# Important announcements

- grades for midterm 1 are out
- example for linear regression

# Gradient Descent

Can we always find a closed form expression for $w^* = \arg\min\limits_{w \in W} g(w)$ ?  $\Rightarrow$ **No !**

Linear regression is closed form because of the squared loss & linear function class

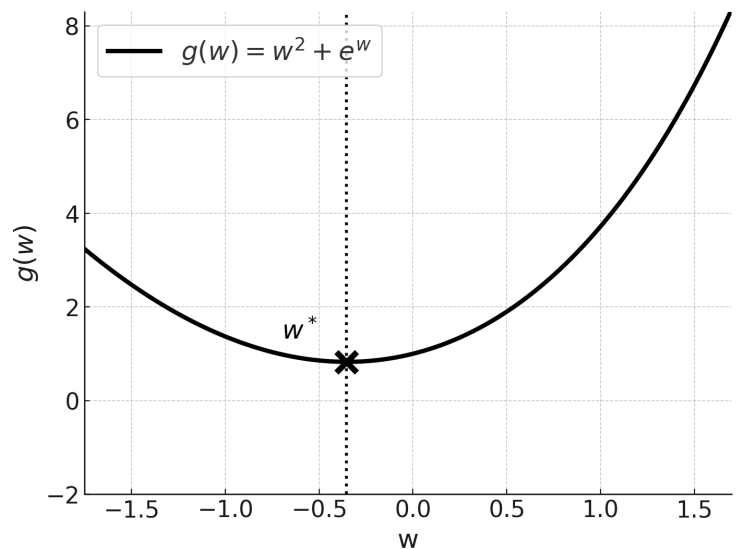Ex: $g(w) = w^2 + e^w$, $g'(w) = 2w + e^w$, $g''(w) = 2 + e^w \geq 0$

g is convex

$g'(w) = 2w + e^w = 0$

$\rightarrow 2w = -e^w$

No way to solve for $w$ !

$\rightarrow$ No closed form solution although $g(w)$ convex.



Gradient descent helps with this problem.

# Second-Order Gradient Descent (Newton-Raphson)

If $g(\omega)$ is a degree 5 polynomial or less

→ then there exists a closed form solution for $g(\omega)=0$

Let's approximate $g(\omega)$ with a convex low degree polynomial (i.e. 2nd degree polynomial)

In general, the Taylor series at a point $\omega^{(0)}$ of $g(\omega)$ is

$$g(\omega) = \sum_{n=0}^{\infty} \frac{g^{(n)}(\omega^{(0)})}{n!} (\omega - \omega^{(0)})^n$$

We use a 2nd order approximation that takes the first 3 terms

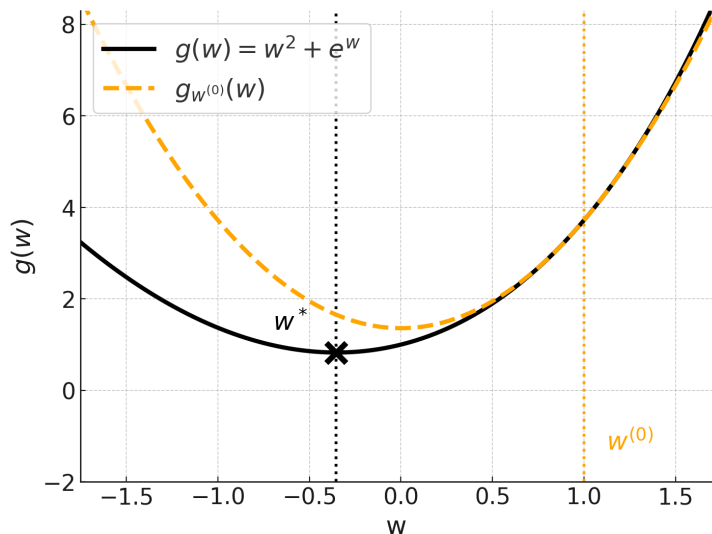→ not exactly $g(\omega)$ but good enough around $\omega^{(0)}$

$$g(\omega) \approx g_{w^{(0)}}(\omega) = \underbrace{g(\omega^{(0)})}_{const} + \underbrace{g'(\omega^{(0)})}_{const}\underbrace{(\omega-\omega^{(0)})}_{const} + \underbrace{\frac{g''(\omega^{(0)})}{2}}_{const} \underbrace{(\omega-\omega^{(0)})^2}_{const}$$

const

→ 2nd order polynomial

Ex: $g(\omega) = \omega^2 + e^\omega$

$\omega^{(0)} = 1$ ← arbitrary choice

$g_{w^{(0)}}(\omega) = g_1(\omega)$

minimize $g_{\omega^{(0)}}(\omega)$:

$$\frac{d}{d\omega} g(\omega) \approx \frac{d}{d\omega} g_{\omega^{(0)}}(\omega) = g'(\omega^{(0)}) + g''(\omega^{(0)})(\omega - \omega^{(0)}) = 0$$

$$\Rightarrow g''(\omega^{(0)}) \omega = g''(\omega^{(0)}) \omega^{(0)} - g'(\omega^{(0)})$$

$$\Rightarrow \boxed{\omega = \omega^{(0)} - \frac{g'(\omega^{(0)})}{g''(\omega^{(0)})}}$$

Ex: $\omega^{(0)} = 1$

$g'(1) = 2 \cdot 1 + e^1$

$g''(1) = 2 + e^1$

$\omega^{(1)} = 1 - \frac{2+e}{2+e} = 0$



$\Rightarrow$ We can approximat $g(\omega)$ again but at $\omega^{(1)}$

$$g_{\omega^{(1)}}(\omega) = g(\omega^{(1)}) + g'(\omega^{(1)})(\omega - \omega^{(1)}) + \frac{g''(\omega^{(1)})}{2}(\omega - \omega^{(1)})^2$$

Ex: $\omega^{(1)} = 0$

$g_{\omega^{(1)}}(\omega) = g_0(\omega)$

And we can minimize $g_{\omega^{(1)}}(\omega)$ to get

$$\omega^{(2)} = \omega^{(1)} - \frac{g'(\omega^{(1)})}{g''(\omega^{(1)})}$$

Ex: $\omega^{(2)} = -\frac{1}{3}$

In general

$$\omega^{(t+1)} = \omega^{(t)} - \frac{g'(\omega^{(t)})}{g''(\omega^{(t)})}$$

where $\omega^{(t+1)}$ approaches $\omega^*$ as $t \to \infty$

# (First-Order) Gradient Descent

Sometimes it is hard to calculate $g''(\omega)$, especially in high dimension.

Instead we can replace it with $\eta^{(t)}$

$$\omega^{(t+1)} = \omega^{(t)} - \eta^{(t)} g'(\omega^{(t)}) \qquad \eta^{(t)} \text{ is the "step size" or}$$

if we know $g''(\omega)$, then we can set $\eta^{(t)} = \frac{1}{g''(\omega^{(t)})}$ "learning rate"
and get back the 2nd order gradient descent.

$\to$ in 1st order gradient descent we still approximate the underlying func with a 2nd order Taylor expansion, but remove the necessity to calculate $g''(\omega)$.

Ex: $\eta^{(t)} = \frac{1}{g''(\omega^{(t)})}$

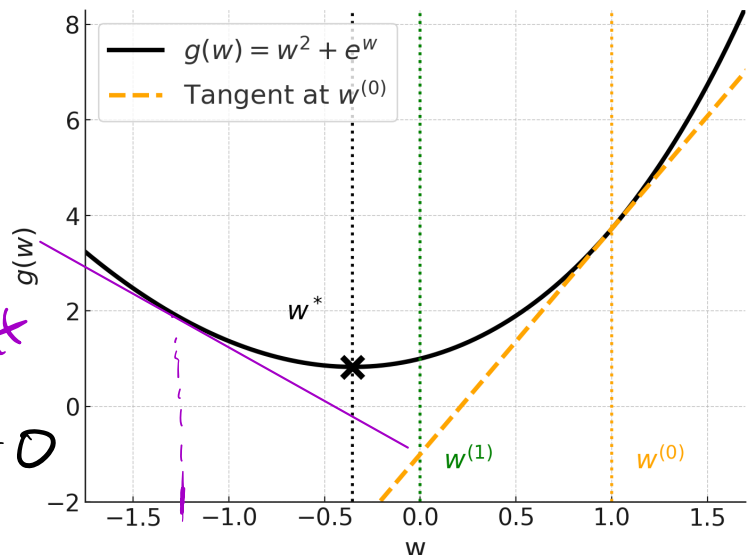$$g(\omega) = \omega^2 + e^\omega$$
$$g'(\omega) = 2\omega + e^\omega$$

$\omega^{(0)} = 1$, $\eta^{(t)} = \frac{1}{2+e}$

$$\omega^{(1)} = \omega^{(0)} - \eta^{(0)} g'(\omega^{(0)}) = 1 - 1 = 0$$

direction of ascent

scale of the step



Legend: $g(w) = w^2 + e^w$ (black), Tangent at $w^{(0)}$ (orange dashed). Axis labels: $g(w)$ vs $w$. Points marked $w^*$, $w^{(1)}$, $w^{(0)}$.

$$\eta^{(1)} = \frac{1}{2 + e^0} = \frac{1}{3}$$

$$\omega^{(2)} = \omega^{(1)} - \eta^{(1)} g'(\omega^{(1)})$$



↑ known

$$g''(\omega)$$

↓ unknown

Ex: $\eta^{(t)} = \frac{1}{10} < \frac{1}{2+e}$
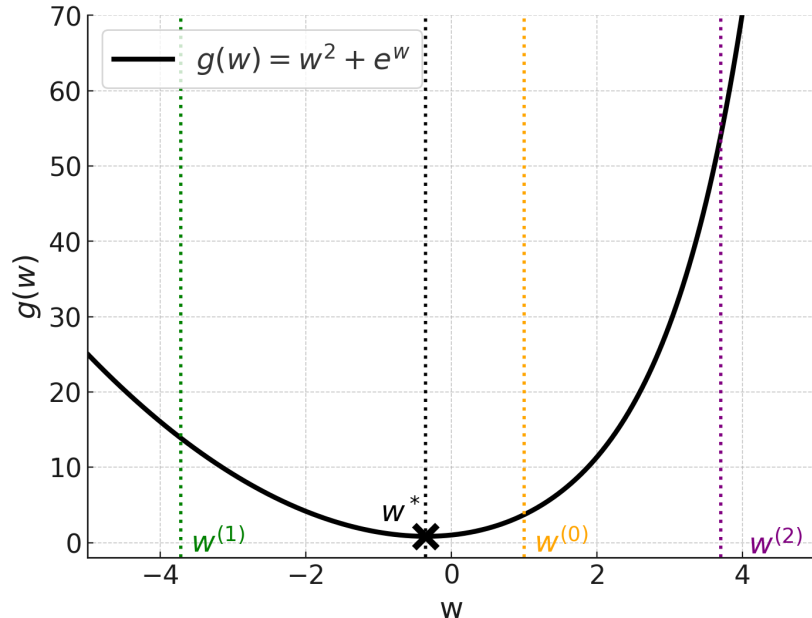
small $\eta^{(t)}$

slowly reaches $\omega^*$



Ex: $\eta^{(t)} = 1 > \frac{1}{2+e}$

large $\eta^{(t)}$

might never reach $\omega^*$ ( diverge )

-> small step size is safe (probably converge to $w^*$) however it might take long

-> large step size get you fast to $w^*$ but might diverge

# Multivariate Gradient Descent

$\vec{w} \in W = \mathbb{R}^d$, $d > 1$, $g(\vec{w})$

Objective: $\vec{w}^* = (w_1^*, \dots, w_d^*) = \arg\min_{\vec{w} \in W} g(\vec{w})$

multivariat gradient descent update rule

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta^{(t)} \nabla g(\vec{w}^{(t)})$$

where $\nabla g(\vec{w}) = \left( \frac{\partial}{\partial w_1} g(\vec{w}), \dots, \frac{\partial}{\partial w_d} g(\vec{w}) \right)^T \in \mathbb{R}^d$

$$\eta^{(t)} \in \mathbb{R}$$

Ex: $W = \mathbb{R}^2$, $g(\vec{w}) = g(w_1, w_2) = w_1^2 + e^{w_1} + w_2^2 + e^{w_2}$

$$\frac{\partial}{\partial w_1} g(\vec{w}) = 2w_1 + e^{w_1} \qquad \frac{\partial}{\partial w_2} g(\vec{w}) = 2w_2 + e^{w_2}$$

$$\vec{w}^{(t+1)} = \begin{pmatrix} w_1^{(t+1)} \\ w_2^{(t+1)} \end{pmatrix} = \begin{pmatrix} w_1^{(t)} \\ w_2^{(t)} \end{pmatrix} - \eta^{(t)} \begin{pmatrix} 2w_1^{(t)} + e^{w_1^{(t)}} \\ 2w_2^{(t)} + e^{w_2^{(t)}} \end{pmatrix}$$

Let $\vec{w}^{(0)} = (1, 1)^T$, $\eta^{(0)} = \frac{1}{2+e}$

$$\vec{w}^{(1)} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} - \frac{1}{2+e} \begin{pmatrix} 2+e \\ 2+e \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} - \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\eta^{(1)} = \frac{1}{2+e} \longrightarrow \vec{w}^{(2)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} - \frac{1}{2+e} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Ex: $g(w_1, w_2) = (1 - w_1)^2 + 100(w_2 - w_1^2)^2$