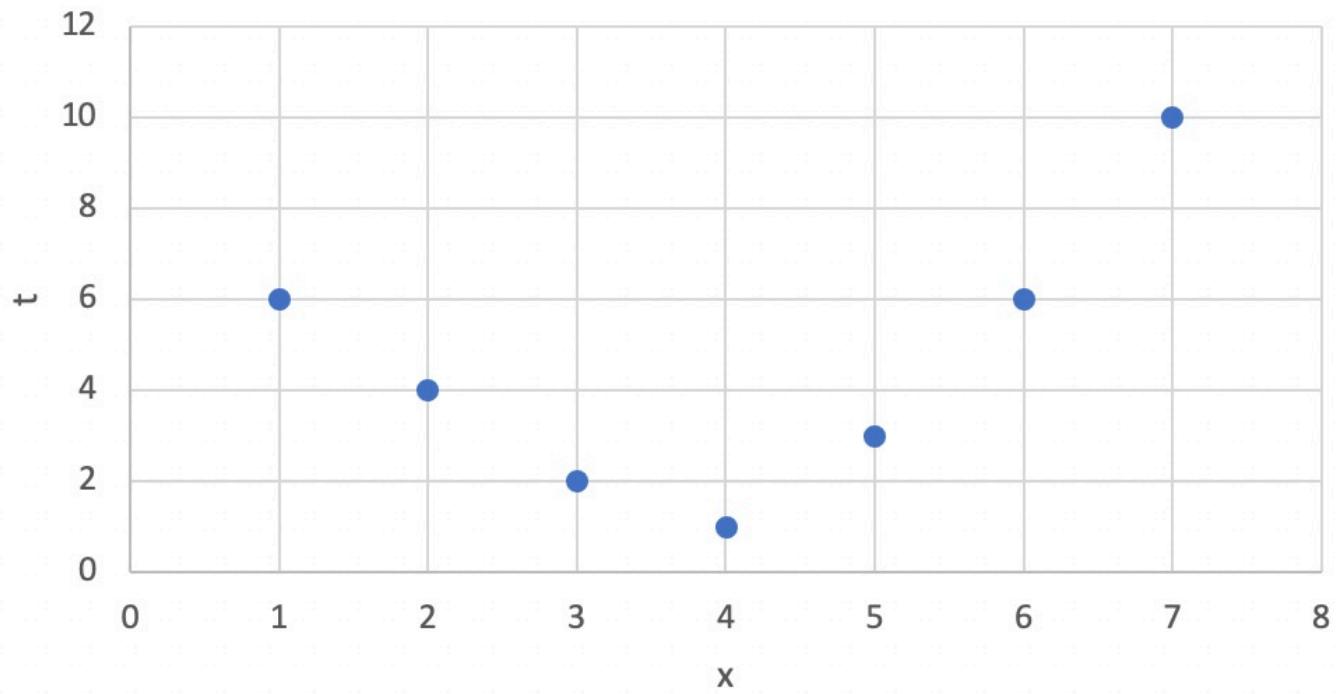


1.1

x	t
1	6
2	4
3	2
4	1
5	3
6	6
7	10

1.1 Graphing



$$1.2. g_{w,b}(x) = wx + b$$

Least squares loss

$$\epsilon(w, b) = \frac{1}{2N} \sum_{i=1}^N (g(i) - t^{(i)})^2 \quad g_{w,b}(wx + b)$$

first term.

$$= \frac{1}{2N} \left[\{(wx^{(1)} + b) - t^{(1)}\}^2 + \{(wx^{(2)} + b) - t^{(2)}\}^2 + \{(wx^{(3)} + b) - t^{(3)}\}^2 + \{(wx^{(4)} + b) - t^{(4)}\}^2 + \{(wx^{(5)} + b) - t^{(5)}\}^2 + \{(wx^{(6)} + b) - t^{(6)}\}^2 + \{(wx^{(7)} + b) - t^{(7)}\}^2 \right]$$

Showing computation only for the first term.

$$\begin{aligned} & [wx^{(1)} + b - t^{(1)}]^2 = [wx^{(1)} + b - t^{(1)}][wx^{(1)} + b - t^{(1)}] \\ & = w^2 x^{(1)2} + b^2 + t^{(1)2} - 2wx^{(1)} - 2bt^{(1)} \\ & \quad + b^2 x^{(1)2} + b^2 - bt^{(1)} \\ & \quad - wt^{(1)2} - t^{(1)2}b + t^{(1)2} \\ & = x^{(1)2} w^2 + b^2 + 2x^{(1)}wb - 2x^{(1)}t^{(1)}w - 2t^{(1)}b + t^{(1)2} \end{aligned}$$

It would be the same for all terms except the $x^{(i)}$ and $t^{(i)}$ will be in its own index.

$$\Rightarrow \frac{1}{2N} \left[\sum_{i=1 \dots N} (x^{(i)})^2 w^2 + (1) b^2 + (2x^{(i)}) wb + (-2x^{(i)}t^{(i)}) w + (-2t^{(i)}) b + (t^{(i)})^2 \right]$$

$$\text{Form wanted} \Rightarrow \frac{1}{2N} \sum_{i=1 \dots N} A_i w^2 + B_i b^2 + C_i wb + D_i w + E_i b + F_i$$

From side by side comparison, we know

$$A_i = [x^{(i)}]^2, \quad B_i = 1, \quad C_i = 2x^{(i)} \quad \text{for } i \in 1 \dots N$$

$$D_i = [-2x^{(i)}t^{(i)}], \quad E_i = -2t^{(i)}, \quad F_i = [t^{(i)}]^2$$

1.3.

From question 2,

$$\mathcal{E}(w, b) = \frac{1}{2N} \sum_{i=1 \dots N} \left[(x^{(i)})^2 w^2 + (1) b^2 + (2x^{(i)}) w b + (-2x^{(i)} t^{(i)}) w + (-2t^{(i)}) b + [t^{(i)}]^2 \right]$$

$$\text{let } A = \sum_i A_i, B = \sum_i B_i \dots \\ = \frac{1}{2N} \{ A(Nw^2) + B(Nb^2) - C(Nwb) + D(Nw) + E(Nb) + F \}$$

$$= \frac{1}{2} Aw^2 + \frac{1}{2} Bb^2 + \frac{1}{2} Cwb + \frac{1}{2} Dw + \frac{1}{2} Eb + \frac{F}{2N}$$

To minimize,

$$\frac{\partial [\mathcal{E}(w, b)]}{\partial w} = 0 ; Aw + \frac{1}{2} Cb + \frac{1}{2} D = 0$$

$$2Aw + Cb + D = 0$$

$$\Rightarrow w = \frac{-D - Cb}{2A} ; \text{sub 'b'}$$

$$\frac{\partial [\mathcal{E}(w, b)]}{\partial b} = 0 ; Bb + \frac{1}{2} Cw + \frac{1}{2} E = 0 \\ \Rightarrow b = \frac{-Cw - E}{2B}$$

To isolate without each dependency, sub in,

Sub 'b' into w,

$$w = \frac{-D - C \left[\frac{-Cw - E}{2B} \right]}{2A}$$

$$2A2Bw = -D(2B) + C^2w + CE$$

$$(2A2B - C^2)w = -2BD + CE$$

$$w = \frac{-2BD + CE}{2A2B - C^2}$$

Sub 'w' into 'b'

$$b = -C \left[\frac{-2BD + CE}{2A2B - C^2} \right] - E$$

$$2Bb [2A2B - C^2] = 2cBD + C^2E - E[2A2B - C^2]$$

$$b = \frac{2cBD + C^2E - E[2A2B - C^2]}{2B[2A2B - C^2]}$$

1.4.

To get numerical values, first evaluate A, B, C, D, E, F with the data-

$$A = \sum [x^{(i)}]^2 = 140$$

$$B = \sum B_i = 7$$

$$C = \sum [2x^{(i)}] = 56$$

$$D = \sum -2x^{(i)}t^{(i)} = -290$$

$$E = \sum -2t^{(i)} = -64$$

$$F = \sum [t^{(i)}]^2 = 202$$

$$\omega = \frac{-2BD + CE}{2A^2B - C^2} = \frac{476}{1184} = \boxed{\frac{17}{28}}$$

$$b = \left(-C \left[\frac{-2BD + CE}{2A^2B - C^2} \right] - E \right) \frac{1}{2B}$$

$$= \left[-56 \left(\frac{476}{1184} \right) + 64 \right] \frac{1}{14}$$

$$= \boxed{\frac{15}{7}}$$

1.5 verified with excel

$$2.2. \bar{x} = \begin{bmatrix} x^{(1)} & 1 \\ \vdots & \vdots \\ x^{(N)} & 1 \end{bmatrix} \quad \bar{x}^T = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(N)} \\ 1 & \dots & \dots & 1 \end{bmatrix}$$

$$\bar{x} \bar{w} = \begin{bmatrix} w x^{(1)} + b \\ \vdots \\ w x^{(N)} + b \end{bmatrix}, \quad \bar{x} \bar{w} - t = \begin{bmatrix} w x^{(1)} + b^{(1)} - t^{(1)} \\ \vdots \\ w x^{(N)} + b^{(N)} - t^{(N)} \end{bmatrix}$$

$$\|\bar{x} \bar{w} - t\|^2 = [w x^{(1)} + b^{(1)} - t^{(1)}]^2 + \dots + [w x^{(N)} + b^{(N)} - t^{(N)}]^2$$

$$\begin{aligned} \nabla_{\bar{w}} \|\bar{x} \bar{w} - t\|^2 &= \left[\frac{\partial}{\partial w} \{ [w x^{(1)} + b^{(1)} - t^{(1)}]^2 + \dots + [w x^{(N)} + b^{(N)} - t^{(N)}]^2 \} \right. \\ &\quad \left. \frac{\partial}{\partial b} \{ [w x^{(1)} + b^{(1)} - t^{(1)}]^2 + \dots + [w x^{(N)} + b^{(N)} - t^{(N)}]^2 \} \right] \\ &= \left[\underbrace{2 [w x^{(1)} + b^{(1)} - t^{(1)}]}_{\sum_{i=1}^N} x^{(1)} + \dots + 2 [w x^{(N)} + b^{(N)} - t^{(N)}] x^{(N)} \right. \\ &\quad \left. - 2 [w x^{(1)} + b^{(1)} - t^{(1)}] + \dots + - 2 [w x^{(N)} + b^{(N)} - t^{(N)}] \right] \end{aligned}$$

$$= \left[\sum_{i=1}^N 2 [w x^{(i)} + b^{(i)} - t^{(i)}] x^{(i)} \right]$$

$$\sum_{i=1}^N 2 [w x^{(i)} + b^{(i)} - t^{(i)}]$$

$$= \left[2 \sum_{i=1}^N [w x^{(i)} + b^{(i)}] x^{(i)} - 2 \sum_{i=1}^N x^{(i)} t^{(i)} \right]$$

$$2 \sum_{i=1}^N [w x^{(i)} + b^{(i)}] - 2 \sum_{i=1}^N t^{(i)}$$

$$= -2 \left[\sum_{i=1}^N [w x^{(i)} + b^{(i)}] - \sum_{i=1}^N t^{(i)} \right]$$

2.2 continued

$$= 2 \left[\sum_{i=1}^N [wx^{(i)} + b^{(i)}] \right] \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(N)} \end{bmatrix} - 2 \begin{bmatrix} x^{(1)} & \dots & x^{(N)} \end{bmatrix} \begin{bmatrix} 1 \\ \dots \\ 1 \\ t^{(1)} \end{bmatrix}$$

$$= 2 \begin{bmatrix} -x^{(1)} \\ \vdots \\ -1 \end{bmatrix} x \bar{w} - 2 x^T \bar{t}$$

$$= 2 x^T x \bar{w} - 2 x^T \bar{t}$$

$$= 2 x^T (x \bar{w} - \bar{t})$$

2.3 From the question before we know

$$A = \sum_{i=1}^n x_i w + b$$

given

$$\nabla_w \|x\bar{w} - \bar{t}\|^2 = 2x^T(x\bar{w}^* - \bar{t}) \Leftrightarrow 2x^T x\bar{w}^* - 2x^T \bar{t} = 0$$

→ To minimize, let the gradient equals zero.

$$2 \begin{bmatrix} x^{(1)} & \dots & x^{(N)} \\ \vdots & \ddots & \vdots \\ x^{(N)} & \dots & 1 \end{bmatrix} \begin{bmatrix} x^{(1)} w + b - t^{(1)} \\ \vdots \\ x^{(N)} w + b - t^{(N)} \end{bmatrix} = 0$$

$$= 2 \begin{bmatrix} \sum (x^{(1)} w + b - t^{(1)}) (x^{(1)}) \\ \vdots \\ \sum x^{(i)} w + b - t^{(i)} \end{bmatrix} = \begin{bmatrix} \sum x^{(1)} w + x^{(1)} b - t^{(1)} x^{(1)} \\ \vdots \\ \sum x^{(i)} w + b - t^{(i)} \end{bmatrix} \stackrel{?}{=} \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

Let's try plugging the numbers from problem 1,

$$\begin{aligned} \{x^{(i)}\}_{i=1 \dots n} &= (1, 2, \dots, 7), \quad \{t^{(i)}\}_{i=1 \dots n} = (6, 4, 2, 1, 3, 6, 10), \quad w = \frac{17}{28}, \quad b = \frac{15}{7} \\ \rightarrow & \begin{bmatrix} (-3.25) + (-1.29) + (5.89) + (14.29) + (10.89) + (-1.29) + (-25.25) \\ (-3.25) + (-0.64) + (1.96) + (3.57) + (2.18) + (-0.214) + (-3.61) \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \text{in QED} \end{aligned}$$

$$\therefore 2x^T x\bar{w}^* - 2x^T \bar{t} = 0$$

∴ QED.

While it can be numerically seen, also analytically,

from 2.2 we've derived:

$$2 \begin{bmatrix} \sum_{i=1}^N [w x^{(i)} + b^{(i)}] \\ \sum_{i=1}^N [w x^{(i)} + b^{(i)}] \end{bmatrix} \begin{bmatrix} x^{(1)} \\ 1 \end{bmatrix} - 2 \begin{bmatrix} x^{(1)} & \dots & x^{(N)} \\ 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} | \\ t^{(1)} \\ | \end{bmatrix}$$

We can put this into a form

$$= 2 \begin{bmatrix} -x^{(1)} \\ -1 \end{bmatrix} x \bar{w} - 2 x^T \bar{t}$$

$$= 2 x^T \bar{w} - 2 x^T \bar{t}$$

Then we know that for \bar{w}^* that minimizes the least squares loss,
the gradient = 0.

$$\Rightarrow 2 x^T \bar{w} - 2 x^T \bar{t} = 0 \text{ for } \bar{w}^*$$

In this is more of a general proof of how we land to the given.
(It required bit of a backtracking from the previous question.)

2.4

$$\nabla_{\bar{w}} \|X\bar{w} - \bar{t}\|^2 = 2X^T X \bar{w}^* - 2X^T \bar{t} = 0 \text{ for } \bar{w}^*$$

$$2X^T X \bar{w}^* = 2X^T \bar{t}$$

$$2[X^T (X^T X)^{-1} \xrightarrow{\rightarrow I} X^T X] \bar{w}^* = 2(X^T X)^{-1} X^T \bar{t}$$

$$\bar{w}^* = (X^T X)^{-1} X^T \bar{t}$$

2.5 verified with Numpy

3.1

$$A = \sum_{i \in 1 \dots N} \bar{x}^{(i)} \bar{x}^{(i)T}; \quad x_j^{(i)}, j \in 1 \dots d$$

$$= \sum_{i \in 1 \dots N} \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_d^{(i)} \end{bmatrix} \begin{bmatrix} x_1^{(i)} & \dots & x_d^{(i)} \end{bmatrix}$$

$$= \begin{bmatrix} x_1^{(1)} x_1^{(1)} & x_1^{(1)} x_2^{(1)} & \dots & x_1^{(1)} x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ x_d^{(1)} x_1^{(1)} & \dots & x_d^{(1)} x_d^{(1)} \end{bmatrix} + \dots + \begin{bmatrix} x_1^{(N)} x_1^{(N)} & \dots & x_1^{(N)} x_d^{(N)} \\ \vdots & \ddots & \vdots \\ x_d^{(N)} x_1^{(N)} & \dots & x_d^{(N)} x_d^{(N)} \end{bmatrix}$$

$$= \sum_{i \in 1 \dots N} \begin{bmatrix} x_1^{(i)} x_1^{(i)} & x_1^{(i)} x_2^{(i)} & \dots & x_1^{(i)} x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ x_d^{(i)} x_1^{(i)} & \dots & x_d^{(i)} x_d^{(i)} \end{bmatrix}$$

$$\bar{x}^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \quad \bar{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_N \end{bmatrix}$$

Date.

No.

3.2

$$b = \sum_{i=1}^n t^{(i)} \bar{x}^{(i)} \text{, show that } \nabla \epsilon(\bar{w}, D) = \frac{1}{N} (A\bar{w} - \bar{b}) + \lambda \bar{w}$$

From the loss function,

$$\epsilon(\bar{w}, D) = \frac{1}{2N} \sum_{i=1 \dots N} (g_{\bar{w}}(\bar{x}^{(i)}) - t^{(i)})^2 + \frac{\lambda}{2} \|\bar{w}\|_2^2$$

Square rooting the whole sum of square of elements = ℓ_2 norm

$$= \frac{1}{2N} \sum_{i=1 \dots N} \underbrace{[(\bar{x}^{(i)} \bar{w}) - t^{(i)}]^2}_{\text{scalar}} + \frac{\lambda}{2} \|\bar{w}\|_2^2$$

$$= \frac{1}{2N} \|\bar{X} \bar{w} - \bar{t}\|_2^2 + \frac{\lambda}{2} \|\bar{w}\|_2^2 \quad ; \text{ capitals absorbed (took account) of the summation.}$$

Now, take the gradient.

$$\nabla \epsilon(\bar{w}, \bar{D})$$

$$= \frac{1}{N} (\bar{X})(\bar{X}^T \bar{w} - \bar{t}) + \lambda \bar{w}$$

known from problem 2.2

④

* Side proof of how $\nabla \|\bar{w}\|_2^2 = 2\bar{w}$

$$\|\bar{w}\|_2^2 = (\sqrt{\sum_{i=1}^N (w_i)^2})^2$$

$$= \sum_{i=1}^N w_i^2$$

taking the derivative,

$$\nabla \sum_{i=1}^N w_i^2 = \sum_{i=1}^N \nabla (w_i^2) \quad (\nabla \text{ operation } \mathbb{R}^n \rightarrow \mathbb{R})$$

$$= 2 \left[\begin{array}{c} \frac{\partial \sum w}{\partial w_1} \\ \frac{\partial \sum w}{\partial w_2} \\ \vdots \\ \frac{\partial \sum w}{\partial w_N} \end{array} \right]$$

$$= 2 \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_N \end{bmatrix} \quad \text{partial derivative makes a column vector}$$

$$= 2\bar{w}$$

Rewriting,

$$\frac{1}{N} (X)(X^T \bar{w} - \bar{t}) + \lambda \bar{w}$$

$$= \frac{1}{N} [\underbrace{X X^T \bar{w}}_{A} - X \bar{t}] + \lambda \bar{w}$$

$$= \frac{1}{N} [A \bar{w} - X \bar{t}] + \lambda \bar{w}$$

$X \bar{t}$ ↴

$$= [\bar{x}^{(1)} \bar{x}^{(2)} \dots \bar{x}^{(N)}] \begin{bmatrix} t^{(1)} \\ \vdots \\ t^{(N)} \end{bmatrix}$$

$$= \bar{x}^{(1)} t^{(1)} + \dots + \bar{x}^{(N)} t^{(N)}$$

$$= \sum \bar{x}^{(i)} t^{(i)}$$

on the document it is
given as $\vec{b} = \sum_{i=1 \dots N} t^{(i)} \bar{x}^{(i)}$

but this is equivalent to

$$\vec{b} = \sum_{i=1 \dots N} \bar{x}^{(i)} t^{(i)}$$

as $t^{(i)}$ are scalars

\Rightarrow this is b by definition.

$$\therefore \nabla \mathcal{E}(\bar{w}, D) = \frac{1}{N} [A \bar{w} - \vec{b}] + \lambda \bar{w} \quad \therefore \text{QED.}$$

3.3

$$\vec{w}^* = \arg \min_{\vec{w}} \mathcal{E}(\vec{w}, 0)$$

Want to minimize the loss

\Rightarrow gradient of loss function = 0.

From the previous question, found the gradient of the loss function.

$$\frac{1}{N} (A\vec{w}^* - \vec{b}) + \lambda \vec{w}^* = 0$$

$$(A\vec{w}^* - \vec{b}) + N\lambda \vec{w}^* = 0$$

$$(A + N\lambda I_b) \vec{w}^* - \vec{b} = 0$$

$$(A + N\lambda I_b) \vec{w}^* = \vec{b} ; \text{ equation be to satisfied for } \vec{w}^*$$

3.4

prove all eigenvalues of A are non-negative.

\Rightarrow prove that $\exists \vec{v}$ s.t. $A\vec{v} = \lambda \vec{v}$ for all $\lambda \geq 0$.

$$\vec{v}^T A \vec{v} = \vec{v}^T \lambda \vec{v}$$

LHS: let \vec{w} be an arbitrary
non zero vector.

RHS:

$$\begin{aligned} \vec{w}^T \lambda \vec{w} &= \lambda \vec{w}^T \vec{w} ; \text{ possible since} \\ &\quad \lambda \text{ is a scalar} \\ &= \lambda \|\vec{w}\|_2^2 \end{aligned}$$

$$\vec{w}^T A \vec{w} = \vec{w}^T (\sum X X^T) \vec{w} ; \text{ from 3.1}$$

$$= \sum \vec{w}^T X X^T \vec{w} ; \text{ took sum outside}$$

$$= \sum (X^T \vec{w})^T (X^T \vec{w})$$

$$\because X^T X = \|X\|_2^2$$

$$= \sum \|X^T \vec{w}\|_2^2 \geq 0$$

\therefore squared, cannot be less than 0.

\vec{w} is a nonzero vector and through the norm square, it will be ≥ 0

\therefore previously shown that LHS ≥ 0

$$\therefore \lambda \geq 0.$$

3.5

$A + \lambda N I_d$ is invertible \Leftrightarrow none of its eigenvalues are zero.

$$(A + \lambda N I_d) \vec{v} = \gamma \vec{v}$$

$$A \vec{v} + \lambda N \vec{v} = \gamma \vec{v}$$

$$A \vec{v} = (\gamma - \lambda N) \vec{v}$$

\because eigenvalue $(A) > 0$ proven in previous question

$$\therefore \gamma - \lambda N > 0 \rightarrow \gamma > \lambda N$$

Knowing $\lambda > 0$ and $N > 0 \Rightarrow \therefore \gamma > 0$

\therefore proven that its eigenvalues are more zero (greater than zero)

$\therefore A + \lambda N I_d$ is invertible.

3.6

From 3.3, we know $(A + \lambda N) \vec{w}^* = \vec{b}$

$$(A + \lambda N I_d)^{-1} \underbrace{(A + \lambda N I_d) \vec{w}^*}_{\text{produces identity matrix}} = (A + \lambda N I_d)^{-1} \vec{b}$$

produces identity matrix

$$\therefore \vec{w}^* = (A + \lambda N I_d)^{-1} \vec{b}$$