

School of Engineering and Applied Science (SEAS), Ahmedabad University

B.Tech (CSE Semester VI)
Machine Learning (CSE 523)

Project Submission #2: Linear Regression

Submission Deadline: March 04, 2020 (11:59 PM)

- Group No.: S_ECC7
- Project Area: Environment and Climate Change
- Project Title: Estimating the effect of climate change on sea level using Machine Learning
- Name of the group members :
 1. Rajvi Patel (AU1741078)

Linear Regression

- Implementation code:

```
1 clc;
2 clear all;
3 close all;
4 tic
5
6 %% Load and process the data
7 frequency = 4;          % 1 for year, 4 for seasons, 12 for months
8 hold_fraction = 0.20;% fraction to Divide data into training and testing part
9 offset =7 ;
10 span = 1;
11 t_xfirst = 1914;
12 t_xlast = 2013;
13 t_xl=2005;
14
15 % define features
16 features = {'globalCO2','globalpop','localprecip','localtemp'};
17
18 % point from where test data is going to be splitted
19 t_xsplit = t_xlast-round((t_xlast-t_xfirst)*hold_fraction);
20 shift = offset+span;
21
22 % Load the data
23 if(frequency==1)
24     appendix = '_y';
25 elseif(frequency==4)
26     appendix = '_s';
27 elseif(frequency==12)
28     appendix = '_m';
29 else error('Invalid data frequency. Select 1, 4 or 12');
30 end
```

```

31 load(['locallevel',appendix]);
32
33 %load all features according to appendix
34 for i = 1:size(features,2)
35     load([char(features(i)),appendix]);
36 end
37
38 % normalize the inputs
39 for i = 1:size(features,2)
40     eval([char(features(i)) appendix '=' char(features(i)) appendix '- mean('
41         char(features(i)) appendix ');' ]]);
42     eval([char(features(i)) appendix '=' char(features(i)) appendix '/ sqrt('
43         var(' char(features(i)) appendix ');' ]]);
44 end
45
46 % Prepare the training and testing set according to type of data(monthly,
47     yearly,seasonly)
48
49 %crop is a function to split
50 y_train = crop(eval(['time_locallevel',appendix]),eval(['locallevel',appendix
51     ]),t_xfirst,t_xsplit);
52 y_test = crop(eval(['time_locallevel',appendix]),eval(['locallevel',appendix])
53     ,t_xsplit,t_xlast);
54 x_train = zeros((t_xsplit-t_xfirst)*frequency,span*frequency*size(features,2))
55     ;
56
57 for i = 0:size(features,2)-1
58     feature = [char(features(i+1)),appendix];
59     for j = 0:span-1
60         for k = 0:frequency-1
61             t = k/frequency;
62             x_train(:,1+i*span*frequency+j*frequency+k) = crop(eval(['time_',
63                 feature]),eval(feature),t_xfirst+j+t,t_xsplit+j+t);
64         end
65     end
66 end
67
68 x_test = zeros((t_xlast-t_xsplit)*frequency,span*frequency*size(features,2));
69 for i = 0:size(features,2)-1
70     feature = [char(features(i+1)),appendix];
71     for j = 0:span-1
72         for k = 0:frequency-1
73             t = k/frequency;
74             x_test(:,1+i*span*frequency+j*frequency+k) = crop(eval(['time_',
75                 feature]),eval(feature),t_xsplit+j+t,t_xlast+j+t);
76         end
77     end
78 end
79
80 x_future = zeros(offset*frequency,span*frequency*size(features,2));
81 for i = 0:size(features,2)-1
82     feature = [char(features(i+1)),appendix];
83     for j = 0:span-1
84         for k = 0:frequency-1
85             t = k/frequency;
86             p=t_xlast+j+t;
87             q=t_xlast+offset+j+t;
88             x_future(:,1+i*span*frequency+j*frequency+k) = crop(eval(['time_',
89                 feature]),eval(feature),t_xl+j+t,t_xl+offset+j+t);
90         end
91     end
92 end

```

```

82     end
83 end
84
85 %create time vector according to training and testing data
86
87 t_traintest = crop(eval(['time_locallevel',appendix]),eval(['time_locallevel',
    appendix]),t_xfirst,t_xlast);
88 t_train=crop(eval(['time_locallevel',appendix]),eval(['time_locallevel',
    appendix]),t_xfirst,t_xsplit);
89 t_test=crop(eval(['time_locallevel',appendix]),eval(['time_locallevel',
    appendix]),t_xsplit,t_xlast);
90 t_future = zeros(offset*frequency,1);
91
92 for i = 0:offset-1
93     for j = 0:frequency-1
94
95         t_future(1+i*frequency+j) = t_xlast+i+j/frequency;
96     end
97 end
98 t = [t_traintest;t_future];
99
100 % add column of 1's to the feature matrix
101 x_train = [ones(size(y_train)) x_train];
102 x_test = [ones(size(y_test)) x_test];
103 x_future = [ones(size(x_future,1),1) x_future];
104
105 % compute theta using normal equation
106 theta = inv(x_train'*x_train)*x_train'*y_train;
107
108 % get results for MLE
109 y_test_MLE = (theta'*x_test)';
110 y_train_MLE = (theta'*x_train)';
111 y_future_MLE = (theta'*x_future)';
112 y_MLE = [y_train_MLE; y_test_MLE; y_future_MLE];
113
114 %Calculate RMSE error for MLE
115 train_Error_MLE = sum((y_train-y_train_MLE).^2)/sum(y_train);
116 test_Error_MLE = sum((y_test_MLE-y_test).^2)/sum(y_test);
117
118 mse_test_ml=((y_test_MLE-y_test).^2)./size(y_test_MLE,1);
119 mse_train_ml=((y_train-y_train_MLE).^2)./size(y_train,1);
120 rmse_train_ml=sqrt(mse_train_ml);
121 rmse_test_ml=sqrt(mse_test_ml);
122 rmse_ml=[rmse_train_ml;rmse_test_ml];
123
124 %calculate noise variance
125 var=0;
126 for i=1:size(mse_train_ml)
127     var=var+mse_train_ml(i);
128 end
129
130 %initialize the variance of the Gaussian prior
131 b_sq=150;
132
133 %calculate theta after applying MAP
134 theta_MAP=inv((x_train'*x_train)+(var/b_sq*eye(size(x_train,2))))*x_train'*
    y_train;
135
136 % get results for MAP
137 y_test_MAP = (theta_MAP'*x_test)';

```

```

138 y_train_MAP = (theta_MAP'*x_train')';
139 y_future_MAP = (theta_MAP'*x_future')';
140 y_MAP = [y_train_MAP; y_test_MAP; y_future_MAP];
141
142 %calculate RMSE error for MAP
143 train_Error_MAP = sum((y_train-y_train_MAP).^2)/sum(y_train);
144 test_Error_MAP =sum((y_test_MAP-y_test).^2)/sum(y_test);
145
146 mse_train_map=((y_train-y_train_MAP).^2)./size(y_train,1);
147 mse_test_map=((y_test_MAP-y_test).^2)./size(y_test_MLE,1);
148 rmse_train_map=sqrt(mse_train_map);
149 rmse_test_map=sqrt(mse_test_map);
150 rmse_map=[rmse_train_map;rmse_test_map];
151
152 disp(['LR RMSE train error: ',num2str(train_Error_MLE)]);
153 disp(['LR RMSE test error: ',num2str(test_Error_MLE)]);
154 disp(['LR RMSE train error after MAP: ',num2str(train_Error_MAP)]);
155 disp(['LR RMSE test error after MAP: ',num2str(test_Error_MAP)]);
156
157 %plot the graphs
158 figure;
159 title('Linear regression (Global)');
160 hold on;
161 plot(t,y_MLE);
162 plot(t,y_MAP);
163 xlabel('Time (years)');
164 ylabel('Sea level (mm)');
165 h1= legend('MLE Prediction','MAP Prediction','Location','northwest');
166
167 figure;
168 title('Linear regression using MLE');
169 hold on;
170 plot(t_traintest,[y_train; y_test]);
171 plot(t_train,y_train_MLE);
172 plot(t_test,y_test_MLE);
173 xlabel('Time (years)');
174 ylabel('Sea level (mm)');
175 h2 = legend('Original','MLE Prediction on train data','MLE Prediction on test
    data','Location','northwest');
176
177 figure;
178 title('Linear regression using MAP');
179 hold on;
180 plot(t_traintest,[y_train; y_test]);
181 plot(t_train,y_train_MAP);
182 plot(t_test,y_test_MAP);
183 xlabel('Time (years)');
184 ylabel('Sea level (mm)');
185 h3 = legend('Original','MAP Prediction on train data','MAP Prediction on test
    data','Location','northwest');
186
187
188 figure
189 title('RMSE Error for MLE');
190 hold on;
191 plot(t_traintest,smooth(rmse_ml));
192 xlabel('Time (years)');
193 ylabel('Root Mean square error');
194
195 figure

```

```

196 title('RMSE Error for MAP');
197 hold on;
198 plot(t_traintest,smooth(rmse_map));
199 xlabel('Time (years)');
200 ylabel('Root Mean square error after MAP');
201
202 figure
203 title('Comparison of RMSE Error for MLE & MAP');
204 hold on;
205 plot(t_traintest,smooth(rmse_ml));
206 plot(t_traintest,smooth(rmse_map));
207 xlabel('Time (years)');
208 ylabel('Root Mean square error');
209 h4 = legend('Root Mean square error for MLE','Root Mean square error for MAP',
    'Location','northwest');
210
211 figure
212 title('Plotting MLE & MAP data');
213 hold on
214 plot([y_train;y_test],[y_train_MLE;y_test_MLE],'ob');
215 plot([y_train;y_test],[y_train_MAP;y_test_MAP],'or');
216 plot([y_train;y_test],[y_train;y_test],'-k','LineWidth',2);
217 l = legend('MLE','MAP','Location','northwest');
218 xlabel('Original data');
219 ylabel('Prediction Data');
220 hold off
221
222 toc
223

```

- **URL links:** Upload your dataset and codes on your drive and paste the short-url links here.

Project: https://drive.google.com/open?id=1DA3vJE_Bru0qtFdTcER8PUL9kwZryIxs

Data: <https://drive.google.com/open?id=1eijmyPbG9kJGlP-YTH4eEkpFTE04-4Y5>

Code: <https://drive.google.com/open?id=1n7n566r8ozbJGRohYKATQh8aABs3ayhy>

- **Inference:**

Aim of my project is to predict Sea level rise for local area(San Fransico) and Global area by applying linear regression algorithm on collected data. As almost half of the world's population lives in coastal regions, the climate change due to increase in sea level is most concerning problem. Having knowledge about increase in sea level can help governments and other interested entities in disaster prevention, real estate evaluation and public safety. For example beach erosion, inundation of low lying areas, salt water intrusion into aquifers and increased flooding.

Initially for the local prediction, I have considered the following features which are affecting the most to rise of sea level: Local precipitation, Local temperature, Local population, Local CO_2 concentration and for the global prediction, I have applied a different set of features:Global precipitation, Global temperature index, Global population, Global CO_2 concentration. If we increased features than prediction will get even worse. So we have to choose features which affects the most to the prediction. The performance of linear regression algorithms under various conditions is derived using cross validation on the latest 20% of the dataset, which was left out from the dataset. This dataset is formed into feature matrix which is Φ and added one column of 1's to meet the dimensions. Using following equation θ_{ML} is

calculated for a hypothesis function to predict the sea level rise,

$$\theta_{ML} = (\Phi^T \Phi)^{-1} \Phi^T y$$

After deriving θ prediction can be done using hypothesis function,

$$y = X^T \theta$$

Observing the value of y for training and testing data,

$$R_{emp}(test) > R_{emp}(train)$$

which leads to overfitting. To minimize the risk regularization parameter is added, which reduces the variance of estimated regression parameters. however, it does this at the expense of adding bias to the estimate. In regularized least squares, we consider the loss function as,

$$\|y - \Phi\theta\|^2 + \lambda \|\theta\|^2$$

On minimizing the loss, we obtain a solution that closely resembles the MAP estimate

$$\theta_{RLS} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$$

where, $\lambda = \frac{\sigma^2}{b^2}$, σ^2 is noise variance and b^2 is the variance of the (isotropic) Gaussian prior $p(\theta) = N(0, b^2 I)$.

To calculate noise variance,

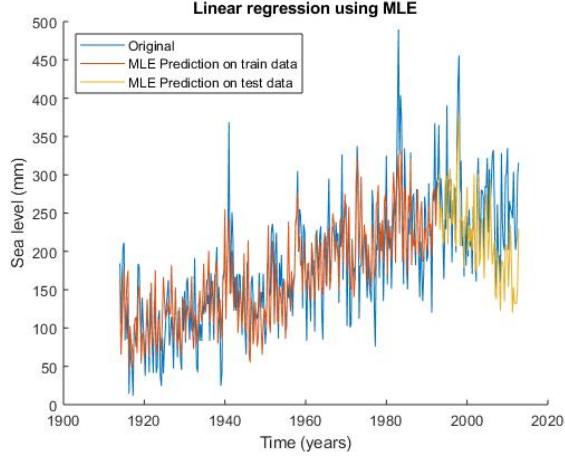
$$\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (y_n - \phi^T(x_n)\theta)^2$$

$$\theta_{MAP} = (\Phi^T \Phi + \frac{\sigma^2}{b^2})^{-1} \Phi^T y$$

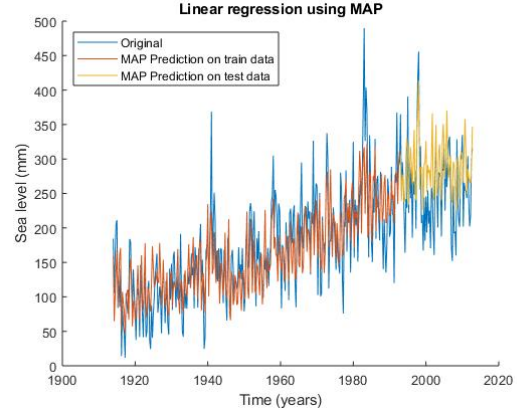
After deriving θ_{MAP} prediction can be done using following hypothesis function,

$$y = X^T \theta_{MAP}$$

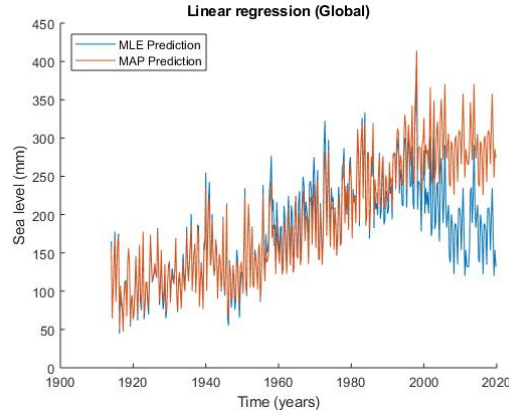
θ 's are actually the weights assigned to a particular feature. Not all parameters affect the hypothesis function equally. Different features contribute to rise differently. The whole purpose of writing the ML algorithm is to come-up with values of θ that will make our hypothesis function fit the training set well.



(a) MLE Prediction



(b) MAP Prediction



(c) Overlapping of MLE and MAP

Figure 1: Linear Regression using MLE and MAP

RMSE Error(mm)		
Model	Training error	Testing Error
MLE	12.8231	19.8554
MAP	13.2227	12.0539

Here in figure 1(a) sea level rise is predicted using maximum likelihood estimation and in figure 1(b) sea level rise is predicted using MAP estimation. Both graphs are plotted along with original data. In figure 1(c) MLE and MAP both are overlapped. All of the above graph and from the first row of table we can say that MLE model models the training data too well but fails to perform well on the testing data which leads to overfitting.

$$R_{emp}(test) > R_{emp}(train)$$

Performing sufficiently good on testing data is considered as a kind of ultimatum in machine learning. Thus, to avoid overfitting λ is used which is called regularization parameter. λ trades off between minimizing the loss on training set and magnitude of the parameter. If λ is higher then it adds high penalty to error term. Regularization parameter λ reduces almost 40% testing error and forced to fit model on unseen data closely.

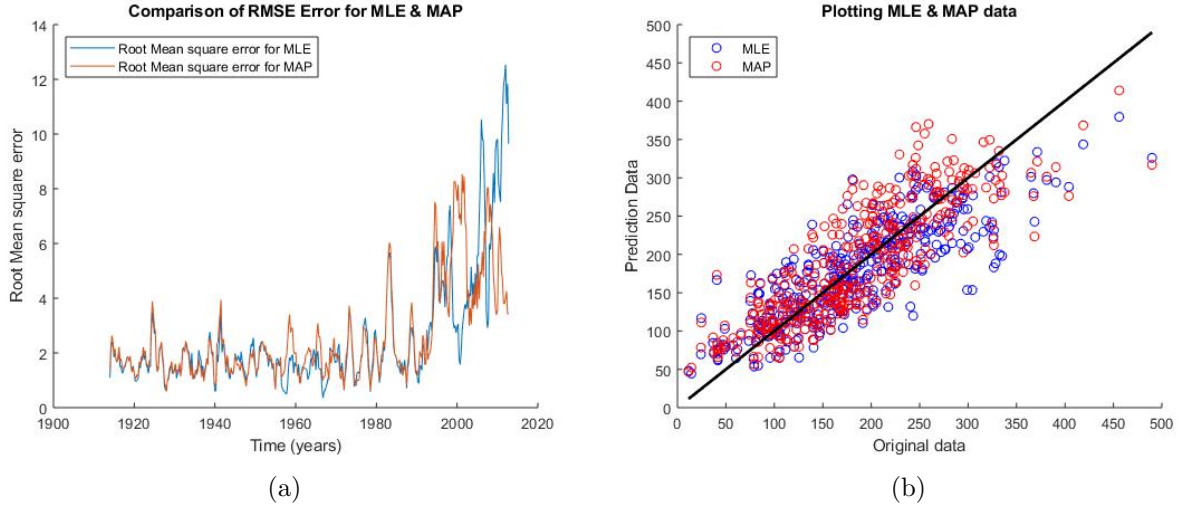


Figure 2: Error graphs in MLE and MAP

Figure 2(a) represents RMSE (root mean square error) of MLE and MAP estimation. For training data MLE model predicts output perfectly but for testing data RMSE error increases exponentially and goes to the peak level. After doing regularization this peak level is forced to go down. In figure 2(b) data of MAP clustered closely to the line than data of MLE which shows better prediction using MAP rather than MLE. Using above observation we can derive following inferences:

	Low training error	High training Error
Low testing error	The model is training	Error in code
High testing error	Overfitting	The model is not training