

EE2211 Lecture 5: Linear Regression & Least Squares.

- Agenda:
- i) Notation for sets & functions / linear & affine functions.
 - ii) Min & Argmin / Differentiation of functions (Grad).
 - iii) Linear Regression
 - Formulation / Objective (Loss) Function / Training / Prediction
 - iv) Predicting multiple outputs.

Reading : Lec_5.pdf Slides.

Office hours: Next Thursday 16 Feb 2pm

Zoom 9839098564

Sets: A set is an unordered collection of objects.

Eg:

$S = \{1, 2, 3, 4, 5, 6\}$ possible outcomes of a dice toss.

$\|S' = \{6, 4, 5, 3, 2, 1\}$

Eg: $S = [a, b] = \{x \in \mathbb{R}; a \leq x \leq b\}$ } intervals of \mathbb{R}
 $S = (a, b) = \{x \in \mathbb{R}; a < x \leq b\}$

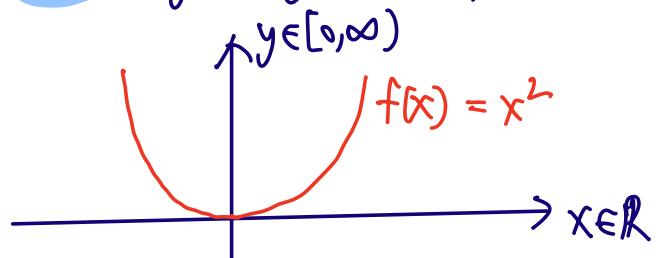
Eg: \mathbb{R} : set of real numbers
 \mathbb{R}^d : set of d -dim real vectors
 $\mathbb{R}^{m \times d}$: set of all size $m \times d$ matrices.

column
rows
 d columns.

Functions: A function f is a map from a set X to another set Y .

$$f: X \rightarrow Y.$$

Eg: $f: \mathbb{R} \rightarrow [0, \infty)$ given by the recipe $f(x) = x^2$



Set of inputs X : domain

Set of possible outputs Y : codomain

Set $\{f(x); x \in X\}$: range or image.

$$f: \underbrace{\{1, 2, 3, 4\}}_{\text{domain}} \rightarrow \underbrace{\{1, 2, 3, 4, 5, 6\}}_{\text{codomain}} \quad \text{range or image.} \quad f(x) = x + 2.$$

Range or image of f is $\{f(x); x \in X\} = \{3, 4, 5, 6\}$

Linear functions: A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is linear if

(i) [Homogeneity] For every vector $x \in \mathbb{R}^d$ & scalar $a \in \mathbb{R}$

$$f(ax) = a f(x),$$

(ii) [Additivity] For every $x, y \in \mathbb{R}^d$,

$$f(x+y) = f(x) + f(y).$$

Question: A linear function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ must pass through the origin $\underline{0} \in \mathbb{R}^d$, i.e.,

$$f(\underline{0}) = 0 \in \mathbb{R}$$

Question: If we have n vectors $\underline{x}_1, \dots, \underline{x}_n \in \mathbb{R}^d$ and n scalars $a_1, \dots, a_n \in \mathbb{R}$, and $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is linear,

$$f\left(\sum_{i=1}^n a_i \underline{x}_i\right) = \sum_{i=1}^n a_i f(\underline{x}_i).$$

Question: $f: \mathbb{R} \rightarrow \mathbb{R}$ is specified by $f(x) = |x|$. Is f linear?

No: Take $x=\underline{x}$, $a=-1$ $f(ax) \stackrel{?}{=} af(x)$

$$f(ax) = |ax| = |x|, \quad af(x) = -|x|$$

Question: Fix $\underline{a} \in \mathbb{R}^d$. Consider the function $f: \mathbb{R}^d \rightarrow \mathbb{R}$

$$f(\underline{x}) = \underline{a}^T \underline{x} = \sum_{i=1}^d a_i x_i, \quad \underline{x} = (x_1, \dots, x_d)$$

f is called the inner product / dot product between \underline{a} & \underline{x} .

Show that the inner product is linear.

Question: Why does a linear f^2 have to pass through $\underline{0} \in \mathbb{R}^d$?

Hence: $\forall a \in \mathbb{R}, \underline{x} \in \mathbb{R}^d, \quad f(a\underline{x}) = af(\underline{x})$

Take $\underline{x} = \underline{x}$, $a=0$ $f(0\underline{x}) = 0 f(\underline{x}) = 0$

$f(0)$

Affine Function: An affine function is a linear function plus a constant.

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ is affine if $f(\underline{x}) = \underline{a}^\top \underline{x} + b$ for some $\underline{a} \in \mathbb{R}^d$, $b \in \mathbb{R}$.

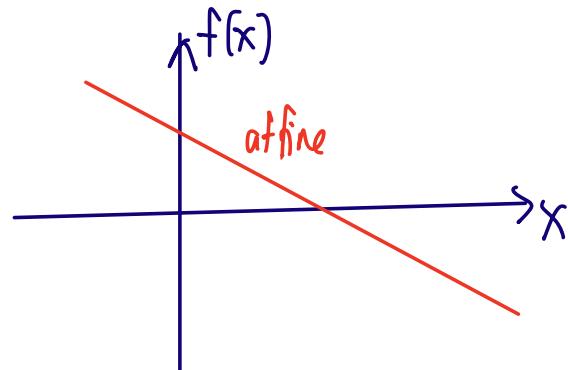
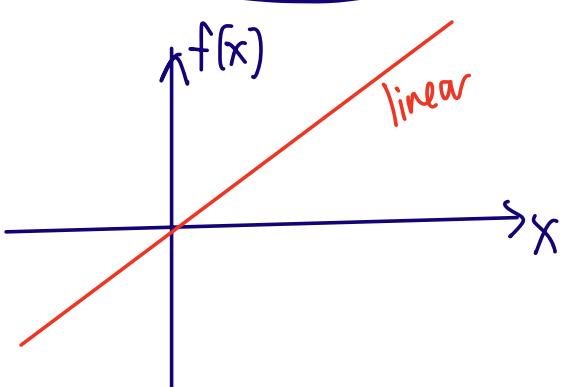
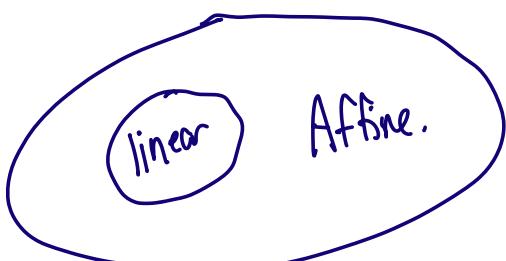
bias offset

Eg: $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ $f(\underline{x}) = f(x_1, x_2) = -x_1 + 3x_2 + 7$ is affine.
 $= \begin{pmatrix} -1 \\ 3 \end{pmatrix}^\top \underline{x} + 7$, $\underline{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

Question: Is a linear f^1 affine? Yes. (Take $b=0$)
 Is an affine f^2 linear? Not necessarily.

A linear f^2 $f(\underline{0}) = 0$.

Look at example $f(x_1, x_2) = -x_1 + 3x_2 + 7$. This function
 $f(0, 0) = 7 \neq 0$. is not linear



Min & Argmin

$$f: [a, b] \rightarrow \mathbb{R}$$

The function f has a local min at $c \in \mathbb{R}$ if

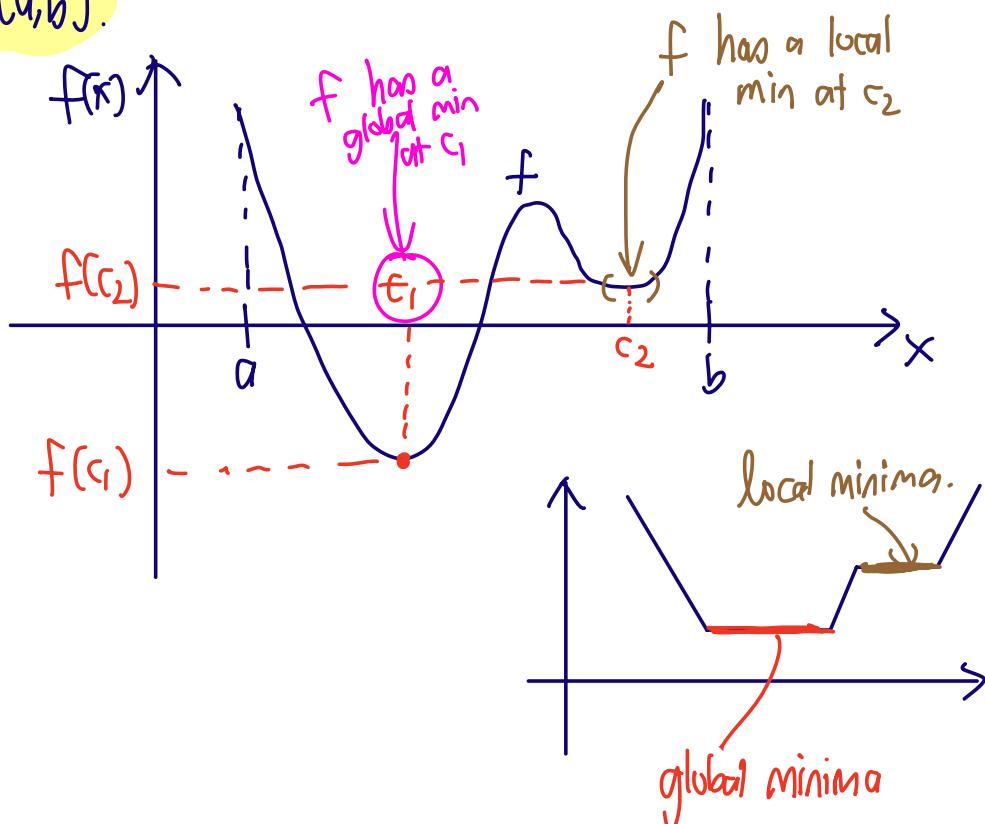
$$f(x) \geq f(c)$$

for all x in some neighborhood of c .

The function f has a global min at $c \in \mathbb{R}$ if

$$f(x) \geq f(c)$$

for all $x \in [a, b]$.



Minimum & argmin.

$f: X \rightarrow Y$, the min $\min_{x \in X} f(x)$ return the smallest value among all elements in the set $\{f(x): x \in X\}$

$f: X \rightarrow Y$, the $\text{argmin } x^* = \underset{x \in X}{\text{argmin}} f(x)$ returns the value $x \in X$ that minimizes $f(x)$.

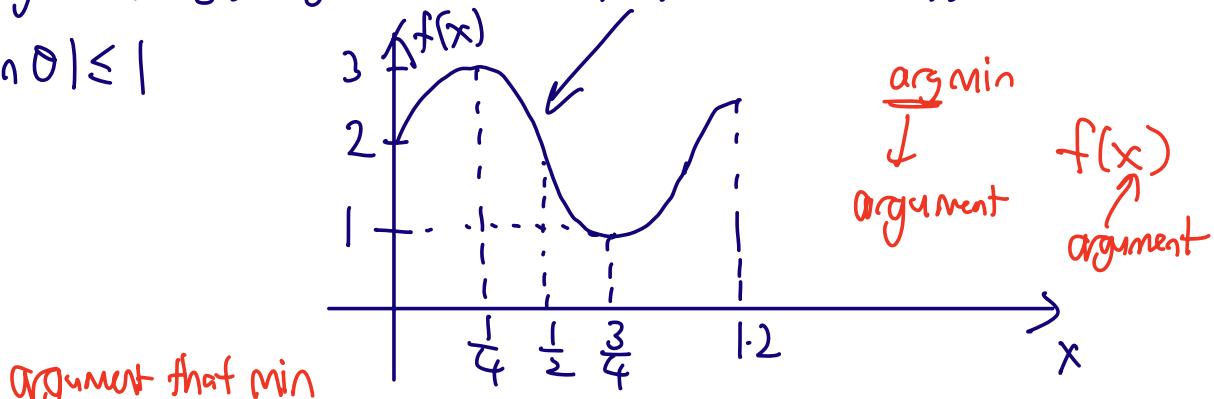
$$f(x^*) = \underset{x \in X}{\min} f(x).$$

Ex: $X = \{0, 1\}$, $f: X \rightarrow \mathbb{R}$, $f(0) = 5$, $f(1) = 7$

$$\underset{x \in X}{\min} f(x) = 5. \quad \underset{x \in X}{\max} f(x) = 7$$

$$\underset{x \in X}{\text{argmin}} f(x) = 0 \quad \underset{x \in X}{\text{argmax}} f(x) = 1.$$

Eg: $f: [0, 1.2] \rightarrow \mathbb{R}$
 $| \sin \theta | \leq 1$



$$\underset{x \in X}{\text{argmin}} f(x) = \frac{3}{4}$$

$$\underset{x \in X}{\min} f(x) = 1.$$

Derivatives or gradient $f: \mathbb{R}^d \rightarrow \mathbb{R}$ $f(\underline{x}) = y$
 d -dim.

Gradient vector / derivative of $\underline{x} \in \mathbb{R}^d$ is the vector

$$\nabla f(\underline{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix} \in \mathbb{R}^d$$

$$f(x_1, x_2) = 2x_1^2 + 5x_1x_2 + 3x_2^2$$

$d=2$

$$\frac{\partial f}{\partial x_1} = 4x_1 + 5x_2 \quad \frac{\partial f}{\partial x_2} = 5x_1 + 9x_2^2$$

$$\nabla f(\underline{x}) = \begin{bmatrix} 4x_1 + 5x_2 \\ 5x_1 + 9x_2^2 \end{bmatrix}$$

Example: Fix $\underline{q} \in \mathbb{R}^d$, $f(\underline{x}) = \underline{q}^\top \underline{x} = \sum_{i=1}^d q_i x_i$
 $\nabla f(\underline{x}) = \underline{q}$

Example: Fix $A \in \mathbb{R}^{d \times d}$, $f(\underline{x}) = \underline{x}^\top A \underline{x} = \sum_{i=1}^d \sum_{j=1}^d x_i x_j A_{ij}$

$$\nabla f(\underline{x}) = (A + A^\top) \underline{x}$$

Most of the time (in this class) A is symmetric ($A = A^\top$)
 $\Rightarrow \nabla f(\underline{x}) = 2A\underline{x}$.

Linear Regression is linear approach for modelling the r/s between a scalar response $y \in \mathbb{R}$ and explanatory variables (attributes, features) $\underline{x} \in \mathbb{R}^d$.

Dataset $\{(x_i, y_i) : 1 \leq i \leq m\}$, $x_i \in \mathbb{R}^d$ (training vector)
 $y_i \in \mathbb{R}$ (target).

Without offset

Form the design matrix

$$\underline{X} = \begin{bmatrix} \underline{x}^T \\ \vdots \\ \underline{x}_m^T \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,d} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m,1} & x_{m,2} & \cdots & x_{m,d} \end{bmatrix}$$

target vector

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^{m \times 1}$$

\cap

$$\mathbb{R}^{m \times d}$$

Goal: Find a solution (or approx solution) $\underline{w} \in \mathbb{R}^d$ to the linear system

$$\underline{X} \underline{w} = \underline{y}$$

With offset: Design a function/model/regressor $f_{\underline{w}, b}$ linear comb. of features in \underline{x} plus possibly a constant.

$$f_{\underline{w}, b}(\underline{x}) = \underline{w}^T \underline{x} + b$$

\leftarrow offset
does not have to be constrained to pass thru origin.

$\underline{w} \in \mathbb{R}^d$: d-dim unknown weight vector

$b \in \mathbb{R}^d$: 1-dim offset, bias scalar unknown.

$f_{\underline{w}, b}$: function f is parameterized by \underline{w}, b .

$$f_{\underline{w}, b}(\underline{x}) = \begin{bmatrix} b \\ \underline{w} \end{bmatrix}^T \begin{bmatrix} 1 \\ \underline{x} \end{bmatrix} = b + \underline{w}^T \underline{x}$$

unknown

$$\underline{w} = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \end{pmatrix}$$

$$\underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

$$\begin{pmatrix} b \\ \underline{w} \end{pmatrix} = \begin{pmatrix} b \\ w_1 \\ w_2 \\ w_3 \end{pmatrix}$$

$$\begin{pmatrix} 1 \\ \underline{x} \end{pmatrix} = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

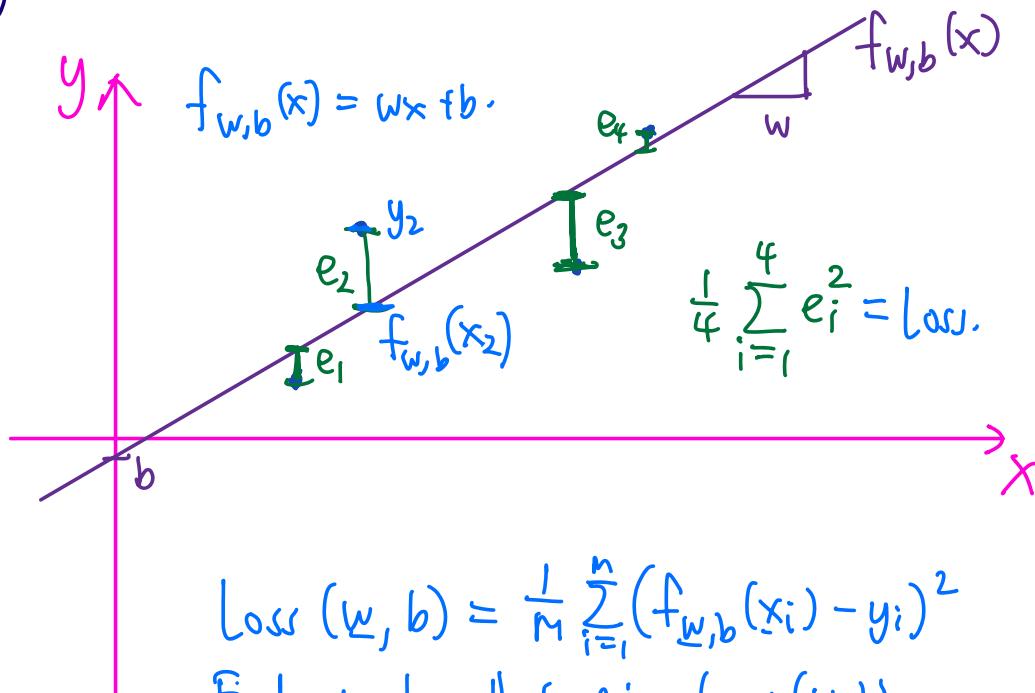
concatenation.

Minimize the error e_i bet. the prediction $f_{\underline{w}, b}(\underline{x}_i)$ and target y_i .

$$e_i = f_{\underline{w}, b}(\underline{x}_i) - y_i$$

$$\text{Loss}(\underline{w}, b) = \frac{1}{m} \sum_{i=1}^m e_i^2 = \frac{1}{m} \sum_{i=1}^m \underbrace{(f_{\underline{w}, b}(\underline{x}_i) - y_i)^2}_{\text{per-sample loss.}}$$

The squared loss (l_2 loss) is differentiable.



Define $\bar{w} \in \mathbb{R}^{d+1}$, $\bar{x}_i \in \mathbb{R}^{d+1}$ (i th training sample)

$$\bar{w} = \begin{bmatrix} b \\ w \end{bmatrix} = \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \quad \bar{x}_i = \begin{bmatrix} 1 \\ x_{i,1} \\ \vdots \\ x_{i,d} \end{bmatrix}$$

Objective (Loss) in Linear Regression

Want to find $\bar{w}^* = \begin{bmatrix} b \\ w^* \end{bmatrix} \in \mathbb{R}^{d+1}$ that min

$$\bar{w}^* = \underset{\bar{w} = [b \ w]}{\operatorname{arg\min}} \text{Loss}(w, b)$$

← is a constant

$$\text{Loss}(w, b) = \frac{1}{m} \sum_{i=1}^m e_i^2 = \cancel{\frac{1}{m}} \sum_{i=1}^m (f_{w,b}(x_i) - y_i)^2$$

$$f_{w,b}(x_i) - y_i = \underbrace{\begin{bmatrix} 1 \\ x_i \end{bmatrix}^T}_{\bar{x}_i^T} \underbrace{\begin{bmatrix} b \\ w \end{bmatrix}}_{\bar{w}} - y_i = \underbrace{x_i^T w + b}_{\text{affine function in } x_i} - y_i$$

So that

$$\sum_{i=1}^m (f_{w,b}(x_i) - y_i)^2 = \sum_{i=1}^m (x_i^T \bar{w} + b - y_i)^2.$$

$$= (\bar{x} \bar{w} - y)^T (\bar{x} \bar{w} - y)$$

X : design matrix $m \times (d+1)$

Unknown $\begin{bmatrix} b \\ w \end{bmatrix} = \bar{w}$

$$X = \begin{bmatrix} | & x_1^T \\ | & x_2^T \\ \vdots & \vdots \\ | & x_m^T \end{bmatrix} = \begin{bmatrix} | & x_{1,1} & \cdots & x_{1,d} \\ | & x_{2,1} & \ddots & x_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ | & x_{m,1} & \cdots & x_{m,d} \end{bmatrix} \in \mathbb{R}^{m \times (d+1)}$$

$(X^T X)^T = X^T X$

Objective:

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

$$\underline{a}^T \underline{b} = \underline{b}^T \underline{a}$$

$$\begin{aligned} &= \bar{w}^T X^T X \bar{w} - \bar{w}^T X^T y - y^T X \bar{w} + y^T y \\ &= \bar{w}^T (X^T X) \bar{w} - 2 \bar{w}^T (X^T y) + y^T y \end{aligned}$$

const wrt \bar{w}

$$\nabla_{\bar{w}} J(\bar{w}) = 2(X^T X) \bar{w} - 2 X^T y = 0$$

$$(X^T X) \bar{w} = X^T y \Rightarrow \bar{w}^* = (X^T X)^{-1} X^T y.$$

Are we allowed to do this?

if $X^T X$ is invertible
 ↓
 X has full column rank.

$\bar{w}^* = (X^T X)^{-1} X^T y$ is the least squares solution.

Given a dataset $\{(x_i, y_i)\}$.

Form the design matrix & target vector

$$X = \begin{bmatrix} 1 & x_{1,1} & \cdots & x_{1,d} \\ 1 & x_{2,1} & \cdots & x_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m,1} & \cdots & x_{m,d} \end{bmatrix} \in \mathbb{R}^{m \times (d+1)}$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m$$

Find the least squares solution.

$$\underline{\underline{w}}^* = \begin{bmatrix} \underline{b}^* \\ \underline{\underline{w}}^* \end{bmatrix} = (X^T X)^{-1} X^T y$$

Training
Learning

if X has full col. rank.

$$\begin{aligned} \underline{x}_{\text{new}} &\in \mathbb{R}^d \\ \text{Prediction: } \hat{y}_{\text{new}} &= \begin{bmatrix} 1 \\ \underline{x}_{\text{new}} \end{bmatrix}^T \begin{bmatrix} \underline{b}^* \\ \underline{\underline{w}}^* \end{bmatrix} \\ &= \underline{b}^* + \underline{x}_{\text{new}}^T \underline{\underline{w}}^*. \end{aligned}$$

Prediction /
Testing

$$\begin{array}{lll} \text{Dataset} & x_1 = -7, x_2 = -5, x_3 = 1, x_4 = 5 & d=1 \\ & y_1 = -6, y_2 = -4, y_3 = -1, y_4 = 4 & m=4. \end{array}$$

$$X = \begin{bmatrix} 1 & -7 \\ 1 & -5 \\ 1 & 1 \\ 1 & 5 \end{bmatrix} \quad y = \begin{bmatrix} -6 \\ -4 \\ -1 \\ 4 \end{bmatrix}$$

$$\begin{array}{ll} X \underline{\underline{w}} = y : & \text{rank}(X) = 2 \\ \cap \\ \mathbb{R}^2 & \text{rank}(\tilde{X}) = \text{rank}([X \ y]) = 3 \\ & \text{rank}(X) < \text{rank}(\tilde{X}) \Rightarrow \text{no solution.} \end{array}$$

$$\bar{w}^* = (X^T X)^{-1} X^T y = \begin{bmatrix} -0.5879 \\ 0.7747 \end{bmatrix}.$$

$x_{\text{new}} = 2$ (Predict its y value)

$$\hat{y}_{\text{new}} = \begin{bmatrix} 1 \\ x_{\text{new}} \end{bmatrix}^T \bar{w}^* = 1(-0.5879) + 2 \times 0.7747 = -2.1374.$$

To get the design matrix from a dataset $\{(x_i, y_i)\}$

$$x_i \in \mathbb{R}^d \quad x_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad x_2 = \begin{bmatrix} 2 \\ 3 \end{bmatrix}, \quad x_3 = \begin{bmatrix} 4 \\ 7 \end{bmatrix}. \quad d=2, \quad m=3.$$

$$X = \begin{bmatrix} 1 & x_1^T \\ 1 & x_2^T \\ 1 & x_3^T \end{bmatrix} = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 3 \\ 1 & 4 & 7 \end{bmatrix} \in \mathbb{R}^{m \times (d+1)}$$

$$x_1^T = [1 \ 2] \quad x_2^T = [2 \ 3] \quad x_3^T = [4 \ 7]$$

Converted the set of training vectors $x_i, 1 \leq i \leq m$ into a design matrix that can be used to learn an affine function

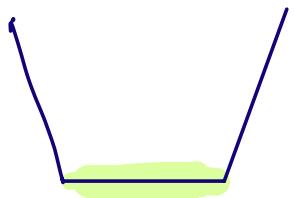
$$(X \bar{w} - y)^T (X \bar{w} - y) = \sum_{i=1}^m (X \bar{w} - y)_i^2$$

$$(g)_i = i^{\text{th}} \text{ comp. of vector } g. \quad = \sum_{i=1}^m (x_i^T \bar{w} + b - y_i)^2$$

X : full column rank $\Rightarrow \bar{w}^* = (X^T X)^{-1} X^T y$ is
 a unique global min for the loss function
TRUE.



X : not full col. rank $\Rightarrow (X^T X) \bar{w}^* = X^T y$
 SVD: singular value decompos. of $X^T X$
 \Rightarrow Multiple global min



$f(\underline{x}) = A\underline{x}$, $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$ $A \in \mathbb{R}^{m \times d}$.
 Check $\forall a, b \in \mathbb{R}$, $\underline{x}, \underline{y} \in \mathbb{R}^d$ $f(a\underline{x} + b\underline{y}) = af(\underline{x}) + bf(\underline{y})$.

$$\begin{aligned} \text{LHS} &= A(a\underline{x} + b\underline{y}) = A(a\underline{x}) + A(b\underline{y}) \quad (\text{additivity}) \\ &= aA\underline{x} + bA\underline{y} \end{aligned}$$

$$\text{RHS} = aA\underline{x} + bA\underline{y} \quad //$$