# Mathematics for AI

## Lecture 28

### Gradient Descent, SGD, Momentum
### Quadratic Approximation, Newton's method

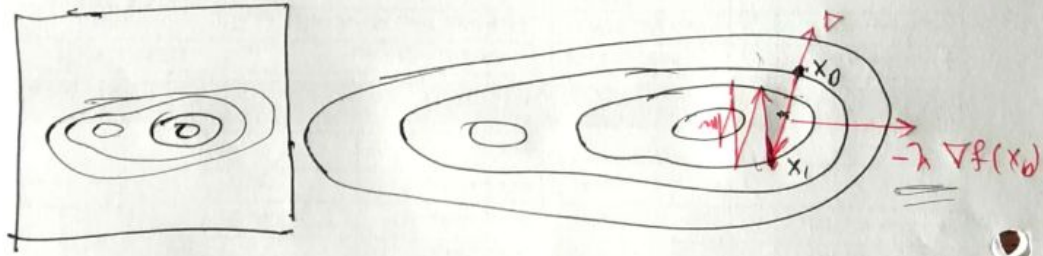# Gradient Descent

# Gradient Descent



we like to have

$f(x_1) < f(x_0)$

line search

$x_0$

$\nabla f(x_0)$

$x_1 \leftarrow x_0 - \lambda \nabla f(x_0)$

step size
learning rate

$-\lambda \nabla f(x_0)$

$x_{t+1} = x_t - \lambda \nabla f(x_t)$   gradient descent

# Stochastic Gradient Descent (SGD)



$$C(\theta) = \sum_{i=1}^{N} \|f(\theta, x_i) - y_i\|^2$$

MA 28 Ⅱ

$$(x_1, y_1), (x_2, y_2), \cdots , (x_n, y_n)$$
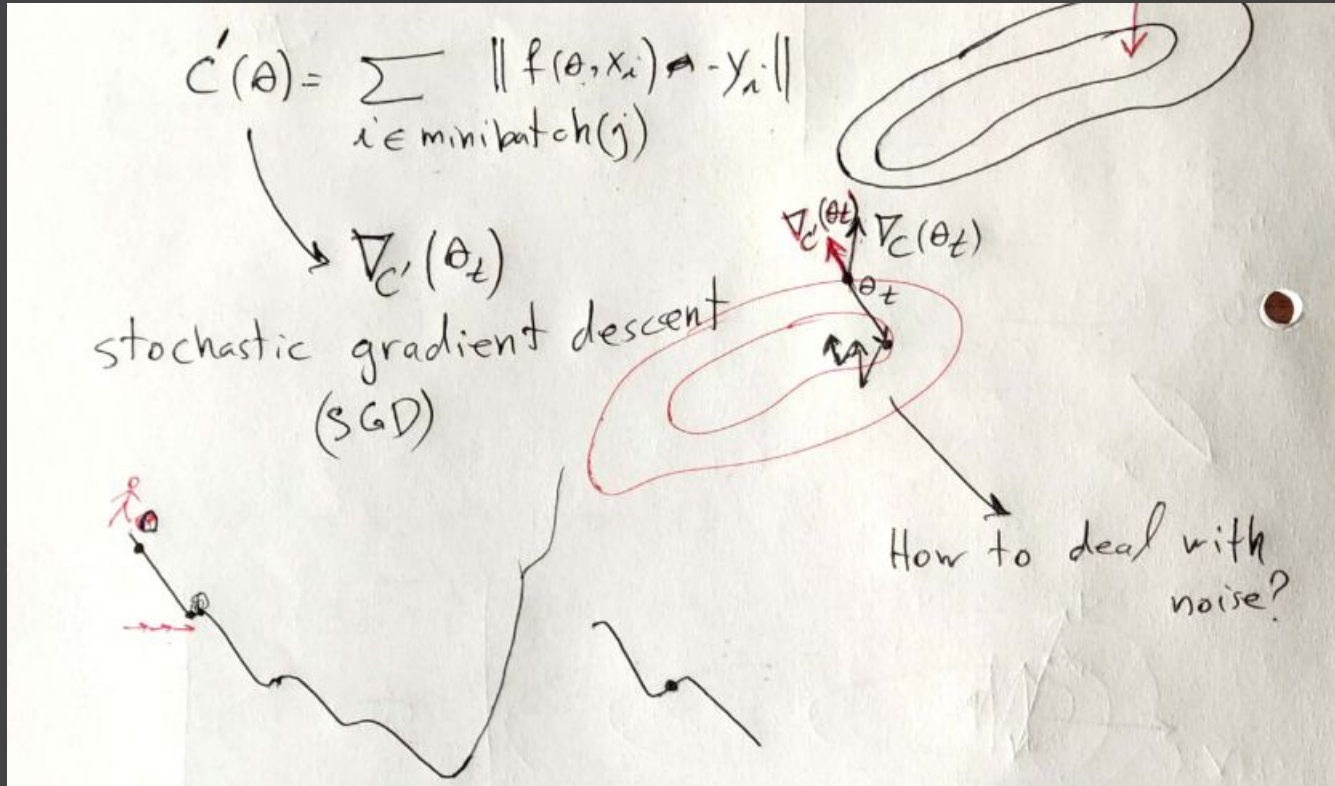
mini-batch

epoch

$$C'(\theta) = \sum_{i \in \text{minibatch}(j)} \|f(\theta, x_i) - y_i\|$$

# Stochastic Gradient Descent (SGD)

# Momentum

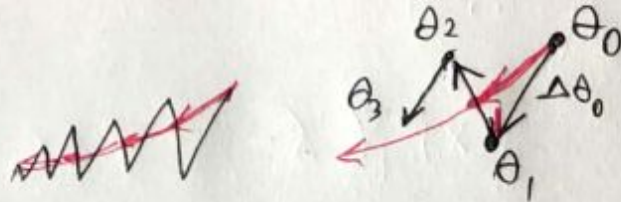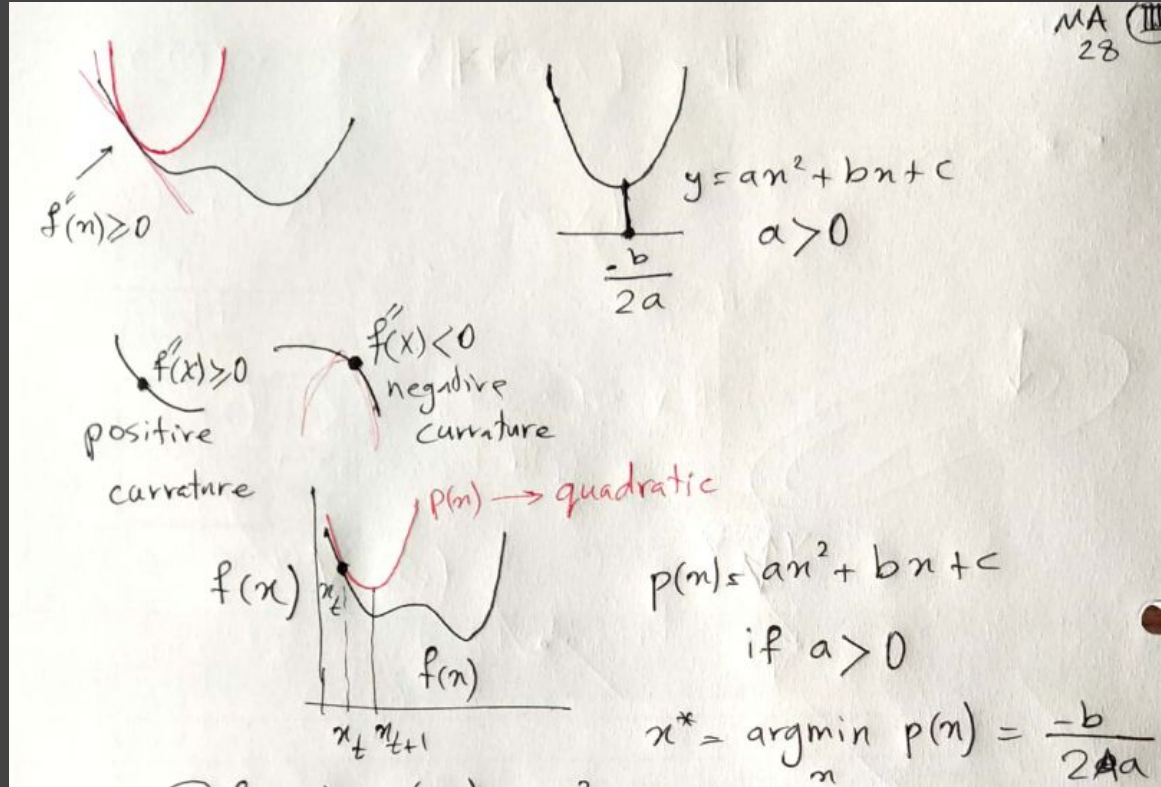$$\theta_{t+1} \leftarrow \theta_t + \Delta\theta_t$$

Gradient descent: $\Delta\theta_t = -\lambda \nabla_c(\theta_t)$

Momentum method: $\Delta\theta_t = -\lambda \nabla_c(\theta_t) + \alpha \Delta\theta_{t-1}$

# Quadratic Approximation

$f''(x) \geq 0$

$y = ax^2 + bx + c$

$a > 0$

$\dfrac{-b}{2a}$

$f''(x) \geq 0$
positive
curvature

$f''(x) < 0$
negative
curvature

$f(x)$ $x_t$

$p(x) \longrightarrow$ quadratic

$f(x)$

$x_t \quad x_{t+1}$

$p(x) = ax^2 + bx + c$

if $a > 0$

$x^* = \underset{x}{\text{argmin }} p(x) = \dfrac{-b}{2a}$

# Quadratic Approximation



carvature

$p(x) \longrightarrow$ quadratic

$f(x)$

$f(x)$

$x_t \quad x_{t+1}$

$p(x) = ax^2 + bx + c$

if $a > 0$

$x^* = \underset{x}{\arg\min} \; p(x) = \dfrac{-b}{2a}$

① $f(x_t) = p(x_t) = ax_t^2 + bx_t + c$

② $f'(x_t) = p'(x_t) = 2ax_t + b$

③ $f''(x_t) = p''(x_t) = 2a \implies a = \dfrac{f''(x_t)}{2}$

② $b = f'(x_t) - 2ax_t = f'(x_t) - f''(x_t)x_t$

$c = f(x_t) - ax_t^2 - bx_t$

# Newton's method



$$\text{(III)} \quad f''(x_t) = p''(x_t) = 2a \implies a = \frac{f''(x_t)}{2}$$

$$\text{(II)} \quad b = f'(x_t) - 2a\, x_t = f'(x_t) - f''(x_t)\, x_t$$

$$c = f(x_t) - a\, x_t^2 - b\, x_t$$

$$x_{t+1} = \frac{-b}{2a} = \frac{-f'(x_t) + f''(x_t)\, x_t}{f''(x_t)}$$

$$x_{t+1} = x_t - \frac{f'(x_t)}{f''(x_t)} \qquad \underline{\text{Newton's method}}$$

روش نیوتن

$$f''(x_t) \; \underset{\text{need to}}{\cancel{\text{most}}} \; \text{be} \quad \text{positive}$$

# Multivariate Quadratic Function

$$P(x) = P\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = 2x_1^2 + x_1 x_2 + 3x_2^2 + 2x_1 - 4x_2 - 10 \quad \text{MA}$$

$$= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 2 & -4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - 10$$

$$= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 2 & -4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - 10$$

$$= \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 2 & 1.5 \\ 1.5 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 2 & -4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - 10$$

symmetric

unique x

$$= x^T A x + \vec{b}^T x + c$$
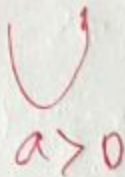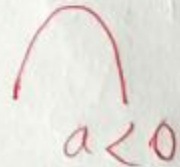
# Multivariate Quadratic Function

Every quadratic function $~~~~~~~~~~ p: \mathbb{R}^2 \longrightarrow \mathbb{R}$

can be uniquely represented as

$$p(x) = x^T A x + b^T x + c$$

where $A \in \mathbb{R}^{n \times n}$ is <u>symmetric</u>, $b \in \mathbb{R}^n$, $c \in \mathbb{R}$

$n = 1$ $~~~~~~~~\cap~~~~~~~~\cup~~~~~~~ a = 0 ~~~~/~~~~ p'(x) = 0$

$~~~~~~~~~~~~~~~~ a < 0 ~~~~~ a > 0 ~~~~~~~~~~~~~~~~~~~~~~ x = -b/2a$

# Stationary point of a Quadratic Function

$$p(x) = x^T A x + b^T x + c$$

$$\nabla_p(x) = 2Ax + b \qquad \nabla_p(x) = 0 \Rightarrow 2Ax = -b$$

$$H(x) = \underline{\underline{2A}} \qquad\qquad x^* = -\frac{A^{-1}b}{2}$$

$$\text{if } A \text{ non singular}$$



$x^*$  $x$

# When is the stationary point the minimum?



$p(x) = \frac{1}{2} x^T A x + b^T x + c$

$\nabla = 0$      $x^* = -\frac{1}{2} A^{-1} b$

$x^*$ is a minimum

directional curvature is positive in all directions.

$\underline{u^T H u > 0}$     for all $u \neq 0$

$\Rightarrow H$ is positive definite

# When is the stationary point the minimum?

$$p(x) = x^T A x + b^T x + c \implies 2A \text{ is positive definite}$$
$$A \text{ is positiv definite}$$

# When do we have a maximum?



$\Rightarrow$ A is negative-definite

# Quadratic Approximation - Multivariate case

$$f(x_t) = P(x_t)$$

$$\nabla f(x_t) = \nabla P = 2Ax_t + b$$

$$H_f(x_t) = \nabla H_p(x_t) = 2A \implies 2A = H_f(x)$$

$$b = \nabla f(x_t) - 2Ax_t = \nabla_f(x_t) - H_f(x_t)\, x_t$$

$H_f(x_t)$ should be positive definite

$$x_{t+1} = -\frac{1}{2}A^{-1}b = -(2A)^{-1}b = -H^{-1}(\nabla - H)x_t = x_t - H^{-1}\nabla$$

# Newton's method - Multivariate case



$$x_{t+1} = x_t - H_f^{-1}(x_t) \nabla_f(x_t)$$

$$\underbrace{H_f^{-1}(x_t)}_{\mathbb{R}^{n \times n}} \quad \underbrace{\nabla_f(x_t)}_{\in \mathbb{R}^n}$$

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

$H_f$ positive-definite

# Quadratic Approximation - Taylor series perspective