



# Practical Optimization

# A Slight Change

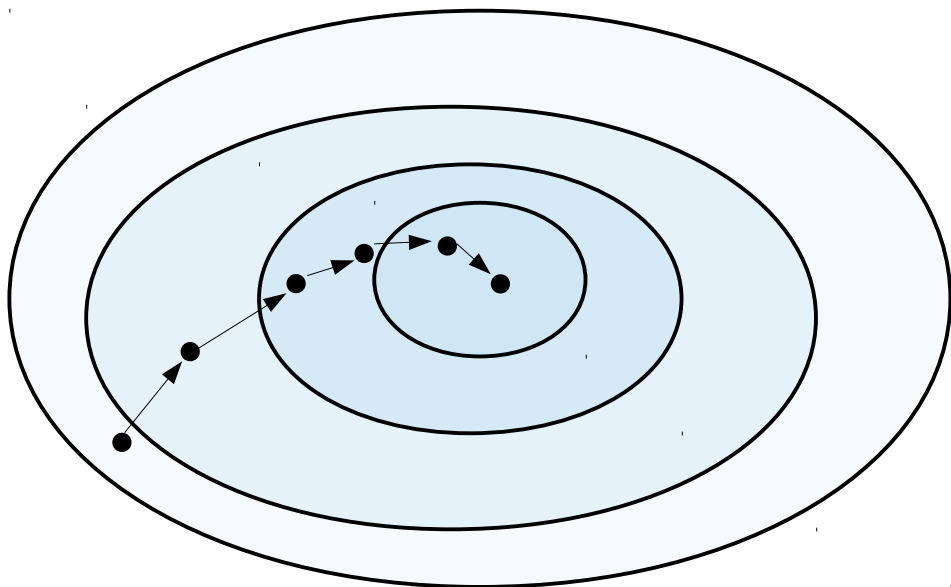
$$D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$$

$$L(\theta) = \sum_i \ell(\theta; \mathbf{x}_i, y_i)$$

$$\begin{aligned} L(\theta) &= \frac{1}{|D|} \sum_i \ell(\theta; \mathbf{x}_i, y_i) \\ &= \mathbb{E}_{\mathbf{x}, y \sim D} [\ell(\theta; \mathbf{x}, y)] \end{aligned}$$

# Gradient Descent

How do we find  $\operatorname{argmin}_{\theta} L(\theta)$



Slow!

Initialize random  $\theta$

N iterations:

Compute  $\nabla_{\theta} L(\theta)$

$\theta \leftarrow \theta - \epsilon \nabla_{\theta} L(\theta)$

$\theta_0 \leftarrow \theta$

For  $\mathbf{x}, y \in D$

$\theta \leftarrow \theta - \frac{\epsilon}{|D|} \nabla_{\theta} \ell(\theta_0; f(\mathbf{x}), y)$

# Stochastic Gradient Descent

## Gradient Descent:

Initialize random  $\theta$

N iterations:

$$\cancel{\theta_0 \leftarrow \theta}$$

For  $x, y \in D$

$$\theta \leftarrow \theta - \frac{\epsilon}{|D|} \nabla_{\theta} \ell(\theta; f(x), y)$$

## Stochastic Gradient Descent:

Initialize random  $\theta$

N iterations:

$$\begin{array}{l} \text{For } x, y \in D \quad \text{Iteration} \\ \theta \leftarrow \theta - \frac{\epsilon}{|D|} \nabla_{\theta} \ell(\theta; f(x), y) \\ \text{Epoch} \end{array}$$

Noisy but much faster than standard GD

# SGD – Variance

$$\mathbb{E}_{\mathbf{x}, y \sim D} [\nabla_{\theta} \ell(\theta; \mathbf{x}, y)] = \nabla_{\theta} L(\theta)$$

$$\nabla_{\theta} \ell(\theta; \mathbf{x}_i, y_i) \neq \nabla_{\theta} L(\theta)$$

$$\begin{aligned} \text{Var} [\nabla_{\theta} \ell(\theta; \mathbf{x}, y)] &= \mathbb{E}_{\mathbf{x}, y \sim D} [\|\nabla_{\theta} \ell(\theta; \mathbf{x}, y) - \nabla_{\theta} L(\theta)\|^2] \\ &= \mathbb{E}_{\mathbf{x}, y \sim D} [\|\nabla_{\theta} \ell(\theta; \mathbf{x}, y)\|^2] - \|\nabla_{\theta} L(\theta)\|^2 \end{aligned}$$

# Mini-batches

## Stochastic Gradient Descent:

Initialize random  $\theta$

N iterations:

For  $x, y \in D$

$$\theta \leftarrow \theta - \frac{\epsilon}{|D|} \nabla_{\theta} \ell(\theta; f(x), y)$$

## Minibatch Gradient Descent:

Initialize random  $\theta$

N iterations:

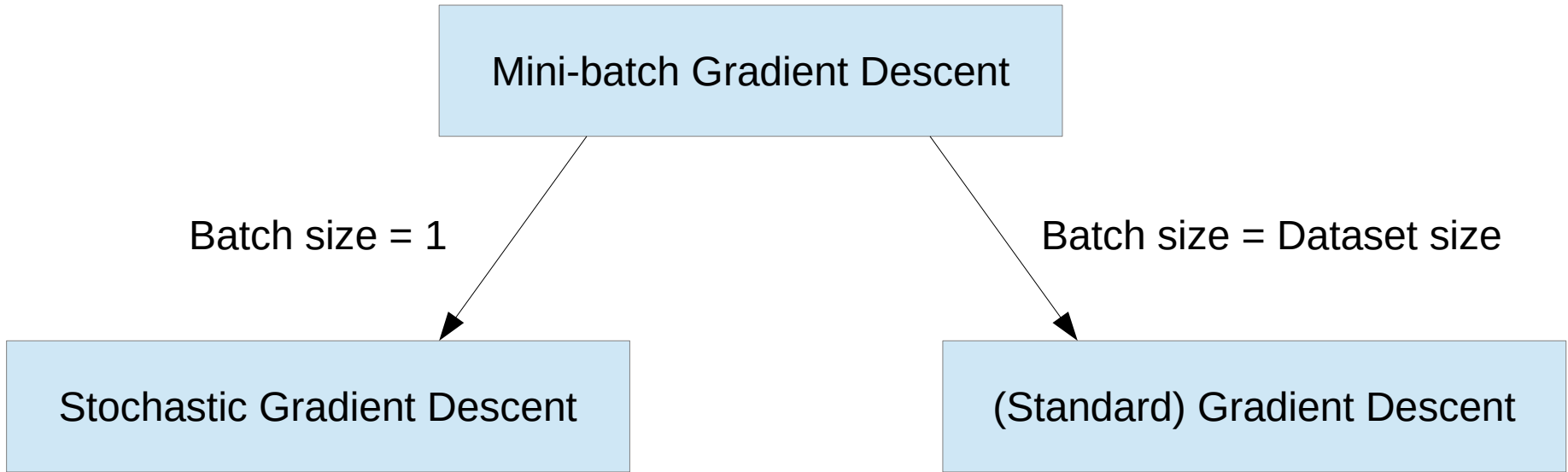
Partition  $D$  into  $D_1, \dots, D_k$

For  $1 \leq i \leq k$

$$\theta \leftarrow \theta - \epsilon \mathbb{E}_{x, y \sim D_i} [\nabla_{\theta} \ell(\theta; f(x), y)]$$

$$\text{Batch size} = |D_1| = |D_2| = \dots = |D_k|$$

# Gradient Descent Algorithms



# Mini-batch Variance

## Stochastic Gradient Descent

$$\text{Var}[\nabla_{\theta} \ell(\theta; \mathbf{x}, y)] = \mathbb{E}_{\mathbf{x}, y \sim D} [\|\nabla_{\theta} \ell(\theta; \mathbf{x}, y)\|^2] - \|\nabla_{\theta} L(\theta)\|^2$$

## Minibatch Gradient Descent

$$\text{Var}[\nabla_{\theta} \ell(\theta; \mathbf{x}, y)] = \mathbb{E}_{D_i} [\mathbb{E}_{\mathbf{x}, y \sim D_i} [\|\nabla_{\theta} \ell(\theta; \mathbf{x}, y)\|^2]] - \|\nabla_{\theta} L(\theta)\|^2$$



Smaller



# Momentum

Initialize random  $\theta$

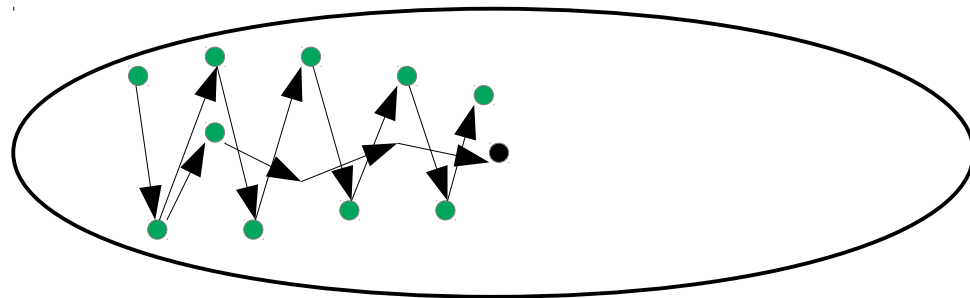
N iterations:

For  $x, y \in D$

~~$$v \leftarrow \nabla_{\theta} \ell(\theta; f(x), y)$$~~

$$v \leftarrow \rho v + \nabla_{\theta} \ell(\theta; f(x), y)$$

$$\theta \leftarrow \theta - \epsilon v$$



$$\rho = 0.9$$