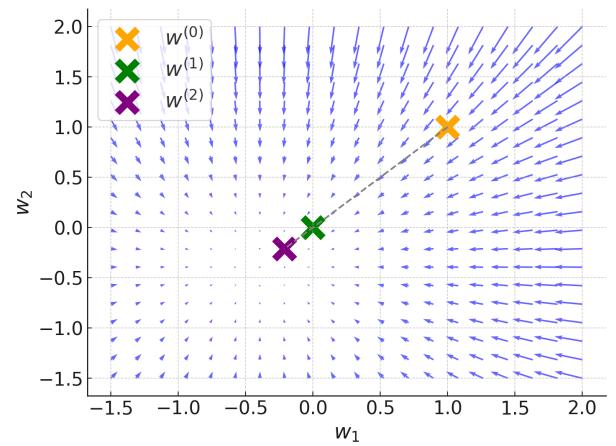
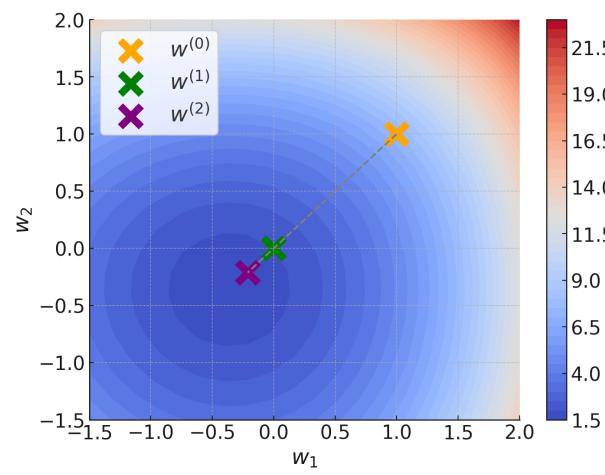
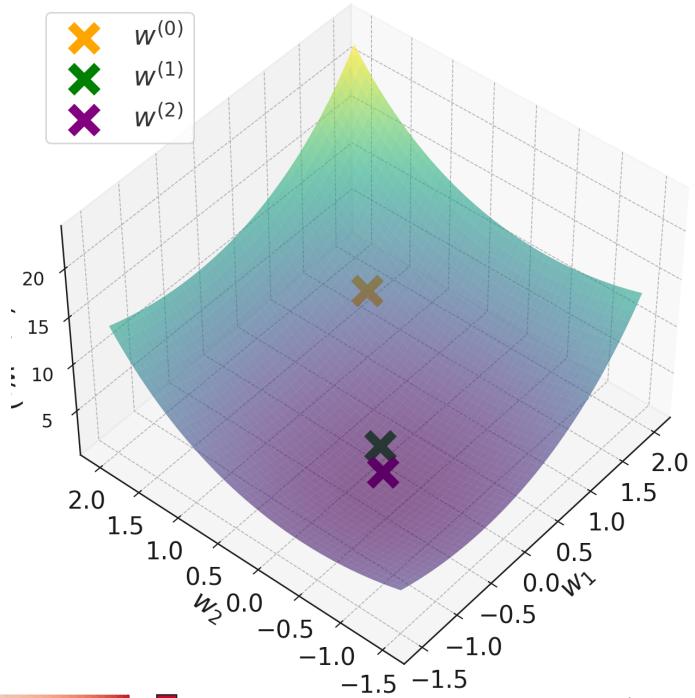


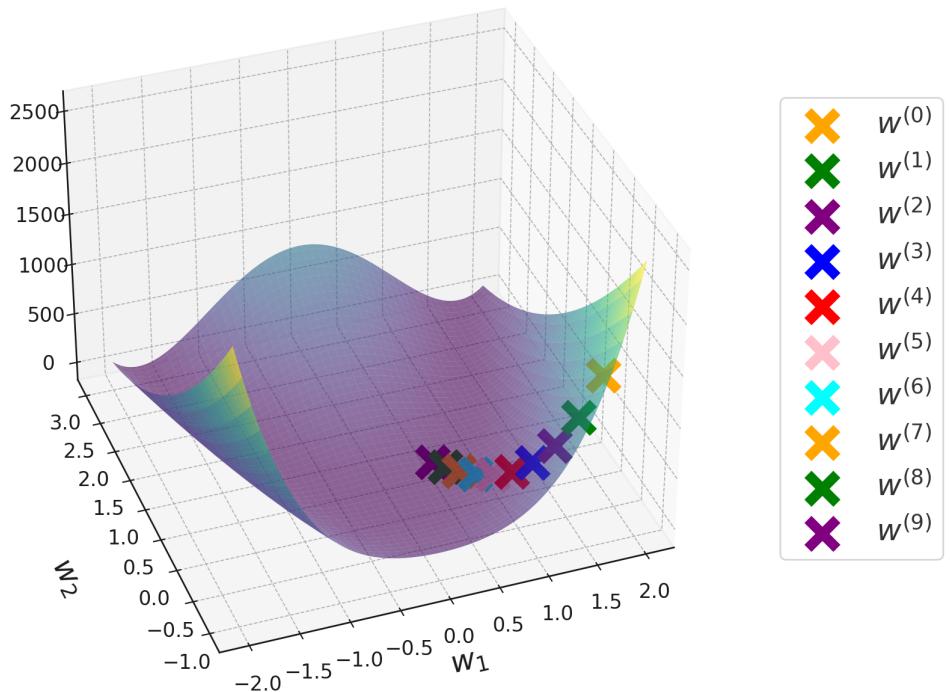
Important announcements

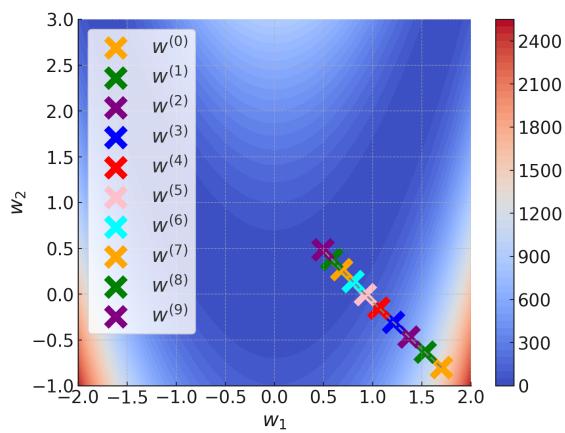
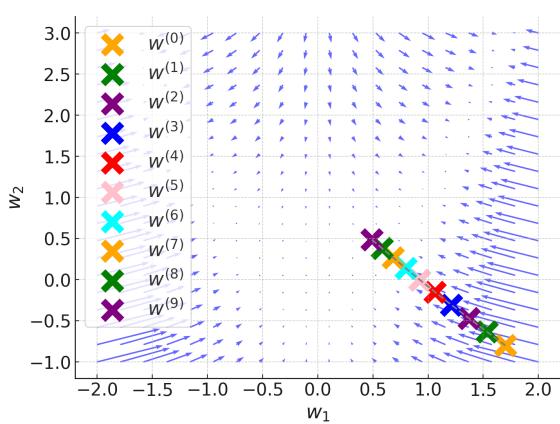
Feb 25

- grades for mid term 1 are out
- example for linear regression



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





Selecting the step size

goal: be close to $\frac{1}{g''(\vec{\omega}^{(t)})}$

Common options (for all $t \in \mathbb{N}$)

1) constant: $\eta^{(t)} = \eta \in (0, \infty)$

2) inverse decaying: $\eta^{(t)} = \frac{\eta}{1 + \lambda t}$ where $\eta, \lambda \in (0, \infty)$

3) exponential decaying: $\eta^{(t)} = \eta e^{-\lambda t}$ where $\eta, \lambda \in (0, \infty)$

4) Normalized gradient decaying $\eta^{(t)} = \frac{\eta}{\varepsilon + \|\nabla g(\vec{\omega}^{(t)})\|}$
 $\text{where } \|\nabla g(\vec{\omega}^{(t)})\| = \sqrt{\sum_{j=1}^d \left(\frac{\partial}{\partial w_j} g(\vec{\omega}^{(t)}) \right)^2}$
 $\varepsilon \text{ is small}$ (ex: $\varepsilon = 10^{-8}$)

Gradient Descent with $\hat{L}(f)$

$$D = \left((\vec{x}_1, y), \dots, (\vec{x}_n, y_n) \right) \in (\mathcal{X} \times \mathcal{Y})^n$$

$\mathcal{X} = \mathbb{R}^{d+1}$, $y = \mathbb{R}$ regression

$$\mathcal{F} = \left\{ f \mid f: \mathcal{X} \rightarrow \mathcal{Y} \text{ and } f(\vec{x}) = \vec{x}^T \vec{w}, \vec{w} \in \mathbb{R}^{d+1} \right\}$$

linear function class

$$l(\hat{y}, y) = (\hat{y} - y)^2$$

squared loss

Objective: find $\vec{\hat{w}} = \underset{\vec{w} \in \mathbb{R}^{d+1}}{\operatorname{argmin}} \hat{L}(\vec{w})$

$$\begin{aligned} \text{where } \hat{L}(\vec{w}) &= \frac{1}{n} \sum_{i=1}^n l(\vec{x}_i^T \vec{w}, y_i) \\ &= \frac{1}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{w} - y_i)^2 \end{aligned}$$

(Batch) Gradient Descent (BGD)

We know the closed form solution $\vec{w} = A^{-1} \vec{b}$ but try to solve the optimization step of ERM with gradient descent anyway

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

$$\begin{aligned} \nabla \hat{L}(\vec{w}) &= \left(\frac{\partial}{\partial w_0} \hat{L}(\vec{w}), \dots, \frac{\partial}{\partial w_d} \hat{L}(\vec{w}) \right)^T \in \mathbb{R}^{d+1} \\ &= \left(\frac{2}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{w} - y_i) \vec{x}_{i,1}, \dots, \frac{2}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{w} - y_i) \vec{x}_{i,d} \right)^T \\ &= \frac{2}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{w} - y_i) \underbrace{\left(\vec{x}_{i,1}, \dots, \vec{x}_{i,d} \right)^T}_{\vec{x}_i^T} \\ &= \frac{2}{n} \sum_{i=1}^1 (\vec{x}_i^T \vec{w} - y_i) \vec{x}_i^T \end{aligned}$$

Algorithm: BGD Linear Regression Learner (with const. step size)

input: $D = \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}, \eta, T$ {number of epochs}

$\vec{w} \in$ random vector in \mathbb{R}^{d+1} {step size}

for $t = 1, \dots, T$

$$\nabla \hat{L}(\vec{w}) \leftarrow \frac{2}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{w} - y_i) \vec{x}_i$$

$$\vec{w} \leftarrow \vec{w} - \eta \nabla \hat{L}(\vec{w})$$

return: $\hat{f}(\vec{x}) = \vec{x}^T \vec{w}$

The reason to still sometimes use GD is computation

computation: closed form solution vs gradient descent

goal: number of computational steps should be minimal

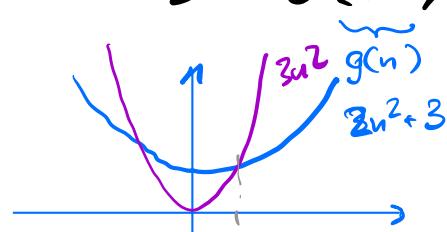
- adding
- subtracting
- multiplications

Big O notation: $f(n) = O(g(n))$ $f(n) < g(n)$ $n \geq n_0$

$$E_F: f(n) = 2n^2 + 3 = O(n^2)$$

$$2n^2 + 3 \leq 3n^2 \text{ for } n \geq 2$$

$$24 + 3 \leq 3 \cdot 4$$



Closed form: $\vec{w} = A^{-1} \vec{b}$ where $A = \sum_{i=1}^n \vec{x}_i \vec{x}_i^T$, $\vec{b} = \sum_{i=1}^n \vec{x}_i y_i$

$$\begin{aligned} \vec{A} &: O(d^2 n) \\ \vec{b} &: O(d n) \end{aligned}$$

$$\begin{aligned} \vec{A}^{-1} &: O(d^3) \\ \vec{A}^{-1} \vec{b} &: O(d^2) \end{aligned}$$

Total computation: $O(d^2 n + d n + d^3 + d^2) = O(d^2 n + d^3)$

$$\text{BGD} \quad \nabla \tilde{L}(\tilde{\omega}^{(t)}) : O(dn) \quad \tilde{w} : O(d)$$

\rightarrow Total computations : $O(dnT)$

Let's compare $O(d^2n + d^3)$ and $O(dnT)$

if $T < d$: $O(dnT) < O(d^2n) \leq O(d^2n + d^3)$
 $\# \text{epochs}$ $\# \text{features}$

\Rightarrow BGD is more computationally efficient if $T < d$

Can we do better in terms of computation?

Mini-batch gradient descent (MBGD)

$$D = \left((\vec{x}_1, y_1), \dots, (\vec{x}_b, y_b), \right. \\ \left(\vec{x}_{b+1}, y_{b+1} \right), \dots, \left(\vec{x}_{2b}, y_{2b} \right), \\ \vdots \\ \left. \left(\vec{x}_{(M-1)b+1}, y_{(M-1)b+1} \right), \dots, (\vec{x}_n, y_n) \right)$$

where b mini-batch size

$$M = \frac{n}{b} \quad \text{number of mini-batches}$$

- good: we want to break up the sum into smaller chunks
- \rightarrow mini-batches
- \oplus sum smaller \rightarrow more efficient
- \ominus estimate is worse

$$\underline{\underline{Ex}}: n = 8, b = 2, M = \frac{8}{2} = 4 \rightarrow \text{integer}$$

$$D = ((\vec{x}_1, y_1), (\vec{x}_2, y_2), \text{MB 1})$$

$$((\vec{x}_3, y_3), (\vec{x}_4, y_4), \text{MB 2})$$

$$((\vec{x}_5, y_5), (\vec{x}_6, y_6), \text{MB 3})$$

$$((\vec{x}_7, y_7), (\vec{x}_8, y_8), \text{MB 4})$$

$M = 4$
mini-batches

$\frac{n}{b}$ is not always an integer. We define

$$M = \text{floor}\left(\frac{n}{b}\right)$$

\rightarrow We will not use $n - Mb < b$ datapoints

New estimate of the expected loss for each mini-batch

$$\hat{L}_m = \frac{1}{b} \sum_{i=(m-1)b+1}^{mb} (\vec{x}_i^T \vec{w} - y_i)^2 \text{ where } m \in \{1, \dots, M\}$$

\rightarrow different estimate for each batch

\hat{L}_m is a sample mean estimate of $L(\vec{w})$ with b datapoints

$\rightarrow \hat{L}_m(\vec{w})$ is a worse estimate than $\hat{L}(\vec{w})$

Algorithm: MBGD Linear Regression Learner (with constant stepsize)

input: $D = ((\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n))$, γ , T , b

$\vec{w} \leftarrow$ random vector in \mathbb{R}^{d+1}

$$M \leftarrow \text{floor}\left(\frac{n}{b}\right)$$

for $t = 1, \dots, T$

randomly shuffle D

for $m = 1, \dots, M$

$$\nabla \hat{L}_m^b(\vec{w}) \leftarrow \frac{2}{b} \sum_{i=mb+1}^{(m+1)b} (\vec{x}_i^\top \vec{w} - y_i) \vec{x}_i$$

$$\vec{w} \leftarrow \vec{w} - \gamma \nabla \hat{L}_m^b(\vec{w})$$

return $\hat{f}(\vec{x}) = \vec{x}^\top \vec{w}$
