

$\qquad\qquad\qquad\qquad\qquad\qquad$ "stochastic gradient"

$$f(x_{n+1}) \leq f(x_n) - \alpha_k \nabla f(x_k)^T \nabla f_{i_k}(x_k) + \frac{L}{2}\alpha_k^2 \|\nabla f_{i_k}(x_k)\|^2$$

Random Variable

No guarantee that this is $> 0$

No guarantee that this is bounded compared to $\nabla f(x_k)^T \nabla f_{i_k}(x_k)$

Taking expected values with respect to $i_k$ on both sides, we get

$$(\star) \quad E_{i_k}\left[f(x_{n+1})\right] \leq E_{i_k}\left[f(x_k)\right] - E_{i_k}\left[\alpha_k \nabla f(x_k)^T \nabla f_{i_k}(x_k)\right] + \frac{L}{2}E_{i_k}\left[\alpha_k^2 \|\nabla f_{i_k}(x_k)\|^2\right]$$

**Assumptions:**

i) For every $k$, $i_k$ is drawn independently of $x_k, \alpha$ and $i_0, i_1, \ldots, i_{k-1}$

ii) For every $k$, $E_{i_k}\left[\nabla f_{i_k}(x_k)\right] = \nabla f(x_k)$

iii) For every $k$, $E_{i_k}\left[\|\nabla f_{i_k}(x_k)\|^2\right] \leq \|\nabla f(x_k)\|^2 + \sigma^2$

for some $\sigma^2 \geq 0$

Under these assumptions, $(\star)$ becomes

i)
$$E_{i_k}\left[f(x_{n+1})\right] \leq f(x_k) - \alpha_k \nabla f(x_k)^T E_{i_k}\left[\nabla f_{i_k}(x_k)\right] + \frac{L}{2}\alpha_k^2 E_{i_k}\left[\|\nabla f_{i_k}(x_k)\|^2\right]$$

ii)
$$= f(x_k) - \alpha_k \|\nabla f(x_k)\|^2 + \frac{L}{2}\alpha_k^2 E_{i_k}\left[\|\nabla f_{i_k}(x_k)\|^2\right]$$

iii)
$$\leq f(x_k) - \alpha_k \|\nabla f(x_k)\|^2 + \frac{L\alpha_k^2}{2}\|\nabla f(x_k)\|^2 + \frac{L}{2}\alpha_k^2 \sigma^2$$

$$\overbrace{\qquad\qquad\qquad\qquad}$$
Right hand side of the GD
inequality

$\downarrow$ "Noise" component

## Remarks

- Assumptions i) and ii) are satisfied by most sampling strategies
  (ex: $i_k$'s iid uniformly distributed in $\{1, \ldots, m\}$)

- Assumption iii) holds when $\|\nabla f_{i_k}(x_k)\|$ are bounded $\forall k$
  and it is implied by $\mathbb{E}_{i_k}\left[\|\nabla f_{i_k}(x_k)\|^2\right] \leq M^2 \quad M \geq 0$

- What we require is:

$$\mathbb{E}_{i_k}\left[x_k - \alpha_k \nabla f_{i_k}(x_k)\right]$$
$$\overset{(i)}{=} x_k - \alpha_k \, \mathbb{E}_{i_k}\left[\nabla f_{i_k}(x_k)\right]$$
$$\overset{(ii)}{=} x_k - \alpha_k \nabla f(x_k)$$

    i). Stochastic gradients are unbiased estimates of the
the true gradient ("on average, a SG step is a
gradient step")

    $(i) + (ii)$

    ii). Stochastic gradients should not deviate too much
from gradients in norm

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

$$\underbrace{\mathbb{E}_{i_k}\left[\|\nabla f_{i_k}(x_k)\|^2\right] - \|\mathbb{E}_{i_k}\left[\nabla f_{i_k}(x_k)\right]\|^2}_{\approx \text{ variance for } \|\nabla f_{i_k}(x_k)\|} \leq \sigma^2 \quad (iii)+(ii)$$

$\uparrow$ "Variance/ Noise parameter"

For a broader class of assumptions, see e.g.

L. Bottou, F. E. Curtis and J. Nocedal, Optimization Methods
for Large-Scale Machine Learning, SIAM Review (2018)

**Theorem:** SG with fixed step size

Suppose we run SG with $\{i_k\}$ satisfying Assumptions (i) − (iii) and

$\forall k, \alpha_k = \alpha \in \left(0, \frac{1}{L}\right]$  ( $f \in C_L^{1,1}$, $\mu$-strongly convex)

Then, for any $K \geq 1$,

$$\mathbb{E}_{i_0, -, i_{K-1}}\left[ f(x_K) - f^* \right] \leq \underbrace{(1-\alpha\mu)^{\textcircled{K}}\left[ f(x_0) - f^* - \frac{\alpha L\sigma^2}{2\mu} \right]}_{\substack{\downarrow > 0 \\ K \to \infty \quad (1-\alpha\mu) \geq 1 - \frac{\mu}{L} \\ \in (0,1)}} + \underbrace{\frac{\alpha L\sigma^2}{2\mu}}_{\substack{\text{Constant} \\ \text{(does not} \\ \text{depend} \\ \text{on } K)}}$$

"Expected convergence result": $\mathbb{E}_{i_0, -, i_{K-1}}\left[ \cdot \right]$

"Last-iterate convergence result": $x_K$

**Recall:** For GD on $C_L^{1,1}$, $\mu$-strongly convex $f$, we have

$$f(x_k) - f^* \leq \underbrace{(1-\alpha\mu)^K (f(x_0) - f^*)}_{\substack{\downarrow 0 \\ K \to \infty}} \qquad \forall \alpha \in \left(0, \frac{1}{L}\right]$$

$\hookrightarrow$ For GD, we can show that $f(x_k) - f^* \underset{K \to \infty}{\longrightarrow} 0$ at a rate $(1-\alpha\mu)^K$

$\hookrightarrow$ For SG, we can show that $\mathbb{E}\left[ f(x_k) - f^* \right] \underset{K \to \infty}{\longrightarrow} \left[ 0, \frac{\alpha L\sigma^2}{2\mu} \right]$ at a rate $(1-\alpha\mu)^K$

- Convergence in expected value (instead of deterministic)
- Convergence to a neighborhood of the optimal value (instead of convergence to the optimal value)

- But same rate of CV!

$\Rightarrow$ From that result, SG seems worse than GD...

$\Rightarrow$ ... however the comparison is conducted with a fixed iteration budget, and iterations of GD and SG have different costs.



Typical plots

$\longrightarrow$ 1 iteration of GD $\equiv$ 1 calculation of $\nabla f(\cdot)$

$\equiv$ $m$ calculation of $\nabla f_i(\cdot)$

$\longrightarrow$ 1 iteration of SG $\equiv$ 1 calculation of $\nabla f_i(\cdot)$

$\Rightarrow$ With that metric, 1 iteration of SG is $m$ times cheaper than 1 iteration of GD.

$\Rightarrow$ Given that an epoch corresponds to $m$ accesses to data points, 1 epoch $\equiv$ 1 iteration of GD

$\qquad$ 1 epoch $\equiv$ $m$ iterations of SG

**Corollary** Suppose we run SG and GD for a fixed number of epochs $N \geq 1$, i.e. $mN$ iterations of SG and $N$ iterations of GD. Then, for the final iterates $x_{mN}^{SG}$ and $x_N^{GD}$, we have

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

$$f(x_N^{GD}) - f^* \leq (1 - \alpha\mu)^{N} \left( f(x_0) - f^* \right)$$

→ SG still converges ==in expectation== to a ==neighborhood of the solution==, however when $n \gg 1$ it does so at a ==faster CV rate than GD==.

---

# Further results

- From expected value to high probability results: $\hat{x}_k \in \{x_0, x_1, -, x_k\}$

  $\Rightarrow$ By drawing an iterate at random and using it as output, you can get guarantees of the form

  $$\mathbb{P}\left( f(\hat{x}_k) - f^* \leq \dots \right) \geq ?$$

- To guarantee CV to a solution rather than a neighborhood,

  - Can use decreasing stepsizes

    $$(\text{Ex}) \quad \alpha_n = O\left(\frac{1}{n}\right) \Rightarrow \mathbb{E}\left[ f(x_k) - f^* \right] \leq O\left(\frac{1}{K}\right)$$

    (NB: still worse than GD in terms of iterations, better in terms of epochs)

  - Can use averaged iterates

    $$\overline{x}_K = \frac{1}{K} \sum_{h=0}^{K-1} x_h$$

    (as in subgradient methods)

  - Can also shrink the neighborhood using batch methods

    $$\left[ 0, \frac{\alpha L \sigma^2}{2\mu} \right] \xrightarrow[\substack{\text{Batch} \\ \text{size} \\ n_b}]{} \left[ 0, \frac{\alpha L \sigma^2}{2\mu \, n_b} \right]$$

- From strongly convex to nonconvex and other
  - There are results for the convex setting and the nonconvex setting
    
    $\Rightarrow$ For nonconvex, look at $\left(\mathbb{E}\left[\dfrac{1}{\sum_{h=0}^{k-1}\alpha_h}\sum_{h=0}^{k-1}\alpha_h\|\nabla f_{i_h}(x)\|^2\right]\right.$
  - Can extend this to nonsmooth functions, typically considering averaged iterates again

---

TAKEAWAYS:
- SG is designed to perform less calculation than GD at every iteration
- Performance depends on the step size and the sampling strategy
- SG has worse iteration CV rates than GD but has better rates in terms of epochs

Exercise: Importance sampling

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

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

Stochastic gradient variant

$\longrightarrow$ Pick $i_k$ such that $\mathbb{P}(i_k = i) = \dfrac{c_i}{\sum_{j=1}^{n} c_j}$

$\longrightarrow$ Set $x_{k+1} = x_k - \dfrac{\alpha_k}{c_{i_k}} \nabla f_{i_k}(x_k), \qquad \alpha_k > 0$

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

b) $f$ is $C_L^{1,1}$ with $L = \dfrac{1}{m} \sum_{i=1}^{m} L_i$

    i) If $\alpha_k = \dfrac{1}{L}$, what is $\dfrac{\alpha_k}{c_{i_k}}$ equal to?

    ii) What can then be the advantage of importance sampling over uniform sampling?