

## IS580 – Knowledge, Discovery and Mining

### Assignment 4

Due Date: 29 May 2018

In this assignment, you are expected to fit a generalized linear model via penalized maximum likelihood (i.e. elastic net), fit a time series model, and make a model selection. For the elastic net task, you are provided with training and test datasets, namely “train.csv” and “test.csv”. For the time series analysis, you are provided with the data set named “timeSeries.txt”.

You can fit the elastic net model with the *glmnet* package. You may refer to the references below for fitting the model:

- <https://cran.r-project.org/web/packages/glmnet/glmnet.pdf>
- [http://web.stanford.edu/~hastie/glmnet/glmnet\\_alpha.html](http://web.stanford.edu/~hastie/glmnet/glmnet_alpha.html)

For the time series analysis, you will use the *tseries* and *forecast* packages.

For the last part (performance analysis), you may refer to the references below for the paired t-test and Wilcoxon signed rank test:

- <https://www.r-bloggers.com/paired-students-t-test/>
- <https://www.r-bloggers.com/wilcoxon-signed-rank-test/>

### Deliverables:

1. RMarkdown file named as “yourID\_Assignment4.Rmd” including R chunks, stats, results of the analysis, graphs and comments.
2. Data file named as “yourID\_Assignment4.Rdata” including elastic net model and predictions of the model in Question A and time series, additive model for the time series and forecasts in the Question B.

### Questions

- A. Consider the data sets “train.csv” and “test.csv”
  1. Generate a multinomial elastic net model with cross-validation option on the “train.csv” data set (**Hint:** Use *cvglmnet* function in *glmnet* package. *cvglmnet* only accepts *matrix* data type, so be sure to convert your data set to matrix. Also, pay attention to the *family* parameter of the function.).
  2. Make predictions for the “test.csv” data set by using the model.
  3. Compare the predictions of the model with the actual classes. Comment on the performance of the model.
  4. Save your model and predictions in the data file “yourID\_Assignment4.RData”
- B. Consider the data set “timeSeries.txt”
  1. Is the time series stationary? Why? Support your reasons with ACF, PACF plots and/or statistical tests. If it is not stationary, convert it to a stationary series. Explain your steps in detail while converting it.

2. Fit an additive model upon the time series. Explain the model fit and why you select it, and show the model is appropriate for the time series.
  3. Make a forecast for the next 12 observations. Plot your forecast.
  4. Save the time series (manipulated if you need), model fit and forecast in the data file "yourID\_assignment4.RData"
- C. Consider the model performance table below. The table shows the AUC (area under curve) values for two classification models on 15 data sets. Which model performed better on the data sets overall? Use paired t-test or Wilcoxon signed rank test in order to show the performance differences of the two models (**Hint:** Consider the distribution of the values before deciding the statistical test). Comment on the results and explain in detail your model selection.

	<i>Model 1</i>	<i>Model 2</i>
<i>Dataset 1</i>	0.018	0.578
<i>Dataset 2</i>	0.022	0.990
<i>Dataset 3</i>	0.455	0.124
<i>Dataset 4</i>	0.153	0.187
<i>Dataset 5</i>	0.373	0.390
<i>Dataset 6</i>	0.555	0.662
<i>Dataset 7</i>	0.158	0.545
<i>Dataset 8</i>	0.777	0.734
<i>Dataset 9</i>	0.308	0.587
<i>Dataset 10</i>	0.081	0.124
<i>Dataset 11</i>	0.369	0.864
<i>Dataset 12</i>	0.816	0.213
<i>Dataset 13</i>	0.247	0.491
<i>Dataset 14</i>	0.185	0.961
<i>Dataset 15</i>	0.312	0.996

### Data set information:

“train.csv” and “test.csv”:

Datasets contain the features extracted from a database of colonoscopic videos showing gastrointestinal lesions. There are 3 types of lesion: hyperplastic, adenoma and serrated adenoma. The last column of the dataset denoted as **Y** corresponds to the lesion name: 3 for adenoma, 1 for hyperplastic, and 2 for serrated. All other columns (denoted as **X1, X2, ..., X698**) are the raw features (without any kind of pre-processing):

#### 422 2D TEXTURAL FEATURES

- First 166 features: AHT: Autocorrelation Homogeneous Texture (Invariant Gabor Texture)
- Next 256: Rotational Invariant LBP

#### 76 2D COLOR FEATURES

- 16 Color Naming
- 13 Discriminative Color
- 7 Hue
- 7 Opponent
- 33 color gray-level co-occurrence matrix

#### 200 3D SHAPE FEATURES

- 100 shapeDNA
- 100 KPCA

“timeSeries.txt”:

This dataset contains 373 instances of daily averaged sensor response for the PT08.S1 (tin oxide) levels in an Air Quality Chemical Multi-sensor Device. The device was located on the field in a significantly polluted area, at road level, within an Italian city.