

**Erasmus+ lecture on  
Machine Learning,  
November 13, 2023**

Simple example

$$f(x) = \frac{1}{2} x^T A x - x^T b$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

$$\begin{aligned} f(x_1, x_2) &= \frac{1}{2} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 2 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &\quad - \begin{bmatrix} 5 \\ 3 \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \end{aligned}$$

$$= x_1^2 + x_1 x_2 + 10 x_2^2 - 5x_1 - 3x_2$$

$$\begin{cases} \frac{\partial f}{\partial x_1} = 0 = 2x_1 + x_2 - 5 \\ \frac{\partial f}{\partial x_2} = 0 = x_1 + 20x_2 - 3 \end{cases}$$

$$\begin{aligned} x_1 &= \frac{97/39}{22.44} \\ x_2 &= \frac{1/39}{20.026} \end{aligned}$$

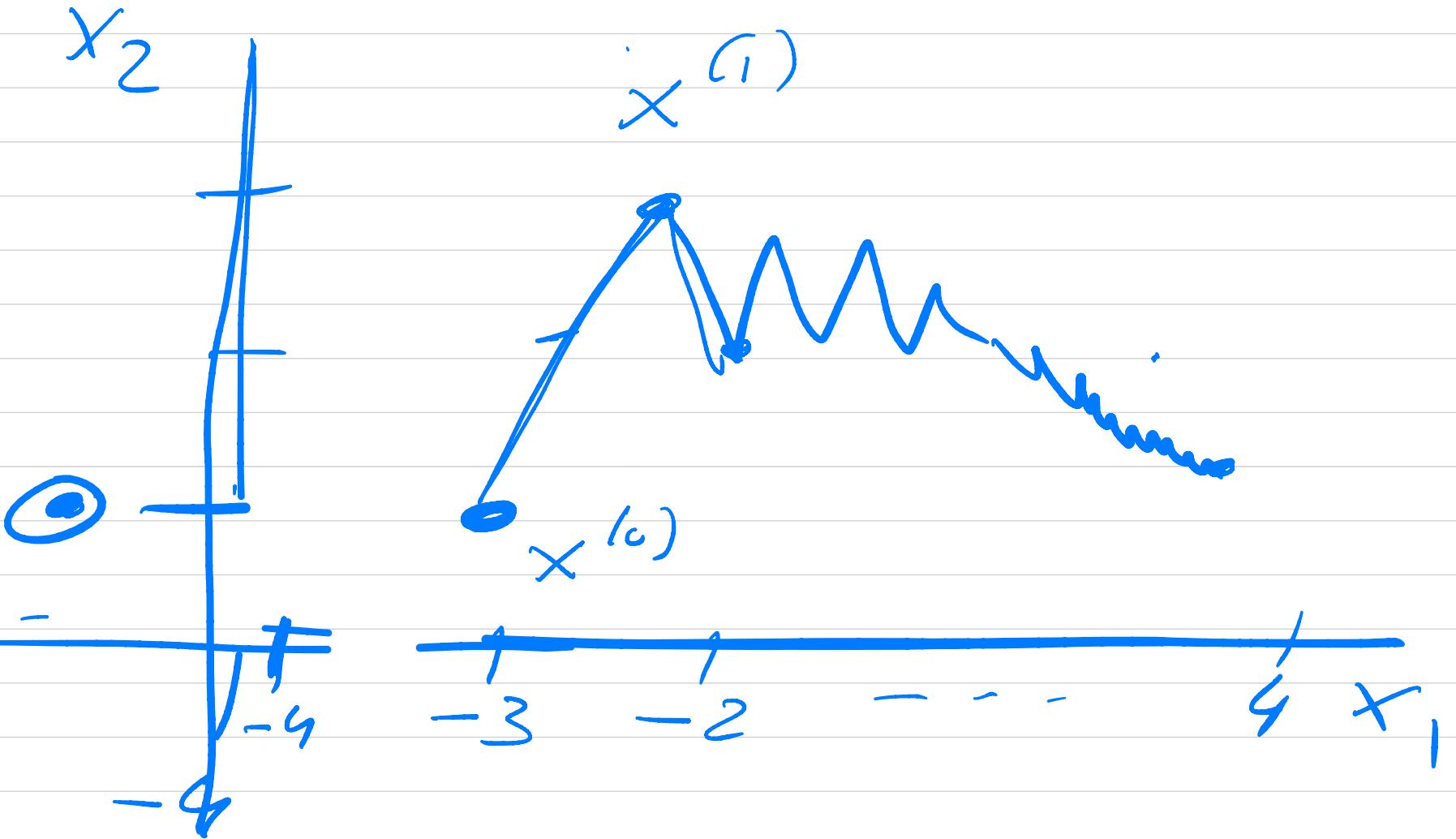
$$f(x_0) \geq f(x_1) \geq f(x_2) \dots \geq f(x^*)$$

$$x^{(0)} = \begin{bmatrix} -3 \\ 1 \end{bmatrix}$$

$$x_i^{(1)} = x_i^{(0)} - \gamma_0 Df(x_i^{(0)})$$

$$i = \{1, 2\}$$

$$x_i^{(k+1)} = x_i^{(k)} - \gamma^{(k)} Df(x_i^{(k)})$$



Gradient with momentum

$$m \frac{d^2x}{dt^2} + \mu \frac{dx}{dt} = -\nabla Q)$$

$$\frac{d^2x}{dt^2} \approx \frac{x(t+\Delta t) + x(t-\Delta t) - 2x(t)}{(\Delta t)^2}$$
$$= \frac{x_{t+\Delta t} + x_{t-\Delta t} - 2x_t}{(\Delta t)^2}$$
$$(x(t+\Delta t), x(t-\Delta t))$$

$$\frac{dx}{dt} \approx \frac{x_{t+\Delta t} - x_t}{\Delta t}$$

$$\frac{m(x_{t+\Delta t} + x_{t-\Delta t} - 2x_t)}{(\Delta t)^2} + \mu \frac{(x_{t+\Delta t} - x_t)}{\Delta t} = -DV(x)$$

Define  $\Delta x_{t+\Delta t} = x_{t+\Delta t} - x_t$

$$\Delta x_t = x_t - x_{t-\Delta t}$$

$$\Delta x_{t+\Delta t} = -\frac{(\Delta t)^2}{m + \mu \Delta t} DV(x)$$

$$S + \frac{m}{m + \mu \Delta t} \Delta x_t$$

$$\Delta x_{t+\Delta t} = -\nabla V(x_t) + \underbrace{\delta \Delta x_t}_{x_t - x_{t-\Delta t}}$$

$$\left\{ \beta^{(i+1)} = \beta^{(i)} - \gamma DC(\beta^{(i)}) \right\}$$

$$x_t \Rightarrow \beta^{(i)} \quad x_{t+\Delta t} \rightarrow \beta^{(i+1)}$$

$$DV \Rightarrow DC(\beta^{(i)}) = g(\beta^{(i)}) = g$$

$$\beta^{(i+1)} = \beta^{(i)} - \gamma g + \underbrace{\delta(\beta^{(i)} - \beta^{(i-1)})}_{\text{hyperparameter}}$$

$$\delta \in [0, 1]$$

## algorithm

regime learning  $\eta(t)$

$\rightarrow$  momentum parameter  
 $\delta$

require initial  $\beta_0, \nu_0$

while stopping criterion not reached

- compute  $g$

- compute velocity update

$$\nu^{(i)} = \delta(\beta^{(i)} - \beta^{(i-1)})$$

-  $\eta g$

- apply update

$$\text{end while } \beta^{(i+1)} = \beta^{(i)} + \eta g^{(i)}$$

$\eta$ :

- constant and choose a grid of values  $\eta = [10^{-5}, 10^{-4}, 10^{-3}, \dots]$

- exponential decay

$$\eta_k = \eta_0 \exp(-K\eta_{k-1})$$

K could be number of iterations

- linear

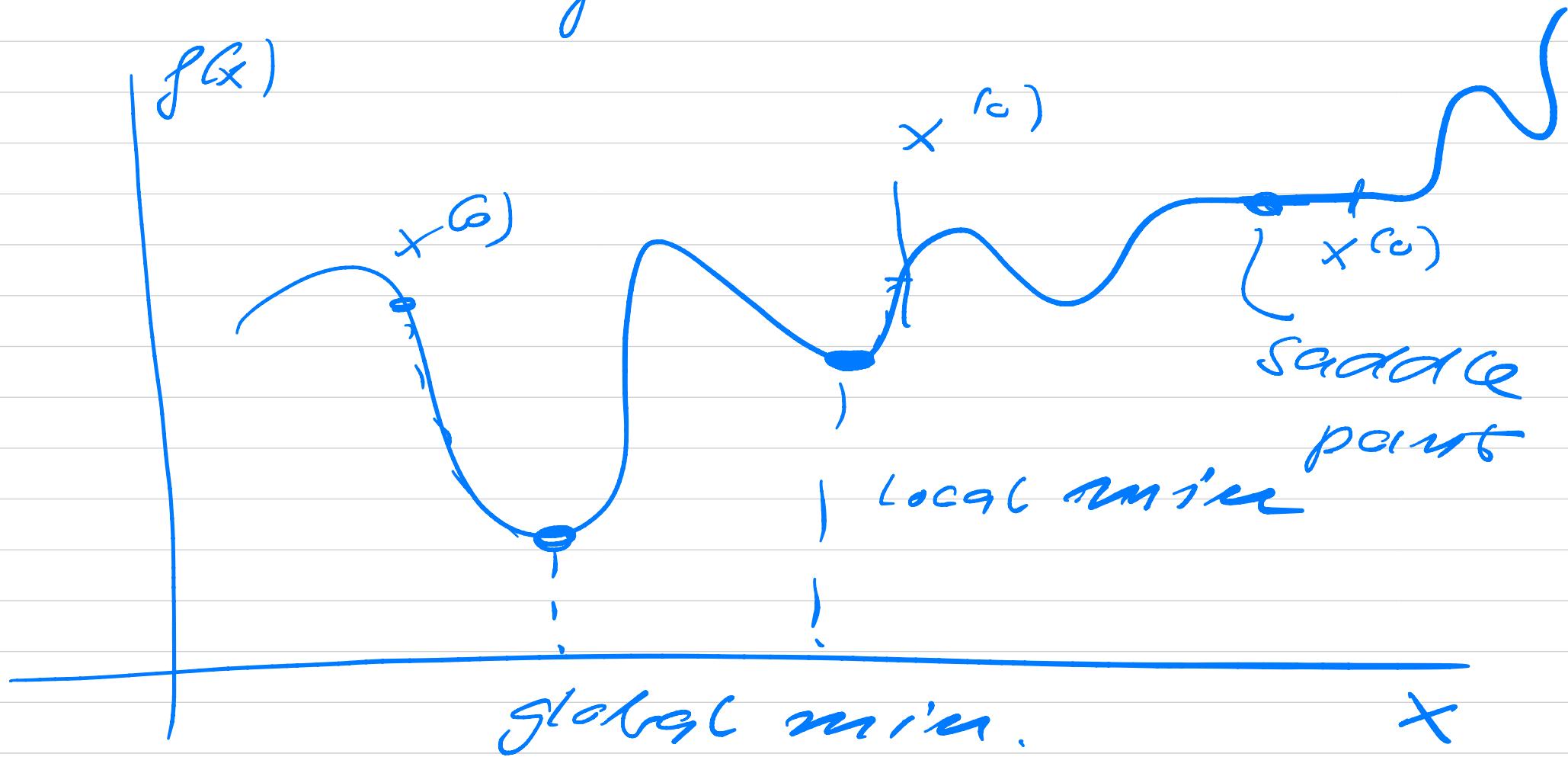
$$\eta_k = (1-\alpha)\eta_0 + \alpha\eta_n$$

$$\alpha = K/\text{const}$$

$$\eta_n \approx \frac{1}{100} \eta_0$$

Adagrad, RMSprop, Adam

stochastic gradient descent.



Split data in minibatches

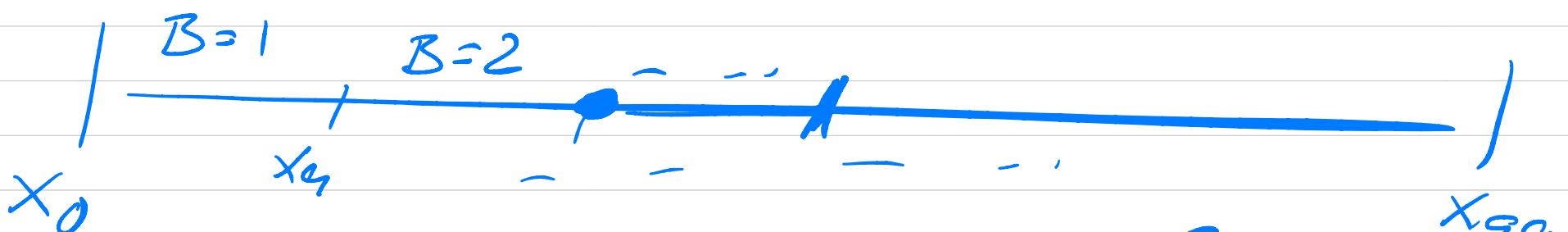
$$n = 100$$

M = 10 minibatches

every minibatch has

$$m = n/M \text{ points (10 points)}$$

$$\nabla_{Bm} = \frac{1}{m} \sum_{i=1}^m \nabla C(x_i)$$



repeat a given number of times = # epochs.

$$n = 1000$$

$$\text{OLS} : g \propto \underbrace{x^T x \beta - x^T y}_{}^{\top}$$

$$x \in \mathbb{R}^{m \times n}$$

$$\text{FLOPs} \propto n^3 = 10^9$$

Batch has only 10 points

$$x^T x \text{ for a minibatch } (10^2)^3$$

$$\approx 1000 \text{ FLOPs}$$

we repeat this 100 times,

$$10^2 \times 10^3 \approx 10^5 \text{ FLOPs},$$