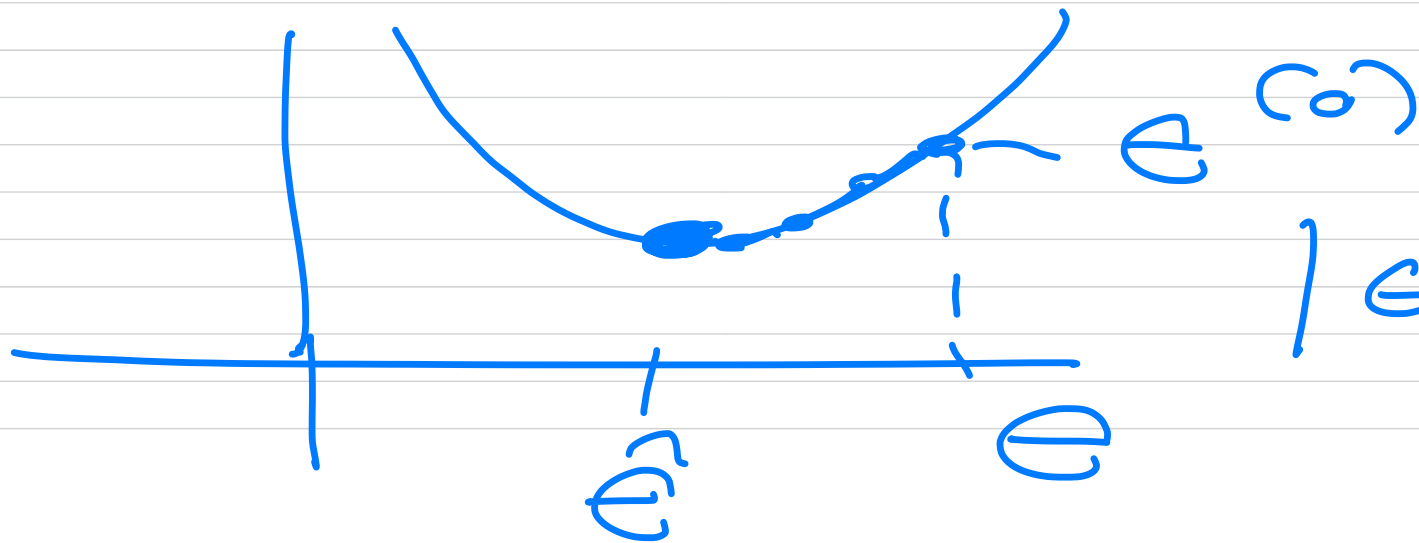# FYS-STK3155/4155 week 37, September 8-12, 2025

Gradient descent methods
- Plain GD
- Momentum GD
- other simple updates of learning rates
- ADAgrad, RMSprop, ADAM
- Stochastic GD
- Examples

what did we obtain last week?

$$\theta^{(m+1)} = \theta^{(m)} - \left(H(\theta^{(m)})\right)^{-1} \nabla_\theta C(\theta^{(m)})$$

– start with a guess $\theta^{(0)}$



$$|\theta^{(m+1)} - \theta^{(m)}| \leq \varepsilon \sim 10^{-5}$$

$$\left( H^{(n)} \right) = \frac{\partial^2 C \left( \theta^{(n)} \right)}{\partial^2 \theta} \implies$$

$$\eta^{(n)}$$

learning rate

$$\theta^{(n+1)} = \theta^{(n)} - \eta \nabla_\theta C \left( \theta^{(n)} \right)$$

plain/simple gradient descent (GD)

# Gradients

## OLS

$$\nabla_\theta C = \frac{2}{n}\left(X^T X \theta - X^T y\right)$$

$$X \in \mathbb{R}^{n \times p} \qquad \theta \in \mathbb{R}^p$$

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

## Ridge

$$\nabla_\theta C = \frac{2}{n}\left(X^T X \theta - X^T y\right) + \lambda \cdot 2\theta$$

# LASSO

$$\frac{2}{n} \left( X^T X \theta - X^T y \right)$$
$$+ \lambda \, sgn(\theta)$$

$$\theta^{(n+1)} = \boxed{\theta^{(n)} - \eta g^{(n)}}$$

$$\nabla_\theta C(\theta^{(n)})$$

Taylor-expand around

keep only terms to 2nd derivative

$$C(\theta^{(n+1)}) = C(\hat{\theta})$$

$$= C(\theta^{(n)}) + g^{T^{(n)}}(\theta^{(n)} - \eta g^{(n)})$$

$$+ \frac{1}{2}(\theta^{(n)} - \eta g^{(n)})^T H^{(n)}$$

$$\times (\theta^{(n)} - \eta g^{(n)})$$

optimal $\eta$ ?

$$\frac{dC}{d\eta} = 0 \qquad \Longrightarrow$$

$$-g^{T(n)} g^{(n)} + \eta \, g^{T(n)} H^{(n)} g^{(n)}$$

$$= 0 \qquad \Longrightarrow$$

$$\eta^{(n)} = \frac{g^{T(n)} g^{(n)}}{\underbrace{g^{T(n)} H^{(n)} g^{(n)}}}$$

$$\eta^{(n)} = \frac{g^T g}{\lambda g^T g} = \frac{1}{\lambda} \qquad H^{(n)} g^{(n)} = \lambda g^{(n)}$$

$\eta$ - requirement

$$\eta < \frac{2}{\lambda_{max}}$$

Largest eigenvalue

of $H^{(n)}$

$\eta \nabla_\theta C$

$\theta^{(0)}$

$\mathcal{M}$ too small

$\mathcal{M}$ too Large

$\hat{\theta}$

$\theta$

$\nabla_\theta C = 0$

GD with momentum

Newton's eq. of motion

$$m \frac{d^2 x}{dt^2} + \underbrace{\mu \frac{dx}{dt}}_{Friction} = -\vec{\nabla} V(x)$$

Discretize

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

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

Define

$$\Delta X_{t+\Delta t} = X_{t+\Delta t} - X_t$$

$$\Delta X_t = X_t - X_{t-\Delta t}$$

$$\frac{m}{\Delta t^2} \Delta X_{t+\Delta t} - \frac{m}{\Delta t^2} \Delta X_t$$

$$+ \mu \frac{\Delta X_{t+\Delta t}}{\Delta t} = - \vec{\nabla} V(x)$$

$$\Delta X_{t+\Delta t} = -\vec{g} \frac{\Delta t^2}{m + \mu \Delta t} + \frac{m \Delta X_t}{m + \mu \Delta t}$$

$$\lim_{\mu \to 0} \delta = 1 \quad \wedge \quad \lim_{\mu \to 8} \delta = 0$$

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

$$\Delta x_{t+\Delta t} = -\eta \vec{g} + \delta \Delta x_t$$

$$x_t \longrightarrow \theta^{(m)}$$

$$x_{t+\Delta t} \Longrightarrow \theta^{(m+1)}$$

$$x_{t-\Delta t} \longrightarrow \theta^{(m-1)}$$

$$e^{(n+1)} = e^{(n)} - \eta \, g\left(e^{(n)}\right)$$

$$+ \, \underset{\uparrow}{s} \left[ e^{(n)} - e^{(n-1)} \right]$$

momentum param

( memory )

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

algorithm:

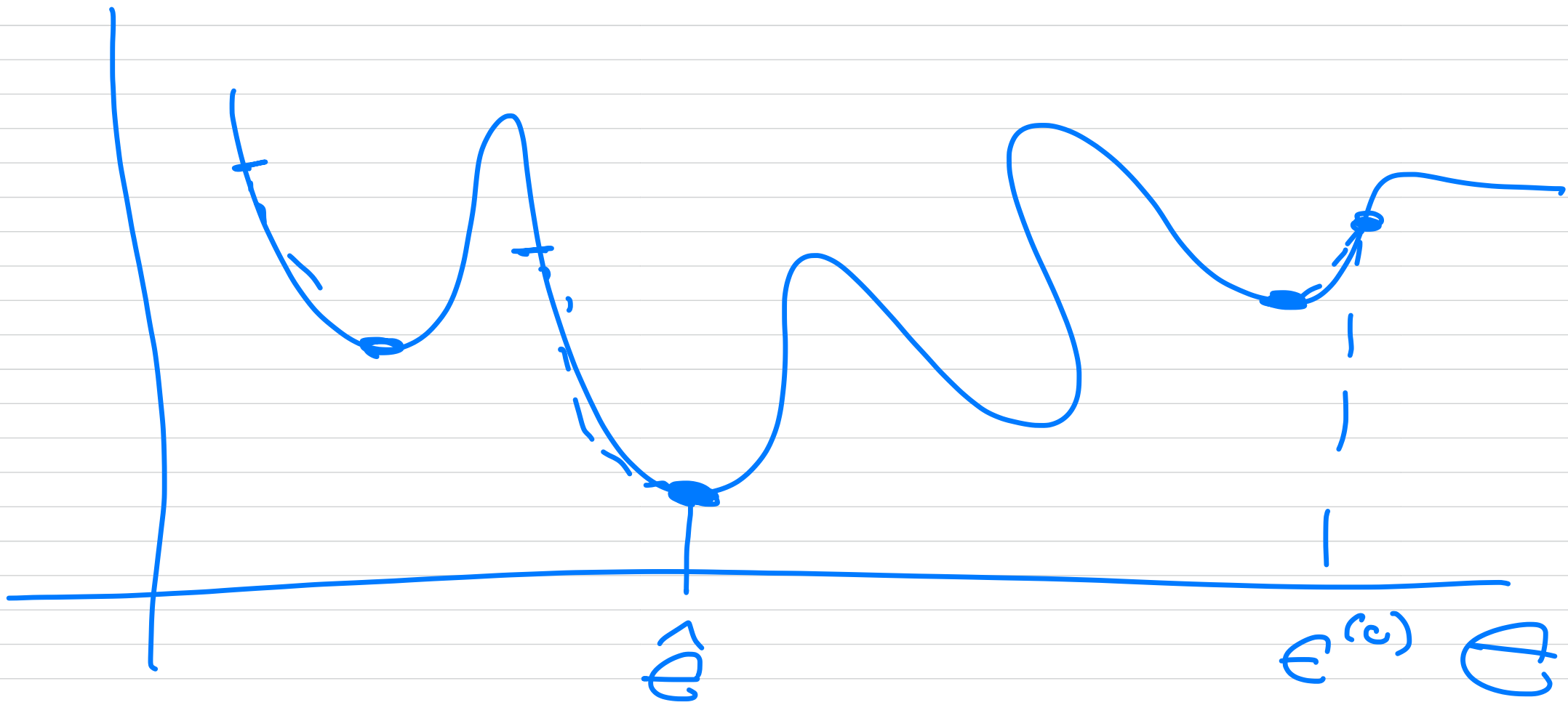fix initial guess $\theta^{(0)}$

fix $\eta^{(0)}$

fix momentum $\delta$

initialize vector $v^{(0)}$

while stopping criterion not met $\quad$ $\textcolor{red}{\theta^{(m)} + \alpha \theta^{(m+1)}}$

$\qquad v^{(m)} = \delta (\theta^{(m)} - \theta^{(m-1)})$

$\qquad \qquad - \eta \, g(\theta^{(m)})$

$\qquad \theta^{(m+1)} = \theta^{(m)} + v^{(m)}$

end while

$$\hat{e} \qquad\qquad\qquad e^{(c)} \quad e$$

cheap ways to update $\eta$

— $\eta$ constant

— exponential decay

$$\eta^{(k)} = \eta^{(0)} \exp(-k \gamma_\eta)$$

$$\gamma_\eta \sim \eta^{(0)}/100 \quad \text{or}$$

similar.

– linear

$$\eta^{(k)} = (1 - \alpha) \eta^{(0)} + \alpha \tilde{\eta}$$

$$\delta \tilde{\eta} \sim \eta^{(0)}/100$$

$$\alpha \in [c, 1]$$

parameters