

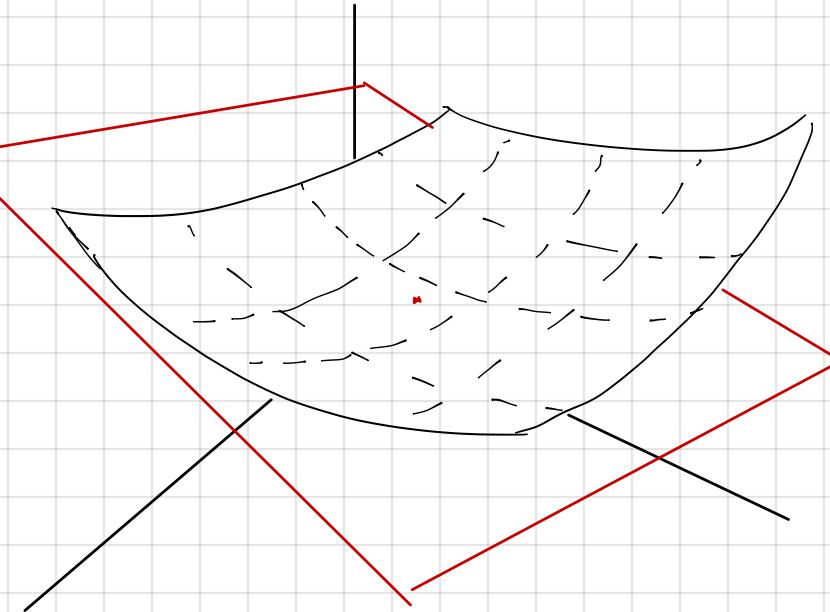
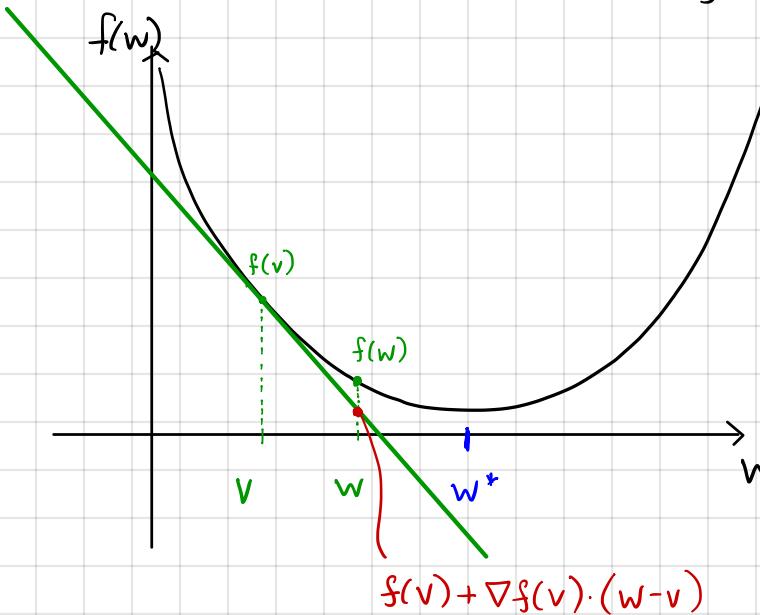
Lecture 14: Basic Convex Optimization

Basic convex optimization

Goal: find $\underline{w}^* = \underset{\underline{w}}{\operatorname{arg\,min}} f(\underline{w})$ when f is a convex function.

A function is convex if $f(\underline{w}) \geq f(\underline{v}) + \nabla f(\underline{v})^\top (\underline{w} - \underline{v})$

- i.e. if it's \geq all its tangents



Ex. $f(\underline{w}) = \|\underline{y} - \underline{X}\underline{w}\|_2^2$. We know $\underline{w}^* = (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \underline{y}$

Gradient descent finds this point iteratively.

- avoids computing matrix inverse
- generalizes to many other problems.

Gradient:

$$\text{if } f(\underline{w}) = \underline{y}^T \underline{y} - 2\underline{w}^T \underline{X}^T \underline{y} + \underline{w}^T \underline{X}^T \underline{X} \underline{w}, \text{ then } \nabla_{\underline{w}} f = 0 - 2\underline{X}^T \underline{y} + 2\underline{X}^T \underline{X} \underline{w}$$

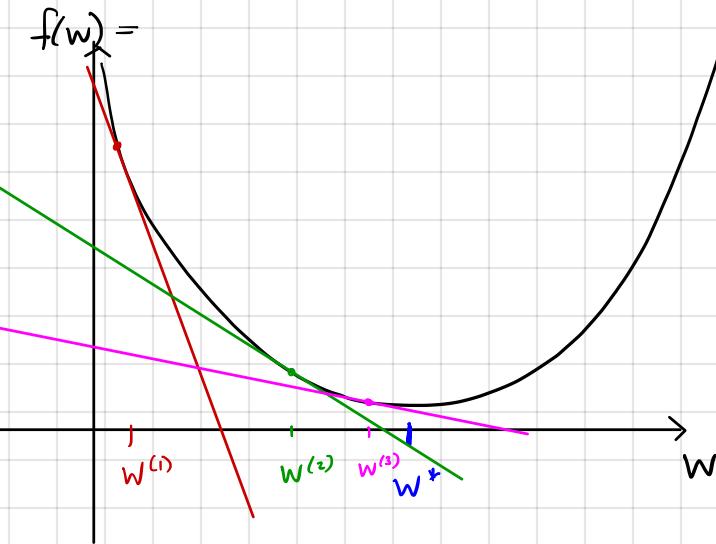
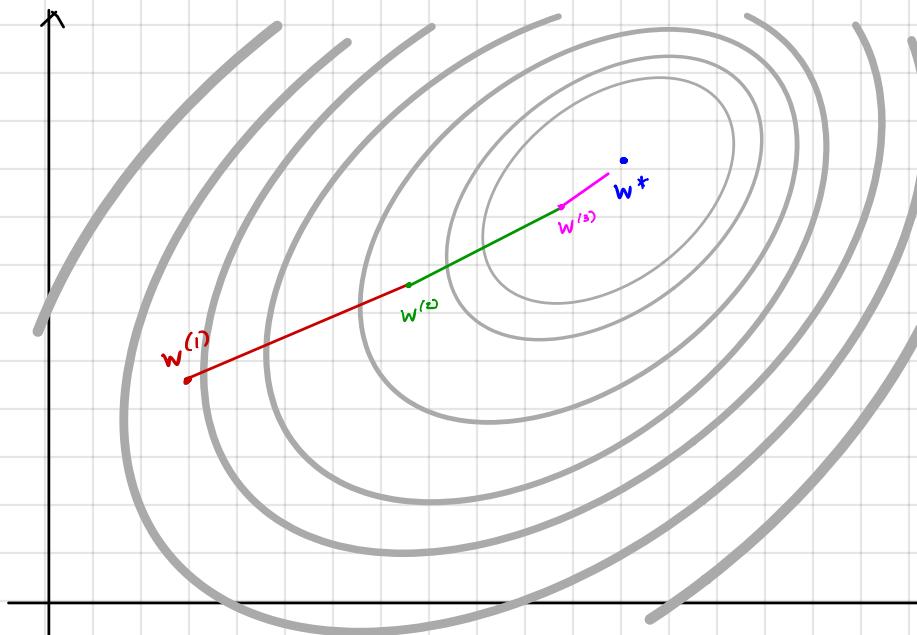
Gradient descent starts with initial guess $\underline{w}^{(1)}$, and then repeatedly takes steps in the direction of the negative gradient.

for $k = 1, 2, 3, \dots$

$$\begin{aligned}\underline{w}^{(k+1)} &= \underline{w}^{(k)} - \gamma (\underline{X}^T \underline{X} \underline{w}^{(k)} - \underline{X}^T \underline{y}) \\ &= \underline{w}^{(k)} - \gamma \underline{X}^T (\underline{X} \underline{w}^{(k)} - \underline{y})\end{aligned}$$

if $\|\underline{w}^{(k+1)} - \underline{w}^{(k)}\|_2 < \varepsilon$, then BREAK

$\gamma > 0$ is step size



More generally:

want to minimize $f(\underline{w})$
initialize with $\underline{w}^{(1)}$
for $k = 1, 2, 3, \dots$

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} - \gamma \nabla_{\underline{w}} f |_{\underline{w}=\underline{w}^{(k)}}$$

if $\|\underline{w}^{(k)} - \underline{w}^{(k+1)}\| < \varepsilon$, then BREAK

Does this work?

Convergence for $f(\underline{w}) = \|\underline{X}\underline{w} - \underline{y}\|_2^2$

$$\text{Want } \|\underline{X}\underline{w}^{(k+1)} - \underline{y}\|_2^2 < \|\underline{X}\underline{w}^{(k)} - \underline{y}\|_2^2$$

$$\text{recall } \underline{w}^{(k+1)} = \underline{w}^{(k)} - \tau \underline{X}^\top (\underline{X}\underline{w}^{(k)} - \underline{y})$$

$$\begin{aligned} \Rightarrow \|\underline{X}\underline{w}^{(k+1)} - \underline{y}\|_2^2 &= \|\underline{X}(\underline{w}^{(k)} - \tau \underline{X}^\top (\underline{X}\underline{w}^{(k)} - \underline{y})) - \underline{y}\|_2^2 \\ &= \|\underline{X}\underline{w}^{(k)} - \underline{y} - \tau \underline{X} \underline{X}^\top (\underline{X}\underline{w}^{(k)} - \underline{y})\|_2^2 \\ &= \|\underline{X}\underline{w}^{(k)} - \underline{y}\|_2^2 - 2\tau (\underline{X}\underline{w}^{(k)} - \underline{y})^\top (\underline{X} \underline{X}^\top (\underline{X}\underline{w}^{(k)} - \underline{y})) + \tau^2 \|\underline{X} \underline{X}^\top (\underline{X}\underline{w}^{(k)} - \underline{y})\|_2^2 \\ &= \|\underline{X}^\top (\underline{X}\underline{w}^{(k)} - \underline{y})\|_2^2 \\ &\leq \|\underline{X}\|_{\text{op}}^2 \|\underline{X}^\top (\underline{X}\underline{w}^{(k)} - \underline{y})\|_2^2 \end{aligned}$$

$$\|\underline{X}\underline{w}^{(k+1)} - \underline{y}\|_2^2 \leq \|\underline{X}\underline{w}^{(k)} - \underline{y}\|_2^2 + 4\tau \left(\tau \|\underline{X}\|_{\text{op}}^2 \|\underline{X}^\top (\underline{X}\underline{w}^{(k)} - \underline{y})\|_2^2 - \|\underline{X}^\top (\underline{X}\underline{w}^{(k)} - \underline{y})\|_2^2 \right)$$

$$= \|\underline{X}\underline{w}^{(k)} - \underline{y}\|_2^2 + 4\tau \|\underline{X}^\top (\underline{X}\underline{w}^{(k)} - \underline{y})\|_2^2 (\tau \|\underline{X}\|_{\text{op}}^2 - 1)$$

$$\Rightarrow \text{if } \tau \|\underline{X}\|_{\text{op}}^2 - 1 < 0 \quad (\tau < \frac{1}{\|\underline{X}\|_{\text{op}}^2}), \text{ then } \|\underline{X}\underline{w}^{(k+1)} - \underline{y}\|_2^2 < \|\underline{X}\underline{w}^{(k)} - \underline{y}\|_2^2$$

$$\text{if } \underline{w}^{(1)} = 0 \text{ and } \tau < \frac{1}{\|\underline{X}\|_{\text{op}}^2}, \text{ then}$$

$$\underline{w}^{(k)} \longrightarrow (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \underline{y} \text{ as } k \rightarrow \infty$$

Another perspective:

$$\begin{aligned}
 w^{(k+1)} &= w^{(k)} + \tau (X^T y - X^T X w^{(k)}) \\
 &= w^{(k)} + \tau X^T X \left[(X^T X)^{-1} X^T y - w^{(k)} \right] \\
 &= w^{(k)} - \tau X^T X (w^{(k)} - w^*)
 \end{aligned}$$

Subtract w^* from both sides:

$$\begin{aligned}
 v^{(k+1)} &= w^{(k+1)} - w^* = w^{(k)} - w^* - \tau X^T X (w^{(k)} - w^*) \\
 &= v^{(k)} - \tau X^T X v^{(k)}
 \end{aligned}$$

$$\begin{aligned}
 \text{or } v^{(k+1)} &= (I - \tau X^T X) v^{(k-1)} \\
 &= (I - \tau X^T X)^{k-1} v^{(0)}
 \end{aligned}$$

for this sequence to go to zero, need all the singular values of $(I - \tau X^T X)$ to be less than 1 in magnitude
this happens if

$$\tau < \frac{1}{\sigma_{\max}^2(X^T X)}$$

$$\begin{aligned}
 I - \tau X^T X &= I - \tau U \Sigma V^T V \Sigma^T U^T \\
 &= U (I - \tau \Sigma^2) U^T
 \end{aligned}$$

⇒ singular values are $1 - \tau \sigma_i^2, 1 - \tau \sigma_i^2, \dots$

⇒ need $\max_i |1 - \tau \sigma_i^2| < 1$

⇒ need $|1 - \tau \sigma_{\max}^2| < 1$

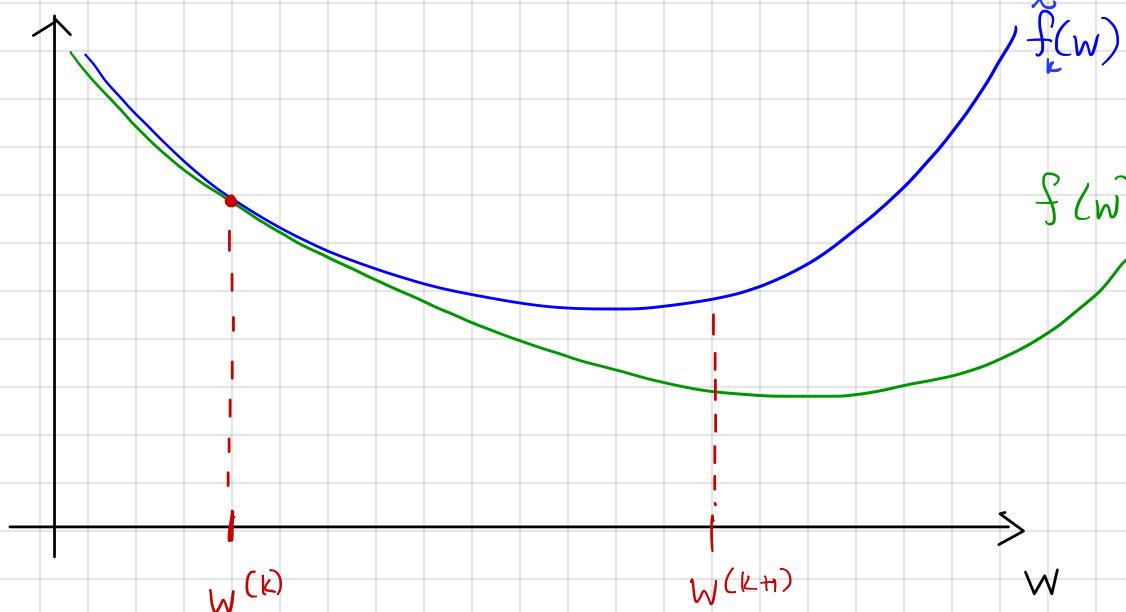
⇒ need $\tau < \frac{1}{\sigma_{\max}^2}$

Where do gradient descent updates come from?

$$\begin{aligned}
 \text{Consider } f(w) &= \|y - Xw\|_2^2 = \|y - Xw^{(k)} + Xw^{(k)} - Xw\|_2^2 \\
 &= \|y - Xw^{(k)}\|_2^2 + 2(y - Xw^{(k)})^\top X(w^{(k)} - w) + \|Xw^{(k)} - Xw\|_2^2 \\
 &\leq \underbrace{\|y - Xw^{(k)}\|_2^2}_{\text{same value regardless of } w} + 2(y - Xw^{(k)})^\top X(w^{(k)} - w) + \|X\|_{\text{op}}^2 \|w^{(k)} - w\|_2^2
 \end{aligned}$$

Let τ be a step size, assume $\tau < \frac{1}{2\|X\|_{\text{op}}^2}$ ($\|X\|_{\text{op}}^2 < \frac{1}{2\tau}$)

$$\Rightarrow f(w) \leq C + 2(y - Xw^{(k)})^\top X(w^{(k)} - w) + \frac{1}{2\tau} \|w^{(k)} - w\|_2^2 =: \tilde{f}_k(w)$$



$$f(w^{(k)}) = \tilde{f}(w^{(k)})$$

Choose $w^{(k+1)}$ to minimize \tilde{f}_k

$$f(w) \leq \tilde{f}_k(w)$$

Aside: $\|Xw\|_2^2 \leq \|X\|_{\text{op}}^2 \|w\|_2^2$
 where $\|X\|_{\text{op}} = \max$ singular value of X .
 because $\|Xw\|_2^2 = \|\Sigma V^T w\|_2^2$
 $= \|\Sigma V^T w\|_2^2$
 $= \sum_i \sigma_i^2 (V^T w)_i^2$
 $\leq \sigma_{\max}^2 \sum_i (V^T w)_i^2$
 $= \sigma_{\max}^2 \|V^T w\|_2^2$
 $= \sigma_{\max}^2 \|w\|_2^2$

$$\hat{w}_{k+1} = \underset{w}{\operatorname{argmin}} \quad 2(y - Xw^{(k)})^T X(w^{(k)} - w) + \frac{1}{2\tau} \|w^{(k)} - w\|_2^2$$

Let $v := 2\tau X^T(y - Xw^{(k)})$

$$= \underset{w}{\operatorname{argmin}} \quad 2 \cdot \underbrace{2\tau(y - Xw^{(k)})^T X(w^{(k)} - w)}_{=: v^T} + \|w^{(k)} - w\|_2^2$$

— independent of w !

$$= \underset{w}{\operatorname{argmin}} \quad 2v^T(w^{(k)} - w) + \|w^{(k)} - w\|_2^2$$

$$= \underset{w}{\operatorname{argmin}} \quad \|v + w^{(k)} - w\|_2^2 - \|v\|_2^2$$

$$= \underset{w}{\operatorname{argmin}} \quad \|v + w^{(k)} - w\|_2^2 = w^{(k)} + v = w^{(k)} + 2\tau X^T(y - Xw^{(k)}) = \text{Gradient Descent Step !!!!}$$

$$= w^{(k)} - 2\tau X^T(Xw^{(k)} - y)$$