

Problem 1

a) We can split the data using K-fold and analyse the data by getting MSE of training folds and MSE of test-folds.

The mean of training folds and MSE mean of test-folds can be used to get learning curves.

We use `validation_curve()` function of sklearn to do the same.

`validation_curve (Rmodel,
training data,
Testing data,
param name,
param range,`

`cross validation
folds,`

`scoring method,
jobs-value)`

This method returns Train-scores and Test-scores which can be plotted using mean of test scores, mean of train scores, along with param-range of penalty parameter ' λ '.

The learning curves show the mean of MSE of training folds and Mean of MSE of ~~test~~ data increase as the penalty parameter increases.

Optimal value of penalty parameter can be found using Randomized SearchCV method which gives best estimator. alpha for the model by iterating over given range of values.

We get $\lambda = 0.00028331735$,

$MSE = 5.6800923134$.

which is best value for range of λ values.

b.

We standardize the input co-variables using `StandardScaler()` method in ~~the~~ scikit learn library using

`StandardScaler().fit(j).transform(j)`

where 'j' is a reshaped array of column of co-variate values.

We once draw resulting learning curves for the Standardized data.

The training data Mean Squared Error curve and testing data MSE curve is set to be closer than in non-standardized learning curve.

The curve after standardization is smoother.

Training score and test scores are set to lesser in standardized version than in the previous learning curve across the range from start to end.

C. Standardization of the data was done in order to decrease the Mean squared error. Hence we consider the

MSE error at beginning and end for lesser penalty parameter to higher penalty parameter for un-standardized Initial data based model and model based on Standardized data.

Non-Standardized based Model

<u>Training data</u> <u>Curve</u>	(At beginning)	<u>Test data</u> <u>Curve</u>
$n = 50$		$n = 50$
$y = -15.98$		$y = -27.24$

<u>Training data</u> <u>Curve</u>	(At End)	<u>Test data</u> <u>Curve</u>
$n = 200$		$n = 200$
$y = -25.66$		$y = -33.49$

Standardized data based Model

Training data Curve (At beginning)

$$n = 50$$

$$y = -12.00$$

Test data Curve

$$n = 50$$

$$y = -20.67$$

Training data Curve (At End)

$$n = 200$$

$$y = -23.29$$

Test data Curve

$$n = 200$$

$$y = -29.77$$

We consider ratios of both Model $MSE(y)$ values At beginning (approximate), $n = 50$.

$$\frac{\text{Ratio}_{\text{train data}}}{MSE} = \left(\frac{15.96}{12.0} \right) \quad \frac{\text{Ratio}_{\text{test data}}}{MSE} = \left(\frac{27.24}{20.67} \right)$$
$$= 1.325 \quad = 1.31785$$

At End, $n = 200$

$$\frac{\text{Ratio}_{\text{train data}}}{MSE} = \left(\frac{25.66}{23.29} \right) \quad \frac{\text{Ratio}_{\text{test data}}}{MSE} = \frac{33.99}{29.77}$$
$$= 1.1017 \quad = 1.1417$$

We can notice that Ratio at beginning point is higher than Ratio at a particular End Point. Hence, standardization has a better or bigger difference in decrease in MSE at beginning of data that is for lower, smaller penalty parameter values.

d. No. We do not enough data.
Because both of the learning curves
never converge across range of penalty
parameters. Hence we are not able to
get optimal value of penalty parameter
' λ ' from the graph.

e. The model performs best for $\lambda =$
 0.00028331735 penalty parameter
for ridge regression.

$$MSE_{\text{test}} = 15.7898752245$$

The model was trained on training data
with MSE on training data $MSE_{\text{tr}} = 5.6800923$
May be the model has got overfitting
with training data hence performs lesser
in case of test data.

$$MSE_{\text{test}} = 15.789875 \text{ is not bad}$$

considering we have small input data
which might lead to overfitting of model.

Problem - 2

a.

In the problem 2)a), we use
`np.random.gamma(1, 1, 1000)` to generate

1000 samples from $\text{Gamma}(1, 1)$ function

Using the list of samples of function Gamma
or x_i distribution of $\text{Gamma}(1, 1)$ function,

we get distribution Z_i by simply

substitute x_i values in the function

or equation $Z = \frac{x^2}{1 + \sin(x)}$.

After getting distribution Z_i of 1000

samples of equation Z , we plot

the histogram of 300 sample to get

the pdf of Z plot.

b. $E[\omega]$ or Expectation of $E[\omega]$ can be found in the following way.

Expectation of a list of N samples is same as the mean of N samples.

In our case, we have a continuous function. 'is'. We have to derive

a ' N ' samples distribution of the function and take the mean of the ' N ' samples so derived.

Here, we choose value of ' z ' to be considered in the function to be 0 to infinity. as for negative value ' z ', $\log z$ has considered which does not exist as \log of negative number is not defined.

Here, we have consider a random number from 0 to 1.79×10^8 as the biggest number allowed in python is 1.8×10^8 .

We can consider a sample of w to be a sample 100000000 numbers

We get mean of 100000000 numbers of samples of distribution of w .

$$\text{Mean } w = 708.778695$$

We get the mean of distribution which is equal to the Expectation of w .

Problem 3

To build the model, we have to replace NA values in training data,

We replace the 'NA' values in training data using simple imputer from sklearn library.

Our strategy is to replace missing value with mean of the rest of the values.

The model is not the best model
Hence again, we use the Randomized
Search CV() function to get best
penalty parameter λ value for the
model. We get $\lambda = 9.29330653033$

~~then~~ we use this best model on
test data after removing 'NA' values
in test data using Simple Impute function of
sklearn. $MSE_{\text{training data}} = 0.16324625$