

# Learning the Optimal Step Size for Gradient Descent on Convex Quadratics

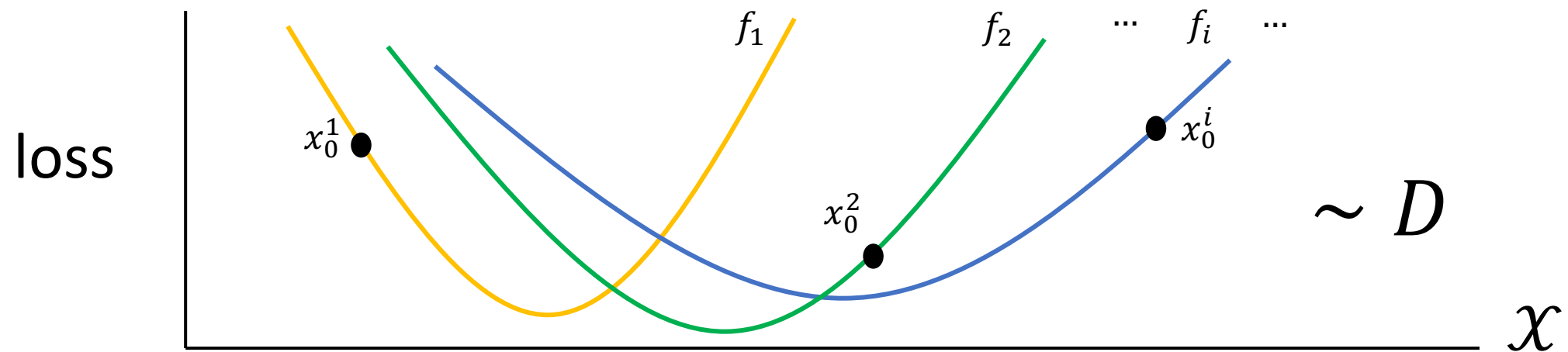
Alexandr Andoni, Daniel Hsu, Tim Roughgarden, **Kiran Vodrahalli**

Columbia University

NYAS ML Symposium, March 2020

# A Distribution on Optimization Problems

**Given:** Distribution  $D$  over  $(f, x_0)$ :



$f: \mathbb{R}^d \rightarrow \mathbb{R}$  is *convex quadratic*

# Learning Gradient Descent Step Size

$$x_{t+1} = x_t - \eta \nabla f(x_t)$$

After  $L$  steps:  $x_L(\eta, x_0)$

**Goal:** Learn optimal *single* step size  $\eta$  for distribution  $D$ .

$$\eta_L^* = \underset{\eta}{\operatorname{argmin}} \mathbb{E}_{f, x_0 \sim D} [f(x_L(\eta, x_0))]$$

# Motivation

- Gupta + Roughgarden '17: Sample complexity of learning step size of GD
- How much does learning the step size help performance?
- Push the limits of performance of a single step size

# Convex Quadratic Loss

$$f(x_L(\eta, x_0)) = \|Ax_L(\eta, x_0) - b\|_2^2$$

# Convex Quadratic Loss

$$\begin{aligned} f(x_L(\eta, x_0)) &= \|Ax_L(\eta, x_0) - b\|_2^2 \\ &= \sum_{i=1}^d \sigma_i^2 z_i^2 (1 - \eta \sigma_i^2)^{2L} \end{aligned}$$

# Convex Quadratic Loss

$$\begin{aligned} f(x_L(\eta, x_0)) &= \|Ax_L(\eta, x_0) - b\|_2^2 \\ &= \sum_{i=1}^d \sigma_i^2 z_i^2 (1 - \eta \sigma_i^2)^{2L} \end{aligned}$$

$$A = U\Sigma V^T \in \mathbb{R}^{n \times d}; n \geq d; z = V^T(x_0 - x^*); Ax^* = b$$

$$\mathbf{spectrum}(AA^T) = \sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_d^2 > 0$$

Worst Case

spectral ratio:  $c = \frac{\sigma_1^2}{\sigma_2^2} > 1$

As  $c \rightarrow \infty$ : No benefit over regular GD!



Worst Case

$$\text{spectral ratio: } c = \frac{\sigma_1^2}{\sigma_2^2} > 1$$

$$\text{condition number: } \kappa = \frac{\sigma_1^2}{\sigma_d^2} > 1$$

As  $c \rightarrow \infty$ : No benefit over regular GD!

**But:** If  $c \rightarrow \infty$ , with  $L = \alpha \cdot c, \alpha \in \mathbb{R}$ :

$$f(x_{L(\eta_L^*, x_0)}) - f(x^*) \leq \exp\left(-L\left(\frac{2}{c} + \frac{1}{\kappa}\right)\right) (f(x_0) - f(x^*))$$

Worst Case

$$\text{spectral ratio: } c = \frac{\sigma_1^2}{\sigma_2^2} > 1$$

$$\text{condition number: } \kappa = \frac{\sigma_1^2}{\sigma_d^2} > 1$$

As  $c \rightarrow \infty$ : No benefit over regular GD!

**But:** If  $c \rightarrow \infty$ , with  $L = \alpha \cdot c, \alpha \in \mathbb{R}$ :

$$f(x_L(\eta_L^*, x_0)) - f(x^*) \leq \exp\left(-L\left(\frac{2}{c} + \frac{1}{\kappa}\right)\right) (f(x_0) - f(x^*))$$

Improvement in the limit for large  $L$ ! ( $c < \kappa$ )

# Main Theorem

spectral ratio:  $c = \frac{\sigma_1^2}{\sigma_2^2} > 1$   
condition number:  $\kappa = \frac{\sigma_1^2}{\sigma_d^2} > 1$

**Theorem (informal):** For  $L$  large enough, for fixed  $f, x_0$  :

$$f(x_L(\eta_L^*, x_0)) - f(x^*) \leq \exp\left(-L\left(2 \log(1 + 1/c) + \frac{1}{\kappa}\right)\right) (f(x_0) - f(x^*))$$

# Main Theorem

spectral ratio:  $c = \frac{\sigma_1^2}{\sigma_2^2} > 1$   
condition number:  $\kappa = \frac{\sigma_1^2}{\sigma_d^2} > 1$

**Theorem (informal):** For  $L$  large enough, for fixed  $f, x_0$  :

$$f(x_L(\eta_L^*, x_0)) - f(x^*) \leq \exp\left(-L\left(2 \log(1 + 1/c) + \frac{1}{\kappa}\right)\right) (f(x_0) - f(x^*))$$

Generalizes to **expectations** over  $(f, x_0) \sim D$

# Main Theorem

spectral ratio:  $c = \frac{\sigma_1^2}{\sigma_2^2} > 1$   
condition number:  $\kappa = \frac{\sigma_1^2}{\sigma_d^2} > 1$

**Theorem (informal):** For  $L$  large enough, for fixed  $f, x_0$  :

$$f(x_L(\eta_L^*, x_0)) - f(x^*) \leq \exp\left(-L\left(2 \log(1 + 1/c) + \frac{1}{\kappa}\right)\right) (f(x_0) - f(x^*))$$

Generalizes to **expectations** over  $(f, x_0) \sim D$

Efficient to learn  $\eta^*$ :

**constant pseudo-dim**

**binary search ERM** to find  $\eta^*$

# Key Bound – Single Problem Instance

$$\left( \frac{Z_{\alpha}^{\frac{1}{2L-1}} c^{\frac{2L+1}{2L-1}}}{1 + Z_{\alpha}^{\frac{1}{2L-1}} c^{\frac{2L+1}{2L-1}}} \right)^{2L-1} \leq \frac{f(x_L(\eta^*, x_0))}{f(x_L(1/\sigma_1^2, x_0))} \leq \left( \frac{Z_{\beta}^{\frac{1}{2L-1}} c^{\frac{2L+1}{2L-1}}}{1 + Z_{\beta}^{\frac{1}{2L-1}} c^{\frac{2L+1}{2L-1}}} \right)^{2L-1}$$

$$Z_{\alpha} = \frac{z_1^2}{\sum_{i>1} z_i^2}; \quad Z_{\beta} = \frac{z_1^2}{z_2^2} \quad \text{Spectral ratio } c = \frac{\sigma_1^2}{\sigma_2^2} > 1$$