

# PES University, Bangalore

Established under Karnataka Act No. 16 of 2013

UE21CS342AA2 - Data Analytics - Worksheet 2b - Ridge, Lasso and Elastic Net Regression  
Designed by Aaditya S Goel, Dept. of CSE - [aadityasgoel@gmail.com](mailto:aadityasgoel@gmail.com)

## Welcome to INADA Ports

Welcome back to your second assignment at the Bangalore Consulting group.

Maritime transport handles 95% of India's trade in volume terms and 70% in value terms and is growing twice as fast as the global maritime trade. The government, therefore, is really pushing hard to give this industry a boost by providing fiscal and non-fiscal initiatives to enterprises to develop and maintain India's coastline. Over 150 initiatives are currently being implemented.

An important part of any coastline development, is the development of breakwaters. A breakwater is a structure built in bodies of water to provide protection and control against the impacts of waves, currents, and tides. It is designed to reduce the energy of incoming waves and create calmer waters behind it.

This allows all kinds of economic activity to occur along the coasts. It enables constructions of ports, recreational areas, protect beach ecosystems etc. A prime example of this would be the tetrapods you see along the Marine Drive in Mumbai.



Figure 1: Image sourced from [www.townmumbai.com](http://www.townmumbai.com)

INADA Ports is the largest operator of ports in India. They are constructing a port along the Western Coast and need your help before they can start building it.

Wave height in a particular region is an important parameter to be considered before construction of a port. There is some historical data available about the climatic conditions prevailing in the region and the height of the waves. Your task is to build a model that can accurately predict the wave height, given these same climatic conditions.

## Overfitting and Regularization

In the realm of predictive modeling, the pursuit of creating a model that perfectly fits the training data can inadvertently lead to a phenomenon known as overfitting. Overfitting occurs when a model becomes

excessively complex, capturing not only the genuine patterns within the data but also the noise and random fluctuations present in the training set.

The hyper-adaptation to the training data renders the model less capable of generalizing to new, unseen data, as it effectively memorizes the training examples rather than discerning meaningful relationships. As a result, an overfitted model may exhibit impressive performance on the training data but performs poorly when faced with real-world scenarios. The delicate balance between capturing essential patterns and avoiding the trap of overfitting underscores the importance of techniques like regularization, which aim to ensure model generalization by restraining excessive complexity.

There are three Regularization techniques we will be dealing with, all of which use the idea of penalizing terms to tackle overfitting.

But before we go any further, let's have a look at the data.

## Data Dictionary

hh: Hour  
WDIR: Wind Direction measured in degrees  
WSPD: Wind Speed measures in metres/second  
GST: Gust speed measures in metres/second  
DPD: Dominant Wave period measures in seconds  
APD: Average Wave period measures in seconds  
MWD: Mean Wind Direction measured in degrees  
PRES: Pressure measured in hectopascal  
ATMP: Atmosphere temperature measured in degrees Celsius  
WTMP: Water temperature measured in degrees Celsius  
DEWP: Dew Point measured in degrees Celsius  
WVHT: Significant Wave Height measured in metres. This the variable we will be predicting

---

## Visualizing the data

Let's visualize this all in the form of a Data Frame

```
data <- read.csv('waves.csv')
# head(data)
summary(data)
```

##	hh	WDIR	WSPD	GST
##	Min. : 0.00	Min. : 1.0	Min. : 0.000	Min. : 0.200
##	1st Qu.: 6.00	1st Qu.: 53.0	1st Qu.: 4.000	1st Qu.: 4.800
##	Median :12.00	Median :166.5	Median : 5.600	Median : 6.600
##	Mean :11.55	Mean :163.1	Mean : 5.738	Mean : 6.852
##	3rd Qu.:17.00	3rd Qu.:257.0	3rd Qu.: 7.300	3rd Qu.: 8.675
##	Max. :23.00	Max. :360.0	Max. :13.500	Max. :16.800
##	DPD	APD	MWD	PRES
##	Min. : 2.150	Min. :2.290	Min. : 0	Min. : 992.1
##	1st Qu.: 3.450	1st Qu.:3.080	1st Qu.: 33	1st Qu.:1010.7
##	Median : 4.000	Median :3.440	Median :138	Median :1014.2
##	Mean : 4.205	Mean :3.529	Mean :137	Mean :1014.2
##	3rd Qu.: 4.760	3rd Qu.:3.890	3rd Qu.:192	3rd Qu.:1018.0
##	Max. :17.390	Max. :5.770	Max. :360	Max. :1030.6
##	ATMP	WTMP	DEWP	WVHT
##	Min. : 3.40	Min. : 3.50	Min. : -4.000	Min. :0.2500
##	1st Qu.: 9.00	1st Qu.: 6.80	1st Qu.: 7.525	1st Qu.:0.3600
##	Median :17.10	Median :16.90	Median :13.500	Median :0.5200

```
## Mean      :15.29   Mean      :14.85   Mean      :12.578   Mean      :0.6414
## 3rd Qu.   :20.90   3rd Qu.   :21.80   3rd Qu.   :17.800   3rd Qu.   :0.8100
## Max.      :27.50   Max.      :25.30   Max.      :23.200   Max.      :2.6300
```

(Ask yourselves, is there merit to thinking of this problem as a time series problem given the distinctly seasonal behaviour of water bodies?)

In this worksheet, we will train the model on 80% and test it on the remaining 20% Therefore, before we build any models, split the data into the Training and testing split

```
# Split the data into training and testing sets using the indices (first 80% for training and next 20% .
```

---

## Correlation plot

You may recall that in the previous worksheet, in the last section with Multiple Linear Regression, we created a correlogram and conducted a fairly informal visual analysis to select the most important factors. Regularization in some ways is a formalization of this same process, where instead of conducting a visual analysis, we have a mathematical basis for selecting (or not selecting or penalizing) certain parameters.

So to form a baseline, let's create a correlogram and repeat our visual analysis process and create a first version of our predictor, an MLR model

```
# Create a correlogram similar to the one created in the previous worksheet
```

Now let's pick the visually best looking parameters for our model

```
# Create the MLR model with the most suitable parameters (There seem to be 4 right? Or is it 6?)
```

Print the R-squared statistic of this model based on its fit over the testing data. (Hint: Recall the value of R-squared statistic as  $1 - \text{RSS}/\text{TSS}$ )

```
# Write your code here
```

---

## Ridge Regression

Ridge regression is a linear regression technique that incorporates L2 regularization to address issues in predictive modeling (overfitting, multi-collinearity etc).

Linear regression, aims to minimize the sum of squared residuals whereas Ridge regression introduces a penalty term proportional to the square of the magnitude of the coefficients. This penalty, controlled by a hyperparameter (often denoted as  $\lambda$ ), discourages large coefficient values, effectively constraining the model's complexity, enhancing its stability and generalization performance.

Perform Ridge Regression on the training data and compare the predictions with the test data to check for the fit of the model. (Hint: Use the glmnet library)

```
# Write your code here
```

Print the R-squared statistic. Does the model seem to be fitting well?

```
# Calculate R^2 here
```

What was your  $\lambda$ ? Is it possible for you to somehow conduct hyperparameter tuning and find the best  $\lambda$  value for the Ridge Regression model? (Hint: use the cv.glmnet function)

```
# Write your code here
```

With the optimal lambda, build the model again and print the coefficients of the various dependent variables. Compare this to coefficients with models that had a higher or lower value of lambda. What can you comment about the relationship between lambda and the strength of regularization?

```
# Write your code here
```

What can you comment about the R-squared statistic value of the best model? Is it higher or lower than in the case of MLR?

Is a higher R-squared statistic always the better model?

What can be done to specifically improve R-squared value for this Ridge Regression model? (Think about the transformations you can do to the data)

---

## Lasso Regression

Lasso regression is similar to Ridge Regression except that instead of L2 regularization, it employs L1 regularization to address the very same issues that Ridge Regression addresses.

There are however, a couple of differences between the two. The first and most obvious being that since Lasso Regression implements L1 regularization, the penalty term in this case is proportional to the absolute value of the coefficient.

Another point to note is that unlike its Ridge counterpart, Lasso Regression can push some coefficients to exactly 0. This effectively drops the feature from the predictive model (Similar to how we drop values through visual analysis). Lasso Regression can thus be used effectively for Feature Selection as well.

The goal here too however is to improve model stability and generalization.

Write code to build a Lasso Regression model similar to how you built the Ridge Regression model. This time incorporate hyperparameter tuning right away. So first print the optimal lambda value.

```
# Write your code here
```

Display the coefficients of all the variables. Do you notice some variables being dropped out? Which ones are they?

```
# Write your code here
```

Let us conclude this one by finding the R-squared statistic for the Lasso Model.

```
#Write your code here
```

---

## Elastic Net Regression

Elastic Net regression, an advanced form of linear regression, combines the benefits of L1 (Lasso) and L2 (Ridge) regularization methods. By integrating both penalty terms, Elastic Net overcomes the limitations of each, offering resilience against multicollinearity, aiding feature selection, and preventing overfitting.

This approach makes Elastic Net a very versatile approach for achieving accurate and efficient models by finding a middle ground between dropping parameters and retaining important predictors.

Build your Elastic Net Regression model incorporating all the steps we previously followed for ridge and lasso regression. (Play around with the alpha value and find out how it affects the model)

```
# Write your code here
```

What can you comment about the coefficients of this model? Is it a simple average of Ridge and Lasso Regression? Or does it vary too? What does that tell you about the number of hyperparameters in Elastic Net Regression compared to the other two models?

---

Amazing work with the analysis. Along with Linear Regression models you now have added Regularization techniques to your repertoire. With each assignment you're making significant progress on your Data Analytics Journey. So, until the next case comes up, Happy Learning!