

$$1) \text{ a. } l_2 \text{ norm : } \min_c \left\{ (c-x_1)^2 + (c-x_2)^2 + (c-x_3)^2 \right\}$$

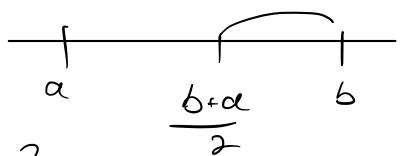
$$\frac{df}{dc} = 2(c-x_1) + 2(c-x_2) + 2(c-x_3) = 6c - 2x_1 - 2x_2 - 2x_3$$

$$6c = 2x_1 + 2x_2 + 2x_3$$

$$c = \frac{x_1 + x_2 + x_3}{3}$$

$$b. l_\infty \text{ norm: } \min_c \left\{ \max(|c-x_1|, |c-x_2|, |c-x_3|) \right\}$$

choose:  $c = \frac{x_3+x_1}{2}$



$$\max \left\{ \left| \frac{x_3+x_1}{2} - x_1 \right|, \left| \frac{x_3+x_1}{2} - x_2 \right|, \left| \frac{x_3+x_1}{2} - x_3 \right| \right\}$$

$$= \left| \frac{x_3+x_1}{2} - x_1 \right| = \frac{x_3 - x_1}{2}$$

$$\frac{x_3 - x_1}{2} > \max \left\{ \dots \right\} \text{ if and only if } \frac{x_3+x_1}{2} \neq c' \quad \text{since } c' \text{ is a fixed value}$$

$$\therefore c' < x_1 \quad \text{then, } c' \notin [x_1, x_3] \quad \text{if } \dots$$

$$\max \{|c'-x_1|, |c'-x_2|, |c'-x_3|\} = |c'-x_3| > |x_1 - x_3| > \frac{x_1 + x_3}{2}$$

$$x_1 \leq C' \angle \frac{x_5 + x_1}{2} : 2 \text{ mgn}$$

$$\max \{ |c' - x_1|, |c' - x_2|, |c' - x_3| \} = |c' - x_3| > \frac{x_1 + x_3}{2}$$

$$\frac{x_3 + x_1}{2} < c' \leq x_3 : 3 \text{ mgn}$$

$$\max \{ |c' - x_1|, |c' - x_2|, |c' - x_3| \} = |c' - x_1| > \frac{x_1 + x_3}{2}$$

$$c. l_1 \text{ norm: } \min_c \{ |c-x_1| + |c-x_2| + |c-x_3| \}$$

choose:  $c = x_2$   $x_1 < x_2 < x_3$

$$\begin{aligned} \text{then } l_1 \text{ norm} &= |x_2 - x_1| + |x_2 - x_2| + |x_2 - x_3| = x_2 - x_1 - (x_2 - x_3) \\ &= x_3 - x_1 \end{aligned}$$

by (3) prove  $\exists c' \in [x_1, x_2]$  s.t.  $|c' - x_1| + |c' - x_2| + |c' - x_3| < x_3 - x_1$

$$\begin{aligned} |c' - x_1| + |c' - x_2| + |c' - x_3| &\stackrel{\text{def}}{=} x_3 - x_1 + (x_2 - c') > x_3 - x_1 \end{aligned}$$

$$\begin{aligned} |c' - x_1| + |c' - x_2| + |c' - x_3| &= x_3 - x_1 + (c' - x_2 - c' + x_3) \stackrel{\text{def}}{=} x_3 - x_1 \end{aligned}$$

$$\begin{aligned} |c' - x_1| + |c' - x_2| + |c' - x_3| &\geq x_3 - x_1 \\ &\stackrel{\text{def}}{=} x_3 - x_1 \end{aligned}$$

$$|c' - x_1| + |c' - x_2| + |c' - x_3| \geq x_3 - x_1$$

② a.  $(A^T A)^T = A^T \cdot (A^T)^T = A^T A$  ✓

$x^T A^T A x = (Ax)^T \cdot (Ax) = \underbrace{\langle Ax, Ax \rangle}_{\text{Do } \mathbb{R}^N} \geq 0$  : אינטראקטיבי

b.  $Cv = \lambda v \Rightarrow (I - C) \cdot v = Iv - Cv = v - \lambda v = (1 - \lambda)v$

$I - C$  אובייקט ווילג'ס ל'ג'ס 1-)

c.  $\text{rank}(A) = n \iff \text{rank}(A^T A) = n$

$Ax = 0$  ווילג'ס,  $x \in \text{Ker}(A)$  :

$$A^T A x = A^T \cdot (Ax) = A^T \cdot 0 = 0 \quad : x \in \text{Ker}(A^T A) \text{ ווילג'ס}$$

בנוסף  $x^T - \circ$  (פ.מ.)  $A^T A x = 0$  ווילג'ס,  $x \in \text{Ker}(A^T A)$  :

$$x^T A^T A x = 0 \Rightarrow \langle Ax, Ax \rangle = 0 \Rightarrow Ax = 0$$

$\text{Ker}(A^T A) = \text{Ker}(A)$  יס'

בנוסף  $A^T A \Leftarrow \dim(\text{rank}(A^T A)) = n \Leftarrow n = \dim(A^T A)$

$\dim \ker(A^T A) = 0 \Leftarrow \text{rank}(A^T A) = n$   $\quad \text{পর্যবেক্ষণ হল } A^T A \text{ এর } n \times n \text{ :} \Rightarrow$

$$\Leftarrow A^T A x = 0 \Leftarrow \begin{matrix} A^T - 0 & \text{সূজি} \\ \text{রেনোন} & \end{matrix} = Ax = 0 \Leftarrow x \in \ker(A) \quad \text{এখন}$$

$$\ker(A) \subseteq \ker(A^T A) \Leftarrow x \in \ker(A^T A)$$

■  $\text{rank}(A) = n \Leftarrow \{0\} = \ker(A) \quad \text{পর্যবেক্ষণ হল } \{0\} = \ker(A^T A) \quad \text{রেনোন}$

- (d) Suppose that  $A \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ . Show that  $A$  is full rank if and only if  $A^T A$  is symmetric and positive definite (you can use the previous section).

$$\text{rank}(A) = n \iff A^T A \text{ উন্নত পদ্ধতি}$$

এখন  $A^T A$  এর প্রযুক্তির সম্বোধন করা হচ্ছে।

$$\ker(A^T A) = \{0\} \quad \text{হল } A^T A x \geq 0 \quad \forall x \in \mathbb{R}^n$$

$$A^T A x > 0 \quad \text{পর্যবেক্ষণ}$$

এখন  $A^T A$  উন্নত পদ্ধতি হল, যে কৌণ

(1)

$$\forall \neq v \in V \quad \text{if} \quad \langle (A^T A + \alpha I) v, v \rangle > 0 \quad \text{then}$$

$$\langle A^T A v, v \rangle \geq 0 \quad \text{by definition} \quad A^T A \geq 0 \quad \text{positive definite} \quad \text{or} \quad \forall v$$

$$\langle (A^T A + \alpha I) v, v \rangle = \langle A^T A v + \alpha I v, v \rangle = \langle A^T A v, v \rangle + \langle \alpha I v, v \rangle$$

$$= \underbrace{\langle A^T A v, v \rangle}_{\geq 0} + \alpha \underbrace{\langle v, v \rangle}_{\geq 0} \Rightarrow > 0$$

$$\therefore \text{positive definite} \quad A^T A + \alpha I \quad - \text{definite}$$

positive definite, positive definite  $A^T A - \text{definite} \quad (\alpha) \text{ proof}$

$$\begin{aligned} \text{positive definite} &\rightarrow \text{positive definite } A^T A + \alpha I \quad \text{proof} \\ \text{positive definite} &\rightarrow \text{positive definite } A^T A + \alpha I \end{aligned}$$

(a) Find the best approximation in a least square sense for the system  $Ax \approx b$  where

$$A = \begin{bmatrix} 2 & 1 & 2 \\ 1 & -2 & 1 \\ 1 & 2 & 3 \\ 1 & 1 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} 6 \\ 1 \\ 5 \\ 2 \end{bmatrix}. \quad (1)$$

Write the normal equations and solve the problem using a computer. You may use built-in functions and provide the code.

$$\text{Find } x \rightarrow \min \|Ax - b\|_2^2, \quad f(x) = \|Ax - b\|_2^2 = (Ax - b)^T(Ax - b)$$

$$= x^T A^T A x - 2b^T A x + b^T b$$

$$\nabla f(x) = 2A^T A x - 2A^T b = 0$$

$$2A^T A x = 2A^T b \quad | : 2$$

$$A^T A x = A^T b \quad \text{Normal Equation}$$

$$x^* = \begin{bmatrix} 1.7 \\ 0.6 \\ 0.7 \end{bmatrix} \quad \begin{array}{l} \text{using } \Rightarrow \text{np.linalg} \\ \text{using } \Rightarrow \text{np.linalg} \end{array}$$

317

$$a = np.array([[[2, 1, 2], [1, -2, 1], [1, 2, 3], [1, 1, 1]]])$$

$$b = np.array([6, 1, 5, 2])$$

$$x = \text{numpy.linalg.lstsq}(a, b)[0]$$

```
[2] import numpy as np
    import pandas as pd
    import matplotlib.pyplot as plt
```

```
a=np.array([[2,1,2],[1,-2,1],[1,2,3],[1,1,1]])
b =np.array([6,1,5,2])
x = numpy.linalg.lstsq(a,b)

print(x)
```

```
Out[2]: (array([1.7, 0.6, 0.7]), array([1.4]), 3, array([4.96978399, 2.53339065, 0.93977598]))
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:3: FutureWarning: `rcond` parameter will change to the default of machine precision times ``max(M, N)``
To use the future default and silence this warning we advise to pass `rcond=None`, to keep using the old, explicitly pass `rcond=-1`.
This is separate from the ipykernel package so we can avoid doing imports until
```

- (b) Is the solution  $\mathbf{x}^*$  that you found in the previous section unique? Explain. What is the minimal objective (loss) value  $\|A\mathbf{x}^* - \mathbf{b}\|_2^2$ ?

(c)  $\text{Proof}$   $A^T A$   $\vdash$   $\text{is full rank}$   $\Rightarrow$   $\text{min}_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|_2$

$$A^T A \mathbf{x} = A^T \mathbf{b} \quad | \quad (A^T A)^{-1} \text{ exists} \quad \text{since } A \text{ is full rank}$$

$$(A^T A)^{-1} (A^T A) \mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$$

$\Rightarrow \mathbf{x}^* = (A^T A)^{-1} A^T \mathbf{b}$   
 $\mathbf{x}^* \in \mathbb{R}^3$ , unique

$\therefore \sqrt{N} \cdot \sqrt{\epsilon} \approx \sqrt{\epsilon}$

$$\|A\mathbf{x}^* - \mathbf{b}\|_2^2 = \left\| \begin{pmatrix} 2 & 1 & 2 \\ 1 & -2 & 1 \\ 1 & 2 & 3 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 1.7 \\ 0.6 \\ 0.7 \end{pmatrix} - \begin{pmatrix} 6 \\ 1 \\ 5 \\ 2 \end{pmatrix} \right\|_2^2 = \left\| \begin{pmatrix} \frac{2}{3} \\ \frac{6}{5} \\ \frac{5}{3} \\ 2 \end{pmatrix} - \begin{pmatrix} 6 \\ 1 \\ 5 \\ 2 \end{pmatrix} \right\|_2^2$$

$$= \left\| \begin{pmatrix} -\frac{3}{5} \\ \frac{1}{5} \\ 0 \\ 1 \end{pmatrix} \right\|_2^2 = \frac{7}{5}$$

- (c) Find the least squares solution of the system in Eq. (1), but now find a solution for which the second equation is almost exactly satisfied (let's say, such that  $|r_2| < 10^{-3}$ ). Hint: use weighted least squares.

:Weighted least Squares  $\rightarrow$  183N120  $\rightarrow$  100  $\rightarrow$  100

$$10^{-3} |r_2| \rightarrow \text{use } \hat{x} = (A^T W A)^{-1} A^T W b \quad \rightarrow \text{IC} \quad (100)$$

Given  $W = \begin{pmatrix} 1 & 0 \\ 10^5 & 1 \\ 0 & 1 \end{pmatrix} \rightarrow 10^5$

$$\hat{x} = \left( \begin{pmatrix} 2 & 1 & 2 \\ 1 & -2 & 1 \\ 1 & 2 & 3 \\ 1 & 1 & 1 \end{pmatrix}^T \begin{pmatrix} 1 & 0 \\ 10^5 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 2 & 1 & 2 \\ 1 & -2 & 1 \\ 1 & 2 & 3 \\ 1 & 1 & 1 \end{pmatrix} \right)^{-1} \cdot \begin{pmatrix} 2 & 1 & 2 \\ 1 & -2 & 1 \\ 1 & 2 & 3 \\ 1 & 1 & 1 \end{pmatrix}^T \begin{pmatrix} 1 & 0 \\ 10^5 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 6 \\ 1 \\ 5 \\ 2 \end{pmatrix}$$

$$= \begin{pmatrix} \frac{1160003}{67002} \\ \frac{229991}{340001} \\ \frac{440005}{67002} \end{pmatrix} \Rightarrow r = Ax - b = \begin{pmatrix} \frac{-210000}{340001} \\ \frac{7}{340001} \\ 0 \\ \frac{350000}{340001} \end{pmatrix} \leftarrow r_2$$

$$\frac{7}{340001} = r_2 < 10^{-3}$$

- (d) Find the least squares solution of the system in Eq. (1), but now add simple Tikhonov regularization term  $\lambda \|\mathbf{x}\|_2^2$  with  $\lambda = 0.5$ .

: Tikhonov regularization  $\rightarrow$  pen

$$\hat{\mathbf{x}} = \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{arg\,min}} \left\| \begin{bmatrix} A \\ \sqrt{\lambda} I \end{bmatrix} \mathbf{x} - \begin{bmatrix} b \\ 0 \end{bmatrix} \right\|_2^2$$

↓

$$\hat{\mathbf{x}} = (A^T A + \lambda I)^{-1} \cdot A^T b$$

$$A = \begin{pmatrix} 2 & 1 & 2 \\ 1 & -2 & 1 \\ 1 & 2 & 3 \\ 1 & 1 & 1 \end{pmatrix}, b = \begin{pmatrix} 6 \\ 1 \\ 5 \\ 2 \end{pmatrix}, \text{ and } \hat{\mathbf{x}} \rightarrow \text{soln}$$

$$\lambda = 0.5$$

$$\hat{\mathbf{x}} = \begin{pmatrix} \frac{2672}{1929} \\ \frac{344}{643} \\ \frac{572}{643} \end{pmatrix}$$

#### 4. Frobenius Norm.

See the definition of Frobenius Norm:

<https://mathworld.wolfram.com/FrobeniusNorm.html>

This norm is often used to compare data sets, or apply manipulations on data sets. It is easy to see that by definition

$$\|A\|_F^2 = \sum_{i,j} (a_{ij})^2 = \sum_j \|\mathbf{a}_j\|_2^2$$

where  $\mathbf{a}_j$  is a column of  $A$ .

(a) Suppose that we want to solve the problem

$$\arg \min_{X \in \mathbb{R}^{n \times n}} \|AX - B\|_F^2$$

where  $A, B \in \mathbb{R}^{n \times n}$ . Find an expression for the solution. When is the solution unique? Hint: this is almost the same as standard least squares.

$$f(x) = \|Ax - B\|_F^2 = \sum_j \|Ax_j - b_j\|_2^2 = \sum_j (Ax_j - b_j)^T (Ax_j - b_j)$$

$x \rightarrow j$  針對  $n$  個  $x_j$

$$= \sum_j (x_j^T A^T - b_j^T)(Ax_j - b_j) = \sum_j x_j^T A^T A x_j - x_j^T A^T b_j - b_j^T A x_j + b_j^T b_j$$

$$\nabla f(x) = \begin{pmatrix} 2A^T A x_1 - 2A^T b_1 \\ 2A^T A x_2 - 2A^T b_2 \\ \vdots \\ 2A^T A x_n - 2A^T b_n \end{pmatrix} = 0$$

( $1 \leq j \leq n$ ) 即  $\sum_j 2A^T A x_j - 2A^T b_j = 0$

$$2A^T A x_j - 2A^T b_j = 0$$

$$A^T A x_j = A^T b_j$$

$$x_j = (A^T A)^{-1} A^T b_j$$

$(A^T A)^{-1}$  存在  $\Leftrightarrow A$  full rank

$\exists (A^T A)^{-1} \Rightarrow A^T b_j$

$$x_j = A^{-1} b_j$$

רוכסן  $A$  רלוֹן בְּגִינָה  $x_j$  כוֹנִינָה, בְּגִינָה  $x_j$  אֲזַמֵּן

- (b) Often, we wish to make data  $A$  look similar to data  $B$  of the same size, to remove artifacts that are not related to the phenomena that we test. For example, when we measure gene expression values, often there's a bias of the particular lab's measurement. To correct that, we want to make the datasets similar, and one option is to solve a minimization of the following form (this is a bit simplified)

$$\arg \min_{D \in \mathbb{R}^{n \times n}, D \text{ is diagonal}} \|DA - B\|_F^2$$

$D$  is a diagonal matrix. In other words, we wish to find a scale  $D_{ii}$  for each row of  $A$  (denoted by  $\mathbf{a}_i$ ) so that  $\|\mathbf{a}_i D_{ii} - \mathbf{b}_i\|_2^2$  is minimized ( $\mathbf{b}_i$  is a row in  $B$ ). Use a computer and find the solution  $D$  for the following matrices:

$$A = \begin{bmatrix} 5 & 6 & 7 & 8 \\ 1 & 3 & 5 & 4 \\ 1 & 0.5 & 4 & 2 \\ 3 & 4 & 3 & 1 \end{bmatrix} \quad B = \begin{bmatrix} 0.57 & 0.56 & 0.8 & 1 \\ 1.5 & 4 & 6.7 & 4.9 \\ 0.2 & 0.1 & 1 & 0.6 \\ 11 & 30 & 26 & 10 \end{bmatrix}$$

$$x_i = \frac{\alpha_i b_i^\top}{\alpha_i \alpha_i^\top}$$

$$D = \begin{pmatrix} x_1 & & & \\ & x_2 & & \\ & & x_3 & \\ & & & x_4 \end{pmatrix} \quad D_{ii} = x_i \text{ ימוי}$$

$$f(x_i) = \|\alpha_i x_i - b_i\|_2^2 = (\alpha_i x_i - b_i) (\alpha_i x_i - b_i)^\top = (\alpha_i x_i - b_i) (x_i^\top \alpha_i^\top - b_i^\top)$$

$$= \alpha_i x_i x_i^\top \alpha_i^\top - \alpha_i x_i b_i^\top - b_i x_i^\top \alpha_i^\top + b_i b_i^\top$$

$$f'(x_i) = 2x_i \alpha_i \alpha_i^\top - 2\alpha_i b_i^\top = 0$$

$$x_i \alpha_i \alpha_i^\top = \alpha_i b_i^\top$$

$$x_i = \frac{\alpha_i b_i^\top}{\alpha_i \alpha_i^\top}$$

$$x_1 = \frac{1981}{\frac{100}{174}} = 0.113$$

$$x_2 = \frac{533}{\frac{5}{51}} = 1.305 \quad x_3 = \frac{241}{\frac{35}{21.25}} = 6.885$$

$$x_3 = \frac{109}{\frac{20}{21.25}} = 0.256$$

$$D = \begin{pmatrix} 0.113 & & & \\ & 1.305 & & \\ & & 0.256 & \\ & & & 6.885 \end{pmatrix}$$

## Code:

```
import matplotlib
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression

path = 'insurData.csv'
df = pd.read_csv(path)
print('\nNumber of rows and columns in the data set: ', df.shape)
print('')
df.head()

# a
df['sex'] = df['sex'].replace(['female'], 0)
df['sex'] = df['sex'].replace(['male'], 1)
df['smoker'] = df['smoker'].replace(['yes'], 1)
df['smoker'] = df['smoker'].replace(['no'], 0)
df['region'] = df['region'].replace(['northeast'], 0)
df['region'] = df['region'].replace(['southeast'], 1)
df['region'] = df['region'].replace(['southwest'], 2)
df['region'] = df['region'].replace(['northwest'], 3)

# b
df['charges'] = df['charges'] / 1000
meanc = df['charges'].mean()
df['charges'] = df['charges'] - meanc
meanb = df['bmi'].mean()
df['bmi'] = df['bmi'] - meanb
meana = df['age'].mean()
df['age'] = df['age'] - meana
```

```

# C
for i in range(5):
    train, test = train_test_split(df, test_size=0.2)
    train_x = train[['age', 'sex', 'bmi', 'children', 'smoker',
'region']].to_numpy()
    train_y = train[['charges']].to_numpy()
    test_x = test[['age', 'sex', 'bmi', 'children', 'smoker',
'region']].to_numpy()
    test_y = test[['charges']].to_numpy()

    ans_train = np.linalg.lstsq(train_x, train_y, rcond=None)[0]
    train_mse = (1 / np.size(train_x, 0)) * (
        np.linalg.norm(train_x @ ans_train - train_y, ord=2,
axis=None, keepdims=False)) * (np.linalg.norm(train_x @ ans_train -
train_y, ord=2, axis=None, keepdims=False))
    print('train MSE is: ')
    print(train_mse)
    ans_test = np.linalg.lstsq(test_x, test_y, rcond=None)[0]
    test_mse = (1 / np.size(test_x, 0)) * (
        np.linalg.norm(test_x @ ans_test -
test_y,ord=2,keepdims=False)) * (np.linalg.norm(test_x @
ans_test - test_y, ord=2, axis=None, keepdims=False))
    print('test MSE is: ')
    print(test_mse)
    compered_mse = train_mse - test_mse
    print('diff is: ')
    print(compered_mse)

#d
matplotlib.pyplot.hist(train_x @ ans_train - train_y,bins=100)
plt.show()

```

**התוצאות לשיער ג':** קיבלנו שהשגיאה הממוצעת היא באיזור ה-40+ (כלומר בפועל פ', 1000), זו תוצאה לא טובת כלומר המודל לא צופה טוב את החינויים.

## **Results :**

Number of rows and columns in the data set: (1338, 7)

///b

	age	sex	bmi	children	smoker	region	charges
0	-20.207025	0	-2.763397	0	1	2	3.614502
1	-21.207025	1	3.106603	1	0	1	-11.544870
2	-11.207025	1	2.336603	3	0	1	-8.820960
3	-6.207025	1	-7.958397	0	0	3	8.714048
4	-7.207025	1	-1.783397	0	0	3	-9.403567
...	...	...	...	...	...	...	...
1333	10.792975	1	0.306603	3	0	3	-2.669874
1334	-21.207025	0	1.256603	0	0	0	-11.064441
1335	-21.207025	0	6.186603	0	0	1	-11.640589
1336	-18.207025	0	-4.863397	0	0	2	-11.262477
1337	21.792975	0	-1.593397	0	1	3	15.870938

///c

train MSE is:

40.782937519716704

test MSE is:

45.84063733557119

diff is:

-5.057699815854484

train MSE is:

41.73079311544873

test MSE is:

42.877532999891436

diff is:

-1.1467398844427024

train MSE is:

41.33098523038746

test MSE is:

44.93500684381962

diff is:

-3.6040216134321597

train MSE is:

40.44832886730718

test MSE is:

47.83190326975075

diff is:

-7.383574402443564

train MSE is:

43.30836210121807

test MSE is:

37.15719576892592

diff is:

6.151166332292156

////d

