

# Advanced Data Mining Homework 4

1. Summarize and derivation linear regression with two regularizers including vanilla one.

A. Derivation vanilla linear regression and Ridge

① vanilla linear regression

$$\mathcal{E}_{LS} = \frac{1}{2} \sum_{t=1}^N \{y_t - w^T \phi(x_t)\}^2 = \frac{1}{2} \|y - \Phi w\|^2$$

$$\frac{\partial \mathcal{E}_{LS}}{\partial w} = \frac{\partial}{\partial w} \left\{ \frac{1}{2} (y - \Phi w)^T (y - \Phi w) \right\} = \frac{\partial}{\partial w} \left\{ \frac{1}{2} (y^T y - w^T \Phi^T y - y^T \Phi w + w^T \Phi \Phi w) \right\}$$

$$= \frac{1}{2} \{ -\Phi^T y - \Phi^T y + 2\Phi^T \Phi w \} = 0$$

$$= -\Phi^T y + \Phi^T \Phi w = 0$$

$$\Phi^T \Phi w = \Phi^T y$$

$$w^* = (\Phi^T \Phi)^{-1} \Phi^T y$$

② Ridge regression

Add regularized term in mean square error

$$\mathcal{E}_{LS}^{\text{Ridge}} = \frac{1}{2} \|y - \Phi w\|^2 + \frac{\lambda}{2} \|w\|^2$$

$$= \frac{1}{2} (y - \Phi w)^T (y - \Phi w) + \frac{\lambda}{2} \|w\|^2$$

$$\frac{\partial \mathcal{E}_{LS}^{\text{Ridge}}}{\partial w} = \frac{\partial}{\partial w} \left\{ \frac{1}{2} (y^T y - w^T \Phi^T y - y^T \Phi w + w^T \Phi \Phi w) + \frac{\lambda}{2} w^T w \right\}$$

$$= \frac{1}{2} \{ -2\Phi^T y + 2\Phi^T \Phi w + 2\lambda w \} = 0$$

$$(\Phi^T \Phi + \lambda I) w = \Phi^T y$$

$$w_{\text{Ridge}}^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$

B. Derivation linear regression using MLE

$$y = \phi w + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I)$$

Then log-likelihood is given by

$$\begin{aligned} \mathcal{L} &= \log P(y | \phi, w) = \sum_{t=1}^N \log P(y_t | \phi(x_t), w) \\ &= \sum_{t=1}^N \log \left\{ \frac{1}{(2\pi)^{\frac{1}{2}} \sigma^2 I} \exp\left(-\frac{(y_t - w^T \phi(x_t))^2}{\sigma^2}\right) \right\} \\ &= -\frac{N}{2} \log \sigma^2 - \frac{N}{2} \log 2\pi - \frac{\sigma^{-2}}{2} \sum_{t=1}^N \{y_t - w^T \phi(x_t)\}^2 \\ &= -\frac{N}{2} \log \sigma^2 - \frac{N}{2} \log 2\pi - \frac{\sigma^{-2}}{2} \|y - \phi w\|^2 \\ &= -\frac{N}{2} \log \sigma^2 - \frac{N}{2} \log 2\pi - \frac{\sigma^{-2}}{2} (y - \phi w)^T (y - \phi w) \end{aligned}$$

MLE result same as least square

$$\frac{\partial \mathcal{L}}{\partial w} = -\frac{\sigma^{-2}}{2} \frac{\partial}{\partial w} (y - \phi w)^T (y - \phi w) = 0$$

$$= \frac{\partial}{\partial w} (y^T y - w^T \phi^T y - y^T \phi w + w^T \phi^T \phi w) = 0$$

$$= -2\phi^T y + 2\phi^T \phi w = 0$$

$$w_{MLE} = (\phi^T \phi)^{-1} \phi^T y$$

### C. Derivation Ridge regression using MAP

$$p(w|y, \phi) = \frac{p(y|\phi, w)p(w)}{\int p(y|\phi, w)p(w)dw} \quad p(w) \sim \mathcal{N}(0, \tau^2)$$

$$\propto p(y|\phi, w)p(w)$$

$$= \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{\|y - \phi w\|^2}{\sigma^2}\right) \cdot \frac{1}{\sqrt{2\pi}\tau^2} \exp\left(-\frac{\|w - 0\|^2}{\tau^2}\right)$$

$$\begin{aligned} \log p(w|y, \phi) &= -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y - \phi w)^T (y - \phi w) \\ &\quad - \frac{D}{2} \log(2\pi) - \frac{D}{2} \log \tau^2 - \frac{1}{2\tau^2} \|w\|^2 \end{aligned}$$

$$\begin{aligned} \frac{\partial \log p(w|y, \phi)}{\partial w} &= -\sigma^{-2} (-2\phi^T y + 2\phi^T \phi w) - 2\tau^{-2} w = 0 \\ &= -2\phi^T y + 2(\phi^T \phi + \frac{\sigma^2}{\tau^2} I) w = 0 \end{aligned}$$

$$2\left(\phi^T \phi + \frac{\sigma^2}{\tau^2} I\right) w = 2\phi^T y$$

$$w_{\text{MAP}} = \left(\phi^T \phi + \frac{\sigma^2}{\tau^2} I\right)^{-1} \phi^T y$$

$$\tau^{-2} \cdot \sigma^2 = \lambda$$

$$= (\phi^T \phi + \lambda I)^{-1} \phi^T y$$

## D. Derivation of LASSO regression

LASSO Regression has a  $l_1$ -norm penalty term so. There is no closed form solution

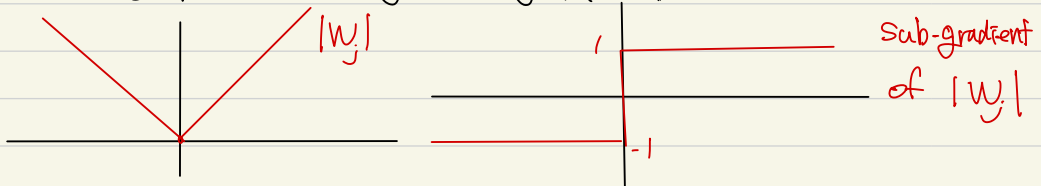
$$\begin{aligned} \mathcal{E}^{\text{LASSO}} &= \|y - \Phi w\|^2 + \lambda \|w\|_1 \quad x \in \mathbb{R}^{n \times p} \quad 12 \ 23 \\ &= \sum_{i=1}^n (y_i - \sum_{j=1}^p w_j \phi_{ij}(x_i))^2 + \lambda \sum_{j=1}^p |w_j| \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{E}^{\text{LASSO}}}{\partial w_j} &= -2 \sum_{i=1}^n \phi_{ij}(x_i) (y_i - \sum_{k=1}^p w_k \phi_{ik}(x_i) - w_j \phi_{ij}(x_i)) + \lambda \frac{\partial}{\partial w_j} \sum_{j=1}^p |w_j| \\ &= -2 \sum_{i=1}^n \phi_{ij}(x_i) (y_i - \sum_{k=1}^p w_k \phi_{ik}(x_i)) + 2 w_j \sum_{i=1}^n \phi_{ij}^2(x_i) + \lambda \frac{\partial}{\partial w_j} \sum_{j=1}^p |w_j| \end{aligned}$$

$\sum_{j=1}^p |w_j|$  is not differentiable, then divide  $\sum_{j=1}^p |w_j|$  into interval to perform differentiation.

$$\lambda \frac{\partial}{\partial w_j} \sum_{j=1}^p |w_j| \begin{cases} \lambda & w_j > 0 \\ [-\lambda, \lambda] & w_j = 0 \\ -\lambda & w_j < 0 \end{cases}$$

$|w|$ 's graph and sub-gradient graph can be shown below



Subgradient is a generalization of the gradient of a differentiable convex function to be applied to a non-differentiable convex function.

$g$  is a subgradient of  $f$  (not necessarily convex) at  $x$  if

$$f(y) \geq f(x) + g^T(y-x)$$

$w_j \neq 0$   $\text{sign}(w_j)$  is subgradient of  $w_j$

$w_j = 0$   $[-1, 1]$  is subgradient of  $w_j$

Therefore differential value when  $w_j=0$  have range between  $-\lambda$  to  $\lambda$

we define  $\ell_j = \frac{1}{n} \sum_{i=1}^n \phi(x_i)(y_i - \sum_{k \neq j} w_k \phi(x_k))$

$$\frac{\partial \mathcal{E}^{\text{LASSO}}}{\partial w_j} = \begin{cases} 2 w_j \frac{1}{n} \sum_{i=1}^n \phi(x_i)^2 - 2 \ell_j - \lambda & \text{when } w_j < 0 \\ [-2 \ell_j - \lambda, -2 \ell_j + \lambda] & \text{when } w_j = 0 \\ 2 w_j \frac{1}{n} \sum_{i=1}^n \phi(x_i)^2 - 2 \ell_j + \lambda & \text{when } w_j > 0 \end{cases}$$

$$2 w_j \frac{1}{n} \sum_{i=1}^n \phi(x_i)^2 - 2 \ell_j - \lambda = 0$$

$$w_j < 0$$

$$\hat{w}_j = \frac{2 \ell_j + \lambda}{2 \frac{1}{n} \sum_{i=1}^n \phi(x_i)^2}$$

$$2 \ell_j + \lambda < 0$$

$$\ell_j < -\frac{\lambda}{2}$$

$$-2 \ell_j - \lambda \leq 0 \leq -2 \ell_j + \lambda$$

$$w_j = 0$$

$$-\frac{\lambda}{2} \leq \ell_j \leq \frac{\lambda}{2}$$

$$2 w_j \frac{1}{n} \sum_{i=1}^n \phi(x_i)^2 - 2 \ell_j + \lambda = 0$$

$$w_j > 0$$

$$\hat{w}_j = \frac{2 \ell_j - \lambda}{2 \frac{1}{n} \sum_{i=1}^n \phi(x_i)^2}$$

$$2 \ell_j - \lambda > 0$$

$$\ell_j > \frac{\lambda}{2}$$

$$\frac{\partial \mathcal{E}^{\text{LASSO}}}{\partial w_j} = \begin{cases} \hat{w}_j = \frac{2 \ell_j - \lambda}{2 \frac{1}{n} \sum_{i=1}^n \phi(x_i)^2} & \ell_j > \frac{\lambda}{2} \\ \hat{w}_j = 0 & \frac{\lambda}{2} \geq \ell_j \geq -\frac{\lambda}{2} \\ \hat{w}_j = \frac{2 \ell_j + \lambda}{2 \frac{1}{n} \sum_{i=1}^n \phi(x_i)^2} & \ell_j < -\frac{\lambda}{2} \end{cases}$$



```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from tqdm import tqdm
from linear_reg import (LinearReg, resid_plot, plot_real_vs_predicted,
                        normality_plot, grid_search_df, plot_r2_rmse)
```

```
In [2]: np.set_printoptions(suppress=True)
data1 = np.loadtxt('./dataset/default_plus_chromatic_features_1059_tracks.txt', delimiter=',')
```

## Data1 Description ¶

- data1 is the data loaded from 'default\_plus\_chromatic\_features\_1059\_tracks'.
- the number of samples : 1059
- the number of features : 116
  - the first 116 columns are audio features of the track
- the number of target : 2
  - the last two columns are the d of the music.
- Training data is allocated as 80 percent of the total data, and the rest is allocated as test data.

```
In [3]: train_sample = round(len(data1) * 0.8)
train_X1, train_y1 = data1[:train_sample, :116], data1[train_sample, 116:]
test_X1, test_y1 = data1[train_sample:, :116], data1[train_sample:, 116:]
```

## 2. Implement three regression models by your own codes to the following dataset.

### 2-A. Vanilla linear regression

First, build a straightforward linear regression of latitude and longitude respectively against features. What is the R-squared of each model? Plot a graph evaluating each regression

```
In [4]: lr = LinearReg()
lr.fit(train_X1, train_y1)
y_tr_pred = lr.predict(train_X1)
y_te_pred = lr.predict(test_X1)
```

vanilla regression is start !

### R-square of linear regression

```
In [5]: lr.report_HW(train_X1, test_X1, train_y1, test_y1, report_type='r2')
```

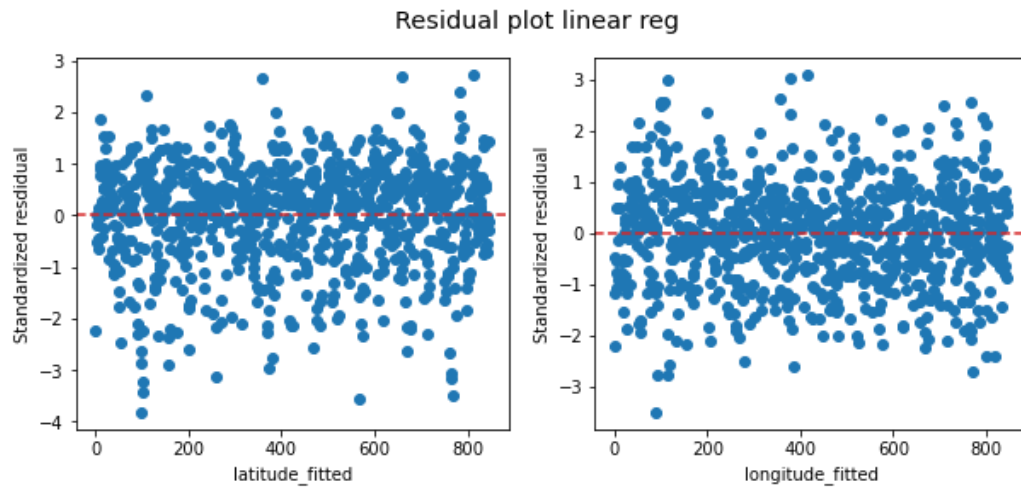
```
Train total R-square : 0.3587
Test total R-square : 0.0907
Train latitude R-square : 0.3106
Test latitude R-square : 0.1131
Train longitude R-square : 0.4068
Test longitude R-square : 0.0683
```

```
In [6]: lr.report_HW(train_X1, test_X1, train_y1, test_y1, report_type='rmse')
```

```
Train total RMSE : 29.479
Test total RMSE : 36.5037
Train latitude RMSE : 15.0804
Test latitude RMSE : 18.4084
Train longitude RMSE : 38.8665
Test longitude RMSE : 48.2304
```

## Residual plot

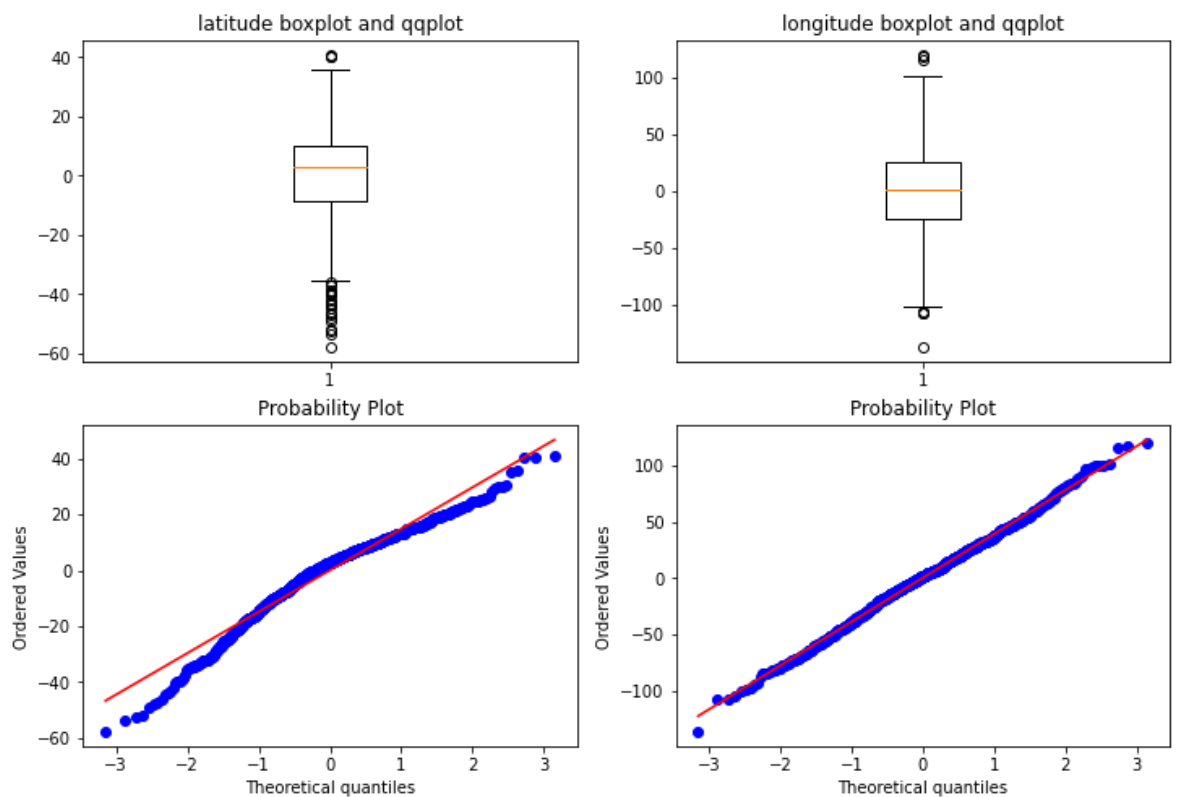
```
In [7]: resid_plot(lr)
```



- When looking at the residual plot, it can be seen that the data are generally linear.

## QQ plot for normality test

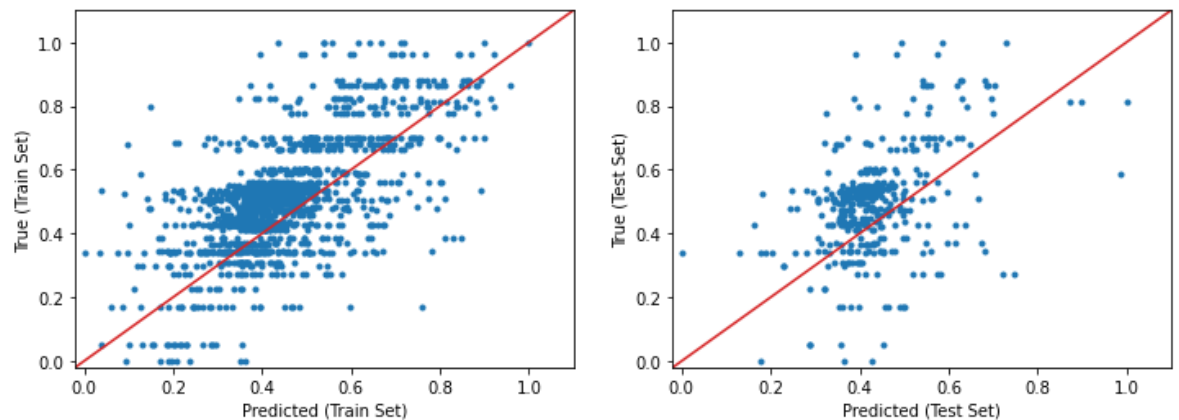
```
In [8]: normality_plot(lr)
```



- Compared to longitude, there are many outliers on the boxplot in latitude, and in Q-Q plot, it can be seen that longitude has normality in residuals compared to latitude.

### predict vs y\_true plot

```
In [9]: fig, axes = plt.subplots(1,2, figsize=(12,4))
plot_real_vs_predicted(train_y1, y_tr_pred, ax=axes[0], mode='Train Set')
plot_real_vs_predicted(test_y1, y_te_pred, ax=axes[1], mode='Test Set')
```



When the predicted value and the actual value match, the blue dots are gathered by a red line. However, in the real graph, you can see blue dots spread around the red line.

## 2-B. Ridge regression

A regression regularized by L2 (equivalently, a ridge regression). You should estimate the regularization coefficient that produces the minimum error. Is the regularized regression better than the unregularized regression?

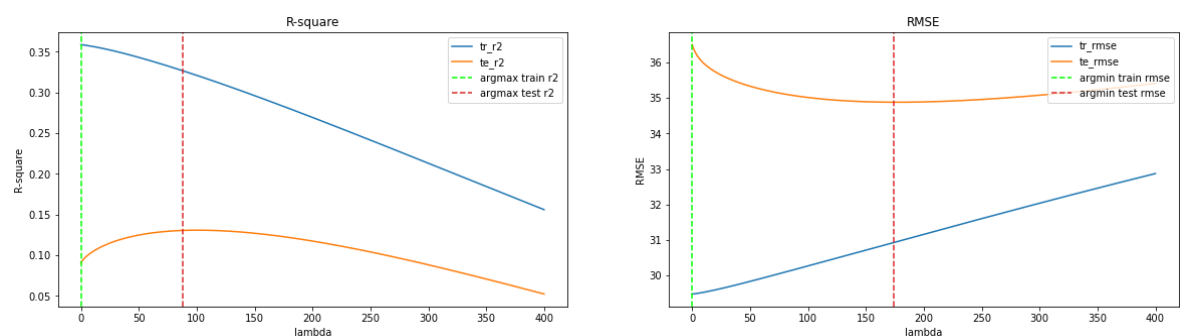
**find best  $\lambda$**

```
In [10]: lambdas_list = [1e-7, 1e-6, 1e-5, 1e-4, 1e-3, 1e-2, 1e-1, 1] + np.linspace(2, 400, 200).tolist()
```

```
In [11]: ridge_score_df = grid_search_df(train_X1, train_y1, test_X1, test_y1, penalty='l2', lambdas_list=lambdas_list)
```

[illegible]

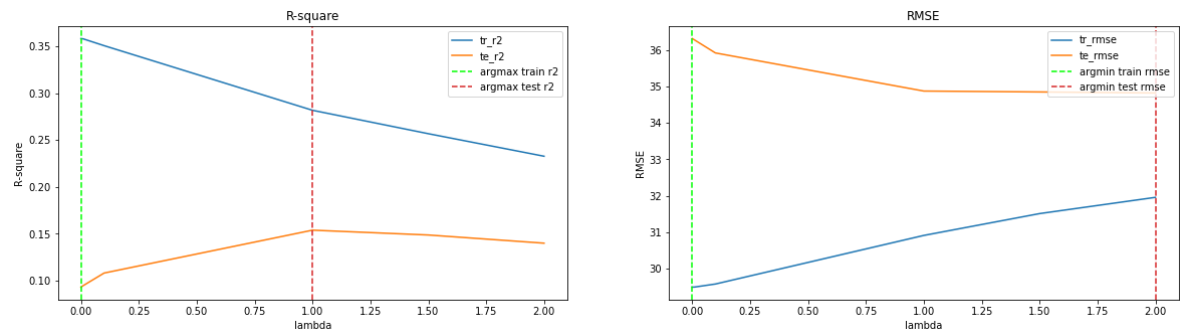
```
In [12]: plot_r2_rmse(ridge_score_df, lambdas_list)
```







```
In [20]: plot_r2_rmse(lasso_score_df, lambdas_list)
```



```
In [21]: best_te_r2_idx = np.argmax(lasso_score_df.values[:,1].round(3), axis=0)
         best_te_rmse_idx = np.argmin(lasso_score_df.values[:,3].round(3), axis=0)
```

```
In [22]: print("best test rmse lambdas : ", lasso_score_df.iloc[best_test_rmse_idx, 2:].name)
         print("best test r2 lambdas : ", lasso_score_df.iloc[best_test_r2_idx, 2:].name)
```

```
best test rmse lambdas : 1.0
best test r2 lambdas : 2.0
```

- LASSO regression also selects the optimal model while adjusting the lambda value.
- If adjust the lambda value, It can be seen that the rsquare of the training set decreases and the rmse increases. However, rsquare and rmse of the test set have been improved.
- If the lambda value exceeds 1.0, the test set also decreases the Rsquare value, so 1.0 was chosen as the optimal lambda value.

## Using best $\lambda$

```
In [23]: lasso = LinearReg(penalty='l1', alpha=1.0, max_iter=1000, learning_rate=1e-4)
```

```
iter : 1000
alpha set 1.0
lasso regression is start !
```

```
In [24]: lasso.fit(train_X1, train_y1)
```

[illegible]

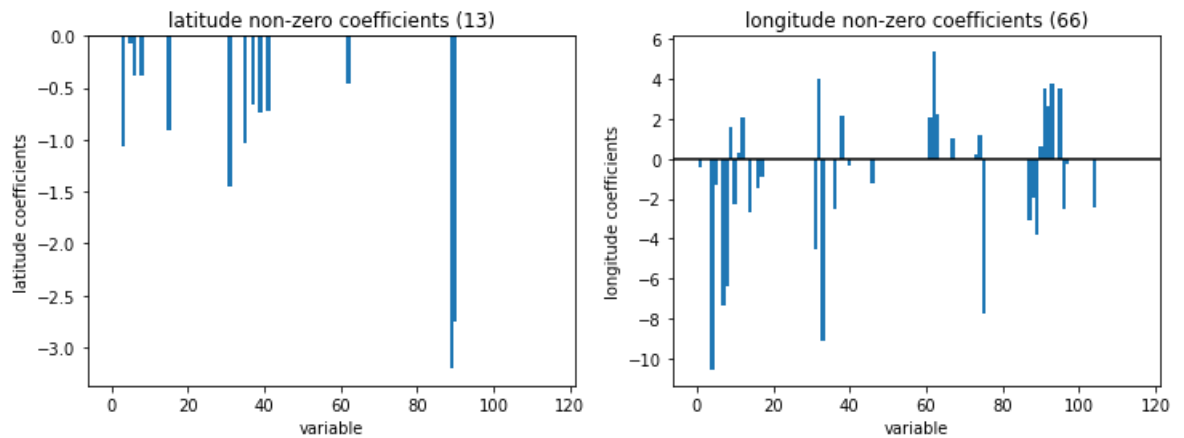
```
In [25]: lasso.report_HW(train_X1, test_X1, train_y1, test_y1)
```

```
Train total R-square : 0.2818
Test  total R-square : 0.1539
Train latitude R-square : 0.2118
Test  latitude R-square : 0.1525
Train longitude R-square : 0.3517
Test  longitude R-square : 0.1552
```

```
In [26]: lasso.report_HW(train_X1, test_X1, train_y1, test_y1, report_type='rmse')
```

```
Train total RMSE : 30.91
Test total RMSE : 34.8781
Train latitude RMSE : 16.1252
Test latitude RMSE : 17.9943
Train longitude RMSE : 40.6304
Test longitude RMSE : 45.9256
```

```
In [27]: fig, axes = plt.subplots(1, 2, figsize=(12,4))
axes[0].bar(range(len(lasso.w[1:,0])), lasso.w[1:,0], width=1)
axes[0].axhline(0, c='black')
axes[0].set_title('latitude non-zero coefficients ({}').format(lasso.w.shape[0] - len(np.where(lasso.w[:,0] == 0)[0]) - 1))
axes[0].set_ylabel('latitude coefficients')
axes[0].set_xlabel('variable')
axes[1].bar(range(len(lasso.w[1:,1])), lasso.w[1:,1], width=1)
axes[1].axhline(0, c='black')
axes[1].set_title('longitude non-zero coefficients ({}').format(lasso.w.shape[0] - len(np.where(lasso.w[:,1] == 0)[0]) - 1))
axes[1].set_ylabel('longitude coefficients')
axes[1].set_xlabel('variable')
plt.show()
```



- When lambda is 1, the number of non-zero coefficients is 13 for longitude and 66 for latitude.

### 3. Summary

In common, in the case of ridge and lasso, it can be seen that the Rsquare of the training set decreases and the RMSE increases. However, it can be seen that the test set is improved. This can be interpreted that Ridge and Lasso alleviate overfitting in the training set.

```
In [40]: result_r2 = pd.DataFrame(raw_data, index=['linear', 'ridge', 'lasso'])
result_r2
```

Out[40]:

|               | Train total R2 | Test total R2 | Train lat R2 | Test lat R2 | Train long R2 | Test long R2 |
|---------------|----------------|---------------|--------------|-------------|---------------|--------------|
| <b>linear</b> | 0.358701       | 0.090679      | 0.310633     | 0.113080    | 0.406768      | 0.068277     |
| <b>ridge</b>  | 0.326652       | 0.130513      | 0.273874     | 0.109973    | 0.379429      | 0.151054     |
| <b>lasso</b>  | 0.281750       | 0.153865      | 0.211797     | 0.152532    | 0.351703      | 0.155198     |

```
In [41]: result_rmse = pd.DataFrame(raw_data2, index=['linear', 'ridge', 'lasso'])
result_rmse
```

Out[41]:

|               | Train total RMSE | Test total RMSE | Train lat RMSE | Test lat RMSE | Train long RMSE | Test long RMSE |
|---------------|------------------|-----------------|----------------|---------------|-----------------|----------------|
| <b>linear</b> | 29.479024        | 36.503714       | 15.080390      | 18.408405     | 38.866535       | 48.230415      |
| <b>ridge</b>  | 30.164292        | 35.068263       | 15.477229      | 18.440624     | 39.752037       | 46.038131      |
| <b>lasso</b>  | 30.909958        | 34.878070       | 16.125237      | 17.994322     | 40.630379       | 45.925635      |