# Stochastic gradient descent
## (contd.)

- Convergence analysis

- mini-batch variant

# Stochastic gradient descent

$$x_{t+1} = x_t - \eta_t \tilde{g}(x_t; \xi)$$

where $\tilde{g}(x_t; \xi)$ is unbiased estimate of $\nabla f(x_t)$, i.e.,

$$E\left[\tilde{g}(x_t; \xi)\right] = \nabla f(x_t)$$

ERM: minimize $f(x) = \dfrac{1}{n} \sum_{i=1}^{n} f_i(x)$

- Sample $i \in [n]$ uniformly at random
- $x_{t+1} = x_t - \eta_t \underbrace{\nabla f_i(x_t)}_{\tilde{g}_t}$ = Stochastic gradient

# Bounded stochastic gradients:

- Same convergence rate as gradient descent method

Claim: Let $f : \mathbb{R}^d \to \mathbb{R}$ be a convex and differentiable function, $\underline{x}^*$ be a global minimum of $f$;

$$\| \underline{x}_0 - \underline{x}^* \| \leq R \quad \text{and that} \quad \mathbb{E}\left[\| \underline{g}_t \|^2\right] \leq B^2 \quad \forall t$$

Then Stochastic gradient descent with

constant step size $\eta = \dfrac{R}{B\sqrt{T}}$ yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\left[f(\underline{x}_t)\right] - f(\underline{x}^*) \leq \frac{RB}{\sqrt{T}}$$

Iteration complexity : $O\left(\dfrac{1}{\varepsilon^2}\right)$          $O\left(\dfrac{1}{\sqrt{T}}\right)$

Recall our vanilla analysis:

$$\underbrace{g_t^T (x_t - x^*)}_{} = \frac{\eta}{2} \| g_t \|^2 + \frac{1}{2\eta} \left( \| x_t - x^* \|^2 - \| x_{t+1} - x^* \|^2 \right)$$

Telescoping sum:

$$\sum_{t=0}^{T-1} g_t^T (x_t - x^*) = \frac{\eta}{2} \sum_{t=0}^{T-1} \| g_t \|^2 + \frac{1}{2\eta} \left( \| x_0 - x^* \|^2 - \| x_T - x^* \|^2 \right)$$

$$\leq \frac{\eta}{2} \sum_{t=0}^{T-1} \| g_t \|^2 + \frac{1}{2\eta} \| x_0 - x^* \|^2$$

Taking expectation on both sides

$$\sum_{t=0}^{T-1} \mathbb{E} \left[ \tilde{g}_t^T (x_t - x^*) \right] \leq \frac{\eta}{2} \sum_{t=0}^{T-1} \underbrace{\mathbb{E} \left[ \| \tilde{g}_t \|^2 \right]}_{\leq B^2} + \frac{1}{2\eta} \underbrace{\| x_0 - x^* \|^2}_{\leq R^2}$$

We have the lower bound:

$$\mathbb{E} \left[ \tilde{g}_t^T (x_t - x^*) \right] \geq \mathbb{E} \left[ f(x_t) - f(\tilde{x}) \right]$$

$$\sum_{t=0}^{T-1} \mathbb{E}\left[f(\underline{x}_t) - f(\underline{x}^*)\right] \leq \frac{\eta}{2}\beta^2 T + \frac{1}{2\eta}R^2$$

$$= g(\eta)$$

Choose $\eta$ that minimize the upper bound:

$$\frac{1}{2}\beta^2 T - \frac{1}{2\eta^2}R^2 = 0$$

$$\eta = \frac{R}{\beta\sqrt{T}}$$

for which we have $O\left(\frac{1}{\sqrt{T}}\right)$

$\Rightarrow$ This can be directly extended to "Projected Stochastic gradient descent"

- Sample $i \in [n]$

- $y_{t+1} = x_t - \eta_t \, \tilde{g}_t$

- $x_{t+1} = P_C(y_{t+1})$

Proj. SGD

$\min. \ f(x)$

$\underline{x} \in C$

# Strong Convexity:

- $f$ is differentiable and $\mu$ strongly convex; with a decreasing stepsize

$$\eta_t = \frac{2}{\mu(t+1)}$$

Stochastic gradient descent yields

$$\mathbb{E}\left[f\left(\frac{2}{T(T+1)}\sum_{t=1}^{T} t \cdot x_t\right) - f(x^*)\right] \leq \frac{2B^2}{\mu(T+1)}$$

$$B = \max_{t=1,\ldots T} \mathbb{E}\left[\|\tilde{g}_t\|\right].$$

- We don't assume smoothness of $f$

  – diminishing step size

( Similar to the analysis of subgradient )

Recall our vanilla analysis:

$$\underline{g}_t^T (\underline{x}_t - \underline{x}^*) = \frac{\eta}{2} \| \underline{g}_t \|^2 + \frac{1}{2\eta} \left( \| \underline{x}_t - \underline{x}^* \|^2 - \| \underline{x}_{t+1} - \underline{x}^* \|^2 \right)$$

$$\mathbb{E} \left[ \underline{\tilde{g}}_t^T (\underline{x}_t - \underline{x}^*) \right] = \frac{\eta_t}{2} \mathbb{E} \left[ \| \underline{\tilde{g}}_t \|^2 \right] + \frac{1}{2\eta_t} \left[ \mathbb{E} \left[ \| \underline{x}_t - \underline{x}^* \|^2 \right] \right.$$

$$\left. - \mathbb{E} \left[ \| \underline{x}_{t+1} - \underline{x}^* \|^2 \right] \right]$$

Use strong convexity lower bound:

$$\mathbb{E} \left[ \underline{\tilde{g}}_t^T (\underline{x}_t - \underline{x}^*) \right] = \mathbb{E} \left[ \nabla f^T (\underline{x}_t) (\underline{x}_t - \underline{x}^*) \right]$$

$$\geq \mathbb{E} \left[ f(\underline{x}_t) - f(\underline{x}^*) \right]$$

$$+ \frac{\mu}{2} \mathbb{E} \left[ \| \underline{x}_t - \underline{x}^* \|^2 \right]$$

$$\Rightarrow \mathbb{E}\left[f(x_t) - f(x^*)\right] \le \frac{B^2 \eta_t}{2} + \frac{1}{2}\left(\eta_t^{-1} - \mu\right) \mathbb{E}\left[\|x_t - x^*\|^2\right]$$

$$- \frac{\eta_t^{-1}}{2} \mathbb{E}\left[\|x_{t+1} - x^*\|^2\right]$$

Substituting $\quad \eta_t = \dfrac{2}{\mu(t+1)}$ :

$$t \cdot \mathbb{E}\left[f(x_t) - f(x^*)\right] \le \frac{B^2 t}{\mu(t+1)} + \frac{\mu}{4}\left[t(t-1)\mathbb{E}\|x_t - x^*\|^2\right]$$

$$- (t+1)t \; \mathbb{E}\left[\|x_{t+1} - x^*\|^2\right]$$

$$\le \frac{B^2}{\mu} + \frac{\mu}{4}\left[t(t-1)\mathbb{E}\|x_t - x^*\|^2\right]$$

$$- (t+1)t \; \mathbb{E}\left[\|x_{t+1} - x^*\|^2\right]$$

Sum from $t = 1 \ldots T$.

$$\sum_{t=1}^{T} t \cdot \mathbb{E}\left[ f(\underline{x}_t) - f(\underline{x}^*) \right] \leq \frac{B^2 T}{\mu} + \frac{\mu}{4}\left[ 0 - T(T+1)\, \mathbb{E}\left[ \|\underline{x}_{t+1} - \underline{x}^*\|^2 \right] \right]$$

$$\leq \frac{B^2 T}{\mu}$$

We have $\quad \dfrac{2}{T(T+1)} \displaystyle\sum_{t=1}^{T} t = 1$ .

$$\mathbb{E}\left[ f\left( \frac{2}{T(T+1)} \sum_{t=1}^{T} t \cdot \underline{x}_t \right) - f(\underline{x}^*) \right] \leq \frac{2B^2}{\mu(T+1)}$$

$$\Rightarrow \quad \varepsilon\text{-accuracy requires } O\!\left( \frac{1}{\varepsilon} \right) \text{ steps.}$$

- Now, natural to ask if $f$ is $L$-smooth and $\mu$-strongly convex, will we get $O\left(\log\left(\frac{1}{\epsilon}\right)\right)$ (linear convergence) similar to the deterministic case.

Answer is <u>No.</u>

- Self-tuning property: $\nabla f(x) \longrightarrow 0$ as $\underline{x} \rightarrow \underline{x}^*$

$\Rightarrow$ Allows a big step size $\left[\frac{1}{L} \text{ or } \frac{2}{\mu + L}\right]$

$\Rightarrow$ So far $\eta \sim \frac{1}{\sqrt{T}}$ or $\eta_t = \frac{2}{\mu(t+1)}$

- No self-tuning for SGD. $\mathbb{E}\left[\|\tilde{g}_x\|_2^2\right] \not\rightarrow 0$ as $\underline{x} \rightarrow \underline{x}^*$

   - SGD responds to every new sample
   - Choose small steps close to the optimal

- $\mu$ - Strongly convex and $L$ - smooth

Suppose $\mathbb{E}\left[\|\tilde{g}_x\|_2^2\right] \leq \sigma_g^2 + C_g \|\nabla F(x)\|_2^2$
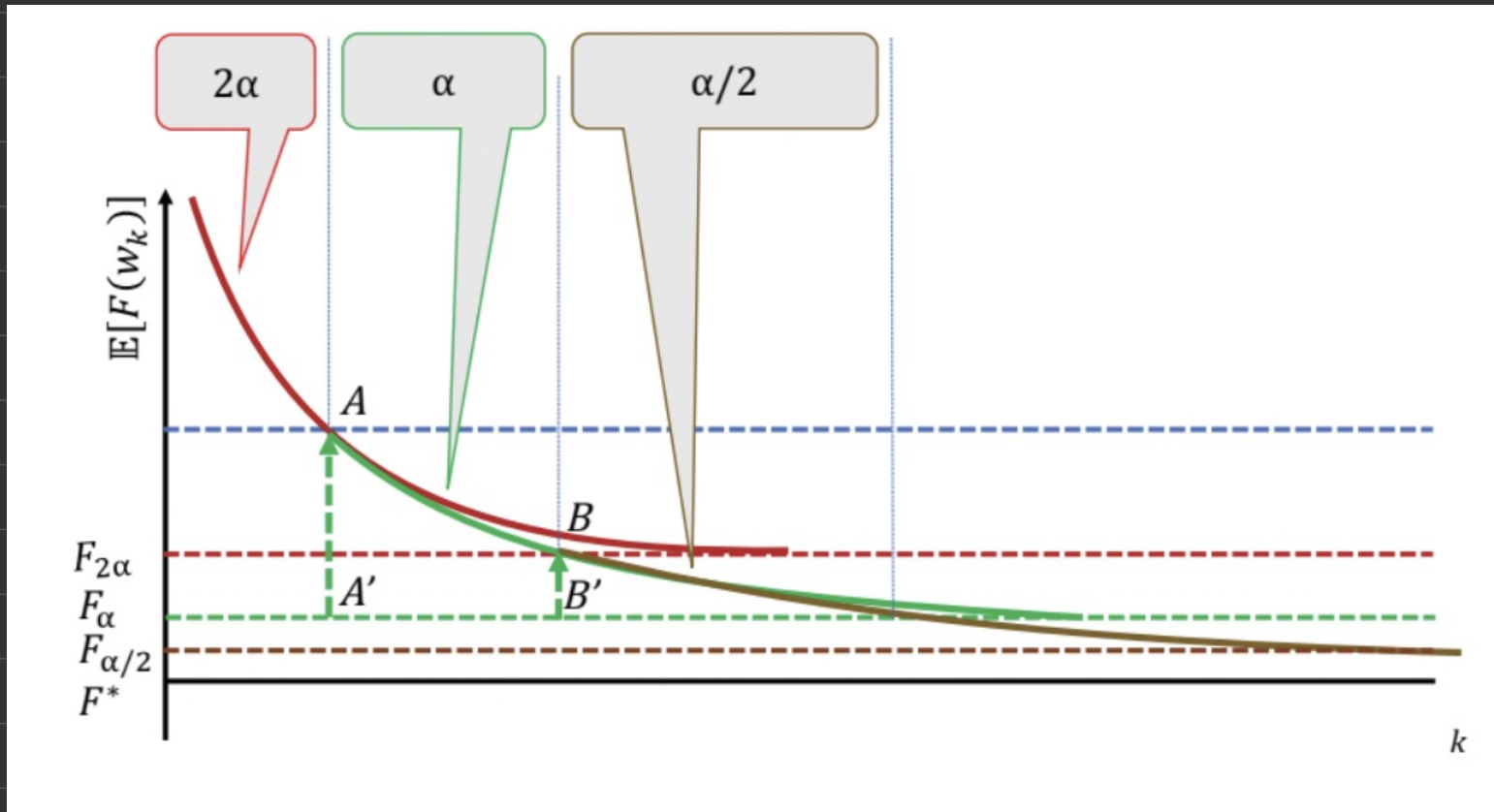
Then, SGD with fixed stepsize $\eta_t = \eta \leq \dfrac{1}{L C_g}$

yields

$$\mathbb{E}\left[f(x_t) - f(x^*)\right] \leq \dfrac{\eta L \sigma_g^2}{2\mu} + (1-\eta\mu)^t \left[f(x_0) - f(x^*)\right]$$

- $\sigma_g = 0$ : linear convergence

- Converges to some neighborhood of $x^*$

# Practical trick:



When progress stalls, half the stepsize & repeat

## Key question:

SGD with big stepsizes poorly suppresses noise. Larger stepsizes are needed for faster convergence.

How to reduce the variance?

Average iterates to reduce variance and improve convergence.

# Mini-batch variants: (Tame the variance)

- Instead of choosing a single $f_i$ from $\frac{1}{n} \sum_{i=1}^{n} f_i(x)$,

  let us pick several of them to form $\tilde{g}_t$

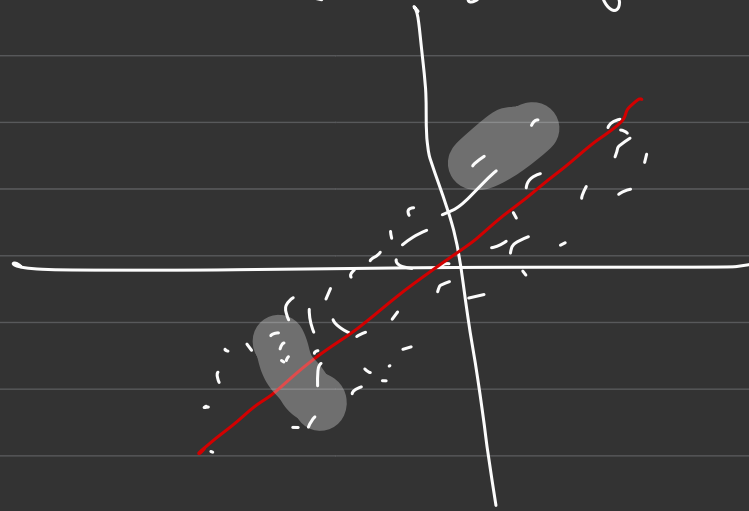- Let us pick $B \{ f_i \}$ : $f_1, f_2 \ldots f_B$

  and average the gradients:

  $$I_1, I_2 \ldots I_j \underset{B}{\sim} \text{uniform}(1, \ldots, n)$$

  $$x_{t+1} = x_t - \eta \frac{1}{B} \sum_{j=1}^{B} \nabla f_{I_j}(x_t)$$

- Stochastic gradient:
  $$\mathbb{E}\left[ \frac{1}{B} \sum_{j=1}^{B} \nabla f_{I_j}(x_t) \right] = \frac{1}{B} \sum_{j=1}^{B} \mathbb{E}\left[ \nabla f_{I_j}(x_t) \right]$$

  $$= \frac{1}{B} \sum_{j=1}^{B} \nabla f(x_t) = \nabla f(x_t)$$

- $B = 1$ , we have SGD $\leftarrow$

- $B = m$ , we have full gradient descent

- Reduces variance: (average of independent r.v. reduces variance)



- parallelization:

$$\tilde{g}_t = \frac{1}{B} \sum_{j=1}^{B} \nabla f_{I_j}(x_t)$$

can be computed independently in parallel