

$$1.1 \quad \textcircled{1} \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix} \quad W = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$$Y^T X = \left[ \frac{\partial}{\partial x_{ij}} y_i \right] \left[ \frac{\partial}{\partial x_{ij}} x_{ij} \right] \cdots \left[ \frac{\partial}{\partial x_{ij}} x_{id} \right] \quad |x|_d$$

$$Y^T X W = \sum_{i=1}^d w_i \sum_{j=1}^h x_{ij} y_i \quad |x|$$

$$\frac{\partial (Y^T X W)}{\partial W} = \left[ \frac{\partial f}{\partial w_i} \right] \cdot \underbrace{W \cdot \frac{\partial f}{\partial w_i}}_{= \frac{\partial (w_i \sum_{j=1}^h x_{ij} y_i)}{\partial w_i}}$$

$$\frac{\partial f}{\partial w_i} = \frac{\partial (w_i \sum_{j=1}^h x_{ij} y_i)}{\partial w_i} + \cancel{\frac{\partial (\sum_{k \neq i} w_k \sum_{j=1}^h x_{kj} y_j)}{\partial w_i}}$$

$$= \sum_{j=1}^h x_{ij} y_i + 0 = \sum_{j=1}^h x_{ij} y_i \quad \text{for } \forall i \in \{1, \dots, d\}$$

$$\frac{\partial (Y^T X W)}{\partial W} = \begin{bmatrix} \sum_{j=1}^h x_{1j} y_i \\ \vdots \\ \sum_{j=1}^h x_{nj} y_i \end{bmatrix} = \cancel{X^T Y} \begin{bmatrix} (X^T)_1 Y \\ (X^T)_2 Y \\ \vdots \\ (X^T)_d Y \end{bmatrix} = X^T Y$$

Q.E.D.

$$2) \quad W^T W = \sum_{i=1}^d w_i^2$$

$$\frac{\partial (W^T W)}{\partial W} = \left[ \frac{\partial (W^T W)}{\partial w_i} \right] = \begin{bmatrix} 2w_1 \\ 2w_2 \\ \vdots \\ 2w_d \end{bmatrix} = 2 \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = 2W$$

Q.E.D.

$$(3) \quad W^T X = \begin{bmatrix} w_1 & w_2 & \cdots & w_d \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1d} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{bmatrix} = \sum_{i=1}^d w_i x_{ii} - \cancel{\sum_{i=1}^d w_i x_{id}}$$

$$\cancel{\sum_{i=1}^d w_i x_{id}} \rightarrow (3)$$

$$(3) \quad W^T X = \sum_{i=1}^d w_i x_{ii} - \sum_{i=1}^d w_i x_{id} \rightarrow$$

$$(3) \quad W^T X = \left[ \sum_{i=1}^d w_i x_{ii}, \sum_{i=1}^d w_i x_{12}, \dots, \sum_{i=1}^d w_i x_{id} \right]$$

$$W^T X W = \sum_{j=1}^h w_j \sum_{i=1}^d x_{ij} w_i \quad \frac{\partial (W^T X W)}{\partial W} = \left[ \frac{\partial f_{W^T X W}}{\partial w_1}, \frac{\partial f_{W^T X W}}{\partial w_2}, \dots \right]$$

$$\frac{\partial (W^T X W)}{\partial w_k} = \frac{\partial (w_k \sum_{i=1}^d x_{ik} w_i)}{\partial w_k} + \frac{\partial (\sum_{j \neq k} w_j \sum_{i=1}^d x_{ij} w_i)}{\partial w_k}$$

$$\begin{aligned} &= \sum_{i=1}^h x_{ik} w_i + w_k x_{kk} + \sum_{j \neq k} w_j x_{kj} \\ &= \sum_{i=1}^h x_{ik} w_i + \sum_{j=1}^d x_{kj} w_j = X^T W + (X^T)_k W \\ &= \sum_{i=1}^h x_{ij} y_i \quad \text{for } \forall k \in \{1, 2, \dots, d\} \end{aligned}$$

$$\begin{aligned} \text{so } \frac{\partial (W^T X W)}{\partial W} &= \begin{bmatrix} X^T W + (X^T)_1 W \\ X^T W + (X^T)_2 W \\ \vdots \\ X^T W + (X^T)_d W \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} W + \cancel{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} W} \\ &= \cancel{X^T W} + X^T W = (X + X^T) W \end{aligned}$$

(2)

(1) Let  $X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix}$   $|N|d$   $Y = \begin{bmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_N^T \end{bmatrix}$   $|N|d$

$\min_{W, b} (f_{W, b}(x_i)) \quad W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$

1.2

$$(1) \quad \text{Let } X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} (N \times (d+1)) \quad Y = \begin{bmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_N^T \end{bmatrix} (N \times 1)$$

$$\bar{W} = \begin{bmatrix} -b \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} (d+1) \times 1 \quad f_{W, b}(x_i) \quad f_{W, b}(X) = X \bar{W}$$

$$J(W) = \frac{1}{N} \sum_{i=1}^N (f_{W, b}(x_i) - y_i)^2 + \lambda \bar{W}^T \bar{W} = \|X \bar{W} - Y\|_2^2 + \lambda \bar{W}^T \bar{W}$$

$$\cancel{X \bar{W}} = (X \bar{W} - Y) + \lambda \bar{W}^T \bar{W}$$

$$\therefore \nabla J(W) = 2X^T(X \bar{W} - Y) + 2\lambda \bar{W} = (2X^T X + 2\lambda I) \bar{W} = 2X^T Y$$

$$\nabla J(W) = 0 \Rightarrow (2X^T X + 2\lambda I) \bar{W} = 2X^T Y$$

$$\bar{W} = (X^T X + 2\lambda I)^{-1} X^T Y$$

(2) ① for each step  $x_t$ , calculate  $\nabla J(W)$

② choose learning rate  $\alpha$  by backtracking

③ update  $\bar{W}_t$  by  $\bar{W}_{t+1} = \bar{W}_t - \alpha \nabla J(W)$

1.3

(1) since  $x \in \mathbb{R}$ , the domain is a convex set.  
 $\forall x, y \in \mathbb{R}$

$$\begin{aligned} f(\lambda x + (1-\lambda)y) &= \lambda^2 x^2 + (1-\lambda)^2 y^2 \\ \text{since } \lambda \in [0, 1] \\ 0 < \lambda^2 < 1, 0 < (1-\lambda)^2 &\Rightarrow \lambda^2 x^2, (1-\lambda)^2 y^2 \\ \text{since } x^2 \geq 0, y^2 \geq 0 \\ \lambda^2 x^2 \leq \lambda x^2, (1-\lambda)^2 y^2 &\\ f(\lambda x + (1-\lambda)y) &= (\lambda x + (1-\lambda)y)^2 \end{aligned}$$

$$\begin{aligned} f'(x) &= 2x & f''(x) &= 2 \\ \forall d \in \mathbb{R} \quad d^T f'(x) d &= 2d^2 \geq 0 \end{aligned}$$

so  $f(x)$  is PSD  $\Rightarrow f(x)$  is convex

(2)  $\forall x, y \in \mathbb{R}$ . domain is  $\mathbb{R}$ , hence convex

$\forall x, y \in \mathbb{R}$

$$\begin{aligned} f(\lambda x + (1-\lambda)y) &\leq \lambda f(x) + (1-\lambda)f(y) \leq \lambda(ax+b) + (1-\lambda)(ay+b) \\ &= a(\lambda x + (1-\lambda)y) + \lambda b + (1-\lambda)b \\ &= a(\lambda x + (1-\lambda)y) + b = f(\lambda x + (1-\lambda)y) \end{aligned}$$

$$f(\lambda x + (1-\lambda)y) = \lambda f(x) + (1-\lambda)f(y)$$

$$f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$$

so  $f(x)$  is convex

since  $\exists x, y$ , s.t.  $f(\lambda x + (1-\lambda)y) = \lambda f(x) + (1-\lambda)f(y)$

$f(x)$  is not strictly convex

(3)  $f(x) \quad x \in \mathbb{R} \Rightarrow$  domain convex

$$\forall x, y \in \mathbb{R} \quad f(\lambda x + (1-\lambda)y) = \lambda|x|/\lambda|x| + \lambda|y|$$

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

$$= \lambda f(x) + (1-\lambda)f(y)$$

$$\text{hence } f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$$

hence  $f(x)$  is convex

when  $x \geq 0, y \geq 0$

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

$$\text{hence } \exists x, y \in \mathbb{R} \quad f(\lambda x + (1-\lambda)y) = \lambda f(x) + (1-\lambda)f(y) \not\Rightarrow \text{strictly convex}$$

1.4

$$x_i \sim \text{Laplace}(u, b) \Rightarrow f(x_i) = \frac{1}{2b} \exp(-\frac{|x_i - u|}{b})$$

$$(u, b)_{\text{MLB}} = \arg \max_{(u, b)} \log \left( \prod_{i=1}^N f(x_i | u, b) \right)$$

$$\log \left( \prod_{i=1}^N f(x_i | u, b) \right) = \log \left( \prod_{i=1}^N \frac{1}{2b} \exp(-\frac{|x_i - u|}{b}) \right)$$

$$\begin{aligned} &= -\log \left( \frac{1}{2b} \right)^N \exp \left( -\frac{|x_i - u|}{b} \right) = \log \left( \frac{1}{2b} \right)^N + \log \exp \left( -\frac{|x_i - u|}{b} \right) \\ &= -N \log \left( \frac{1}{2b} \right) - \frac{|x_i - u|}{b} \\ &= \arg \max_{(u, b)} \left( -N \log \left( \frac{1}{2b} \right) - \frac{|x_i - u|}{b} \right) \\ &= \text{argmax} \end{aligned}$$

$$\begin{aligned} &= \log \left( \frac{1}{2b} \right)^N \exp \left( -\frac{|x_i - u|}{b} \right) = \log \left( \frac{1}{2b} \right)^N + \log \exp \left( -\frac{|x_i - u|}{b} \right) \\ &= -N \log \left( \frac{1}{2b} \right) + \frac{\sum_{i=1}^N |x_i - u|}{b} \end{aligned}$$

hence

$$(u, b)_{\text{MLB}} = \arg \max_{(u, b)} -N \log \left( \frac{1}{2b} \right) - \sum_{i=1}^N \frac{|x_i - u|}{b}$$

$$= \arg \min_{(u, b)} N \log \left( \frac{1}{2b} \right) + \sum_{i=1}^N \frac{|x_i - u|}{b}$$



# Report-Assignment1

## Linear Model

In this assignment, I use ridge linear regression to fit the data. Specifically, given an input selected features  $x \in R^d$ , where d is the number of features, the predicted housing price  $\hat{y}$  can be predicted by

$$\hat{y} = \bar{x}^T \bar{w}$$

where

$$\bar{x} = [1, x_1, x_2, \dots, x_d]^T \in R^{d+1}$$

$$\bar{w} = [b, w_1, w_2, \dots, w_d]^T \in R^{d+1}$$

## Loss Function

Since I use ridge linear regression, the parameters  $\bar{w}$  is trained by minimizing the loss function:

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_{\bar{w}}(x_i) - y_i)^2 + \frac{\lambda}{2} \bar{w}^T \bar{w}$$

where m is the number of training samples.

The loss function can also be written as

$$J(w) = \frac{1}{2m} \|X\bar{w} - Y\|^2 + \frac{\lambda}{2} \hat{w}^T \hat{w}$$

where

$$X = \begin{matrix} \bar{x}_1^T \\ \bar{x}_2^T \\ \dots \\ \bar{x}_m^T \end{matrix} \in R^{m \times (d+1)}$$

$$Y = \begin{matrix} y_1 \\ y_2 \\ \dots \\ y_m \end{matrix} \in R^{m \times (d+1)}$$

$$\hat{w} = [0, w_1, w_2, \dots, w_d]^T \in R^{d+1}$$

## Gradient Descend

The procedure of gradient descend can be summarized as:

- Calculate the gradient of loss function w.r.t. weight

$$\nabla J(\bar{w}) = \frac{\partial J}{\partial \bar{w}}$$

- Calculate learning rate  $\alpha$  by Armijo Line Search
- Update  $w$  by  $\bar{w} = \bar{w} - \alpha \nabla J(\bar{w})$
- Check whether the procedure should be terminated
- Go back to step 1 if not terminated

The gradient can be calculated as:

$$\nabla J(\bar{w}) = \frac{1}{m} X^T (X \bar{w} - Y) + \lambda \hat{w}$$

The learning rate  $\alpha$  for each step is selected from a set of weights  $\{1, \sigma, \sigma^2, \sigma^3, \dots\}$  by trying through the set starting from 1, until

$$J(w^k - \alpha_k \nabla J(w^k)) - J(w^k) \leq -\gamma \alpha_k \nabla J(w^k)^T \nabla J(w^k)$$

Whether to terminate the procedure is determined by the loss difference before and after the weight is updated. The procedure is terminated when

$$||J(w^k) - J(w^k - \alpha_k \nabla J(w^k))||^2 \leq threshold$$

Therefore, in total, there are 4 hyperparameters:  $\lambda$ ,  $\sigma$ ,  $\gamma$  and  $threshold$ .

## Features Selection

In order to avoid overfitting, instead of simply using all features, I only select several features for price prediction. In this report, only part of analysis plots will be shown. To see all plots, please refer to code.ipynb. Features that have an obvious linear relationship with the price will be chosen. Specifically, features rm, age, dis and lstat are chosen. As is shown in the following plot, these features have strong relationship with price.

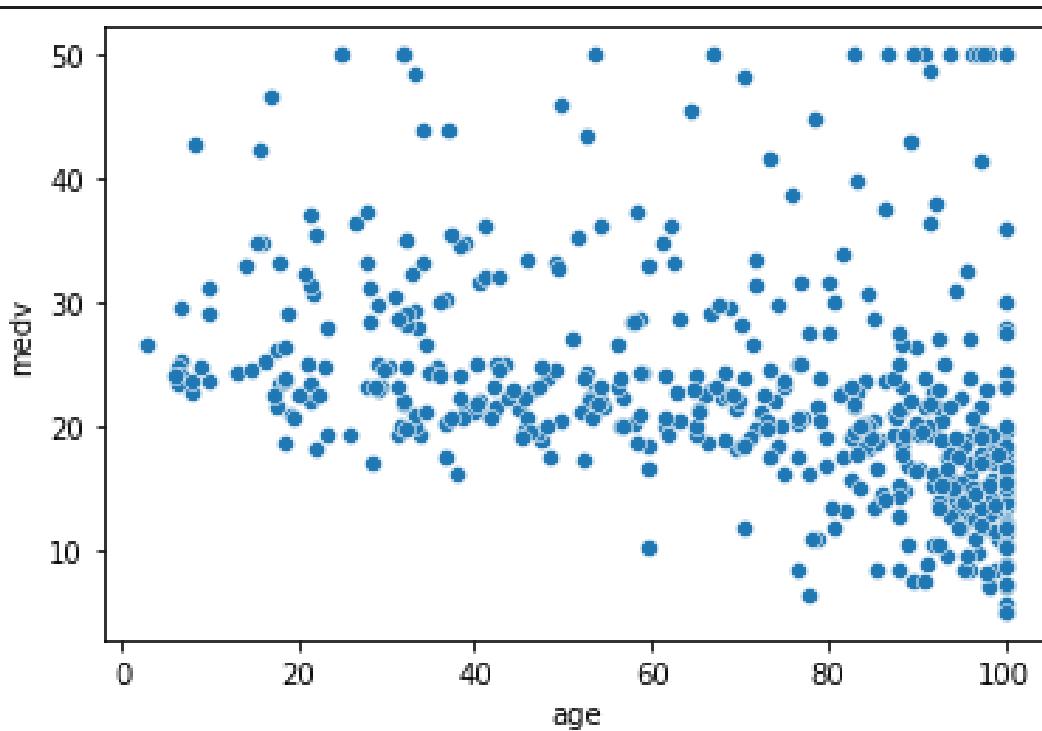


Figure 1: Relationship between age and medv

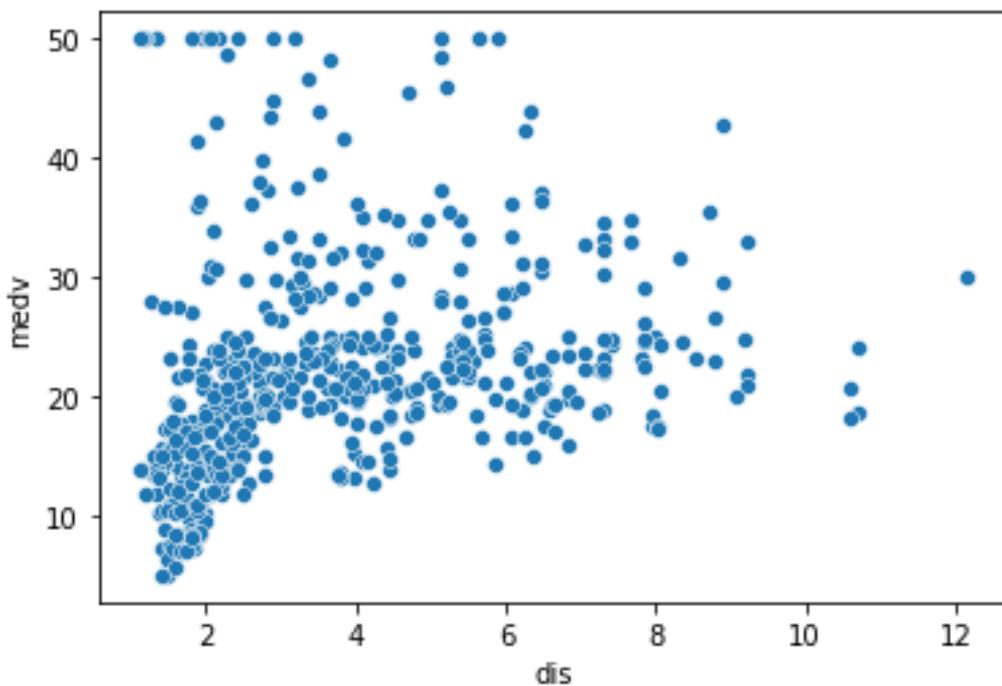


Figure 2: Relationship between `dis` and `medv`

---

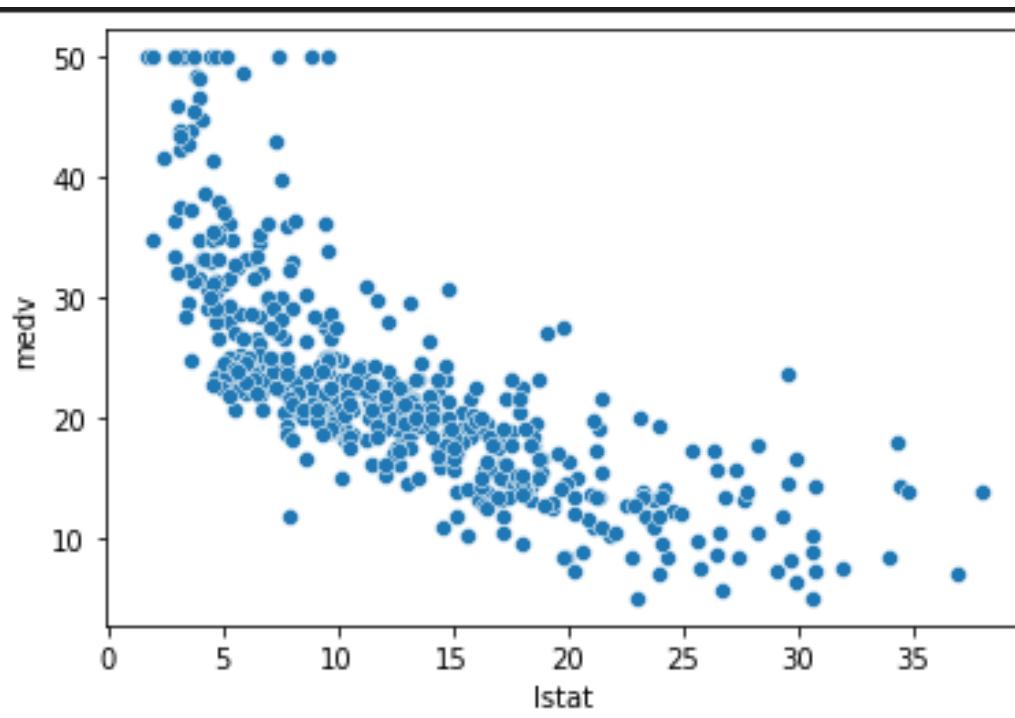


Figure 3: Relationship between `Istat` and `medv`

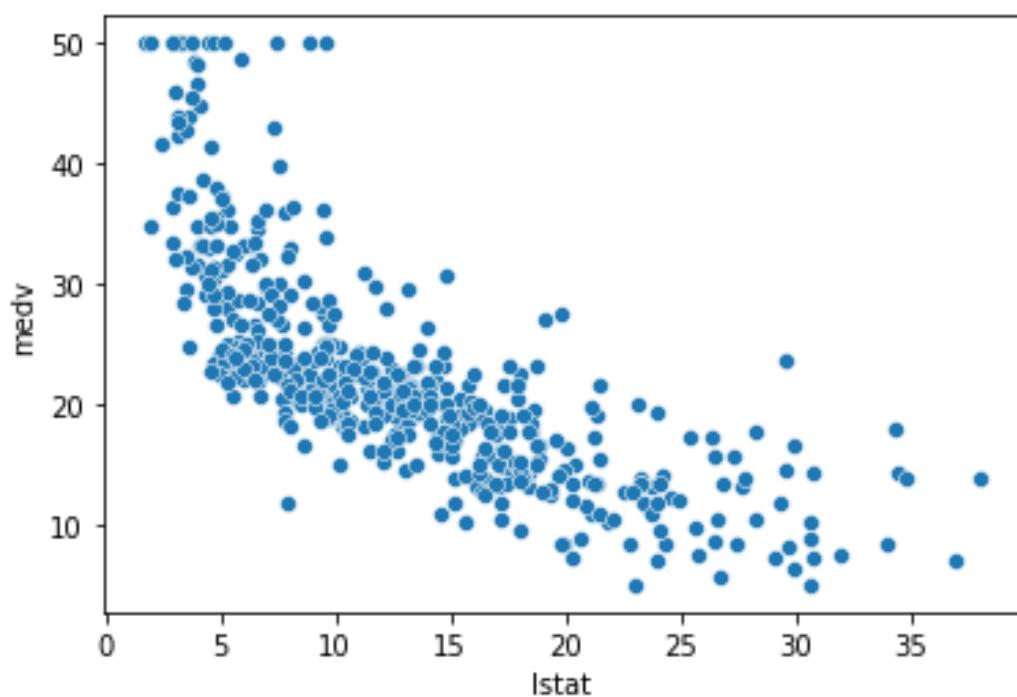


Figure 4: Relationship between rm and medv

As is shown in the heatmap, these features are not depenednt on each other. None of the two features are of the same distribution.

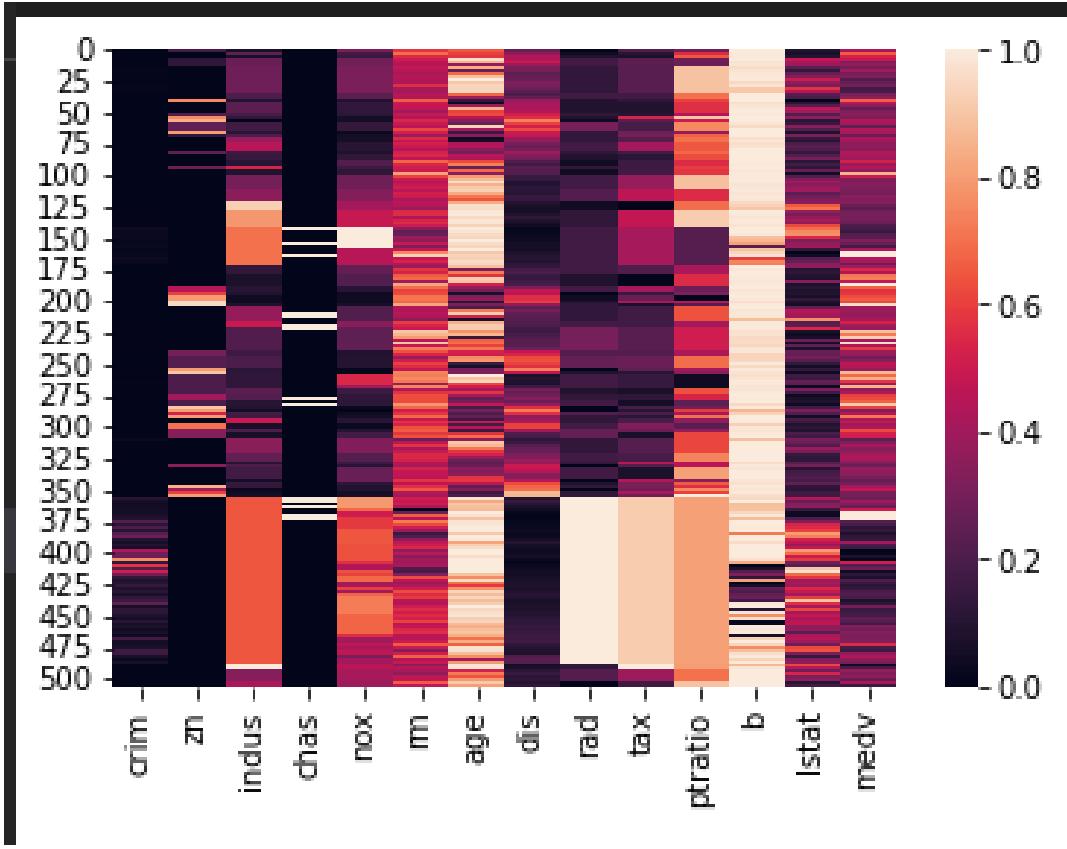


Figure 5:Heatmap

However, there exists a weak correlation between some pairs, as is shown in relevance plot. Here only several of these plots will be shown. To see all plots, please refer to code.ipynb.

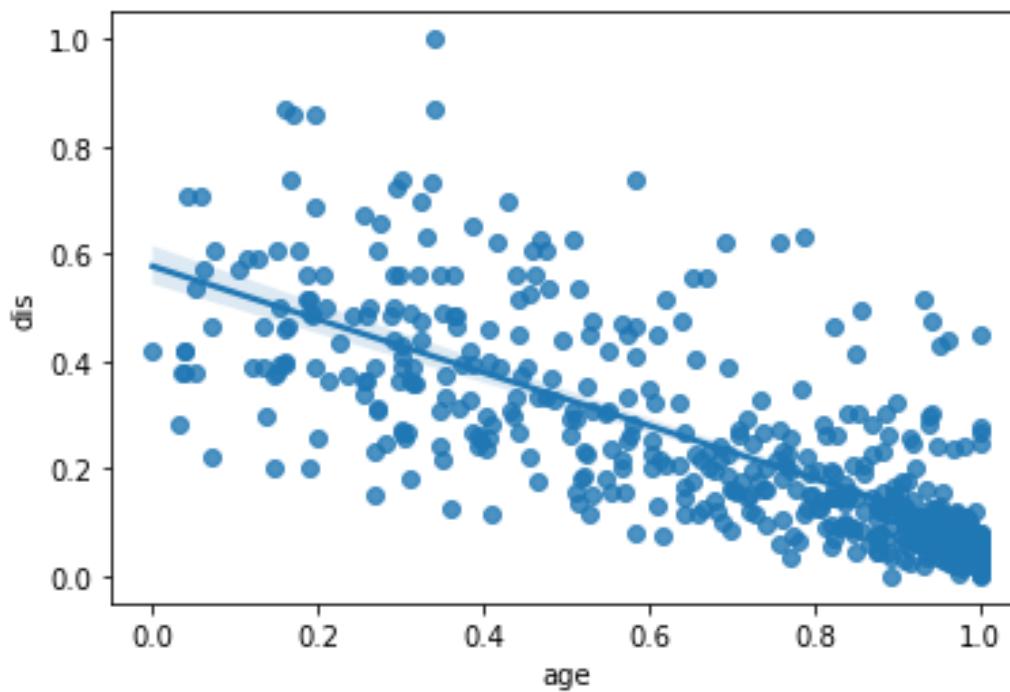


Figure 6: Relation between age and dis

---

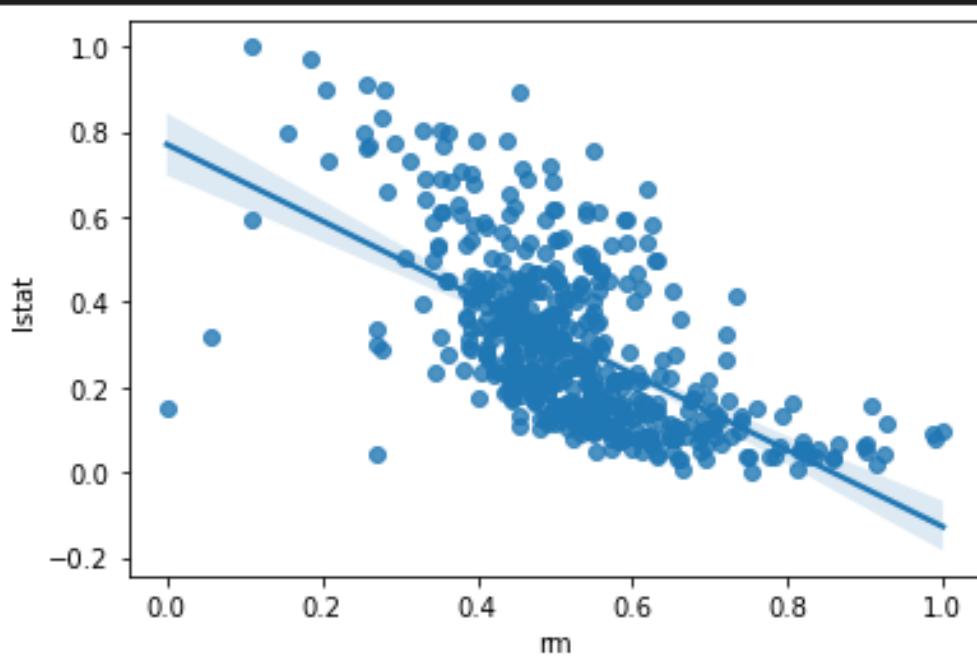


Figure 7: Relation between rm and lstat

---

# Training and Testing

10 experiments are conducted. In each experiment, the data is randomly split into training and testing set in a ratio of 8:2. In each experiment, the linear regression model is trained from random initialization, and the tested by the testing set.

## Training

In training, the hyperparameters are set as

$$\lambda = 0.01$$

$$\sigma = 0.5$$

$$\gamma = 0.1$$

$$threshold = 10^{-6}$$

In 10 experiments, the training is terminated at around 170 to 200 iterations. Two plots of loss function value over iterations are shown here. To see all plots, please refer to code.ipynb.

Experiment-10

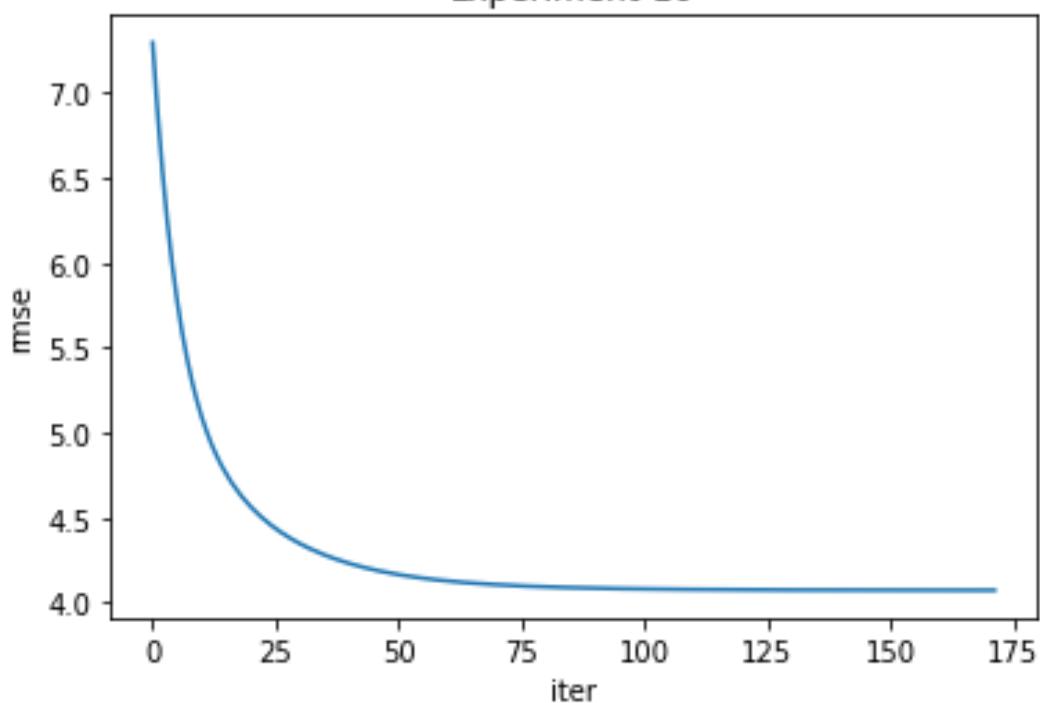


Figure 8:Loss over iteration at experiment1

Experiment-1

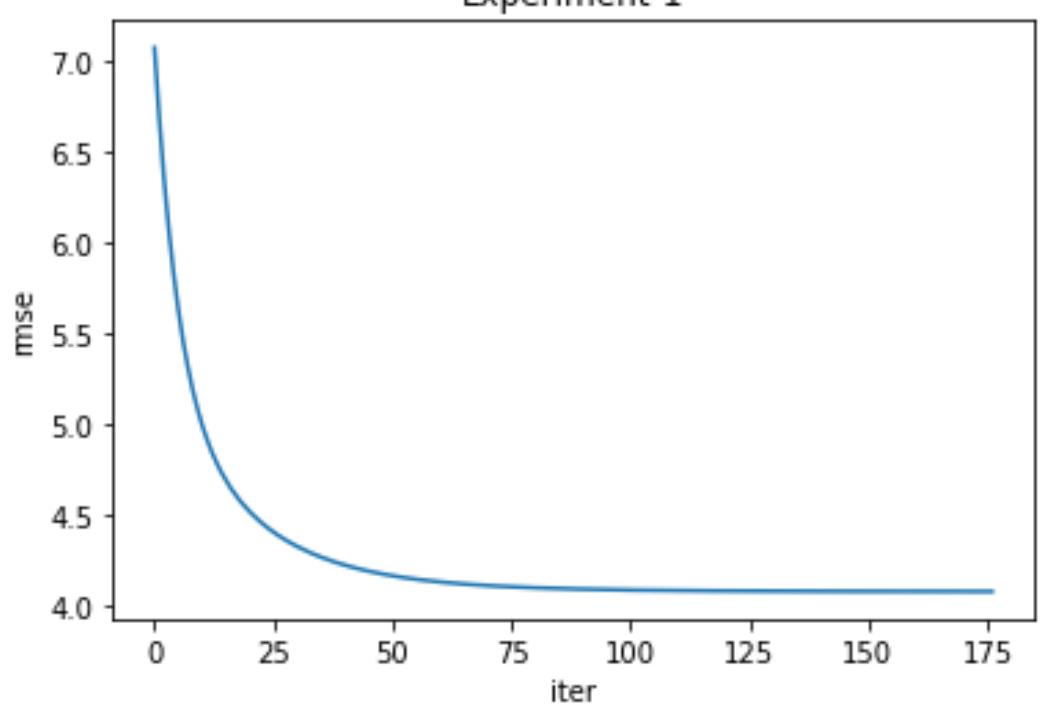


Figure 9:Loss over iteration at experiment10

---

## Testing

RMSE is used in testing the model.

$$rmse(w) = \sqrt{\frac{1}{k} \sum_{i=1}^k (f_{\bar{w}}(x_i) - y_i)^2}$$

where k is the number of testing sample. The mean testing rmse over 10 experiments is 3.807362106194172. To see the testing results of all 10 experiments, please refer to code.ipynb.