# Stochastic Gradient Descent

Axel Böhm

October 19, 2021

## Finite sum structure

Many optimization problems in Data science are sum structured:

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

◇ known as empirical risk (minimization)

◇ $f_i$ corresponds to the loss of the $i$-th observation

◇ for example: linear regression

$$f(x) = \|Ax - b\|^2 = \sum_{i=1}^{n} \left(a_i^T x - b_i\right)^2$$

◇ evaluating $\nabla f$ can be expensive if $n$ is large

## Risk minimization

In theory we would even like to minimize the population risk

$$f(x) = \mathbb{E}_\xi[f(x, \xi)]$$

$\diamond$ Typically no access to $f$

# (vanilla) Stochastic gradient descent

> sample $i \in 1, \ldots, n$ uniformly at random
>
> $x_{k+1} = x_k - \alpha \nabla f_i(x_k)$.

$\diamond$ requires only **one** gradient instead of $n$ per iteration.

$\diamond$ we call $g_t := \nabla f_i(x_k)$ a stochastic gradient (estimator)

## Unbiased

◇ Can't really use convexity as before since

$$f(x_k) - f(x^*) \leq \langle \nabla f_i(x_k), x^* - x_k \rangle$$

might not hold in general.

◇ But holds in expectation!

◇ For this we need that $\nabla f_i(x)$ is unbiased estimator of $\nabla f(x)$

$$\mathbb{E}[\nabla f_i(x)] = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(x) = \nabla f(x)$$

◇ We would like to conclude that

$$\mathbb{E}\left[\langle g_k, x^* - x_k\rangle\right] = \langle\mathbb{E}[g_k], \mathbb{E}[x^* - x_k]\rangle$$

but this is not so clear since $x_k$ is also stochastic and in general $\mathbb{E}[XY] \neq \mathbb{E}[X]\mathbb{E}[Y]$.

◇ We use the **conditional Expectation** $\mathbb{E}[\cdot|x_k]$ (read as expectation of $\cdot$ given $x_k$). Then

$$\mathbb{E}\left[\langle g_k, x^* - x_k\rangle|x_k\right] = \langle\mathbb{E}[g_k|x_k], x^* - x_k\rangle = \langle\nabla f(x_k), x^* - x_k\rangle.$$

◇ Together with the tower property $\mathbb{E}[\mathbb{E}[X|Y]] = \mathbb{E}[X]$:

$$\begin{aligned}
\mathbb{E}\left[\langle g_k, x^* - x_k\rangle\right] &= \mathbb{E}\left[\mathbb{E}\left[\langle g_k, x^* - x_k\rangle|x_k\right]\right] \\
&= \mathbb{E}\left[\langle\nabla f(x_k), x^* - x_k\rangle\right] \leq f(x^*) - f(x_k).
\end{aligned}$$

# Convergence statement: $\mathcal{O}(\epsilon^{-2})$ steps

### assumptions

◇ $f$ is convex and differentiable

◇ $\|x_0 - x^*\| \leq R$

◇ stochastic gradient are bounded in expectation $\mathbb{E}[\|g_k\|^2] \leq B^2$

### Theorem

*With the assumptions above and stepsize*

$$\alpha = \frac{R}{B\sqrt{k}}$$

*yields*

$$\mathbb{E}\left[f(\bar{x}_i) - f^*\right] \leq \frac{RB}{\sqrt{k}}.$$

error bound holds in expectation

# Proof

### Proof.

We start as usual ($g_k$ is a stochastic gradient)

$$\|x_{k+1} - x^*\|^2 \leq \|x_k - \alpha g_k - x^*\|^2$$
$$= \|x_k - x^*\|^2 + 2\alpha\langle g_k, x^* - x_k\rangle + \alpha^2\|g_k\|^2.$$

Now take expectation

$$\mathbb{E}\left[\|x_{k+1} - x^*\|^2\right] \leq \mathbb{E}\left[\|x_k - x^*\|^2\right] + 2\alpha\mathbb{E}[f^* - f(x_k)] + \alpha^2\mathbb{E}[\|g_k\|^2].$$

Bound gradients and telescope to finish the proof. □

## Comparing constants: SGD vs. GD

◇ GD: In the bounded (sub-)gradient analysis we assumed $\|\nabla f(x)\|^2 \leq B_{BG}^2$. For finite-sum this gives

$$\left\| \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(x) \right\|^2 \leq B_{BG}^2$$

◇ SGD: We assumed that the expected squared norm are bounded, i.e.

$$\mathbb{E}[\|\nabla f_i(x)\|^2] = \frac{1}{n} \sum_{i=1}^{n} \|\nabla f_i(x)\|^2 \leq B_{SGD}^2$$

By convexity we have that

◇ $B_{GD}^2 \approx \left\| \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(x) \right\|^2 \leq = \frac{1}{n} \sum_{i=1}^{n} \|\nabla f_i(x)\|^2 \approx B_{SGD}^2$

◇ but usually comparable

## Minibatch SGD

Instead of just using a single element $f_i$ we can use several $S \subset \{1, \ldots, n\}$

$$g_k := \frac{1}{|S|} \sum_{j \in S} \nabla f_j(x_k)$$

Interpolates between

◇ $|S| = 1 \Leftrightarrow$ (vanilla) SGD, as defined earlier

◇ $|S| = n \Leftrightarrow$ (batch) GD