

Lecture FYS-
STK3155/4155,
September 28, 2023

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^P}{\arg \min} - \underbrace{\log P(D|\beta)}_{c(\beta)}$$

$$\frac{\partial c}{\partial \beta} \Big| = g$$

$$\frac{\partial^2 c}{\partial \beta \partial \beta^T} = H$$

Perform Taylor expansion
around $\hat{\beta}$
 $\hat{\beta} - \beta^{(0)}$

$$C(\hat{\beta}) = C(\beta^{(n)}) + (g^{(n)})^T (\hat{\beta} - \beta^{(n)}) \\ + \frac{1}{2} (\hat{\beta} - \beta^{(n)})^T H^{(n)} (\hat{\beta} - \beta^{(n)}) \\ + \dots$$

$$b = \hat{\beta} - \beta^{(n)}$$

$$C(\hat{\beta}) \approx C(\beta^{(n)}) + (g^{(n)})^T b \\ + \frac{1}{2} b^T H^{(n)} b$$

$$\frac{\partial C}{\partial b} = 0 = g^{(n)} + H^{(n)} \cdot b$$

$$\hat{\beta} = \beta^{(m+1)} = \beta^{(m)} -$$

$$[H(\beta^{(m)})]^{-1} g^{(m)}(\beta^{(m)})$$

in more general terms

$$C(p) \Rightarrow f(x) = C + g^T x + \frac{1}{2} x^T A x$$

$$\frac{\partial f}{\partial x} = 0 \quad Ax + g \Rightarrow \\ x = -A^{-1}g$$

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

↑
learning rate

Taylor expand $C(\hat{\beta})$ again

$$C(\beta^{(n)} - \gamma g^{(n)}) \approx C(\beta^{(n)})$$

$$- \gamma (g^{(n)})^T g^{(n)} +$$

$$\frac{1}{2} \gamma^2 (g^{(n)})^T H^{(n)} g^{(n)} + \dots$$

Take derivative wrt γ

$$(g^{(n)} \rightarrow g; H^{(n)} \rightarrow H)$$

$$\gamma = \frac{g^T g}{g^T H g}$$

$$Hg = \lambda g$$

$$\gamma = \frac{g^T g}{\lambda g^T g} = \frac{1}{\lambda}$$

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

Adagrad } \rightarrow Momentum

RMSprop }

ADAM

Plain gradient descent

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

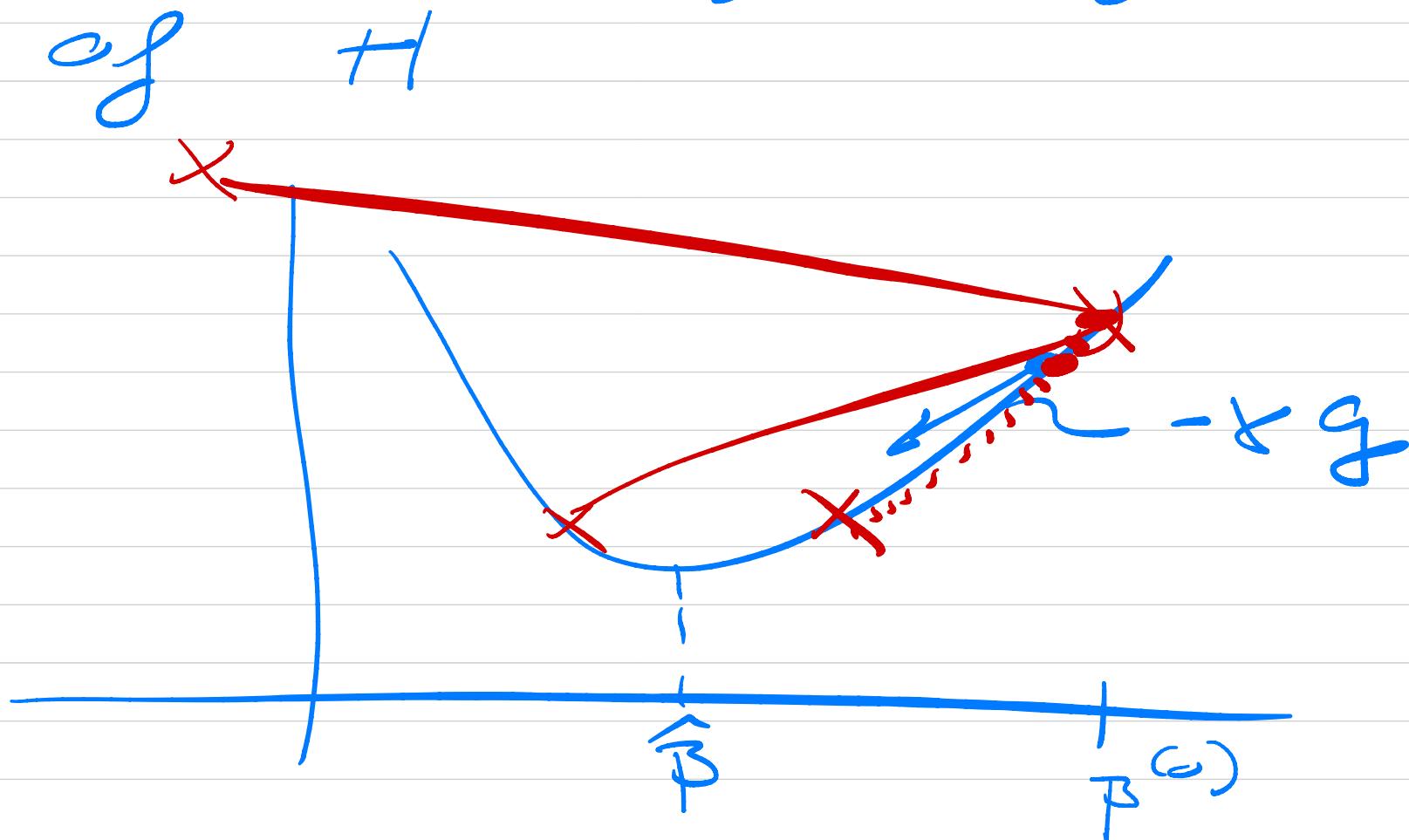
Convergence test

$$\|\beta^{(n+1)} - \beta^{(n)}\|_1 \leq \Sigma \sim 10^{-8}$$

Explicit convergence criterion

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

λ_{\max} largest eigenvalue
of H



Gradient descent with
momentum term.

Newton's eq of motion

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

discretize w.r.t $t, \Delta t$

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

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

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

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

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

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

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

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

$$\delta = \frac{m}{m + \mu \Delta t} \quad \gamma = \frac{(\Delta t)^2}{m + \mu \Delta t}$$

$$\Delta x_t + \Delta t = -\gamma \nabla V(x) + \delta \Delta x_t$$

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

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

$\delta(x_t - x_{t-\Delta t})$

memory
term

$$x_t \Rightarrow \beta_i \quad x_{t+\Delta t} \Rightarrow \beta_{i+1}$$

$$\nabla V \Rightarrow g(\beta_i) \quad \Delta \beta_{i+1} = \beta_{i+1} - \beta_i$$

$$\beta_{i+1} = \beta_i - \gamma_i g(\beta_i) + \delta(\beta_i - \beta_{i-1})$$

algorithm :

Require : learning rate γ
and δ

Require : initial β_i value
and N_i

$$N_i = \delta(\beta_i - \beta_{i-1}) - \gamma g(\beta_i)$$

$$\beta_{i+1} = \beta_i + N_i$$

while stopping criterion not met
- compute gradient

- compute

$$\nu_i = \delta(\beta_i - \beta_{i-1}) - \gamma g_i'$$

- update β_{i+1}

$$\beta_{i+1} = \beta_i + \nu_i'$$

end while .