# Session 13: Stochastic gradient descent

Optimization and Computational Linear Algebra for Data Science

Léo Miolane

# Contents

# Stochastic gradient descent

# Setting

In machine learning, one often has to minimize functions of the form

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

where $f_i : \mathbb{R}^n \to \mathbb{R}$.

# Setting

# Stochastic gradient descent

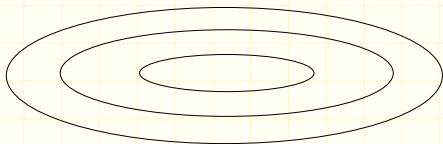$$f(x) = \frac{1}{N} \sum_{i=1}^{N} f_i(x).$$

Starting at some $x_0 \in \mathbb{R}^n$, perform the updates:
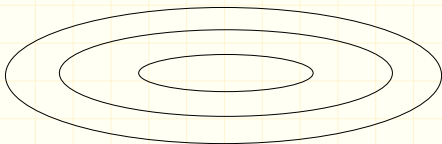
Pick $i$ uniformly at random in $\{1, \ldots, N\}$,

Update $x_{t+1} = x_t - \alpha_t \nabla f_i(x_t)$,

# Tradeoffs in SGD

**Rapidly decaying step sizes**
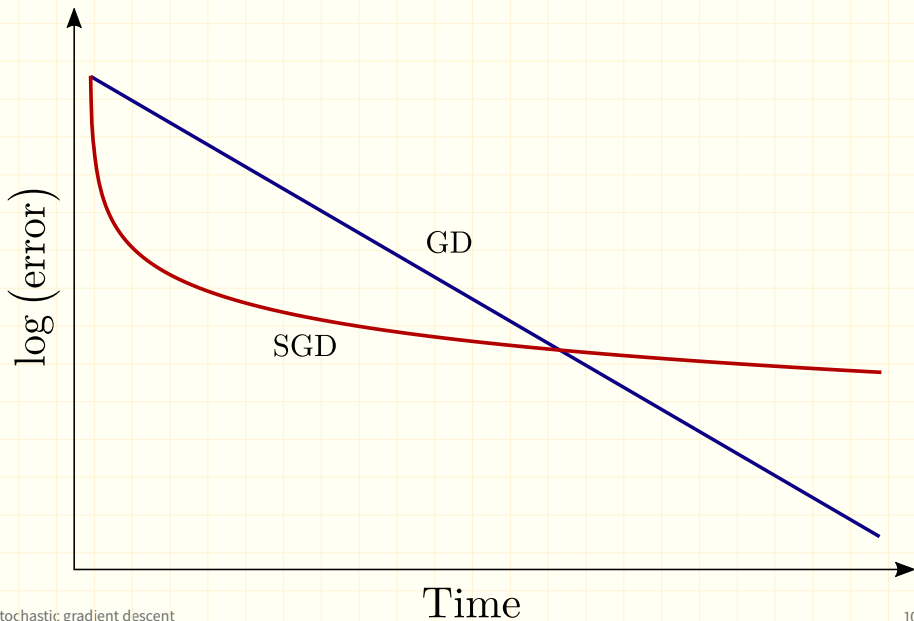
**Slowly decaying step sizes**

# Convergence analysis

# GD vs SGD

**Gradient descent**

**Stochastic gradient descent**

# GD vs SGD

# Questions?

# Questions?