

Lecture 4: Least Squares and Optimization

Example

Least Squares & Classification

Setup: n training samples, $(\underline{x}_i, y_i) \in \mathbb{R}^p \times \{-1, +1\}$ for $i=1, \dots, n$

$$\text{let } X = \begin{bmatrix} - \underline{x}_1^T - \\ - \underline{x}_2^T - \\ \vdots \\ - \underline{x}_n^T - \end{bmatrix}, \quad \underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Assume $n \geq p$, X is rank- p . Then

compute $\hat{\underline{w}} = \underset{\underline{w}}{\operatorname{argmin}} \|\underline{y} - X\underline{w}\|_2^2$,

$\tilde{\underline{y}} = X\hat{\underline{w}} \in \mathbb{R}^n$ (i.e. not ± 1 or -1 labels)

Classification rule: predict $+1$ if $\tilde{y}_i > 0$ and -1 if $\tilde{y}_i < 0$

$$\Rightarrow \hat{y}_i = \operatorname{sign}(\tilde{y}_i)$$

$$\hat{\underline{y}} = \operatorname{sign}(\tilde{\underline{y}})$$

for a new sample $\underline{x}_{\text{new}} \in \mathbb{R}^p$, want to predict new, unknown label y_{new} .

$$\tilde{y}_{\text{new}} = \langle \underline{x}_{\text{new}}, \hat{\underline{w}} \rangle \in \mathbb{R}$$

$$\hat{y}_{\text{new}} = \operatorname{sign}(\tilde{y}_{\text{new}})$$

alternatively, we might be tempted to consider

$$\hat{\underline{w}} = \underset{\underline{w}}{\operatorname{argmin}} \|\underline{y} - \operatorname{sign}(X\underline{w})\|_2^2$$

but this is very hard to solve computationally

Optimization approach

$$\hat{\underline{w}} = \text{"argument } \underline{w} \text{ that minimizes"} \sum_{i=1}^n r_i^2(\underline{w})$$

$$= \operatorname{argmin}_{\underline{w}} \sum_{i=1}^n r_i^2(\underline{w})$$

$$= \operatorname{argmin}_{\underline{w}} \sum_{i=1}^n \left(y_i - \sum_{j=1}^p w_j x_{i,j} \right)^2$$

$$= \operatorname{argmin}_{\underline{w}} \|\underline{r}\|_2^2$$

$$= \operatorname{argmin}_{\underline{w}} \underbrace{\|y - X\underline{w}\|_2^2}_{f(\underline{w})}$$

$$= \operatorname{argmin}_{\underline{w}} \underline{y}^T \underline{y} - \underline{y}^T X \underline{w} - \underline{w}^T X^T \underline{y} + \underline{w}^T X^T X \underline{w}$$

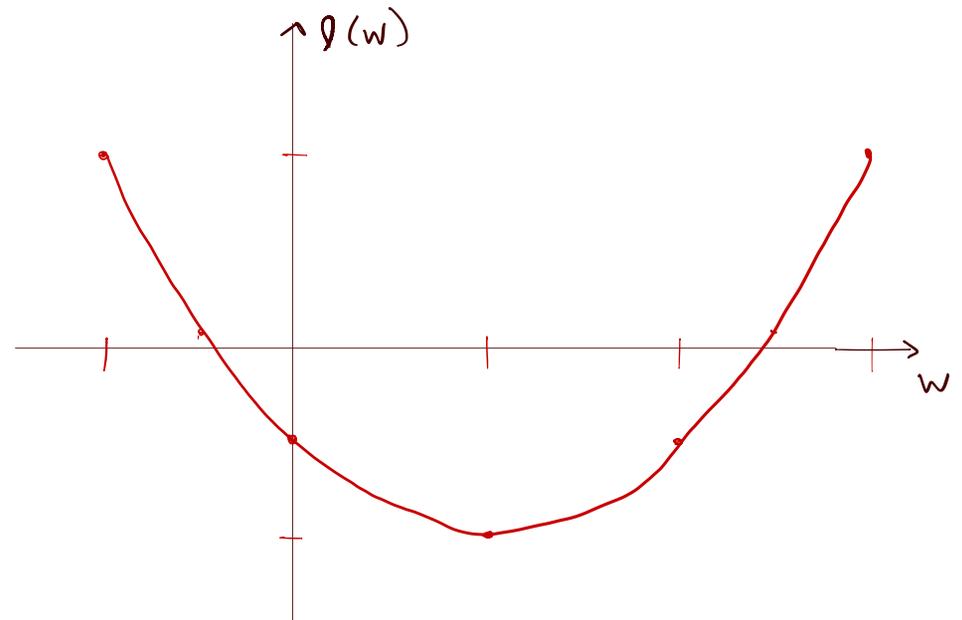
2-norm or
Euclidean norm:
 $\|\underline{a}\|_2 = \left(\sum_{i=1}^n a_i^2 \right)^{1/2}$

Warmup.

$$l(w) = \frac{1}{2} w^2 - w - \frac{1}{2}$$

$$\hat{w} = \operatorname{argmin}_w l(w)$$

$$\frac{dl}{dw} = 2 \cdot \frac{1}{2} w - 1 = 0 \Rightarrow \hat{w} = 1$$



Positive definite matrices

From the geometric perspective, we saw that it was important for finding a unique least squares solution $\hat{\underline{w}}$ that $X^T X$ be invertible.

Is this important in the optimization setting as well? YES!

The following two things are equivalent for $X \in \mathbb{R}^{n \times p}$ with $n \geq p$, $\text{rank}(X) = p$
(X has p linearly independent columns)

① $X^T X \in \mathbb{R}^{p \times p}$ is invertible ($(X^T X)^{-1}$ exists)

② $X^T X$ is positive definite

A matrix Q is positive definite (p.d.) if

$$\underline{w}^T Q \underline{w} > 0 \quad \text{for all } \underline{w} \neq 0$$

Short hand: $Q \succ 0$

A matrix Q is positive semi-definite (p.s.d.) if

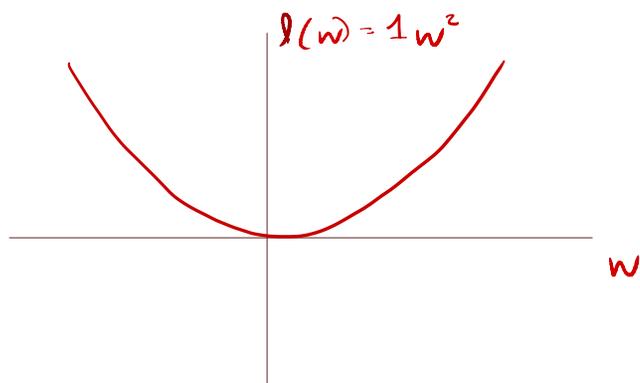
$$\underline{w}^T Q \underline{w} \geq 0 \quad \text{for all } \underline{w} \neq 0$$

Short hand: $Q \succeq 0$

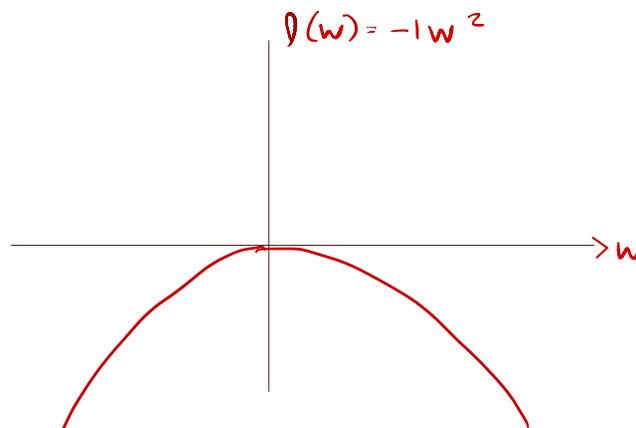
This does not
mean all
elements of
 Q are positive
or non-negative!

Ex: $w \in \mathbb{R}, Q \in \mathbb{R} \Rightarrow w^T Q w = Q w^2 > 0$ if $Q > 0$

imagine trying to minimize $J(w) = Q w^2$



easy to minimize
(convex)

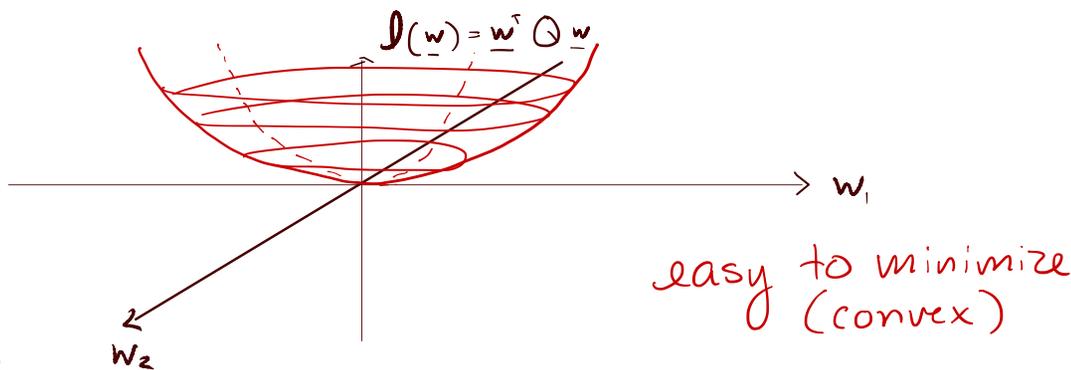


hard to minimize

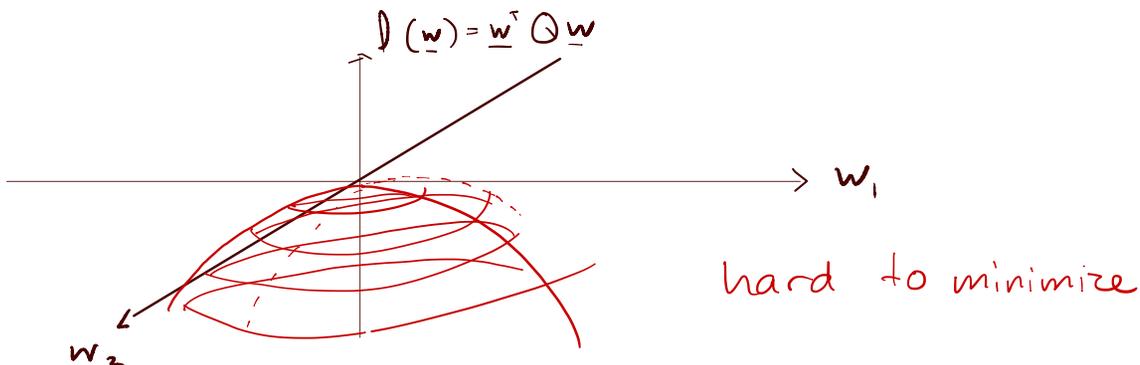
Ex: $\underline{w} \in \mathbb{R}^2, Q \in \mathbb{R}^{2 \times 2}$

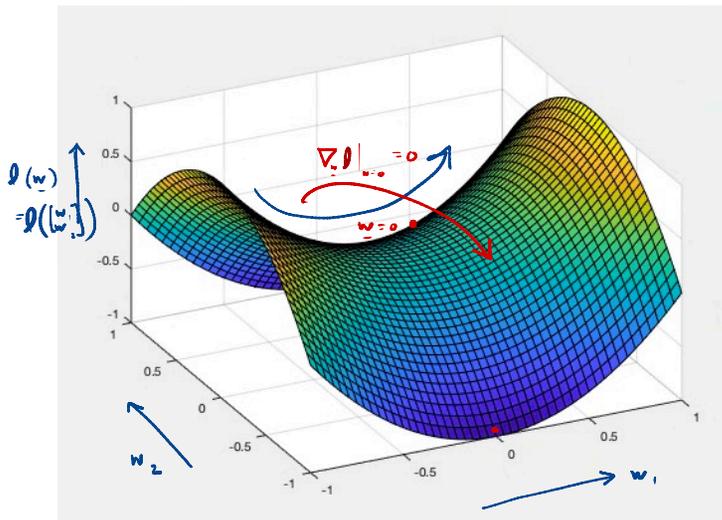
$$Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \Rightarrow$$

$$Q > 0 \text{ b/c } \underline{w}^T Q \underline{w} = w_1^2 + w_2^2 > 0$$



$$Q = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix} \Rightarrow$$



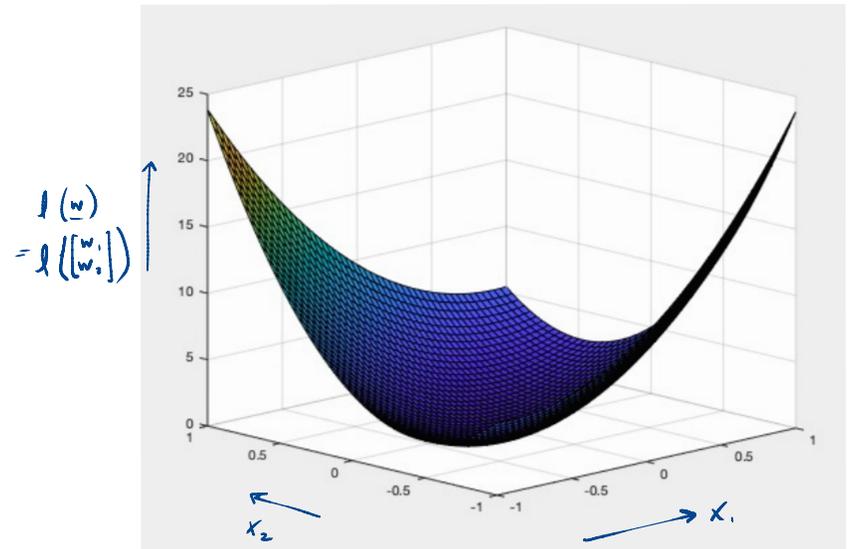


$$Q = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

$$l(\underline{w}) = \underline{w}^T Q \underline{w}$$

this an example of a "saddle point"

Q is NOT positive definite



$$Q = \begin{bmatrix} 5 & -5 \\ -5 & 1 \end{bmatrix}$$

$$l(\underline{w}) = \underline{w}^T Q \underline{w}$$

$$l(\underline{w}) > 0 \text{ for } \underline{w} \neq 0$$

$\Rightarrow Q$ is positive definite

Recall Q is pos. def.
if $\underline{w}^T Q \underline{w} > 0$ for
all $\underline{w} \neq 0$

Properties of Positive Definite Matrices

1) if $P > 0$ and $Q > 0$, then $P+Q > 0$

$$\underline{w}^T P \underline{w} > 0 \text{ \& } \underline{w}^T Q \underline{w} > 0 \Rightarrow \underline{w}^T (P+Q) \underline{w} = \underline{x}^T P \underline{x} + \underline{x}^T Q \underline{x} > 0$$

2) if $Q > 0$ and $a > 0$, then $aQ > 0$

$$\underline{w}^T Q \underline{w} > 0 \Rightarrow \underline{w}^T (aQ) \underline{w} = a \underline{w}^T Q \underline{w} > 0$$

3) for any A , $A^T A \geq 0$ and $AA^T \geq 0$

if the columns of A are linearly independent, then $A^T A > 0$

Note $\underline{w}^T \underline{w} \geq 0$ always, and $\underline{w}^T \underline{w} = 0$ only if $\underline{w} = 0$.

$$\text{Let } \tilde{\underline{w}} := A \underline{w}$$

$$\underline{w}^T A^T A \underline{w} = \tilde{\underline{w}}^T \tilde{\underline{w}} \geq 0. \quad \text{now } \tilde{\underline{w}}^T \tilde{\underline{w}} = 0 \text{ only if } \tilde{\underline{w}} = A \underline{w} = 0. \quad A \underline{w} = 0 \text{ if either (a) } \underline{w} = 0$$

or (b) columns of A are linearly dependent.

4) if $Q > 0$, then Q^{-1} exists (consider special case where Q is diagonal $Q = \begin{bmatrix} q_1 & & 0 \\ & q_2 & \\ 0 & & \ddots \\ & & & q_n \end{bmatrix}$)
 $\Rightarrow Q > 0 \iff \text{all } q_i > 0$

5) $Q > P$ means $Q - P > 0$

If $Q = X^T X$, is Q positive definite?

For Q to be P.D., need $\underline{w}^T Q \underline{w} > 0$ for all $\underline{w} \neq 0$

$$\text{now } \underline{w}^T Q \underline{w} = \underline{w}^T X^T X \underline{w} = \tilde{\underline{w}}^T \tilde{\underline{w}} \geq 0 \quad \text{where we define } \tilde{\underline{w}} := X \underline{w}$$

For Q to be positive definite (not positive semi definite, where $\underline{w}^T Q \underline{w}$ might = 0)

we need to ensure $\tilde{\underline{w}}^T \tilde{\underline{w}} > 0$

Now $\tilde{\underline{w}}^T \tilde{\underline{w}} = 0$ only if $\tilde{\underline{w}} = 0$ (because $\tilde{\underline{w}}^T \tilde{\underline{w}} = \tilde{w}_1^2 + \tilde{w}_2^2 + \dots + \tilde{w}_p^2$)

$\tilde{\underline{w}} = 0$ if $X \underline{w} = 0$. Can $X \underline{w} = 0$ for some $\underline{w} \neq 0$?

Recall $X \underline{w}$ = weighted sum of columns of X . I.e. $X \underline{w} = w_1 \underline{X}_1 + w_2 \underline{X}_2 + \dots + w_p \underline{X}_p$.

Recall defn of linear independence: the columns of X are L.I. if no weighted sum = 0 (unless all weights = 0)

If cols. of X are L.I., then $X \underline{w} \neq 0$ unless $\underline{w} = 0$

$\Rightarrow \underline{w}^T X^T X \underline{w} > 0 \Rightarrow X^T X$ is pos definite

Least squares optimization problem

$$\hat{\underline{w}} = \underset{\underline{w}}{\operatorname{argmin}} l(\underline{w}) \quad \text{where } l(\underline{w}) = \underline{y}^T \underline{y} - \underline{y}^T X \underline{w} - \underline{w}^T X^T \underline{y} + \underline{w}^T X^T X \underline{w}$$

$l(\underline{w})$ maps $\underline{w} \in \mathbb{R}^p$ to \mathbb{R}

Assume $l(\underline{w})$ is convex (more on this later):

When w is a scalar, we set derivative $\frac{dl}{dw}$ to zero and solve for w .

When \underline{w} is a vector, we set gradient $\nabla_{\underline{w}} l$ to zero and solve for \underline{w}

$$\nabla_{\underline{w}} l := \begin{bmatrix} df/dw_1 \\ df/dw_2 \\ \vdots \\ df/dw_p \end{bmatrix}$$

$$\text{Ex. } l(\underline{w}) = \underline{w}^T \underline{c} = c_1 w_1 + c_2 w_2 + \dots + c_p w_p$$

$$\nabla_{\underline{w}} l = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_p \end{bmatrix} = \underline{c}$$

"Rule 1"

$$\text{Ex. } l(\underline{w}) = \|\underline{w}\|^2 = \underline{w}^T \underline{w} = w_1^2 + w_2^2 + \dots + w_p^2$$

$$\nabla_{\underline{w}} l = \begin{bmatrix} 2w_1 \\ 2w_2 \\ \vdots \\ 2w_p \end{bmatrix} = 2\underline{w}$$

$$\text{Ex. } l(\underline{w}) = \underline{w}^T Q \underline{w} = \sum_{i=1}^P \sum_{j=1}^P w_i Q_{ij} w_j$$

$$\frac{d(w_i Q_{ij} w_j)}{dw_k} = \begin{cases} 2Q_{kk} w_k & \text{if } k=i=j \\ Q_{kj} w_j & \text{if } i=k \neq j \\ Q_{ik} w_i & \text{if } j=k \neq i \end{cases}$$

$$\frac{dl}{dw_k} = \sum_{i=1}^P \sum_{j=1}^P \frac{d(w_i Q_{ij} w_j)}{dw_k}$$

$$\Rightarrow \nabla_{\underline{w}} l = Q \underline{w} + Q^T \underline{w}.$$

if Q is symmetric (ie. $Q = Q^T$), then

$$\nabla_{\underline{w}} \underline{w}^T Q \underline{w} = 2Q \underline{w}$$

"Rule 2"

$$\text{Ex. Least squares: } l(\underline{w}) = \underline{y}^T \underline{y} - 2 \underline{w}^T X^T \underline{y} + \underline{w}^T X^T X \underline{w}$$

$$l(\underline{w}) = \sum_{i=1}^n (y_i - \underline{x}_i^T \underline{w})^2 = \left\| \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} -x_1^T \\ -x_2^T \\ \vdots \\ -x_n^T \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_r \end{bmatrix} \right\|_2^2 = \| \underline{y} - X \underline{w} \|_2^2$$

$$= (\underline{y} - X \underline{w})^T (\underline{y} - X \underline{w})$$

$$= \underline{y}^T \underline{y} - \underline{y}^T X \underline{w} - \underline{w}^T X^T \underline{y} + \underline{w}^T X^T X \underline{w}$$

$$= \underline{y}^T \underline{y} - \underbrace{2 \underline{w}^T X^T \underline{y}}_{\text{use Rule 1}} + \underbrace{\underline{w}^T X^T X \underline{w}}_{\text{use Rule 2}} \rightarrow \underbrace{X^T X}_{\text{symmetric}} \text{ here is like } Q \text{ in Rule 2}$$

$$\nabla_{\underline{w}} l = 0 - 2X^T \underline{y} + 2X^T X \underline{w}$$

$$\nabla_{\underline{w}} \ell = 0 - 2X^T y + 2X^T X \underline{w}$$

If $Q = X^T X > 0$, then we can compute gradient and set it to zero because f is convex

$$\nabla_{\underline{w}} \ell = 0$$

$$\Rightarrow X^T X \hat{\underline{w}} = X^T y$$

$$\Rightarrow \hat{\underline{w}} = (X^T X)^{-1} X^T y \quad (\text{what we got from geometric perspective!})$$

So far, we've assumed $X \in \mathbb{R}^{n \times p}$ has $n \geq p$ and p cols of X are L.I.

- Columns of X are L.I. $\Rightarrow Q = X^T X > 0$ (i.e. $X^T X$ is positive definite)
- $X^T X$ is pos definite \Rightarrow least square loss $y^T y - 2\underline{w}^T X^T y + \underline{w}^T X^T X \underline{w}$ is convex
- if $X^T X > 0$, then $X^T X$ has inverse
- $X^T X$ has inverse $\Rightarrow \hat{\underline{w}} = (X^T X)^{-1} X^T y$ exists and is unique