

Lecture 4: Least squares and optimization

Mathematical Foundations of Machine Learning
University of Chicago

Given: vector of labels $y \in \mathbb{R}^n$

matrix of features $X \in \mathbb{R}^{n \times p}$

Want: vector of weights $w \in \mathbb{R}^p$

Assume: $n \geq p$

$\text{Rank}(X) = p$ (X has p

linearly independent columns)

If $y = Xw$, then we have a system of n linear equations;

i^{th} equation: $y_i = w_1 x_{i1} + w_2 x_{i2} + \dots + w_p x_{ip}$

$$= \sum_{j=1}^p w_j x_{ij} = \langle w, x_i \rangle$$

↑ i^{th} row of X (transposed)

In general, $y \neq Xw$ for any w (because of modeling errors, noise)

Define residual $r_i = r_i(w) = y_i - \langle w, x_i \rangle$

$$\left. \begin{array}{c} n \\ \left\{ \begin{array}{c} X \\ \underbrace{\hspace{1cm}}_p \end{array} \right. \end{array} \right\} = \left. \begin{array}{c} w \\ y \end{array} \right\} \quad \leftarrow \begin{array}{l} n \text{ equations} \\ p \text{ unknowns} \end{array}$$

let $x_i \in \mathbb{R}^p$ be feature vector of i^{th} sample;
 $= (i^{\text{th}} \text{ row of } X)^T$

let $x_j \in \mathbb{R}^n$ be j^{th} feature for all samples
 $= j^{\text{th}}$ col of X

LEAST SQUARES ESTIMATION: if columns of X are linearly independent

find w to minimize $\sum_{i=1}^n |r_i(w)|^2 \Rightarrow \hat{w} = (X^T X)^{-1} X^T y = X^\# y \quad (X^\# = (X^T X)^{-1} X^T = \text{"pseudoinverse"})$

$\underbrace{\hspace{1cm}}_{= \|r\|^2} \quad \hat{y} = X\hat{w} = X(X^T X)^{-1} X^T y = P_X y = \text{projection of } y \text{ onto range}(X)$

Least Squares & Classification

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

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

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

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

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

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}

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

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

Optimization approach

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

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

$$= \underset{\underline{w}}{\operatorname{argmin}} \|\Sigma\|_2^2$$

$$= \underset{\underline{w}}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \sum_{j=1}^p w_j x_{ij})^2$$

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

$f(\underline{w})$

$$= \underset{\underline{w}}{\operatorname{argmin}} \underline{y}^\top \underline{y} - \underline{y}^\top \underline{X}\underline{w} - \underline{w}^\top \underline{X}^\top \underline{y} + \underline{w}^\top \underline{X}^\top \underline{X}\underline{w}$$

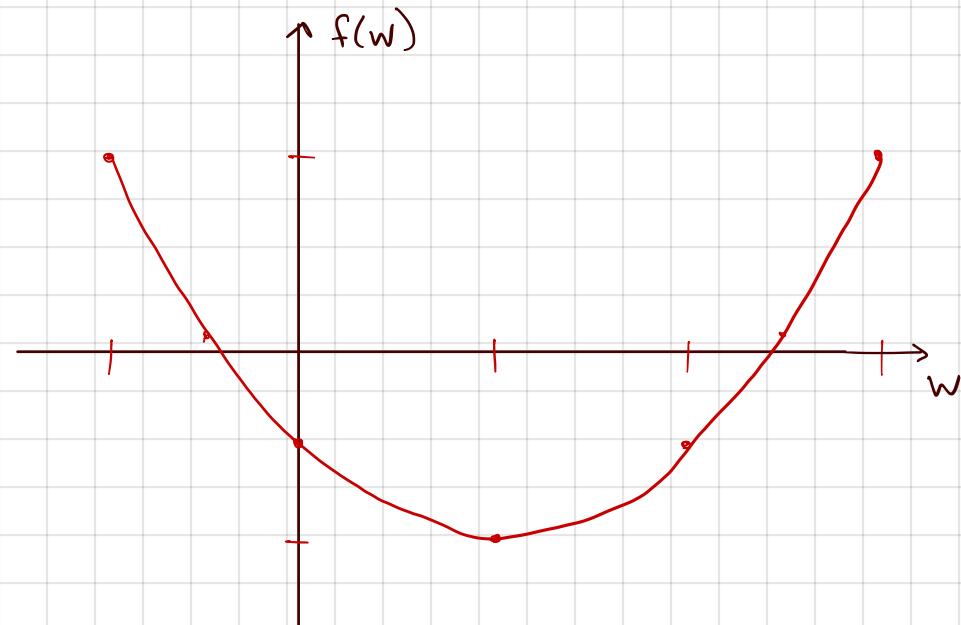
2-norm or
Euclidean norm:
 $\|\underline{a}\|_2 := \left(\sum_{i=1}^n a_i^2 \right)^{1/2}$

Warmup.

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

$$\hat{w} = \underset{w}{\operatorname{argmin}} f(w)$$

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



Positive definite matrices

From the geometric perspective, we saw that it was important for finding a unique least squares solution \hat{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

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

Shorthand: $Q > 0$

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

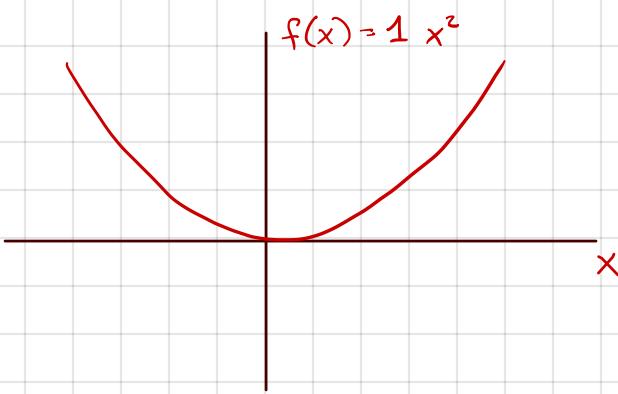
This does not
mean all
elements of
 Q are positive
or non-negative!

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

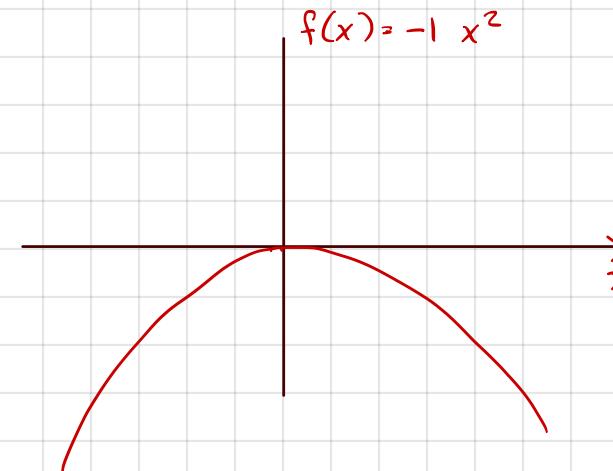
Shorthand: $Q \succeq 0$

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

imagine trying to minimize $f(\underline{x}) = Q \underline{x}^2$



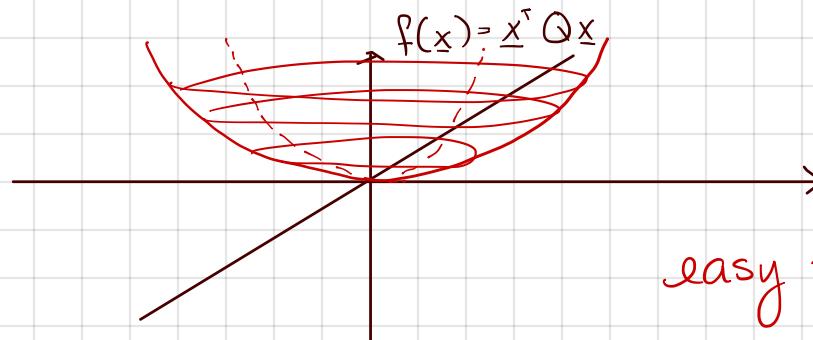
easy to minimize
(convex)



hard to minimize

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

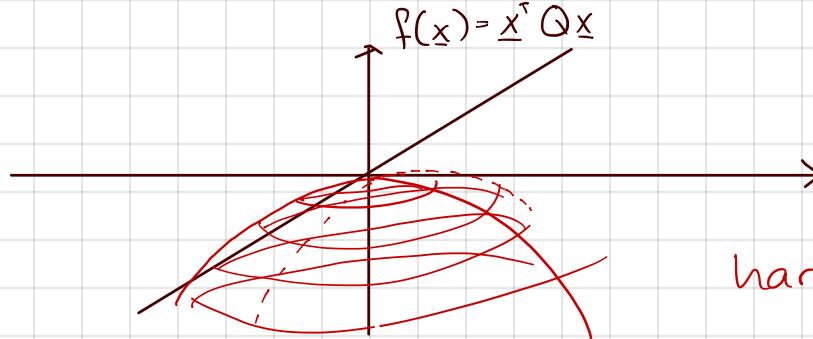
$$Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \Rightarrow$$



easy to minimize
(convex)

$$Q > 0 \text{ b/c } \underline{x}^T Q \underline{x} = x_1^2 + x_2^2 > 0$$

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



hard to minimize

Properties of Positive Definite Matrices

1) if $P \succ 0$ and $Q \succ 0$, then $P+Q \succ 0$

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

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

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

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

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

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

Let $\tilde{\underline{x}} := A \underline{x}$

$$\underline{x}^T A^T A \underline{x} = \tilde{\underline{x}}^T \tilde{\underline{x}} \geq 0. \quad \text{now } \tilde{\underline{x}}^T \tilde{\underline{x}} = 0 \text{ only if } \tilde{\underline{x}} = A \underline{x} = 0. \quad A \underline{x} = 0 \text{ if either } @ \underline{x} = 0$$

or (b) columns of A are linearly dependent.

4) if $A \succ 0$, then A^{-1} exists

5) $Q \succ P$ means $Q - P \succ 0$

Least squares optimization problem

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

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

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

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

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

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

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

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

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

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

$$\text{Ex. } f(\underline{w}) = \underline{w}^\top Q \underline{w} = \sum_{i=1}^P \sum_{j=1}^P w_i Q_{ij} w_j \Rightarrow \nabla_{\underline{w}} f = Qx + Q^\top x.$$

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

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

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

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

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

if $Q > 0$, then we can compute gradient and set it to zero because f is convex

$$\nabla_{\underline{w}} f = \underline{0} - 2X^\top \underline{y} + 2X^\top X \underline{w} = 0$$

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

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