

Regularized Greedy Forest Version 1: User Guide

1 Introduction

Regularized greedy forest (RGF) [2] is a tree ensemble learning method. The purpose of this document is to guide you through the software package of RGF implementation. [2] should be referred to in order to understand RGF. With the code provided in the package, you can do the following:

Using the stand-alone executable:

- Conduct RGF training for regression tasks.
- Conduct RGF training for binary classification tasks.
- Apply trained models to new data.
- Reproduce the experimental results reported in [2].

Possible future work items include an interface for multiclass classification. It is possible to call RGF functions from your code, but documentation for this purpose is not included in this guide.

This guide is organized as follows. Section 2 quickly goes through the functionality. Data formats are described in Section 3. Section 4 explains all the parameters. In the Appendix, batch processing tools to run the experiments reported in [2] are introduced.

2 Get started: A quick tour of functionality

The first thing to do is to download the package, extract the content, and create the executable. Please follow the online instructions and README.

The top directory of the extracted content is `rgf_v1`. Make sure that the executable `rgf` (or `rgf.exe` on Windows) is at `rgf_v1/bin`. To go through the examples in this section, always set the current directory to `rgf_v1/test` because the resource files configured for using sample toy data assume that the current directory is `rgf_v1/test`. For the same reason, the path expressions in this section should be understood to be relative to `rgf_v1/test` unless it begins with `rgf_v1`.

2.1 Get started: `train` – training

To get started, set the current directory to `rgf_v1/test` and enter in the command line:

```
perl call_exe.pl ../bin/rgf train sample/train
```

If successful, the last several lines of the screen should look like:

```
Generated 5 model file(s):
output/sample.model-01
output/sample.model-02
output/sample.model-03
output/sample.model-04
output/sample.model-05

Sat Dec 10 12:17:39 2011: Done ...
elapsed: 0.115
```

What happened is: the Perl script `call_exe.pl` read parameters in the configuration file `sample/train.inp` and passed them to the executable `../bin/rgf` along with the action to take (`train`); as a result, training was conducted and five models were saved to files. Here is the configuration file `sample/train.inp` that `call_exe.pl` was given:

```
#### sample input to "train" ####
train_x_fn=sample/train.data.x # Training data points
train_y_fn=sample/train.data.y # Training targets
model_fn_prefix=output/sample.model
reg_L2=1 # Regularization parameter
algorithm=RGF
loss=LS # Square loss
test_interval=100 # Save models every time 100 leaves are added.
max_leaf_forest=500 # Stop training when #leaf reaches 500
Verbose
```

It essentially says: conduct training with the training data points in `sample/train.data.x` and the training targets in `sample/train.data.y` and save the trained models to the files. Any texts from “#” to the end of line are comments. The line “`model_fn_prefix...`” indicates that the system should generate model path names using the string “`output/sample.model`” as a prefix and attaching sequential numbers “-01”, “-02”, ... to it. It also says: training should proceed until the number of leaf nodes in the forest reaches 500 (`max_leaf_forest=500`); and every time another 100 leaf nodes are added to the forest, the system should simulate the end of training and save the model for later testing (`test_interval=100`). As a result, you should obtain five models each of which contains approximately 100, 200, 300, 400, or 500 leaf nodes. We delay explanation of other parameters until Section 4 where complete lists of parameters are given. The format of training data files is described in Section 3.1.

2.2 Models: Why one call of training results in several model files

As seen above, one call of RGF training typically results in several model files. This fact may deserve some explanation as it differs from typical boosting methods.

For example, suppose that you conducted Gradient Boosting [1] training using regression trees as base learners; stopped training when you obtained 500 trees; and saved them to a file for later testing. Then using these 500 trees, you can test 500 distinct models, each of which consists of the first k trees where $k = 1, 2, \dots, 500$, by simply changing the number of trees to be used for making predictions. You do not have to save 500 models individually. This is because Gradient Boosting does not change the previously-generated trees; it only adds new trees as it proceeds. The same can be said about other typical boosting methods such as AdaBoost.

By contrast, RGF performs *fully-corrective update of weights*, which updates the weights of *all* the leaf nodes of *all* the trees, in the designated interval and at the end of training. For this reason, if we save the model of, for example, 500 trees, then these 500 trees can be used *only* for testing the additive model of 500 trees. Unlike the Gradient Boosting example above, the first k trees of these 500 trees do not constitute a meaningful model. If we stopped training when k trees were obtained, the weights assigned to the nodes of the k trees would be totally different from those of the first k trees of the 500 trees.

It might be simpler if the system let a user only specify when to stop training and only return one model, but it would be very inefficient to train several models of different sizes this way. For efficiency, our implementation trains several models of different sizes in one call by *simulating the end of training* in the interval designated by `test_interval`. More precisely, training branches into two in the designated interval. One continues training as if nothing happened, and the other ends training, which triggers weight optimization (if it has not been triggered by the designated optimization interval), and tests or saves the model. That is how one call of training produces several models.

2.3 predict: To apply a model to new data

The next example reads a model from one of the five model files generated in Section 2.1 and applies the model to new data. Set the current directory to `rgf_v1/test` and enter:

```
perl call_exe.pl ../bin/rgf predict sample/predict
```

If successful, after parameters are displayed, something similar to the following should be displayed:

```
output/sample.pred: output/sample.model-03,#leaf=301,#tree=73
Sat Dec 10 13:20:54 2011: Done ...
```

which indicates that the prediction values were saved to `output/sample.pred`; the model was read from the file `output/sample.model-03` and it contained 301 leaf nodes and 73 trees.

The configuration file `sample/predict.inp` we used is:

```
#### sample input to "predict"
test_x_fn=sample/test.data.x      # Test data points
model_fn=output/sample.model-03   # Model file
prediction_fn=output/sample.pred   # Where to write prediction values
```

It says: read the model from `output/sample.model-03`; apply it to the data points in `sample/test.data.x`; and save the prediction values to `output/sample.pred`. The format of the prediction file is described in Section 3.3.

2.4 Executable `rgf` and Perl script `call_exe.pl`

The executable `rgf`, called through the Perl script in the examples above, takes two arguments:

`rgf action parameters`

<i>action</i>	train predict train_test
	train Conduct training and save the trained models to files. Input: training data; Output: models.
	predict Apply a model to new data. Input: a model and test data; Output: predictions
	train_test Train and test the models in one call. Input: training data and test data; Output: performance results. Optional output: models.
<i>parameters</i>	Parameters are in the form of: <i>keyword1=value1, keyword2=value2, Option1,...</i> Example: <i>algorithm=RGF,train_x_fn=data.x,train_y_fn=data.y,...</i>

To get help on parameters, call `rgf` with *action* but without *parameters*, for example:

```
rgf train
rgf predict
```

Since parameters could be long and tedious to type in, the Perl script `call_exe.pl` introduced above is provided to ease the job. It essentially reads parameters from a configuration file and concatenates them with delimiter “,” to pass to `rgf`. The syntax is as follows:

```
perl call_exe.pl executable action config_pathname
```

<i>executable</i>	Typically, <code>../bin/rgf</code> , i.e., <code>rgf_v1/bin/rgf</code> .
<i>action</i>	train predict train_test
<i>config_pathname</i>	Path name to the configuration file without extension. The extension of configuration files must be “.inp”.

In the configuration files, any text from “#” to the end of line is considered to be a comment.

Additionally, `call_exe.pl` provides an interface to perform several runs in one call with one configuration file. This is convenient, for example, for testing different degrees of regularization with other parameters fixed. `sample/regress_train_test.inp` provides a self-explaining example.

2.5 train_test: train, apply, and evaluate models

`train_test` performs training and test in one call. What `train_test` does can also be done by combining `train` and `predict` and writing an evaluation routine by yourself. The advantage of `train_test` other than convenience is faster testing. As explained in Section 2.2, during training, all the weights are updated whenever fully-corrective update of weights is performed. However, the structure of the previously-generated trees mostly remains the same. By training and testing simultaneously, such invariant information can be reused, which saves testing time.

To try the example configuration for `train_test`, set the current directory to `rgf_v1/test`, and enter:

```
perl call_exe.pl ../bin/rgf train_test sample/train_test
```

If successful, the last several lines of the screen should look like:

Generated 5 model file(s):

output/m-01
output/m-02
output/m-03
output/m-04
output/m-05

Sat Dec 10 10:17:50 2011: Done ...
elapsed: 0.135

The configuration file `sample/train_test.inp` is:

```
#### sample input to "train_test" ####
train_x_fn=sample/train.data.x # Training data points
train_y_fn=sample/train.data.y # Training target
test_x_fn=sample/test.data.x # Test data points
test_y_fn=sample/test.data.y # Test target
evaluation_fn=output/sample.evaluation
                                # Where to write evaluation results
model_fn_prefix=output/m        # Save models. This is optional.
algorithm=RGF
reg_L2=1                         # Regularization parameter
loss=LS                          # Square loss
test_interval=100                # Test at every 100 leaves
max_leaf_forest=500             # Stop training when 500 leaves are added
Verbose
```

It is mostly the same as the configuration file for `train` in Section 2.1 except that test data is specified by `test_x_fn` (data points) and `test_y_fn` (targets) and `evaluation_fn` indicates where the performance evaluation results should be written. In this example, model files are saved to files, as `model_fn_prefix` is specified. If `model_fn_prefix` is omitted, the models are not saved.

Now check the evaluation file (`output/sample.evaluation`) that was just generated. It should look like the following except that the items following `cfg` are omitted here:

```
#tree,29,#leaf,100,acc,0.61,rmse,0.9886,sqerr,0.9773,#test,100,cfg,...,output/m-01
#tree,52,#leaf,200,acc,0.66,rmse,0.9757,sqerr,0.952,#test,100,cfg,...,output/m-02
#tree,73,#leaf,301,acc,0.66,rmse,0.9824,sqerr,0.9651,#test,100,cfg,...,output/m-03
#tree,94,#leaf,400,acc,0.69,rmse,0.9767,sqerr,0.9539,#test,100,cfg,...,output/m-04
#tree,115,#leaf,501,acc,0.67,rmse,0.985,sqerr,0.9702,#test,100,cfg,...,output/m-05
```

Five lines indicate that five models were trained and tested. For example, the first line says: a model with 29 trees and 100 leaf nodes was applied to 100 data points and classification accuracy was found to be 61%, and the model was saved to `output/m-01`.

The evaluation file format is described in Section 3.4. The format of training data and test data files is described in Section 3.1.

3 Input/output file format

This section describes the format of input/output files.

3.1 Data file format

3.1.1 Data points

The data points (or feature vectors) should be given in a plain text file of the following format. Each line represents one data point. In each line, values should be separated by one or more white space characters. All the lines should have exactly the same number of values. The values should be in the format that is recognized as valid floating-point number expressions by `atof` of C libraries. The following example represents three data points of five dimensions.

```
0.3    -0.5  1  0  2
1.555  0    0  2.8 0
0      0    0  3  0
```

(NOTE) Currently, there is no support for categorical values. All the values must be numbers. This means that categorical attributes, if any, need to be converted to indicator vectors in advance.

Alternative data format for sparse data points For *sparse* data which has many zero components (e.g., bag-of-word data), the following format can be used instead. The first line should be “sparse d ” where d is the feature dimensionality. Starting from the second line, each line represents one data point. In each line, non-zero components should be specified as *feature#*:*value* where *feature#* begins from 0 and goes up to $d - 1$. For example, the three data points above can be expressed as:

```
sparse  5
0 : 0.3   1 : -0.5  2 : 1   4 : 2
0 : 1.555 3 : 2.8
3 : 3
```

3.1.2 Targets

The target values should be given in a plain text file of the following format. Each line contains the target value of one data point. The order must be in sync with the data point file. If the data is for the classification task, the values must be in $\{1, -1\}$, for example:

```
+1
-1
-1
```

If paired with the data point file example above, this means that the target value of the first data point $[0.3, -0.5, 1, 0, 2]$ is 1 and the target value of the second data point $[1.555, 0, 0, 2.8, 0]$ is -1 , and so on.

For regression tasks, the target values could be any real values, for example:

```
0.35
1.23
-0.0028
```

3.2 Data point weight file format

As introduced later, `train` and `train_test` optionally take the user-specified weights of data points as input. The data point weights should be given in a plain text file of the same format as the target file. That is, each line should contain the user-specified weight of one data point, and the order must be in sync with the data point file of training data.

3.3 Prediction file format

`predict` outputs prediction values to a file. The prediction file is a plain text file that contains one prediction value per line. The order of the values is in sync with the data point file of test data.

3.4 Evaluation file format

`train_test` outputs performance evaluation results to a file in the CSV format. Here is an example:

```
#tree,115,#leaf,500,acc,0.64,rmse,0.9802,sqerr,0.9607,#test,100,cfg,...
#tree,213,#leaf,1000,acc,0.65,rmse,0.9721,sqerr,0.945,#test,100,cfg,...
```

In the evaluation file each line represents the evaluation results of one model. In each line, each value is preceded by its descriptor; e.g., “`#tree,115`” indicates that the number of trees is 115 in the tested model. In the following, y_i and p_i are the target value and prediction value of the i -th data point, respectively; $\mathcal{I}(x)$ is the indicator function so that $\mathcal{I}(x) = 1$ if x is true and 0 otherwise; and m is the number of test data points.

Descriptor	Meaning
<code>#tree</code>	Number of trees in the model
<code>#leaf</code>	Number of leaf nodes in the model
<code>acc</code>	Accuracy regarding the task as a classification task. $\sum_{i=1}^m \mathcal{I}(y_i \cdot p_i > 0) / m$.
<code>rmse</code>	RMSE regarding the task as a regression task. $\sqrt{\sum_{i=1}^m (y_i - p_i)^2 / m}$
<code>sqerr</code>	Square error. $\text{RMSE} \times \text{RMSE}$.
<code>#test</code>	Number of test data points m .
<code>cfg</code>	Some of training parameters.

In addition, if models were saved to files, the last item of each line will be the model path name.

(NOTE) Although performances are shown in several metrics, depending on the task some are obviously meaningless and should be ignored, e.g., accuracy should be ignored on the regression task; RMSE and square error should be ignored on the classification task especially when exponential loss is used.

3.5 Model files

The model files generated by `train` or `train_test` are binary files. Caution is needed *if* you wish to share model files between the environments with different *endianness*. By default the code assumes *little-endian*. To share model files between environments with different endians the executable used in the *big-endian* environment needs to be compiled in a certain way; see README for detail.

4 Parameters

4.1 Overview of RGF training

Since many of the parameters are for controlling training, let us first give a brief overview of RGF training, focusing on the things that can be controlled via parameters. [2] should be referred to for more precise and complete definition.

Suppose that we are given n training data points $\mathbf{x}_1, \dots, \mathbf{x}_n$ and targets y_1, \dots, y_n . The additive model obtained by RGF training is in the form of: $h_{\mathcal{F}}(\mathbf{x}) = \sum_v \alpha_v \cdot g_v(\mathbf{x})$, where v goes through all the leaf

nodes in the forest \mathcal{F} , $g_v(\mathbf{x})$ is the *basis function* associated with node v , and α_v is its *weight* or coefficient. Initially, we have an empty forest with $h_{\mathcal{F}}(\mathbf{x}) = 0$. As training proceeds, the forest \mathcal{F} obtains more and more nodes so the model $h_{\mathcal{F}}(\mathbf{x})$ obtains more and more basis functions. The training objective of RGF is to find the model that minimizes regularized loss, which is the sum of loss and a regularization penalty term:

$$\frac{1}{n} \sum_{i=1}^n \ell(h_{\mathcal{F}}(\mathbf{x}_i), y_i) + \mathcal{G}(\mathcal{F}), \quad (1)$$

where ℓ is a loss function; and $\mathcal{G}(\mathcal{F})$ is the regularization penalty term. RGF grows the forest with greedy search so that regularized loss is minimized, while it performs fully-corrective update of weights to minimize the regularized loss in the designated interval. The loss function ℓ and the interval of weight optimization can be specified by parameters.

There are three methods of regularization discussed in [2]. One is L_2 regularization on leaf-only models in which the regularization penalty term $\mathcal{G}(\mathcal{F})$ is:

$$\lambda \cdot \sum_v \alpha_v^2 / 2,$$

where λ is a constant. This is equivalent to standard L_2 regularization and penalizes larger weights. The other two are called *min-penalty regularizers*. Their definition of the regularization penalty term over each tree is in the form of:

$$\lambda \cdot \min_{\{\beta_v\}} \left\{ \sum_v \gamma^{d_v} \beta_v^2 / 2 : \text{some conditions on } \{\beta_v\} \right\},$$

where d_v is the depth of node v ; and λ and γ are constants. While [2] should be consulted for precise definition of min-penalty regularizers, one thing to note here is that a larger $\gamma > 1$ penalizes deeper nodes (corresponding to more complex basis functions) more severely. Parameters are provided to choose the regularizer or to specify the degree of regularization through λ or γ .

On the regression tasks, it is sensible to normalize targets so that the average becomes zero since regularization shrinks weights towards zero. An option switch `NormalizeTarget` is provided for this purpose. When it is turned on, the model is fitted to the normalized targets $[y_i - \bar{y}]_{i=1}^n$ where $\bar{y} = \sum_{i=1}^n y_i / n$ and the final model is set to $h_{\mathcal{F}}(\mathbf{x}) + \bar{y}$.

The regularized loss in (1) can be customized not only by specifying a loss function but also by specifying user-defined weights. Let $w_i > 0$ be the user-defined weight assigned to the i -th data point. Then instead of (1) the system will minimize the following:

$$\frac{1}{\sum_{i=1}^n w_i} \sum_{i=1}^n w_i \cdot \ell(h_{\mathcal{F}}(\mathbf{x}_i), y_i) + \mathcal{G}(\mathcal{F}). \quad (2)$$

Finally, in our implementation, fully-corrective weight update is done by coordinate descent as follows:

```

for  $j = 1$  to  $J$  do
  | foreach leaf node  $v$  do
  | | // Update  $\alpha_v$  by the Newton step with step size  $\eta$  to minimize regularized loss  $\mathcal{Q}$ .
  | |  $\alpha_v \leftarrow \alpha_v - \eta \cdot \frac{\partial \mathcal{Q} / \partial \delta_v |_{\delta_v=0}}{\partial^2 \mathcal{Q} / \partial \delta_v^2 |_{\delta_v=0}}$  //  $\delta_v$  is the additive change to  $\alpha_v$ .
  | end
end

```

The constants J and η above can be changed using the parameters described below, although, in our experiments, we never used them as the default values worked fine on a number of datasets.

4.2 Naming conventions and notation

There are two types of parameters: keyword-value pairs and options. The keywords begin with a lower-case letter and should be followed by *=value*, e.g., “`loss=LS`”. The options begin with an upper-case letter, e.g., “`Verbose`”, and indicate to turn on corresponding option switches, which are off by default.

In the lists below, ‘*’ in front of the keyword indicates that the designated keyword-value pair is required and cannot be omitted.

4.3 Parameters for `train`

The parameters described in this section are used by the `train` function. `train` trains models and saves them to files.

4.3.1 Parameters to control input and output for `train`

A pair of training data files (one contains the data points and the other contains the target values) are required. Another required parameter is the prefix of model path names, which is used to generate model path names by attaching to it sequential numbers “-01”, “-02”, and so on. The reason why one call of training typically produces multiple model files is explained in Section 2.2.

Optionally, training can be resumed from the point where training was ended last time, which we call *warm-start*. To do warm-start, the model file from which training should be resumed needs to be specified. Also optionally, user-defined weights of training data points can be specified through `train_w_fn`. They are used as in (2).

Required parameters to control input and output for <code>train</code>	
* <code>train_x_fn=</code>	Path to the data point file of training data.
* <code>train_y_fn=</code>	Path to the target file of training data.
* <code>model_fn_prefix=</code>	To save models to files, path names are generated by attaching “-01”, “-02”,... to this value.
Optional parameters to control input and output for <code>train</code>	
<code>train_w_fn</code>	Path to the file of user-defined weights assigned to training data points.
<code>model_fn_for_warmstart=</code>	Path to the model file from which training should do warm-start.

4.3.2 Parameters to control training

In the list below, the first group of parameters are most important in the sense that they would affect either accuracy of the models or speed of training directly, and they were actually used in the experiments reported in [2]. The second group of parameters never needed to be specified in our experiments, as the default values worked fine on a number of datasets, but they may be useful in some situations. The third group is for displaying information and specifying the memory allocation policy.

The variables below refer to the corresponding variables in the overview in Section 4.1.

Parameters to control training	
algorithm=	RGF RGF_Opt RGF_Sib (Default: RGF) RGF: RGF with L_2 regularization on leaf-only models RGF_Opt: RGF with min-penalty regularization RGF_Sib: RGF with min-penalty regularization with the sum-to-zero sibling constraints.
loss=	Loss function $\ell(p, y)$. LS Expo (Default: LS) LS: square loss $(p - y)^2/2$; Expo: exponential loss $\exp(-py)$.
max_leaf_forest=	Training will be terminated when the number of leaf nodes in the forest reaches this value. It should be large enough so that a good model can be obtained at some point of training, whereas a smaller value makes training time shorter. Appropriate values are data-dependent and in [2] varied from 1000 to 10000. (Default:10000)
NormalizeTarget	If turned on, training targets are normalized so that the average becomes zero. It was turned on in all the regression experiments in [2].
* reg_L2=	λ . Used to control the degree of L_2 regularization. Crucial for good performance. Appropriate values are data-dependent. In [2], either 1, 0.1, or 0.01 often produced good results though with exponential loss (loss=Expo) some data required smaller values such as 1e-10 or 1e-20.
reg_depth=	γ . Must be no smaller than 1. Meant for being used with algorithm=RGF_Opt RGF_Sib. A larger value penalizes deeper nodes more severely. In [2] it was set to 2 or 4 with algorithm=RGF_Sib. (Default: 1)
test_interval=	Test interval in terms of the number of leaf nodes. For example, if test_interval=500, every time 500 leaf nodes are newly added to the forest, end of training is simulated and the model is tested or saved for later testing. For efficiency, it must be either multiple or divisor of the optimization interval (opt_interval: default 100). If not, it may be changed by the system automatically. (Default:500)
Parameters that are probably rarely used	
min_pop=	Minimum number of training data points in each leaf node. Smaller values may slow down training. Too large values may degrade model accuracy. (Default:10)
num_iteration_opt=	J . Number of iterations of coordinate descent to optimize weights. (Default:10 for square loss; 5 for exponential loss and the likes)
num_tree_search=	Number of trees to be searched for the nodes to split. The most recently-grown trees are searched first. (Default:1)
opt_interval=	Weight optimization interval in terms of the number of leaf nodes. For example, if opt_interval=100, weight optimization is performed every time approximately 100 leaf nodes are newly added to the forest. (Default:100)
opt_stepsize=	η . Step size of Newton updates used in coordinate descent to optimize weights. (Default:0.5)
Other parameters	
Verbose	Print information during training.
Time	Measure and display elapsed time for node search and weight optimization.
memory_policy=	Conservative Generous. (Default:Generous)

memory_policy=	Conservative Generous. (Default:Generous)
----------------	---

4.4 Parameters for predict

predict reads a model saved by train or train_test, applies it to new data, and saves prediction values to a file.

Parameters for predict	
* test_x_fn	Path to the data point file of test data.
* model_fn	Path to the model file.
* prediction_fn	Path to the prediction file to write prediction values to.

4.5 Parameters for train_test

train_test trains models with training data and evaluates them on test data in one call.

4.5.1 Parameters to control input and output for train_test

train_test requires a pair of training data files (one contains the data points and the other contains the target values) and a pair of test data files.

Optionally, the models can be saved to files by specifying model_fn_prefix. The value specified with model_fn_prefix is used to generate model path names by attaching to it sequential numbers “-01”, “-02”, and so on. The reason why one call of training typically produces multiple model files is explained in Section 2.2. Other things that can be done optionally are the same as train. That is, optionally, training can be resumed from the point where training was ended last time (*warm-start*). Also optionally, user-defined weights of training data points can be specified through train_w_fn; see Section 4.1 for how they are used.

Parameters to control input and output for train_test	
* train_x_fn=	Path to the data point file of training data.
* train_y_fn=	Path to the target file of training data.
* test_x_fn=	Path to the data point file of test data.
* test_y_fn=	Path to the target file of test data.
evaluation_fn	Path to the file to write performance evaluation results to. If omitted, the results are written to stdout.
Append_evaluation	Open the file to write evaluation results to with the append mode.
model_fn_prefix=	If omitted, the models are not saved to files. Model path names are generated by attaching “-01”, “-02”,... to this value to save models.
train_w_fn	Path to the file of user-defined weights assigned to training data points.
model_fn_for_warmstart=	Path to the input model file from which training should do warm-start.

4.5.2 Parameters to control training

The parameters to control training for train_test are the same as those for train; see Section 4.3.2.

References

- [1] Jerome Friedman. Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, 29, 2001.
- [2] Rie Johnson and Tong Zhang. Learning nonlinear functions using regularized greedy forest. Technical report, Tech Report: arXiv:1109.0887, 2011.
- [3] Robert E. Schapire. The boosting approach to machine learning: An overview. *Nonlinear Estimation and Classification*, 2003.

A Appendix

A.1 To run the experiments in [2]

To run the experiments described in [2], you need to download the data archive file `data_for_ex.tar.gz` and extract the data files to `rgf_v1/exp_data`. Once this is done, 392 files whose names end with either “.x” or “.y” should be at `rgf_v1/exp_data`. These files are the exact training portions and test portions of the data used in [2].

A.1.1 To run the RGF experiments

The RGF results in the paper can be obtained by running the following three scripts:

<code>exp_regress</code>	to run the regression experiments
<code>exp_classif_square</code>	to run the classification experiments with square loss
<code>exp_classif_expo</code>	to run the classification experiments with exponential loss

The scripts need to be executed at `rgf_v1/exp` as they use relative path to access resources. They take the name of the executable as an argument. The files ending with `.cmd` are Windows command files for Windows users. Set the current directory to `rgf_v1/exp` and enter the following at the Windows command prompt, to start the experiments.

```
exp_regress      ..\bin\rgf
exp_classif_square ..\bin\rgf
exp_classif_expo ..\bin\rgf
```

The files ending with `.ss` are shell scripts for the users of Unix-like systems. Enter the following in the command line to start the experiments:

```
./exp_regress.ss      ../bin/rgf
./exp_classif_square.ss ../bin/rgf
./exp_classif_expo.ss  ../bin/rgf
```

(NOTE1) Each script could take some hours as it goes through a number of configurations over a number of datasets. If you are interested in running only part of the experiments, comment out the lines you do not want to run.

(NOTE2) You may need to customize the shell scripts for your environment.

A.1.2 To run the GBDT and AdaBoost experiments

In [2], Gradient Boosting Decision Tree (GBDT) [1] and AdaBoost (see for example [3]). were also tested for comparison. To run these experiments, you need to have an additional executable `others` (or `others.exe` on Windows) at `rgf_v1/bin`. See README for how to build `others`.

The GBDT and AdaBoost results in the paper can be obtained by running the following two scripts:

```
exp_square_gb    to run the experiments with GBDT with square loss
exp_expo_gb_ada  to run the experiments with GBDT with exponential loss and AdaBoost
```

The scripts need to be executed at `rgf_v1/exp` as they use relative path to access resources. They take the name of the executable as an argument. The files ending with `.cmd` are Windows command files for Windows users. Set the current directory to `rgf_v1/exp` and enter the following at the Windows command prompt, to start the experiments.

```
exp_square_gb    ..\bin\others
exp_expo_gb_ada  ..\bin\others
```

The files ending with `.ss` are shell scripts for the users of Unix-like systems. Enter the following in the command line to start the experiments:

```
./exp_square_gb.ss    ../bin/others
./exp_expo_gb_ada.ss  ../bin/others
```

(NOTE1) Since it takes many hours to run the entire set of experiments with GBDT described in the paper, the scripts above are configured to run selected best-performing configurations. Still, each script could take some hours. If you are interested in running only part of the experiments, comment out the lines you do not want to run.

(NOTE2) You may need to customize the shell scripts for your environment.

A.1.3 Interpretation of the results

After running the experiments, the results are found at `rgf_v1/exp/output`. The `*.perf` files contain the performance results; see Section 3.4 for the format. The `*.perf.log` files contain the parameter settings used for the corresponding `*.perf` files. The first part of the filename indicates the dataset and configuration as in:

dataname_run_config.perf[.log]

<i>dataname</i>	Dataset name. “syn...” and “rsyn...” are synthesized data for classification and regression, respectively. It also indicates that <code>rgf_v1/exp/input/dataname.ds</code> was used to configure the experiments.
<i>run</i>	Run identifier. 01–10 for synthesized data and 01–03 for others.
<i>config</i>	<code>rgf_v1/exp/input/config.alg</code> was used to configure the experiments.

A.1.4 Averaging the results of multiple runs

A tool to compute the average performance of several runs, as is presented in [2], is also provided:

```

avg_runs_leaf   to plot performance against the number of leaf nodes
avg_runs_tree   to plot performance against the number of trees

```

These take the executable (rgf) as an argument. For example, set the current directory to `rgf_v1/exp` and enter in the Windows command prompt:

```

avg_runs_leaf ..\bin\rgf

```

The scripts generate CSV files, one for each dataset, e.g., `syn164t100.csv` contains all the results on the 64-leaf synthesized data classification results. Essentially, the performances can be plotted similarly to [2] by using: the first column as the x -axis values; the fourth and later columns as the y -axis values; and the first row as legends. More precisely, the format of the CSV files is as follows.

Cell positions	Meaning of values
row#=1, col# \geq 4	Short description of the configuration associated with this column; can be used as legends to draw graphs.
row#=2, col# \geq 4	Long description of the configuration associated with this column, which ends with “(#run= k)” indicating the average was taken over k runs.
row# \geq 3, col#=1	The number of leaf nodes if produced by <code>avg_runs_leaf</code> or the number of trees if produced by <code>avg_runs_tree</code> ; associated with the results in this row.
row# \geq 3, col#= 4, 5, 6	Performance results of the best-performing non-RGF configurations, judged by the peak performance. Accuracy(%) for classification; square error for regression.
row# \geq 3, col# \geq 7	Performance results of all the configurations. Accuracy(%) for classification; square error for regression.

(NOTE1) To obtain the same results as [2], it is important to make sure that the number of runs averaged is correct – which is 10 on the synthesized data and 3 on the others. Check the long descriptions in the second line of the CSV files. They should all end with “(#run=10)” on the synthesized data and “(#run=3)” on the other datasets. If not, something is wrong; check the performance files and log files in `rgf_v1/exp/output`.

(NOTE2) The results in [2] can be exactly reproduced by the provided Windows executable, which was compiled by Microsoft Visual C++ 2005 Express Edition. The results may slightly vary when other compilers are used due to the differences in compiler-dependent floating-point processing and optimization.