

hw1 - writing.

杨景兰 121090699

1.1 Given the following denominator layout derivatives, (13 points)

- **Differentiation of a scalar function w.r.t. a vector:** If  $f(\mathbf{w})$  is a scalar function of  $d$  variables,  $\mathbf{w}$  is a  $d \times 1$  vector, then differentiation of  $f(\mathbf{w})$  w.r.t.  $\mathbf{w}$  results in a  $d \times 1$  vector

$$\frac{df(\mathbf{w})}{d\mathbf{w}} = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix}$$

- **Differentiation of a vector function w.r.t. a vector:** If  $\mathbf{f}(\mathbf{w})$  is a vector function of size  $h \times 1$  and  $\mathbf{w}$  is a  $d \times 1$  vector, then differentiation of  $\mathbf{f}(\mathbf{w})$  w.r.t.  $\mathbf{w}$  results in a  $d \times h$  matrix

$$\frac{d\mathbf{f}(\mathbf{w})}{d\mathbf{w}} = \begin{bmatrix} \frac{\partial f_1}{\partial w_1} & \cdots & \frac{\partial f_h}{\partial w_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_1}{\partial w_d} & \cdots & \frac{\partial f_h}{\partial w_d} \end{bmatrix}$$

Please prove the following derivatives:

Consider  $\mathbf{X} \in \mathbb{R}^{h \times d}$  and  $\mathbf{y} \in \mathbb{R}^{h \times 1}$ , which are not functions of  $\mathbf{w}$ :

$$\frac{d(\mathbf{y}^T \mathbf{X} \mathbf{w})}{d\mathbf{w}} = \mathbf{X}^T \mathbf{y}, \quad (4 \text{ points})$$

$$\frac{d(\mathbf{w}^T \mathbf{w})}{d\mathbf{w}} = 2\mathbf{w}, \quad (4 \text{ points})$$

Consider  $\mathbf{X} \in \mathbb{R}^{d \times d}$  and  $\mathbf{w} \in \mathbb{R}^{d \times 1}$  (5 points):

$$\frac{d(\mathbf{w}^T \mathbf{X} \mathbf{w})}{d\mathbf{w}} = (\mathbf{X} + \mathbf{X}^T) \mathbf{w}$$

proof = ①  $\mathbf{y} \in \mathbb{R}^{h \times 1}$ ,  $\mathbf{X} \in \mathbb{R}^{h \times d}$ ,  $\mathbf{w} \in \mathbb{R}^{d \times 1}$

$$\mathbf{y}^T \mathbf{X} = [\mathbf{y}_1, \dots, \mathbf{y}_h] \cdot \begin{bmatrix} \mathbf{x}_{1,1} & \dots & \mathbf{x}_{1,d} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{h,1} & \dots & \mathbf{x}_{h,d} \end{bmatrix} = [\mathbf{y}_1 \mathbf{x}_{1,1} + \dots + \mathbf{y}_h \mathbf{x}_{h,1}, \dots, \mathbf{y}_1 \mathbf{x}_{1,d} + \dots + \mathbf{y}_h \mathbf{x}_{h,d}] \in \mathbb{R}^{1 \times d}$$

$$(\mathbf{y}^T \mathbf{X}) \mathbf{w} = [\mathbf{y}_1 \mathbf{x}_{1,1} + \dots + \mathbf{y}_h \mathbf{x}_{h,1}, \dots, \mathbf{y}_1 \mathbf{x}_{1,d} + \dots + \mathbf{y}_h \mathbf{x}_{h,d}] \cdot \begin{bmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_d \end{bmatrix} = (\mathbf{y}_1 \mathbf{x}_{1,1} + \dots + \mathbf{y}_h \mathbf{x}_{h,1}) \mathbf{w}_1 + \dots + (\mathbf{y}_1 \mathbf{x}_{1,d} + \dots + \mathbf{y}_h \mathbf{x}_{h,d}) \mathbf{w}_d \in \mathbb{R}^{1 \times 1}$$

$\Rightarrow$  differentiation of a scalar function w.r.t. a vector =

$$\frac{d(\mathbf{y}^T \mathbf{X} \mathbf{w})}{d\mathbf{w}} = \begin{bmatrix} \frac{\partial [(\mathbf{y}_1 \mathbf{x}_{1,1} + \dots + \mathbf{y}_h \mathbf{x}_{h,1}) \mathbf{w}_1 + \dots + (\mathbf{y}_1 \mathbf{x}_{1,d} + \dots + \mathbf{y}_h \mathbf{x}_{h,d}) \mathbf{w}_d]}{\partial \mathbf{w}_1} \\ \vdots \\ \frac{\partial [(\mathbf{y}_1 \mathbf{x}_{1,1} + \dots + \mathbf{y}_h \mathbf{x}_{h,1}) \mathbf{w}_1 + \dots + (\mathbf{y}_1 \mathbf{x}_{1,d} + \dots + \mathbf{y}_h \mathbf{x}_{h,d}) \mathbf{w}_d]}{\partial \mathbf{w}_d} \end{bmatrix} = \begin{bmatrix} \mathbf{y}_1 \mathbf{x}_{1,1} + \dots + \mathbf{y}_h \mathbf{x}_{h,1} \\ \vdots \\ \mathbf{y}_1 \mathbf{x}_{1,d} + \dots + \mathbf{y}_h \mathbf{x}_{h,d} \end{bmatrix}$$

$$\mathbf{X}^T \mathbf{y} = \begin{bmatrix} \mathbf{x}_{1,1} & \dots & \mathbf{x}_{h,1} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{1,d} & \dots & \mathbf{x}_{h,d} \end{bmatrix} \mathbf{X} \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_h \end{bmatrix} = \begin{bmatrix} \mathbf{y}_1 \mathbf{x}_{1,1} + \dots + \mathbf{y}_h \mathbf{x}_{h,1} \\ \vdots \\ \mathbf{y}_1 \mathbf{x}_{1,d} + \dots + \mathbf{y}_h \mathbf{x}_{h,d} \end{bmatrix} = \frac{d(\mathbf{y}^T \mathbf{X} \mathbf{w})}{d\mathbf{w}}$$

$$\text{Hence, } \frac{d(\mathbf{y}^T \mathbf{X} \mathbf{w})}{d\mathbf{w}} = \mathbf{X}^T \mathbf{y}$$

②  $\mathbf{w} \in \mathbb{R}^{d \times 1}$  =

$$\mathbf{w}^T \mathbf{w} = [\mathbf{w}_1, \dots, \mathbf{w}_d] \cdot \begin{bmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_d \end{bmatrix} = \mathbf{w}_1^2 + \dots + \mathbf{w}_d^2 \in \mathbb{R}^{1 \times 1}$$

$\Rightarrow$  Differentiation of a scalar function w.r.t. a vector =

$$\frac{d(\mathbf{w}^T \mathbf{w})}{d\mathbf{w}} = \begin{bmatrix} \frac{\partial (\mathbf{w}_1^2 + \dots + \mathbf{w}_d^2)}{\partial \mathbf{w}_1} \\ \vdots \\ \frac{\partial (\mathbf{w}_1^2 + \dots + \mathbf{w}_d^2)}{\partial \mathbf{w}_d} \end{bmatrix} = \begin{bmatrix} 2\mathbf{w}_1 \\ \vdots \\ 2\mathbf{w}_d \end{bmatrix} = 2\mathbf{w}$$

$$\text{Hence, } \frac{d(\mathbf{w}^T \mathbf{w})}{d\mathbf{w}} = 2\mathbf{w}$$

③  $\mathbf{w} \in \mathbb{R}^{d \times 1}$ ,  $\mathbf{X} \in \mathbb{R}^{d \times d}$  =

$$\mathbf{w}^T \mathbf{X} = [\mathbf{w}_1, \dots, \mathbf{w}_d] \begin{bmatrix} \mathbf{x}_{1,1} & \dots & \mathbf{x}_{1,d} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{d,1} & \dots & \mathbf{x}_{d,d} \end{bmatrix} = [\mathbf{w}_1 \mathbf{x}_{1,1} + \dots + \mathbf{w}_d \mathbf{x}_{d,1}, \dots, \mathbf{w}_1 \mathbf{x}_{1,d} + \dots + \mathbf{w}_d \mathbf{x}_{d,d}] \in \mathbb{R}^{1 \times d}$$

$$(\mathbf{w}^T \mathbf{X}) \mathbf{w} = [\mathbf{w}_1 \mathbf{x}_{1,1} + \dots + \mathbf{w}_d \mathbf{x}_{d,1}, \dots, \mathbf{w}_1 \mathbf{x}_{1,d} + \dots + \mathbf{w}_d \mathbf{x}_{d,d}] \cdot \begin{bmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_d \end{bmatrix} = [(\mathbf{w}_1 (\mathbf{w}_1 \mathbf{x}_{1,1} + \dots + \mathbf{w}_d \mathbf{x}_{d,1}) + \dots + \mathbf{w}_d (\mathbf{w}_1 \mathbf{x}_{1,d} + \dots + \mathbf{w}_d \mathbf{x}_{d,d}))] \in \mathbb{R}^{1 \times 1}$$

$\Rightarrow$  Differentiation of a scalar function w.r.t. a vector =

$$\frac{d(\mathbf{w}^T \mathbf{X} \mathbf{w})}{d\mathbf{w}} = \begin{bmatrix} \frac{\partial [(\mathbf{w}_1 (\mathbf{w}_1 \mathbf{x}_{1,1} + \dots + \mathbf{w}_d \mathbf{x}_{d,1}) + \dots + \mathbf{w}_d (\mathbf{w}_1 \mathbf{x}_{1,d} + \dots + \mathbf{w}_d \mathbf{x}_{d,d}))]}{\partial \mathbf{w}_1} \\ \vdots \\ \frac{\partial [(\mathbf{w}_1 (\mathbf{w}_1 \mathbf{x}_{1,1} + \dots + \mathbf{w}_d \mathbf{x}_{d,1}) + \dots + \mathbf{w}_d (\mathbf{w}_1 \mathbf{x}_{1,d} + \dots + \mathbf{w}_d \mathbf{x}_{d,d}))]}{\partial \mathbf{w}_d} \end{bmatrix} = \begin{bmatrix} (\mathbf{w}_1 \mathbf{x}_{1,1} + \dots + \mathbf{w}_d \mathbf{x}_{d,1}) + (\mathbf{w}_1 + \dots + \mathbf{w}_d) \mathbf{x}_{1,d} \\ \vdots \\ (\mathbf{w}_1 \mathbf{x}_{1,d} + \dots + \mathbf{w}_d \mathbf{x}_{d,d}) + (\mathbf{w}_1 + \dots + \mathbf{w}_d) \mathbf{x}_{d,d} \end{bmatrix} = \begin{bmatrix} 2\mathbf{w}_1 \mathbf{x}_{1,1} + \dots + (\mathbf{x}_{1,d} + \mathbf{x}_{d,1}) \mathbf{w}_d \\ \vdots \\ (\mathbf{x}_{1,d} + \mathbf{x}_{d,1}) \mathbf{w}_1 + \dots + 2\mathbf{x}_{d,d} \mathbf{w}_d \end{bmatrix}$$

$$(X+X^T)W = \left( \begin{bmatrix} x_{1,1} & \dots & x_{1,d} \\ \vdots & & \vdots \\ x_{d,1} & \dots & x_{d,d} \end{bmatrix} + \begin{bmatrix} x_{1,1} & \dots & x_{d,1} \\ \vdots & & \vdots \\ x_{d,1} & \dots & x_{d,d} \end{bmatrix} \right) \cdot \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} 2x_{1,1} & \dots & x_{1,d} + x_{d,1} \\ \vdots & & \vdots \\ x_{1,d} + x_{d,1} & \dots & 2x_{d,d} \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} 2x_{1,1}w_1 + \dots + (x_{1,d} + x_{d,1})w_d \\ \vdots \\ (x_{1,d} + x_{d,1})w_1 + \dots + 2x_{d,d}w_d \end{bmatrix} = \frac{d(W^T X W)}{dW}$$

Hence,  $\frac{d(W^T X W)}{dW} = (X+X^T)W$

1.2 Suppose we have training data  $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ , where  $x_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}^d, i = 1, 2, \dots, N$ . Consider  $f_{w,b}(x) = x^T w + b$ . (12 points)

(1) Find the closed-form solution of the following problem

$$\min_{w,b} \sum_{i=1}^N (f_{w,b}(x_i) - y_i)^2 + \lambda \bar{w}^T \bar{w}, \quad (1)$$

where  $\bar{w} = \begin{bmatrix} 0 \\ w \end{bmatrix} = [0, w_1, w_2, \dots, w_d]^T$ . (6 points)

1)  $f_{w,b}(x) = x^T w + b = XW = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} W = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & & \vdots \\ 1 & x_{N,1} & \dots & x_{N,d} \end{bmatrix} \begin{bmatrix} w_0 & \dots & w_d \\ w_1 & \dots & w_d \\ \vdots & & \vdots \\ w_d & \dots & w_d \end{bmatrix} = \begin{bmatrix} f_{w,0} & \dots & f_{w,d} \\ \vdots & & \vdots \\ f_{w,1} & \dots & f_{w,d} \end{bmatrix}$

*Annotations:  $\begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & & \vdots \\ 1 & x_{N,1} & \dots & x_{N,d} \end{bmatrix} \in \mathbb{R}^{N \times (d+1)}$ ,  $\begin{bmatrix} w_0 & \dots & w_d \\ w_1 & \dots & w_d \\ \vdots & & \vdots \\ w_d & \dots & w_d \end{bmatrix} \in \mathbb{R}^{(d+1) \times d}$ ,  $f_{w,0} \dots f_{w,d}$  are sample j's output.*

$$E = XW - Y$$

$$J(W) = \text{trace}(E^T E)$$

$$= \text{trace}[(XW - Y)^T (XW - Y)]$$

$$\min_{W,b} \sum_{i=1}^N (f_{w,b}(x_i) - y_i)^2 + \lambda \bar{w}^T \bar{w}$$

$$= \min_W (XW - Y)^T (XW - Y) + \lambda \bar{w}^T \bar{w}$$

$$\frac{\partial}{\partial W} (XW - Y)^T (XW - Y) + \lambda \bar{w}^T \bar{w} = 0$$

$$\Rightarrow 2X^T XW - 2X^T Y + 2\lambda \hat{I}_d W = 0$$

$$\Rightarrow (X^T X + \lambda \hat{I}_d) W = X^T Y$$

$$\Rightarrow W = (X^T X + \lambda \hat{I}_d)^{-1} X^T Y$$

Note that  $(X^T X + \lambda \hat{I}_d)^{-1}$  is guaranteed to be invertible given  $\lambda > 0$ .

(2) Show how to use gradient descent to solve the problem. (6 points)

$$\text{Assume } f(W) = \sum_{i=1}^N (f_{w,b}(x_i) - y_i)^2 + \lambda \bar{w}^T \bar{w} \\ = (XW - Y)^T (XW - Y) + \lambda \bar{w}^T \bar{w}$$

start from  $W_0$ , set  $\epsilon > 0, k=0$ .

1) If  $\|2X^T XW_k - 2X^T Y + 2\lambda \hat{I}_d W_k\| < \epsilon$ , output  $W_k$ .

Otherwise, continue.

2) compute search direction  $d^k = -(2X^T XW_k - 2X^T Y + 2\lambda \hat{I}_d W_k)$

3) choose step size  $\alpha_k$  via exact line search or backtracking line search.

4) set  $W_{k+1} = W_k + \alpha_k d^k, k=k+1$

Go back to step 1.

1.3 Prove that:

(1)  $f(x) = x^2$  is convex. (4 points)

(2) every affine function  $f(x) = ax + b$  is convex, but not strictly convex; (4 points)

(3)  $f(x) = |x|$  is convex, but not strictly convex; (5 points)

Proof: (1)  $f(x) = x^2, x \in \mathbb{R}$ .

$$\Rightarrow f'(x) = 2x$$

$$\Rightarrow f''(x) = 2 > 0$$

$\forall x \in \mathbb{R}, f''(x) > 0$ . Therefore,  $f(x)$  is convex.

(2)  $f(x) = ax + b$

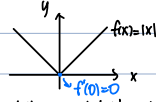
$$\Rightarrow f'(x) = a$$

$$\Rightarrow f''(x) = 0$$

$\forall x \in \mathbb{R}, f'(x) = 0$ . Therefore,  $f(x)$  is convex but not strictly convex.

(3)  $f(x) = |x|$

Suppose  $\forall x, y \in \mathbb{R}, \forall \alpha \in [0, 1]$ .



$$f(\alpha x + (1-\alpha)y) = |\alpha x + (1-\alpha)y| \leq |\alpha x| + |\alpha(1-\alpha)y| = \alpha|x| + (1-\alpha)|y| = \alpha f(x) + (1-\alpha)f(y)$$

When  $xy > 0$ , the equality holds.

Therefore,  $f(x)$  is convex but not strictly convex

1.4 Suppose  $x_1, x_2, \dots, x_N$  are drawn from  $\text{Laplace}(\mu, b)$ . Calculate the MLE (maximum likelihood estimation) of  $\mu$  and  $b$ . Hint: Use the logarithmic trick to process multiple exponential items. (12 points)

$$p(x) = \frac{1}{2b} \exp\left(-\frac{|x-\mu|}{b}\right)$$

$$L(x) = \prod_{i=1}^N p(x_i) = \left(\frac{1}{2b}\right)^N \exp\left(-\frac{\sum_{i=1}^N |x_i - \mu|}{b}\right)$$

$$\ell(x) = -N \log(2b) - \frac{1}{b} \sum_{i=1}^N |x_i - \mu|$$

$$\frac{\partial \ell(x)}{\partial b} = -\frac{N}{b} + \frac{1}{b^2} \sum_{i=1}^N |x_i - \mu| = 0 \Rightarrow \hat{b} = \frac{1}{N} \sum_{i=1}^N |x_i - \mu|$$

$$\frac{\partial \ell(x)}{\partial \mu} = -\frac{1}{b} \sum_{i=1}^N \frac{\partial |x_i - \mu|}{\partial \mu} \quad \frac{\partial |x|}{\partial x} = \frac{\partial \sqrt{x^2}}{\partial x} = x(x^2)^{-\frac{1}{2}} = \frac{x}{|x|} = \text{sgn}(x)$$

$$= -\frac{1}{b} \sum_{i=1}^N \text{sgn}(x_i - \mu)$$

$$-\frac{1}{b} \sum_{i=1}^N \text{sgn}(x_i - \mu) = 0$$

①  $N$  is odd. we choose  $\hat{\mu} = \text{median}(x_1, \dots, x_N)$

②  $N$  is even. Ranking the observations as  $x_{(1)} \leq x_{(2)} \leq \dots, x_{(N)}$  and then choose  $\hat{\mu} = x_{(N/2)}$  or  $\hat{\mu} = x_{(N+1)/2}$

In summary,  $\hat{\mu} = \text{median}(x_1, \dots, x_N)$  is the MLE of  $\mu$ , and  $\hat{b} = \frac{1}{N} \sum_{i=1}^N |x_i - \mu|$  is the MLE of  $b$ .

## DDA3020 Assignment 1 report (coding part)

### 2. Programming Question

#### 2.1 Load data and summarize characteristics

Use `pd.read_csv()` function to read the data file `boston.csv`. Notice that the loading of the data is using relative path instead of absolute path. Hence, when executing the code of the model, one should place the train data `txt` file and test data `test` file under the same directory of the model python file.

Below is the general content of the boston data. (generated by `df.head()`)

	crim	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	b	lstat	medv
0	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242	17.8	392.83	4.03	34.7
3	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222	18.7	396.90	5.33	36.2

We guess that 'medv' and 'rm' have a strong relationship, since as 'rm' becomes larger, 'medv' also becomes larger.

I use 3 methods to test whether there exists 'NaN' or 'Null' in this dataset.:

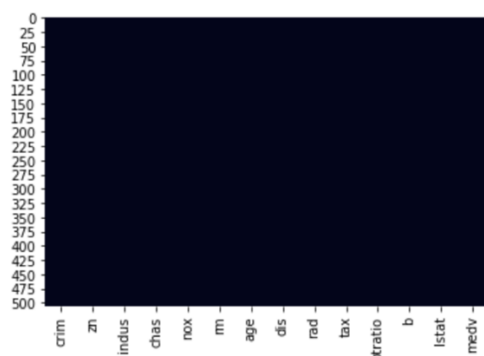
#### 1. `df.isnull()`

	crim	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	b	lstat	medv
0	False	False	False	False	False	False	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False	False	False	False	False	False	False
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
501	False	False	False	False	False	False	False	False	False	False	False	False	False	False
502	False	False	False	False	False	False	False	False	False	False	False	False	False	False
503	False	False	False	False	False	False	False	False	False	False	False	False	False	False
504	False	False	False	False	False	False	False	False	False	False	False	False	False	False
505	False	False	False	False	False	False	False	False	False	False	False	False	False	False

506 rows x 14 columns

NaN (Null) corresponds to True, and otherwise.

#### 2. `sns.heatmap(df.isnull(), cbar=False)`



The picture is all black, no white space in it → no NaN (Null) in this dataset.

#### 3. `df.isna().sum()`

```

crim      0
zn        0
indus     0
chas      0
nox       0
rm        0
age       0
dis       0
rad       0
tax       0
ptratio   0
b         0
lstat     0
medv      0
dtype: int64

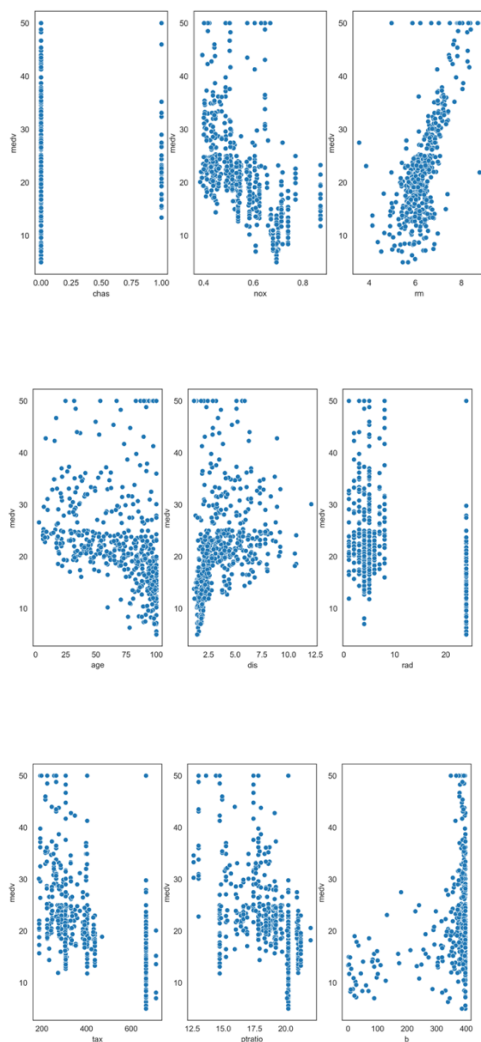
```

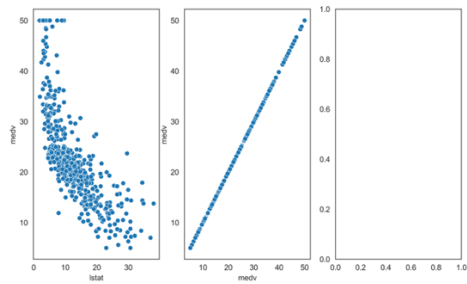
As the function name said, count the sum of the number of NaN or Null in every column of this dataset, respectively.

Hence, we conclude that there is no Nan or Null in this boston dataset.

2.2 visualize the dataset

Use `sns.scatterplot()` to plot the MEDV distributions over each attribute. And below are the figures we get.





From the figures above, we insist the assumption for the most relevant attribute for MEDV is RM.

### 2.3 Pairwise correlation on data

Use `sns.heatmap(df.corr(),annot=True)` to draw the covariance figures.



Assume if covariance  $>0.4$  or  $<-0.4$ , then the attribute is good predictor

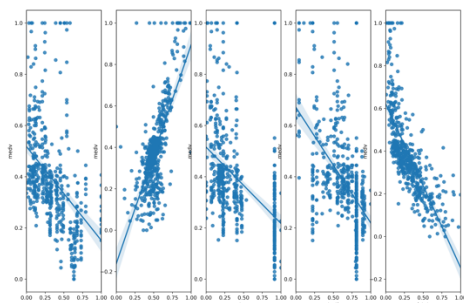
The covariances between medv and nox,rm,tax,ptratio,lstat are -0.43,0.7,-0.47,-0.51,and -0.74, respectively.

Good predictors: nox,rm,tax,ptratio,and lstat

### 2.4 scale the columns and draw the relevance

Use `MinMaxScaler()` to scale the attributes.

Then use `seaborn.regplot()` to plot the relevance of these columns against MEDV with 95% confidence interval.



### 2.5 train linear regression model and test on it

(1). linear model description

$$f_{w,b} = x^T w + b$$

The set of feature vector  $x$  is the columns of “nox,rm,tax,ptratio,and lstat” in the boston

dataset.

The target output  $y$  is the column of "medv" in the boston dataset.

We solve this linear regression by gradient descent.

We split the data into two parts: training (80%) and test (20%).

(2). loss function

$$\text{loss function: } J(\bar{w}) = \frac{1}{2} \sum_{i=1}^m (x_i^T w + b - y_i)^2 = \frac{1}{2} (X\bar{w} - y)^2$$

$$\bar{w}^* = \arg_{\bar{w}} \min J(\bar{w})$$

$w$  can be updated by gradient descent algorithm,

$$w \leftarrow w - \alpha \frac{\partial J(w)}{\partial w}, \frac{\partial J(w)}{\partial w} = X^T (Xw - y)$$

Where  $X = [(1, x_1^T); (1, x_2^T); \dots; (1, x_m^T)] \in \mathcal{R}^{m \times (d+1)}$ , and  $\bar{w} = [b; w] \in \mathcal{R}^{(d+1) \times 1}$

We transform  $x_i^T w + b$  into  $X\bar{w}$  by "get\_full\_sample\_matrix" this self-define function.

(3). hyperparameter settings

After transforming into  $X\bar{w}$ , we only need to train one hyperparameter  $\bar{w}$ .

(4). RMSE equation

Mathematic equation of RMSE is:

$$RMSE = \sqrt{\sum_{i=1}^m \frac{(\bar{y}_i - y_i)^2}{m}}$$

In python, we use `RMSE=np.sqrt(metrics.mean_squared_error(y, \bar{y}))` to compute.

(5). outputs (errors, plots, figures)

Assume step size=0.0001, iteration number=50;

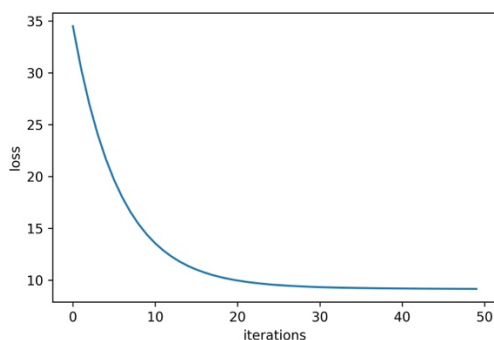
we get that:

`W= [[0.19354444],[0.05286568],[0.10555365],[0.05506099],[0.0970676 ],[0.04649039]]`

RMSE of train= 0.20864370794310755

RMSE of test= 0.22787361921161076

the loss curve in the training process is as below:



2.6 repeat

We choose step size from list [0.01,0.001,0.0001,0.00001], iteration number from list

[50,100,500].

Repeat 2.5 and get solutions as below:

1. step size= 0.01 iteration steps= 50

RMSE of train= 1.176899655389494e+43

RMSE of test= 1.2011007219694868e+43

2. step size= 0.001 iteration steps= 50

RMSE of train= 0.2058516518776543

RMSE of test= 0.19632119342529225

3. step size= 0.001 iteration steps= 100

RMSE of train= 0.20649137801162187

RMSE of test= 0.1929998937307037

4. step size= 0.001 iteration steps= 500

RMSE of train= 0.19124148666161425

RMSE of test= 0.24220397879116876

5. step size= 0.0001 iteration steps= 50

RMSE of train= 0.2140196628121716

RMSE of test= 0.18928550674499142

6. step size= 0.0001 iteration steps= 100

RMSE of train= 0.2098229912994151

RMSE of test= 0.20121029612621325

7. step size= 0.0001 iteration steps= 500

RMSE of train= 0.20552547567407597

RMSE of test= 0.20527230239717645

8. step size= 1e-05 iteration steps= 50

RMSE of train= 0.3284879924913493

RMSE of test= 0.3325026354759041

9. step size= 1e-05 iteration steps= 100

RMSE of train= 0.2708984994062212

RMSE of test= 0.2611387972606977

10. step size= 1e-05 iteration steps= 500

RMSE of train= 0.2114803182586758

RMSE of test= 0.20545886408108904

Conclusion, as step size increases, it takes less iterations to get close to the optimal solution.

But if step size is too big, the result may diverge, and loss may be too large.