

Lecture FYS-  
Stk3155/4155,  
September 30,  
2024

Basic algo : GD, and GD with momentum

GD with momentum

Define learning rate  $\gamma$   
initialize  $\beta$ -parameters

Set # of iterations, set  $S_{[t_0, t]}$   
for # iterations and while  
not reached stop criterion

$$\beta_{(t)} = \beta_{(t-1)} - \gamma g(\beta_{(t)})$$

(no momentum)  $- \gamma g(\beta_{(t)})$

$$\beta_{(i+1)} = \beta_{(i)} + \delta_i^{-1}$$

stop if  $|\beta_{(i+1)} - \beta_{(i)}| \leq \epsilon$

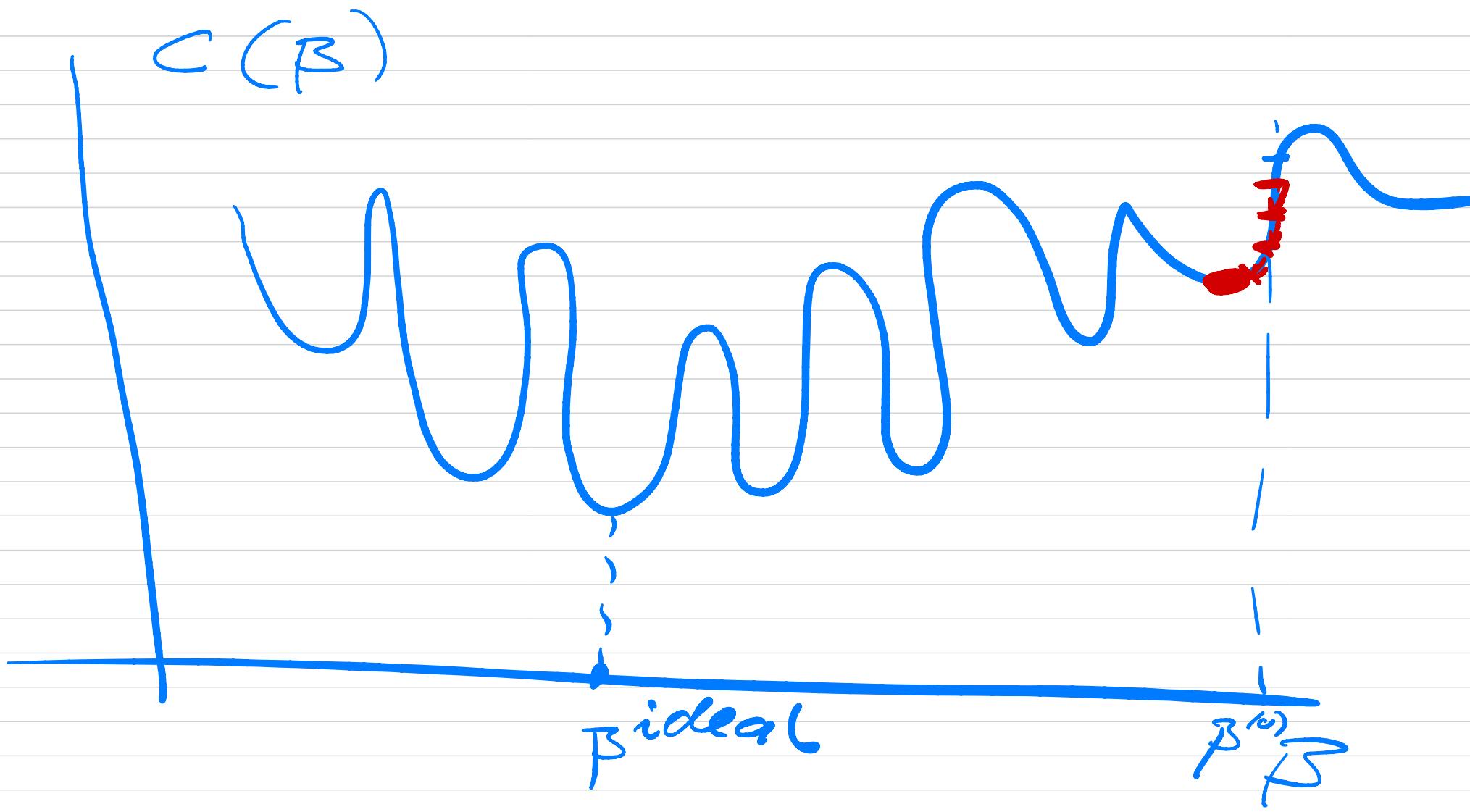
end while.

$10^{-5}$

AND-GATE



$$\begin{matrix} & x & y & z \\ \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} & \times & \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} & \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \end{matrix}$$



## Updating learning rates

- constant  $\gamma$
- exponential decay

$$\gamma^{(k)} = \gamma^{(0)} \exp(-\kappa \gamma_p)$$

$\gamma_p$  is a parameter

- linear

↓ initial value

$$\gamma^{(k)} = (1 - \alpha) \gamma^{(0)} + \alpha \gamma_p$$

$\gamma_p$  is a parameter

$$\gamma_p \approx \frac{1}{100} \gamma^{(0)}$$

select an initial  $\gamma^{(0)}$

$$\gamma^{(0)} = [10^{-5}, 10^{-4}, \dots, 10^{-1}]$$

- ADAgrad

- RMS prop

- ADAM

## ADAGRAD

require learning rate ( $\eta$ ) &  $\beta^k$   
initial guess  $\beta^{(0)}$

constant  $S$  (small) for  
numerical stability.

while stopping criterion

not met

(with SGD or without)

- compute gradient  $g(\beta^{(i)})$

- accumulate squared  
gradient

$$r^{(i+1)} = r^{(i)} + g \cdot g(\beta^{(i)})$$

compute update of

$$\gamma^{(i)} = \frac{\gamma^{(0)}}{\delta + \sqrt{\eta^{(i)}}}$$

update

$$\beta^{(i+1)} = \beta^{(i)} - \gamma^{(i)} g(\beta^{(i)})$$

end while

# Automatic differentiation

Example

$$f(x) = \exp(x^2)$$

$$= b = \exp(a)$$

$$a = x^2$$



$$\frac{df}{dx} = \underbrace{\frac{df}{db}}_1 \frac{db}{da} \frac{da}{dx}$$
$$\exp(a) \quad 2x$$

aside

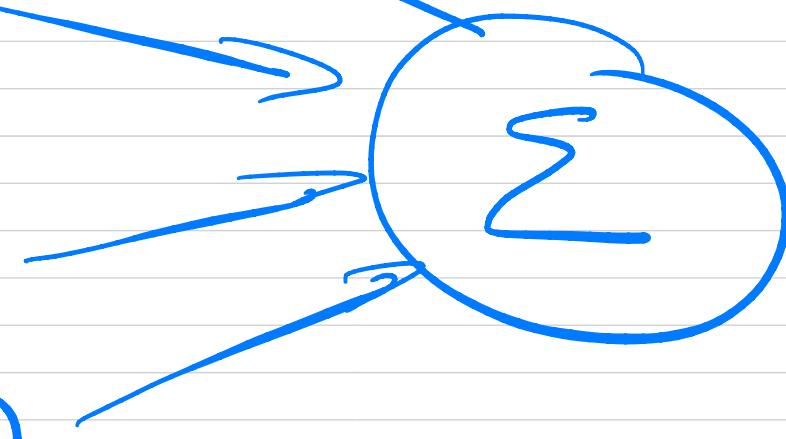
$(x_1, y_1)$

$(x_2, y_2)$

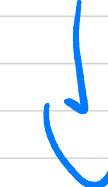
$(x_3, y_3)$

:

$(x_n, y_n)$



output  
 $z$



para  
metra

$C(z_1, t, \epsilon)$

observation

$$\frac{df}{dx} = \begin{bmatrix} \frac{df}{dr} & \frac{df}{da} \end{bmatrix} \frac{da}{dx}$$

reverse mode

or

$$\frac{df}{dr} \begin{bmatrix} \frac{dr}{da} & \frac{da}{dx} \end{bmatrix}$$

Forward mode