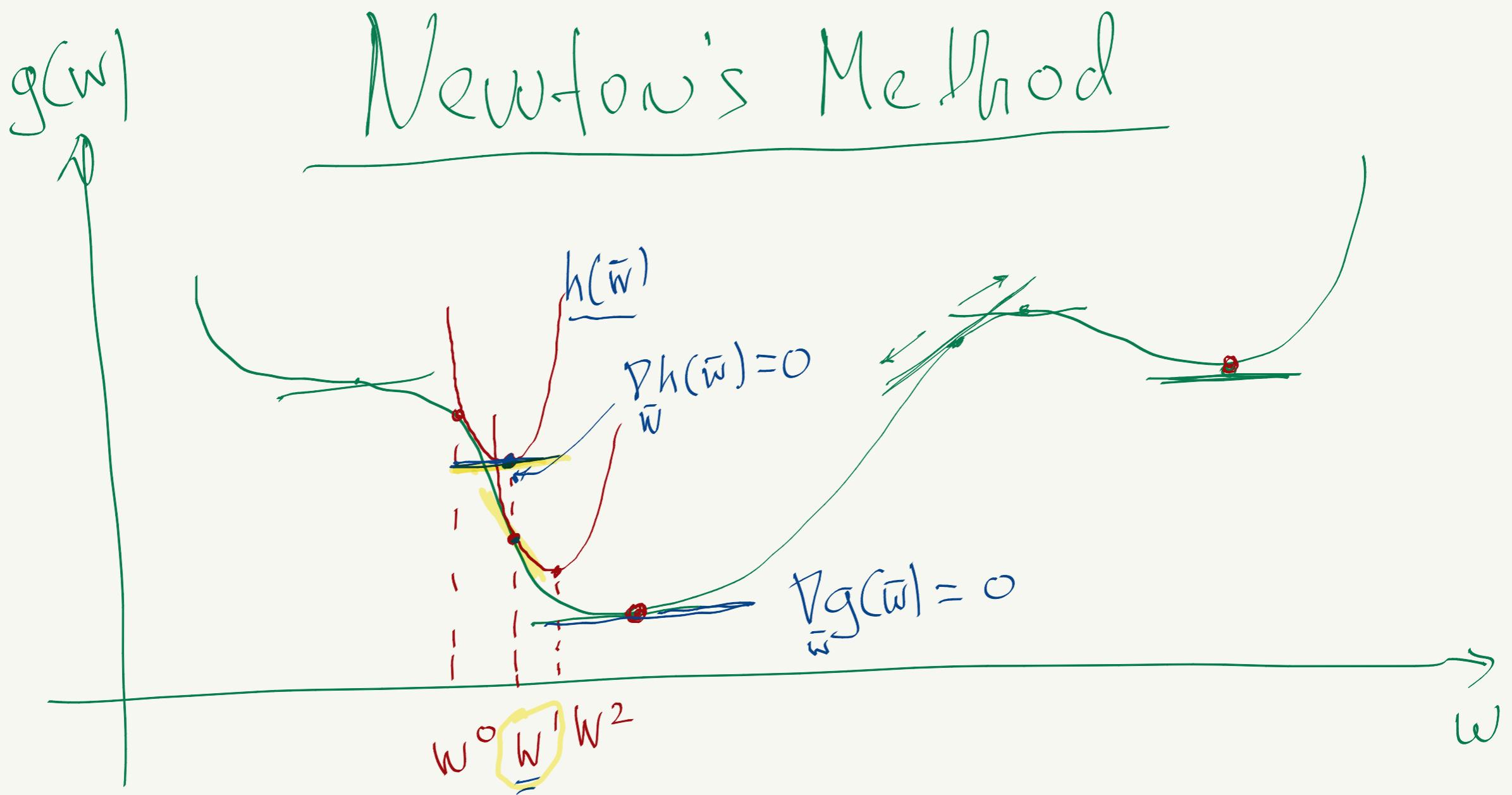


4/15/20



\mathcal{L} nd order Taylor series approximation

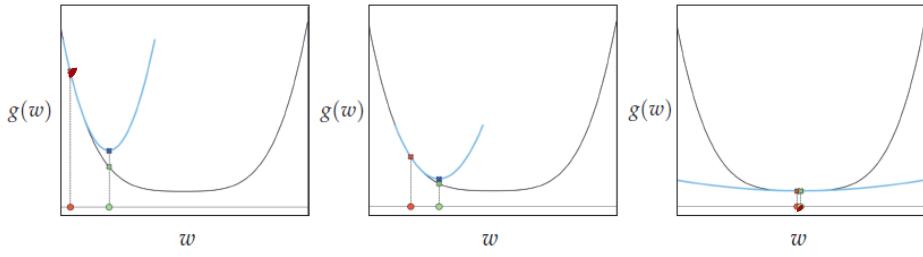
$$\nabla h(\bar{w}) = \nabla g(\bar{w}^0) + \underbrace{\nabla (\nabla g(\bar{w}^0)^T (\bar{w} - \bar{w}^0))}_{\textcircled{1}} + \frac{1}{2} \underbrace{\nabla ((\bar{w} - \bar{w}^0)^T \cdot \nabla g(\bar{w}^0) \cdot (\bar{w} - \bar{w}^0))}_{\textcircled{2}} + \textcircled{3}$$

Find stationary point

$$\boxed{\nabla_{\bar{w}} h(\bar{w}) = 0}$$

① $\nabla_w g(\bar{w}^0) = 0$

(T)



(B)

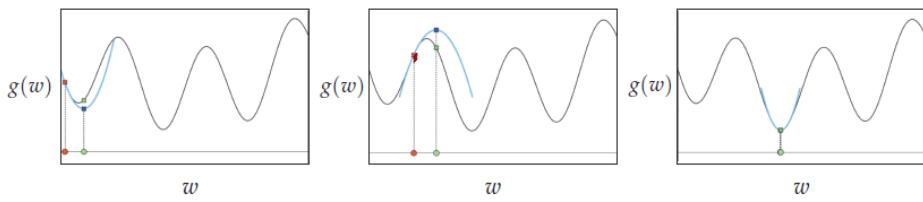


Figure 4.5 Figure associated with Example 4.4. See text for details.

⇒ (T): $g(w) = \frac{1}{50} (w^4 + w^2) + 0.5$

⇒ (B): $g(w) = \sin(3w) + 0.1w^2 + 1.5$

$$\begin{aligned}
 ② \quad & \nabla \left(\nabla g(\bar{w}^o)^T (\bar{w} - \bar{w}^o) \right) \\
 &= \underbrace{\nabla_w \left(\nabla g(\bar{w}^o)^T \bar{w} \right)}_{\text{O}} - \nabla \left(\nabla g(\bar{w}^o)^T \bar{w}^o \right) \\
 &= \nabla \left(\bar{w}^T \cdot \nabla g(\bar{w}^o) \right) \\
 &= \nabla g(\bar{w}^o)
 \end{aligned}$$

$$\left[\begin{array}{l} \nabla(\bar{w}^T \bar{c}) = \bar{c} \\ \nabla(\bar{c}^T \bar{w}) = \bar{c} \\ \bar{w}^T \bar{c} = \bar{c}^T \bar{w} \end{array} \right]$$

$$\begin{aligned}
 ③ \quad & \nabla \left(\underbrace{(\bar{w} - \bar{w}^o)^T}_{\bar{x}^T} \underbrace{\nabla^2 g(\bar{w}^o)}_A (\bar{w} - \bar{w}^o) \right) \\
 & \quad \left[\begin{array}{l} \nabla_{\bar{x}} (\bar{x}^T A \bar{x}) \\ = 2A\bar{x} \\ \frac{d}{dx} (\bar{x}^2) = 2 \cdot \bar{x} \end{array} \right] \\
 & \nabla_w \left(f^T(\bar{w}) A f(\bar{w}) \right) \\
 &= \nabla_{f(\bar{w})} \left(f^T(\bar{w}) A f(\bar{w}) \right) \nabla_w f(\bar{w}) \\
 &\quad \Rightarrow \boxed{2 \nabla^2 g(\bar{w}^o) (\bar{w} - \bar{w}^o) \cdot \nabla_w (\bar{w} - \bar{w}^o)^1}
 \end{aligned}$$

$$\rightarrow \nabla h(\bar{w}) = 0$$

$$\Rightarrow \boxed{\nabla g(\bar{w}^o) + \nabla^2 g(\bar{w}^o) \cdot (\bar{w} - \bar{w}^o) = 0}$$

$$\Rightarrow \nabla^2 g(\bar{w}^o) \cdot \bar{w} = \nabla^2 g(\bar{w}^o) \cdot \bar{w}^o - \nabla g(\bar{w}^o)$$

constant

if Hessian is invertible

$$\Rightarrow \bar{w} = (\nabla^2 g(\bar{w}^o))^{-1} \left(\nabla^2 g(\bar{w}^o) \bar{w}^o - \nabla g(\bar{w}^o) \right)$$

$$\Rightarrow \bar{w} = \bar{w}^o - (\nabla^2 g(\bar{w}^o))^{-1} \cdot \nabla g(\bar{w}^o)$$

$\downarrow \bar{w}'$

General Form of Newton's Method

Start \bar{w}^0

$$\rightarrow \bar{w}^k = \bar{w}^{k-1} - \left(\nabla^2 g(\bar{w}^{k-1}) \right)^{-1} \cdot \nabla g(\bar{w}^{k-1})$$

Gradient Descent

$$\bar{w}^k = \bar{w}^{k-1} - \alpha^k \cdot \nabla g(\bar{w}^{k-1})$$

$$\bar{w}^k = \bar{w}^{k-1} + \alpha^k \bar{d}^{k-1}$$

$$GD: \bar{d}^{k-1} = -\nabla g(\bar{w}^{k-1}) ; \alpha^k = \alpha^k$$

$$NM: \bar{d}^{k-1} = -\left(\nabla^2 g(\bar{w}^{k-1})\right)^{-1} \cdot \nabla g(\bar{w}^{k-1}) ; \alpha^k = 1$$

$$w^k = w^{k-1} - \frac{\frac{d}{dw} g(w^{k-1})}{\frac{d^2}{dw^2} g(w^{k-1}) + \epsilon}$$

Newton's method \approx GD

- ⊕ no α to be chosen
- ⊖ need inverse of Hessian
(lack of scalability
w/ input dimensions)
- ⊖ does not handle well non-convexity
(α local min)
(α local max)
- ⊕ it does not suffer from zig-zagging
- ⊕ typically converges much faster
(for convex functions)

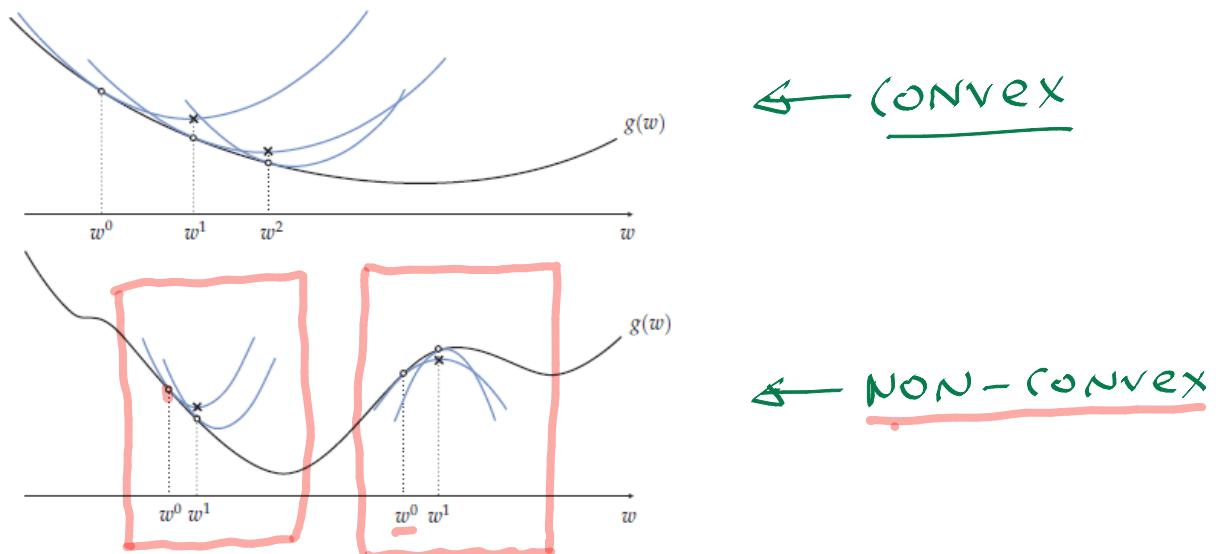


Figure 4.6 Newton's method illustrated. To find a minimum of g , Newton's method hops down the stationary points of quadratic approximations generated by its second-order Taylor series. (top panel) For convex functions these quadratic approximations are themselves always convex (whose only stationary points are minima), and the sequence leads to a minimum of the original function. (bottom panel) For nonconvex functions quadratic approximations can be concave or convex depending on where they are constructed, leading the algorithm to possibly converge to a maximum.

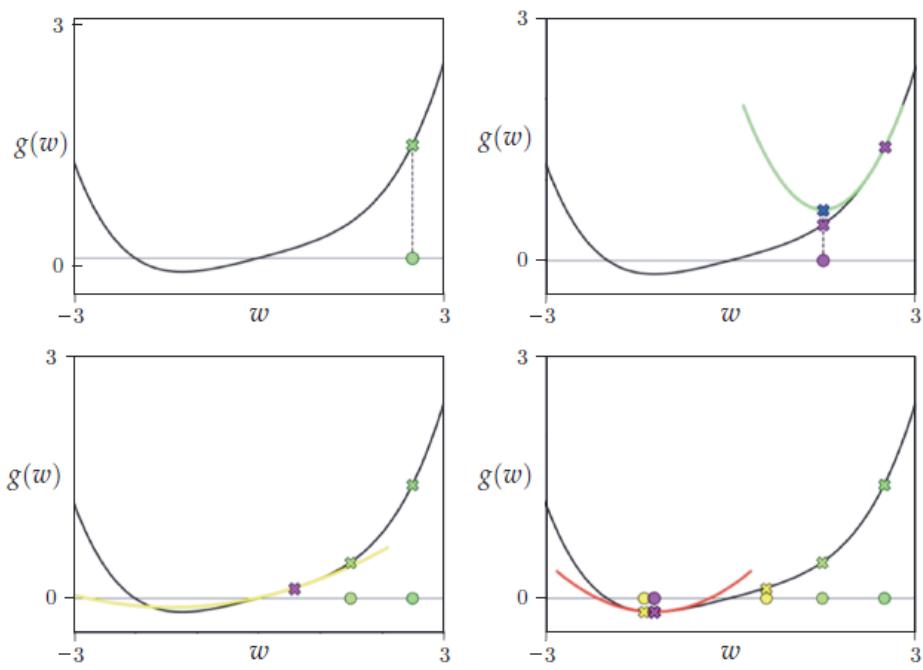


Figure 4.7 Figure associated with Example 4.5, animating a run of Newton's method applied to the function in Equation (4.18). See text for further details.

Dealing w/ NON-convexity \rightarrow regularization

to avoid "division by zero"

$$\bar{w}^k = \bar{w}^{k-1} - \underbrace{\left(\nabla^2 g(\bar{w}^{k-1}) + \epsilon I \right)^{-1}}_{\text{red}} \nabla g(\bar{w}^{k-1})$$

$$\Rightarrow (\nabla^2 g(\bar{w}^{k-1}) + \epsilon I) \bar{w}^k = (\nabla^2 g(\bar{w}^{k-1}) + \epsilon I) \bar{w}^{k-1} - \nabla g(\bar{w}^{k-1})$$

this equation can be interpreted as the stationary point of a slightly adjusted 2nd order Taylor series approximation

$$h(\bar{w}) = \begin{cases} h_1(\bar{w}) & \\ \quad g(\bar{w}^{k-1}) + \nabla g(\bar{w}^{k-1})^T (\bar{w} - \bar{w}^{k-1}) \\ \quad + \frac{1}{2} (\bar{w} - \bar{w}^{k-1})^T \nabla^2 g(\bar{w}^{k-1}) (\bar{w} - \bar{w}^{k-1}) \\ + \frac{\epsilon}{2} \|\bar{w} - \bar{w}^{k-1}\|_2^2 & \text{regularizer} \\ & h_2(\bar{w}) \end{cases}$$

$$(\bar{w} - \bar{w}^{k-1})^T \left(\frac{\epsilon}{2} I \right) (\bar{w} - \bar{w}^{k-1})$$

N positive eigenvalues = $\epsilon/2$

$$h(\bar{w}) = h_1(\bar{w}) + \frac{\epsilon}{2} h_2(\bar{w})$$

(ϵ -values)

as $\epsilon \uparrow$

we convexify $h(\bar{w})$

flat or non-convex always convex

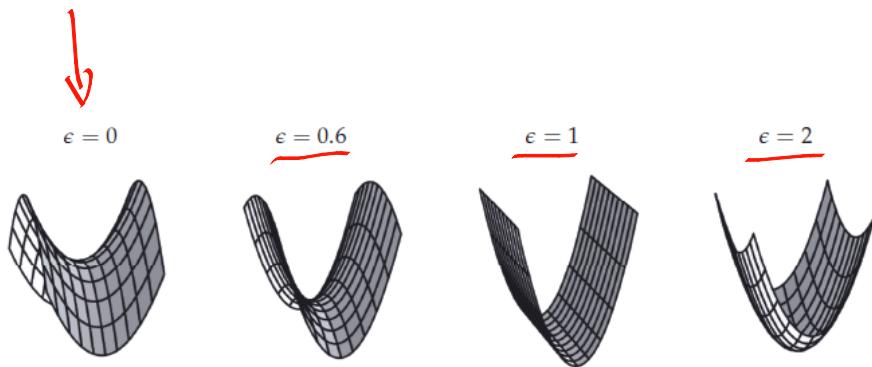


Figure A.12 Figure associated with Example A.9. From left to right, a nonconvex quadratic function is slowly turned into a convex function via the weighted addition of a convex quadratic. See text for further details.

$$h(\bar{w}) = h_1(w_1, w_2) + \epsilon \cdot h_2(w_1, w_2)$$

↓ ↓

$$\rightarrow h_1(w_1, w_2) = w_1^2 - w_2^2$$

$$\rightarrow h_2(w_1, w_2) = w_1^2 + w_2^2 = \|\bar{w}\|_2^2$$

Dealing w/ scaling w/ input dimension

\Rightarrow Hessian-Free Methods
(appendix A)

A. Subsampling the Hessian

- keep only diagonal value

decouple along each coordinate

$$w_n = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ \underline{\underline{w_n}} \\ \vdots \\ w_N \end{bmatrix} \quad w_n^k = w_n^{k-1} - \frac{\frac{\partial}{\partial w_n} g(\bar{w}^{k-1})}{\frac{\partial^2}{\partial w_n^2} g(\bar{w}^{k-1})}$$

B. Quasi-Newton methods

Replace by Secant Matrix
Hessian
(use of low rank matrices)

$$A = uu^T$$

$$A = \begin{bmatrix} & & \\ & & \end{bmatrix} \begin{bmatrix} & \\ & \end{bmatrix}$$

↑ rank 1