

MATH 565: Lecture II (02/17/2026)

Today:

- * Newton's update for nonquadratic J
- * line search in Newton
- * Newton in regression, L_2 -SVM

Recall Newton's method update: $\bar{w} \leftarrow \bar{w} - H^{-1} \nabla J$

Newton's method works perfectly for quadratic J . What about other J ?

Consider $J(w) = -w^3 + 4w^2 + 1$

Consider Newton update at $w_0 = 2$ ($J(\bar{w}_0) = 9$).

$w^\circ = w_0$ iteration index

$$\nabla J = -3w^2 + 8w \Rightarrow \nabla J(w_0) = 4$$

$$HJ = -6w + 8 \Rightarrow HJ(w_0) = -4$$

$$w' \leftarrow w^\circ - H^{-1} \nabla J = 2 - \left(-\frac{1}{4}\right) \cdot 4 = 3$$

$$J(w') = -27 + 36 + 1 = 10$$

So, the Newton update increases J !

$H < 0$ here! So,
Newton's method
pushes w (from $w_0=2$)
to $w'=3$, which is
a local max for the
quadratic approximation.

Quadratic approximation of $J(w)$ at $w_0=2$

$$J(w) = J(w_0) + (w-2) \cdot 4 + \frac{1}{2} (w-2)^2 (-4)$$

$$= 9 + 4w - 8 - 2(w^2 - 4w + 4)$$

$$= -2w^2 + 12w - 7$$

this parabola is
opening down (\cap)

$\rightarrow w=3$ is a local maximum for
this quadratic approximation

Line Search

Modify Newton update to

$$\bar{\omega}^{k+1} \leftarrow \bar{\omega}^k - \alpha H_k^{-1} \nabla J_k$$

accept, i.e., update $\bar{\omega}$ as described above, only if $J(\bar{\omega}^{k+1}) < J(\bar{\omega}^k)$
 else, change α , or start over from a different $\bar{\omega}_0$.

Example

$$J(\omega) = \omega^2 - \ln(\omega)$$

$$\nabla J = 2\omega - \frac{1}{\omega}$$

$$HJ = 2 + \frac{1}{\omega^2}$$

(log barrier function)

used in interior point methods.

keeps the algorithm from getting too close to $\omega=0$.

$$\text{generalization: } J(\bar{\omega}) = \sum_{i=1}^{d_1} J_i(\omega_i)$$

$$\text{where } J_i(\omega_i) = \omega_i^2 - \ln(\omega_i)$$

$\Rightarrow J$ has unique global minimum
 at $\omega^* = \frac{1}{\sqrt{2}} = 0.707$.

$$HJ = \begin{bmatrix} \ddots & & 0 \\ \vdots & 2 + \frac{1}{\omega_i^2} & \ddots \\ 0 & \ddots & \ddots \end{bmatrix}$$

Starting @ $\omega_0 = 2$, apply Newton's update:

$$\nabla J = \frac{7}{2} (= 3.5)$$

$$HJ = \frac{9}{4} (= 2.25)$$

$$\omega \leftarrow \omega_0 - \left(\frac{4}{9} \right) \cdot \left(\frac{7}{2} \right) = \frac{4}{9} \approx 0.44 \left(< \frac{1}{\sqrt{2}} \right)$$

\hookrightarrow We've overshot (passed over) the global minimum!

With $\alpha \approx 0.83$,

$$\omega' = \omega_0 - \alpha H^{-1} \nabla \approx 0.707.$$

Can adapt various line search selection approaches introduced for gradient descent for Newton's method as well.

Newton's Method in ML

Newton in Regression

$$J = \frac{1}{2} \|D\bar{w} - \bar{y}\|^2 = \frac{1}{2} (D\bar{w} - \bar{y})^T (D\bar{w} - \bar{y})$$

$$\nabla J = D^T D\bar{w} - D^T \bar{y}$$

$$HJ = D^T D$$

$$HJ = \sum_{i=1}^n \bar{x}_i \bar{x}_i^T \quad (\text{sum of outer products})$$

$$= \sum_{i=1}^n HJ_i \quad \xrightarrow{\text{Hessian of } J = \text{sum of point-specific Hessians}}$$

$$D = \begin{bmatrix} & & & & d \\ 1 & 2 & \dots & n \\ \vdots & & & \bar{x}_i^T \\ n & & & \end{bmatrix}_{n \times d}$$

$$\text{Newton update} \quad \bar{w} \leftarrow \bar{w} - \underbrace{H^{-1} \nabla J}_{\text{Hessian of } J}$$

$$\begin{aligned} \bar{w} &\leftarrow \bar{w} - (D^T D)^{-1} (D^T D \bar{w} - D^T \bar{y}) \\ &= \bar{w} - \bar{w} + (D^T D)^{-1} D^T \bar{y} \end{aligned}$$

$$\bar{w} = (D^T D)^{-1} D^T \bar{y} \rightarrow \text{optimal solution in a single update.}$$

Recall that we obtained the same solution using gradient descent (GD) approaches. $\nabla J = \bar{0}$ gave $\bar{w} = (D^T D)^{-1} D^T \bar{y}$.

Newton in SVM

We look at L_2 -SVM. (The hinge loss J is not smooth)

$$\begin{aligned} J = J_{L_2\text{-SVM}}(\bar{w}) &= \frac{1}{2} \sum_{i=1}^n \max \left\{ 0, (1 - y_i (\bar{w}^\top \bar{x}_i)) \right\}^2 \\ &= \sum_{i=1}^n J_i \quad \text{for} \quad J_i = \frac{1}{2} \max \left\{ 0, 1 - y_i (\bar{w}^\top \bar{x}_i) \right\}^2. \\ \text{Here, } J_i &= f_i(z) = \frac{1}{2} \max \left\{ 0, 1 - y_i z \right\}^2 \text{ for } z = \bar{w}^\top \bar{x}_i. \end{aligned}$$

$$\frac{\partial f_i(z)}{\partial z} = -y_i \max \left\{ 0, 1 - y_i z \right\}$$

$$\Rightarrow \nabla J_i = \frac{\partial J_i}{\partial \bar{w}} = -y_i \underbrace{\max \left\{ 0, 1 - y_i (\bar{w}^\top \bar{x}_i) \right\}}_{\text{scalar}} \bar{x}_i.$$

If we were using least squares classification (instead of L_2 -SVM), we get $\nabla J_{LS} = -y_i (1 - y_i (\bar{w}^\top \bar{x}_i)) \bar{x}_i$

We can rewrite ∇J_i using the indicator function $\delta(\cdot)$:

$$\nabla J_i = \underbrace{(\bar{w}^\top \bar{x}_i - y_i) \delta(1 - y_i (\bar{w}^\top \bar{x}_i) > 0)}_{\text{scalar}} \cdot \bar{x}_i$$

$$\begin{aligned} \Rightarrow \nabla J(\bar{w}) &= \sum_{i=1}^n \nabla J_i \\ &= D^\top \Delta_{\bar{w}} (D\bar{w} - \bar{y}) \end{aligned}$$

$$\text{where } \Delta_{\bar{w}} = \left[\text{diag}(\delta(1 - y_i (\bar{w}^\top \bar{x}_i) > 0)) \right]$$

↳ diagonal matrix whose i^{th} entry is $\delta(1 - y_i (\bar{w}^\top \bar{x}_i) > 0)$.

Hessian?

$$HJ = \sum_{i=1}^n H_i = \sum_{i=1}^n HJ_i$$

Note: $\nabla J_i = (\bar{w}^\top \bar{x}_i - y_i) \underbrace{\delta(1 - y_i(\bar{w}^\top \bar{x}_i)) > 0}_{\text{scalar}} \cdot \bar{x}_i$

$$= s_i(\bar{w}) \cdot \bar{x}_i \quad \text{where } s_i(\bar{w}) = -y_i \max\{0, 1 - y_i(\bar{w}^\top \bar{x}_i)\}$$

$$\begin{aligned} \Rightarrow HJ_i &= \bar{x}_i \left[\frac{\partial s_i}{\partial \bar{w}} \right]^\top \\ &= \bar{x}_i \left(y_i^2 \delta(1 - y_i(\bar{w}^\top \bar{x}_i)) > 0 \right) \bar{x}_i^\top \quad y_i^2 = 1 \\ &= \delta(1 - y_i(\bar{w}^\top \bar{x}_i)) \bar{x}_i \bar{x}_i^\top \end{aligned}$$

$$\Rightarrow HJ = \sum_{i=1}^n HJ_i = \sum_{i=1}^n \delta(1 - y_i(\bar{w}^\top \bar{x}_i)) \bar{x}_i \bar{x}_i^\top$$

For J_{LS} , we get $H = D^\top D$.

Here, for L_2 -SUM, we get

$$HJ_{L_2\text{-SVM}} = D^\top \Delta_{\bar{w}} D$$

$$\text{where } \Delta_{\bar{w}} = [\text{diag}(\delta(1 - y_i(\bar{w}^\top \bar{x}_i)) > 0)].$$