# EXAMPLE &TUTORIAL

HOW TO USE «*SEMI-AUTOMATIC FEATURE ENGINEERING TOOL*» FOR PRODUCING LSM
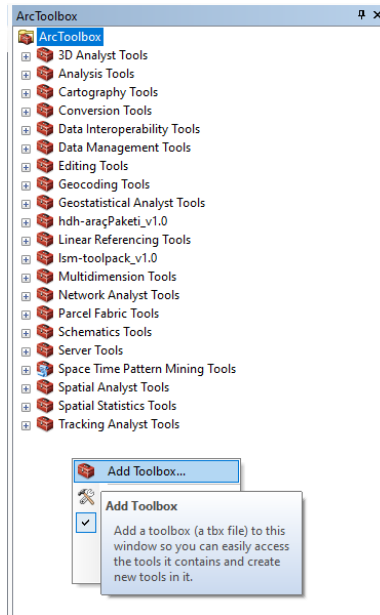
# Requirements

- ArcGIS 10.3.1 or later or ArcGIS Pro 1.1 or later (don't have it? try trial edition)

- R Statistical Computing Software, 3.3.2 or later (What is R?). If you're experiencing issues when using a new version of R (i.e.g, R 4.0 and upper) Recommended Version v3.6.3

- ArcGIS R-Bridge -- Recommended Version v1.0.1.239

- Java Runtime Environment for FSelector package

- If you are setting up modules for the first time, you will need an internet connection to install the R-packages on the depository.
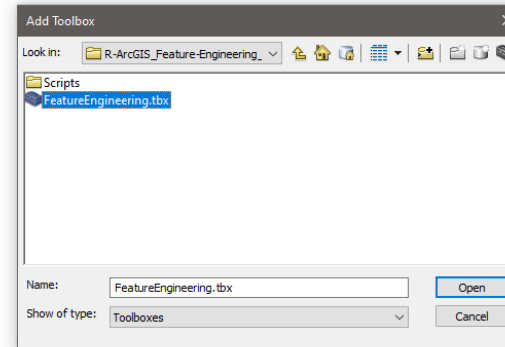
# Installation

- Download this repository and unzip 'R-ArcGIS_Feature-Engineering_ToolPack-main.zip' into the folder C:\\*targetfolder*.

  - https://github.com/emrehanks/R-ArcGIS_Feature-Engineering_ToolPack
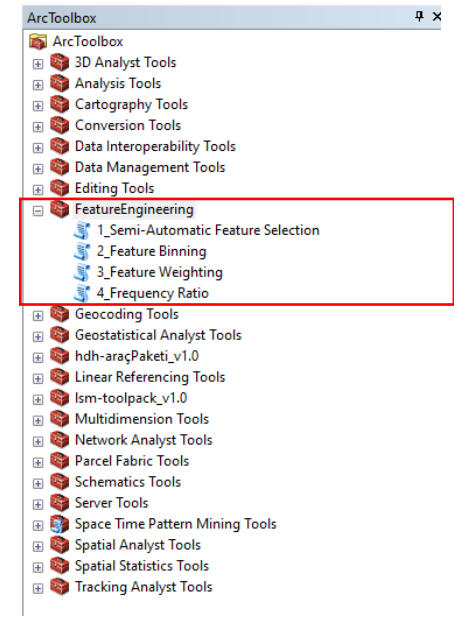
- Open your ArcGIS application:

1. Navigate ArcToolBox panel
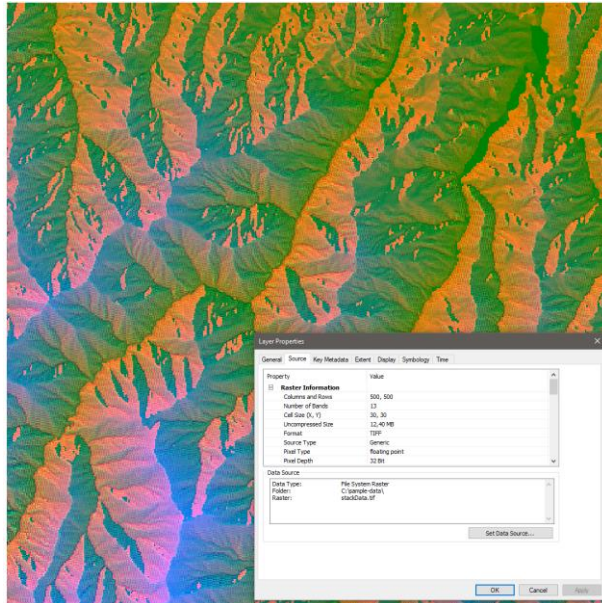
2. Add Toolbox on ArcMap
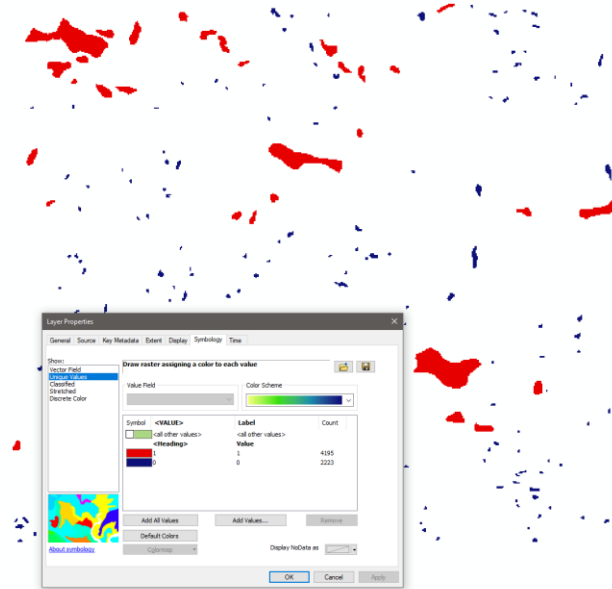
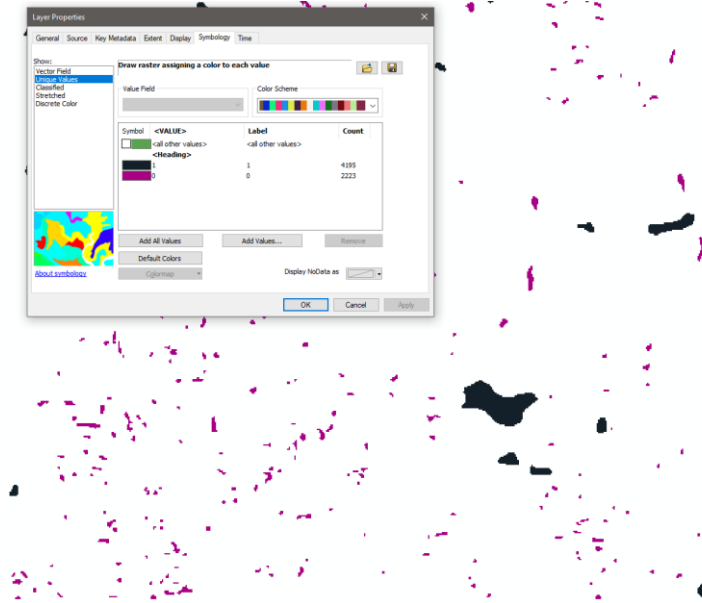3. Be sure to toolbox added on ArcToolBox panel.

# Sample Data

- Download and and unzip 'sample-data.zip' into the folder C:\\*targetfolder*.

- This data contains:



**stackData.tif** -> Multi-bands geographic image (13 bands, 500*500 and 30m pixel size)
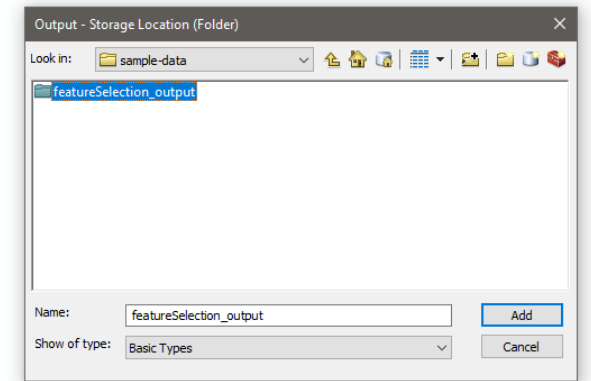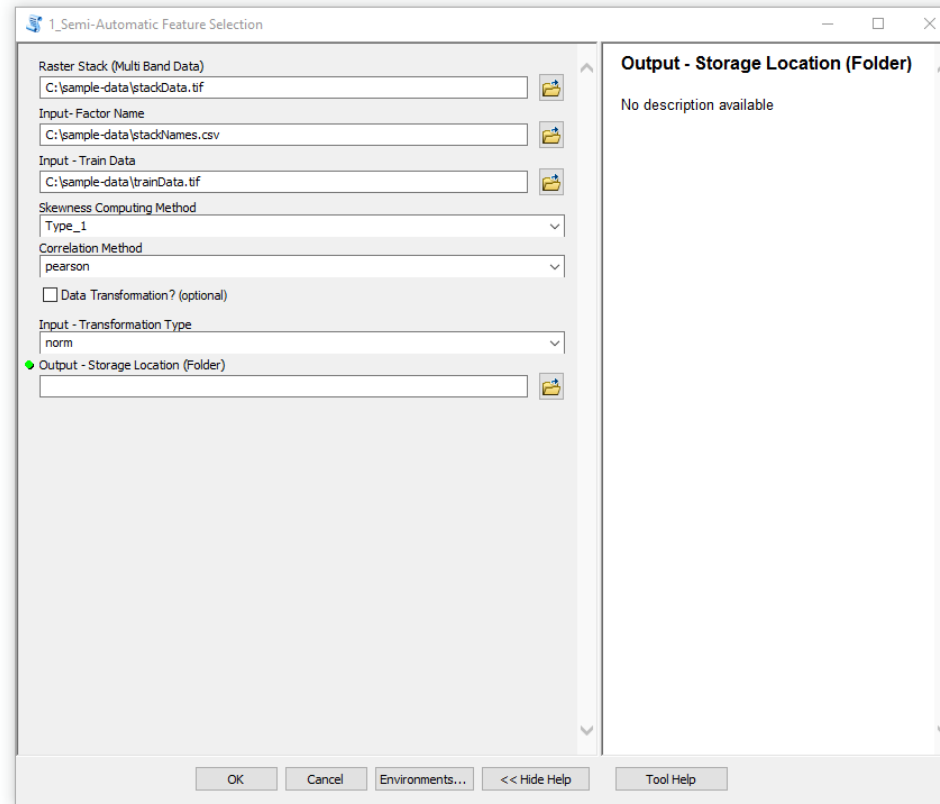**Not**: band or factor names was included in this folder as name as «stackNames.csv»

**trainData.tif**-> Single-bands geographic image (30m pixel size)

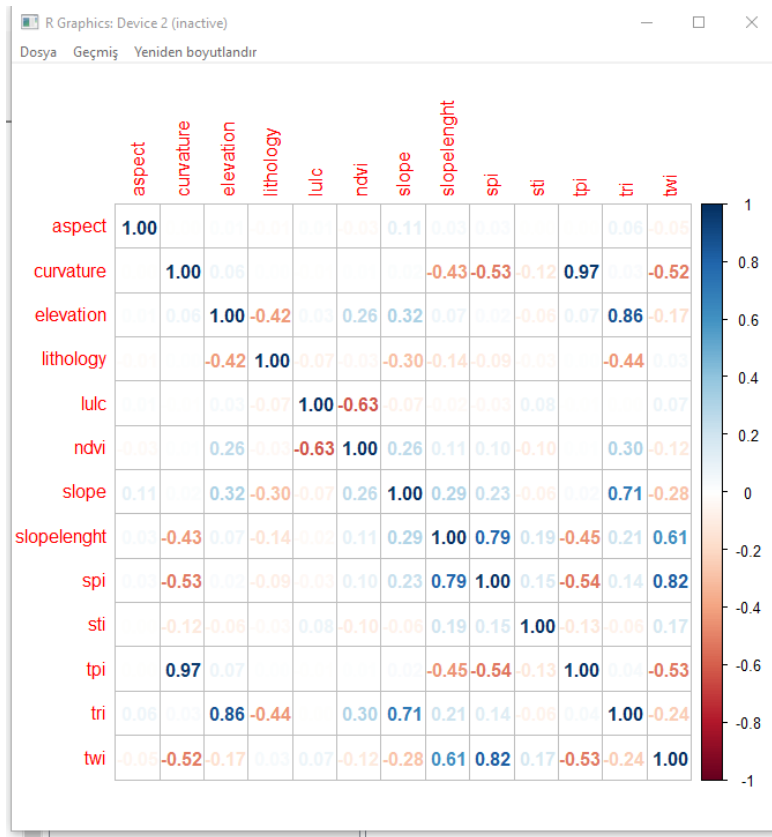**validationData.tif**-> Single-bands geographic image (30m pixel size)

# 1_Semi-Automatic Feature Selection

1. Import stackData.tif

2. If available, import stackNames.csv.

3. Import train image file **Not:** Landslide area must be labeled as '1' and non-landslide areas must be labeled as '0'.

4. Select type of skewness function.

5. Select type of Correlation method

6. Give the output location.

7. Select a folder and click Add

8. Click «OK»

# 1_Semi-Automatic Feature Selection

- Note that the data with correlation and click Ok on information window
- In this tutorial; The data namely **curvature, tri, and spi** were identified as the most correlated features.

# 1_Semi-Automatic Feature Selection

1. You can select multiple features by «ctrl and click».

2. The second information graph about correlation statistic was shown as below:

3. Pearson Correlation reported that the NDVI feature is not statistically significant. The tool is awaiting confirmation from us whether to delete this factor or not.

# 1_Semi-Automatic Feature Selection

1. The before and after analyzed of the skewness score report is also shown in the dialog panel: **Not**: *Do not fill the check bar (i.e., Close this dialog when completed successfully) if you want to note and keep all scores for future use.*

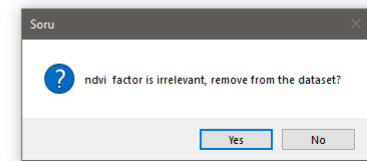2. Histogram plots of features after skewness computing method processing.

2. All feature selection processes successfully completed. Users also can access all process output from the selected folder.

# 2_Feature Binning

1. Import only countinuos feature in (Factors-Raster) tab.

2. Select Classifier method.

3. Enter number of class

4. If the user wants to process the feature weighting for future work, the pick as "**Split Sub-Classes**" must be approved

5. Import only features that categorical data type such as lithology and LULC. *However, since both lithology and lulc log are transformed by the Log Transformation process in this study, therefore both data should be selected in the first input area.*

6. Give the output location for the classified features.

7. Select a folder and click Add

8. Give the output location for split features

9. Select a folder and click Add

10. Click «OK»

# 2_Feature Binning

- Feature Binning process is completed.
- Classified and split features are stored in the selected folder after the process.

# 3_Feature Weighting

1. Import only split feature in (Factors-Raster) tab. Not: For this study data, the SPI factor was eliminated because of some problem about spi data type. Therefore, only 8 factors were selected for the weighing process.

2. Select Weight method. Suer selected several method for this process. However, not guaranteed to work every method including in this tab. **Recommended to user find the most suitable method by trial and error.**

3. Import train data.

4. Give the output location for the weighted features as a raster stack data.

5. Select a folder and give raster data and clik OK

6. Give the output location for weight score of features

7. Select a folder and give XLS or XLSX document name

8. Click «OK»

# 3_Feature Weighting

- Raster stack data and lists of weight scores are located in the selected folder.

- Finally, the necessary raster stack data was obtained to generate the Landslide Susceptibility Mapping.

# 4_Frequency Ratio

1. Import only feature produced with the Feature Selection module. Not: For this study data, the SPI factor was eliminated because of some problem about spi data type. Therefore, only 8 factors were selected for the weighing process.

2. Select the Classifier method.

3. Give number of class.

4. If «*msd*» method was selected as a classifier method, the user should be given a standard variation score.

5. Import only Landslide data with a polygon vector.

6. Select a folder and give XLS or XLSX document name

7. Give the output location for the raster stack data.

8. Click «OK»

# 4_Frequency Ratio

- Raster stack data and lists of frequency scores for each feature are located in the selected folder.
- Finally, the necessary raster stack data was obtained to generate the Landslide Susceptibility Mapping.

# Producing LSM using LSM Tool Pack v1.0

- Three basic ML methods namely Extreme Gradient Boosting, Random Forest, and Support Vector Machine can be used for modeling.

- Users easily can be reached at LSM Tool Pack v1.0 on the given Github webpage. https://github.com/emrehanks/R-ArcGIS-LSM_ToolPack

# Producing LSM using LSM Tool Pack v1.0

- Output of _4_Random Forest_ module results

LS Map

Scores of feature importance

AUC-ROC

|               | IncNodePurity |
|---------------|---------------|
| splitaspect   | 46,69777      |
| splitelevation | 104,6985     |
| splitlithology | 158,4259     |
| splitlulc     | 43,22014      |
| splitslope    | 317,4535      |
| splitslopelenght | 27,36402   |
| splittpi      | 36,46575      |
| splittwi      | 94,80457      |

Random Forest AUC : 98.8617

# Producing LSM using LSM Tool Pack v1.0

- Output of *7_Support Vector Machine* (SVM) module results

LS Map

AUC-ROC

# Producing LSM using LSM Tool Pack v1.0

- Output of *8_Extreme Gradient Boosting (XGBoost)* module results

LS Map

Scores of feature importance

AUC-ROC

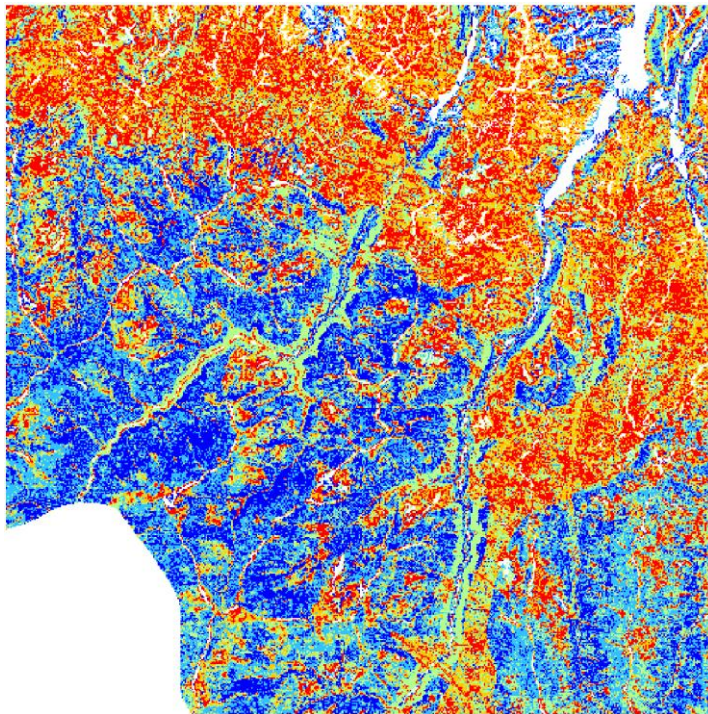| | Feature | Gain | Cover | Frequency |
|---|---|---|---|---|
| 1 | splitslope | 0,314797967 | 0,179756685 | 0,125998771 |
| 2 | splitlithology | 0,175911772 | 0,114584788 | 0,077443147 |
| 3 | splitelevation | 0,162048064 | 0,170299626 | 0,175169023 |
| 4 | splittwi | 0,127587305 | 0,152884603 | 0,118623233 |
| 5 | splitaspect | 0,06697848 | 0,121317078 | 0,170866626 |
| 6 | splitlulc | 0,059982692 | 0,105898053 | 0,094652735 |
| 7 | splittpi | 0,049694773 | 0,091197581 | 0,145666872 |
| 8 | splitslopelenght | 0,042998946 | 0,064061585 | 0,091579594 |

XgBoost AUC : 98.7356

# Producing LSM using LSM Tool Pack v1.0

- Output of *5_Performance Evaluation* module results

5_Performance Evaluation Module GUI

*Performance results of three models*

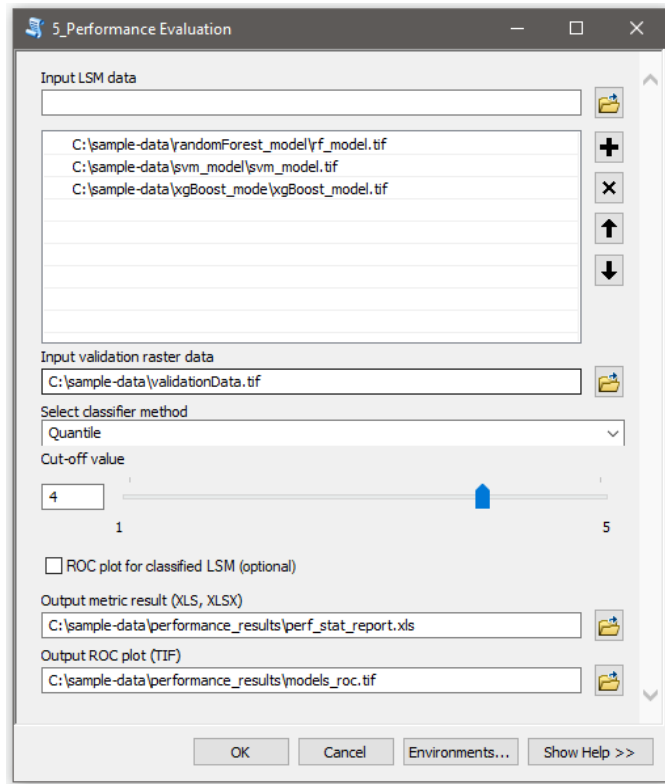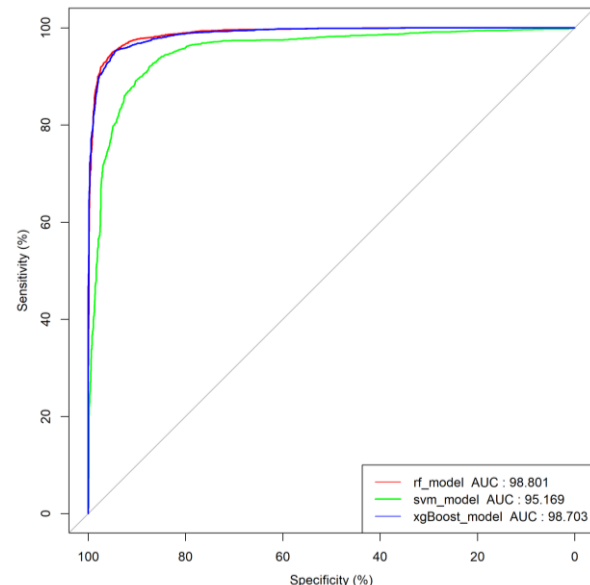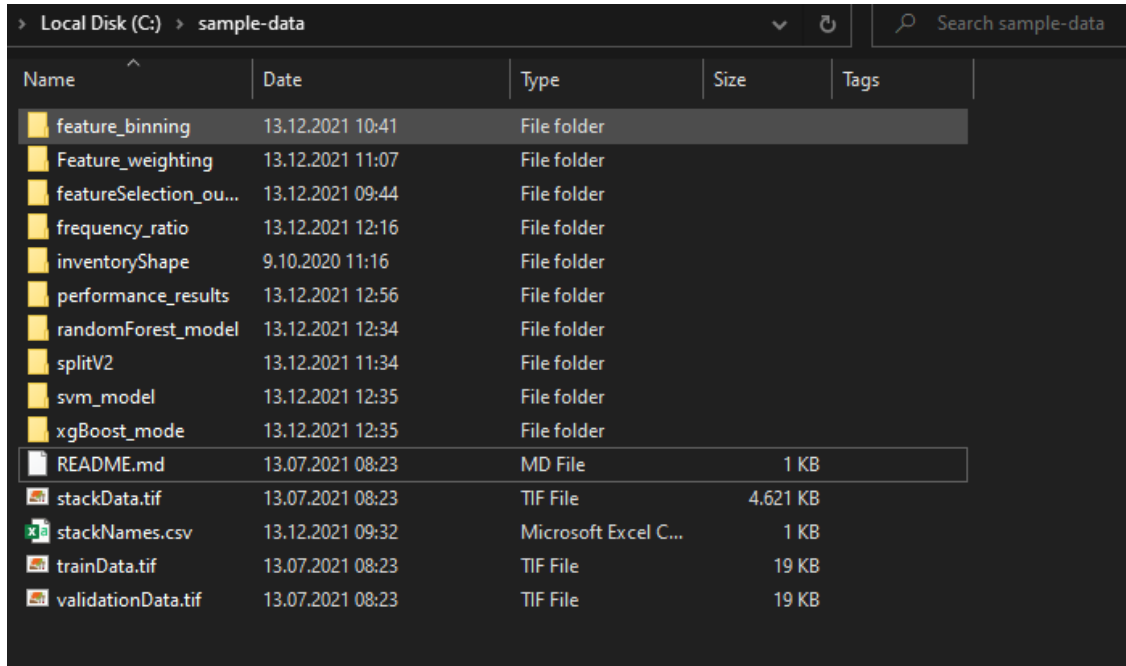| | Accuracy | AUC.Class | AUC.NonC | MAE | RMSE | Kappa | Precision | Recall | F1 |
|---|---|---|---|---|---|---|---|---|---|
| xgBoost_model | **0,909594** | 0,925583 | 0,987035 | 0,090406 | 0,300677 | 0,811683 | 0,987911 | 0,870542 | 0,92552 |
| rf_model | **0,890765** | 0,91232 | 0,988009 | 0,109235 | 0,330507 | 0,775588 | 0,991235 | 0,838119 | 0,90827 |
| svm_model | **0,8552** | 0,873667 | 0,951688 | 0,1448 | 0,380526 | 0,702629 | 0,959145 | 0,810097 | 0,878343 |

*Comparison ROC graph of models*

**5_Performance Evaluation**

Input LSM data

- C:\sample-data\randomForest_model\rf_model.tif
- C:\sample-data\svm_model\svm_model.tif
- C:\sample-data\xgBoost_mode\xgBoost_model.tif

Input validation raster data

C:\sample-data\validationData.tif

Select classifier method

Quantile

Cut-off value

4

1                                                          5

☐ ROC plot for classified LSM (optional)

Output metric result (XLS, XLSX)

C:\sample-data\performance_results\perf_stat_report.xls

Output ROC plot (TIF)

C:\sample-data\performance_results\models_roc.tif

OK    Cancel    Environments...    Show Help >>

rf_model  AUC : 98.801
svm_model  AUC : 95.169
xgBoost_model  AUC : 98.703

Sensitivity (%)
Specificity (%)

# Results

- For this sample study, XGBoost model constructed with RF-selected subset was superior to other models.

- All ouputput was published on article supplementry data folder.

Emrehan Kutlug Sahin, Ismail Colkesen, Suheda Semih Acmali, Aykut Akgun, Arif Cagdas Aydinoglu, Developing comprehensive geocomputation tools for landslide susceptibility mapping: LSM tool pack, Computers & Geosciences, 2020, 104592, ISSN 0098-3004, https://doi.org/10.1016/j.cageo.2020.104592. (http://www.sciencedirect.com/science/article/pii/S009830042030577X)