

BoostTree and BoostForest for Ensemble Learning: Supplementary Materials

Changming Zhao¹, Dongrui Wu^{1,*}, Jian Huang¹, Ye Yuan¹, Hai-Tao Zhang¹, Ruimin Peng¹, Zhenhua Shi¹ and Chenfeng Guo¹

¹Key Laboratory of the Ministry of Education for Image Processing and Intelligent Control, School of Artificial Intelligence and Automation, Huazhong University of Science and Technology, Wuhan, China. *e-mail: drwu@hust.edu.cn

1 Supplementary Algorithms

This section presents pseudo-code of all algorithms introduced in this article.

Supplementary Algorithm 1: LogitBoost for classification¹.

Input: K , the maximum number of iterations;

$\{\mathbf{x}_n, y_n\}_{n=1}^N$, N labeled training samples.

Output: The ensemble $F(\mathbf{x})$.

if $J == 2$ **then**

Initialize $F(\mathbf{x}) = 0$, $w_n = \frac{1}{N}$, and $p(\mathbf{x}_n) = \frac{1}{2}$, $n = 1, 2, \dots, N$

for $k = 1 : K$ **do**

$z_n = \frac{y_n - p(\mathbf{x}_n)}{p(\mathbf{x}_n)(1 - p(\mathbf{x}_n))}$, $n = 1, \dots, N$;

$w_n = p(\mathbf{x}_n)(1 - p(\mathbf{x}_n))$, $n = 1, \dots, N$;

Fit a function $f_k(\mathbf{x})$ by weighted least-squares regression, using $\{\mathbf{x}_n, z_n\}_{n=1}^N$ and weights $\{w_n\}_{n=1}^N$;

$F(\mathbf{x}) \leftarrow F(\mathbf{x}) + f_k(\mathbf{x})$;

$p(\mathbf{x}_n) = \frac{1}{1 + e^{-F(\mathbf{x}_n)}}$, $n = 1, \dots, N$;

end

else

Initialize $F_j(\mathbf{x}) = 0$, $w_{nj} = \frac{1}{N}$, and $p_j(\mathbf{x}_n) = \frac{1}{J}$, $n = 1, 2, \dots, N$, $j = 1, 2, \dots, J$;

for $k = 1 : K$ **do**

for $j = 1 : J$ **do**

$z_{nj} = \frac{y_{nj} - p_j(\mathbf{x}_n)}{p_j(\mathbf{x}_n)(1 - p_j(\mathbf{x}_n))}$, $n = 1, \dots, N$;

$w_{nj} = p_j(\mathbf{x}_n)(1 - p_j(\mathbf{x}_n))$, $n = 1, \dots, N$;

Fit a function $f_{kj}(\mathbf{x})$ by weighted least-squares regression, using $\{\mathbf{x}_n, z_{nj}\}_{n=1}^N$ and weights $\{w_{nj}\}_{n=1}^N$;

end

$f_{kj}(\mathbf{x}) \leftarrow \frac{J-1}{J} \left[f_{kj}(\mathbf{x}) - \frac{1}{J} \sum_{i=1}^J f_{ki}(\mathbf{x}) \right]$;

$F_j(\mathbf{x}) \leftarrow F_j(\mathbf{x}) + f_{kj}(\mathbf{x})$;

$p_j(\mathbf{x}_n) = \frac{e^{F_j(\mathbf{x}_n)}}{\sum_{i=1}^J e^{F_i(\mathbf{x}_n)}}$, $n = 1, \dots, N$;

end

end

Supplementary Algorithm 2: BoostTree using ridge regression (RR) as the node function.

Input: Data = $\{\mathbf{x}_n, y_n\}_{n=1}^N$, N training samples, where $\mathbf{x}_n \in \mathbb{R}^D$;

Pool_{MinSamplesLeaf}, candidate value pool of the minimum number of samples at a leaf;

Pool _{λ} , candidate value pool of the ℓ_2 regularization parameter λ ;

(optional) MaxNumLeaf, the maximum number of leaves.

Output: A BoostTree.

NumLeaf = 1;

$f(\mathbf{x}) \equiv 0$;

root \leftarrow {data = Data, model = f , leftChild = None, rightChild = None};

BoostTree \leftarrow split(root);

split(node){

Let t be the index of the current node;

$\{\mathbf{x}_n, y_n\}_{n \in I} \leftarrow$ node.data;

$F_I(\mathbf{x}) \leftarrow \sum_{m \in \text{Path}(t)} \text{node}_m.\text{model}(\mathbf{x})$;

Randomly select N_{\min}^L and N_{\min}^R from Pool_{MinSamplesLeaf};

$\delta_{\mathcal{L}}^{\max} = 0$;

Randomly select \sqrt{D} features to create a feature subset \mathcal{F} ;

for $d \in \mathcal{F}$ **do**

$S_d = \{x_{n,d} | n \in I\}$;

for s in S_d **do**

$I_L = \{n | x_{n,d} \leq s, n \in I\}$;

$I_R = \{n | x_{n,d} > s, n \in I\}$;

if $|I_L| \geq N_{\min}^L$ and $|I_R| \geq N_{\min}^R$ **then**

 Randomly select λ_L and λ_R from Pool _{λ} ;

 Calculate $\delta_{\mathcal{L}}$ in (12);

if $\delta_{\mathcal{L}} > \delta_{\mathcal{L}}^{\max}$ **then**

$\delta_{\mathcal{L}}^{\max} = \delta_{\mathcal{L}}, \quad d^* = d, \quad s^* = s, \quad \lambda_L^* = \lambda_L, \quad \lambda_R^* = \lambda_R$;

end

end

end

end

if $\delta_{\mathcal{L}}^{\max} > 0$ **then**

$I_L^* = \{n | x_{n,d^*} \leq s^*, n \in I\}$;

$I_R^* = \{n | x_{n,d^*} > s^*, n \in I\}$;

$f_L^* = \text{FitModel}(\{\mathbf{x}_n, y_n\}_{n \in I_L^*}, F_I, \lambda_L^*)$;

$f_R^* = \text{FitModel}(\{\mathbf{x}_n, y_n\}_{n \in I_R^*}, F_I, \lambda_R^*)$;

 Calculate $\delta_{\mathcal{L}}$ in (9);

if $\delta_{\mathcal{L}} > 0$ **then**

 node.leftChild = {data = $\{\mathbf{x}_n, y_n\}_{n \in I_L^*}$, model = f_L^* , leftChild = None, rightChild = None};

 node.rightChild = {data = $\{\mathbf{x}_n, y_n\}_{n \in I_R^*}$, model = f_R^* , leftChild = None, rightChild = None};

 NumLeaf = NumLeaf + 1;

end

if MaxNumLeaf is not supplied, or NumLeaf \leq MaxNumLeaf **then**

 For each leaf, the cross entropy loss (classification) or RMSE (regression) multiplied by the number of samples at that leaf is used as the impurity score;

 Identify node*, the leaf node with the highest impurity score;

 split(node*);

end

end

Supplementary Algorithm 3: FitModel for linear regression.

Input: $\{\mathbf{x}_n, y_n\}_{n \in I}$, sample set of the current node;

F_I , ensemble of the models along the path from the root node to the parent node of the current node;

λ , the ℓ_2 regularization parameter.

Output: Linear regression model f_m for the current node.

$$\tilde{y}_n = y_n - F_I(\mathbf{x}_n), n \in I;$$

Fit $f_m = \text{RidgeRegression}(\{\mathbf{x}_n, \tilde{y}_n\}_{n \in I}, \lambda)$ using ridge regression with regularization parameter λ .

Supplementary Algorithm 4: FitModel for binary classification.

Input: $\{\mathbf{x}_n, y_n\}_{n \in I}$, sample set of the current node;

F_I , ensemble of the models along the path from the root node to the parent node of the current node;

λ , the ℓ_2 regularization parameter.

Output: Linear classifier f_m for the current node.

$$p(\mathbf{x}_n) = \text{sigmod}[F_I(\mathbf{x}_n)], n \in I;$$

$$\tilde{y}_n = \frac{y_n - p(\mathbf{x}_n)}{p(\mathbf{x}_n)(1 - p(\mathbf{x}_n))}, n \in I;$$

$$\tilde{y}_n = \text{Clip}(\tilde{y}_n) \text{ in (18)}, n \in I;$$

$$w_n = p(\mathbf{x}_n)(1 - p(\mathbf{x}_n)), n \in I;$$

$$q_{5\%} = \text{Quantile}(\{w_n \mid n \in I\}, 5);$$

$$\mathcal{D}' = \{(\mathbf{x}_n, \tilde{y}_n, w_n) \mid w_n > q_{5\%}, n \in I\};$$

Fit $f_m = \text{WeightedRidgeRegression}(\mathcal{D}', \lambda)$ using weighted ridge regression on \mathcal{D}' with regularization parameter λ .

Supplementary Algorithm 5: FitModel for J -class ($J > 2$) classification.

Input: $\{\mathbf{x}_n, y_n\}_{n \in I}$, sample set of the current node;

F_I , ensemble of the models along the path from the root node to the parent node of the current node;

λ , the ℓ_2 regularization parameter.

Output: The linear classifier set \mathbf{f}_m for the current node.

Compute $p_j(\mathbf{x}_n) = \text{softmax}^j(\hat{\mathbf{y}}_n)$, $n \in I$, $j = 1, \dots, J$;

for $j = 1 : J$ **do**

$$\tilde{y}_{nj} = \frac{y_{nj} - p_j(\mathbf{x}_n)}{p_j(\mathbf{x}_n)(1 - p_j(\mathbf{x}_n))}, n \in I;$$

$$\tilde{y}_{nj} = \text{Clip}(\tilde{y}_{nj}) \text{ in (18)}, n \in I;$$

$$w_{nj} = p_j(\mathbf{x}_n)(1 - p_j(\mathbf{x}_n)), n \in I;$$

$$q_{5\%} = \text{Quantile}(\{w_{nj} \mid n \in I\}, 5);$$

$$\mathcal{D}' = \{(\mathbf{x}_n, \tilde{y}_{nj}, w_{nj}) \mid w_{nj} > q_{5\%}, n \in I\};$$

Fit $f_j = \text{WeightedRidgeRegression}(\mathcal{D}', \lambda)$ using weighted ridge regression on \mathcal{D}' with regularization parameter λ ;

end

$$f_j(\mathbf{x}) \leftarrow \frac{J-1}{J} \left[f_j(\mathbf{x}) - \frac{1}{J} \sum_{i=1}^J f_i(\mathbf{x}) \right], \quad j = 1, \dots, J;$$

$$\mathbf{f}_m = \{f_1, f_2, \dots, f_J\}.$$

Supplementary Algorithm 6: BoostForest training algorithm.

Input: Data = $\{\mathbf{x}_n, y_n\}_{n=1}^N$, N training samples, where $\mathbf{x}_n \in \mathbb{R}^D$;

NumEstimators, the number of BoostTrees;

Pool_{MinSamplesLeaf}, candidate value pool of the minimum number of samples at a leaf;

Pool _{λ} , candidate value pool of the regularization parameter λ .

Output: A BoostForest

BoostForest = {};

for $i = 1 : \text{NumEstimators}$ **do**

 Bootstrap Data' from Data;

 Train BoostTree _{i} on Data' using Supplementary Algorithm 2;

 Add BoostTree _{i} to BoostForest;

end

Supplementary Algorithm 7: BoostTree-ELM for regression.

Input: Data = $\{\mathbf{x}_n, y_n\}_{n=1}^N$, N training samples, where $\mathbf{x}_n \in \mathbb{R}^D$;

Pool_{MinSamplesLeaf}, candidate value pool of the minimum number of samples at a leaf;

Pool_{NumHiddenNodes}, candidate value pool of the number of hidden nodes;

Pool _{λ} , candidate value pool of the regularization parameter λ ;

(optional) MaxNumLeaf, the maximum number of leaves.

Output: A BoostTree-ELM.

NumLeaf = 1;

$f(\mathbf{x}) \equiv 0$;

root \leftarrow {data = Data, model = f , leftChild = None, rightChild = None};

BoostTree-ELM \leftarrow split(root);

split(node){

Let t be the index of the current node;

$\{\mathbf{x}_n, y_n\}_{n \in I} \leftarrow$ node.data;

$F_I(\mathbf{x}) \leftarrow \sum_{m \in \text{Path}(t)} \text{node}_m.\text{model}(\mathbf{x})$;

Randomly select N_{\min}^L and N_{\min}^R from Pool_{MinSamplesLeaf};

$\delta_{\mathcal{L}}^{\max} = 0$;

Randomly select \sqrt{D} features to create a feature subset \mathcal{F} ;

for $d \in \mathcal{F}$ **do**

$S_d = \{x_{n,d} | n \in I\}$;

for s in S_d **do**

$I_L = \{n | x_{n,d} \leq s, n \in I\}$;

$I_R = \{n | x_{n,d} > s, n \in I\}$;

if $|I_L| \geq N_{\min}^L$ and $|I_R| \geq N_{\min}^R$ **then**

 Randomly select λ_L and λ_R from Pool _{λ} ;

 Calculate $\delta_{\mathcal{L}}$ in (12);

if $\delta_{\mathcal{L}} > \delta_{\mathcal{L}}^{\max}$ **then**

$\delta_{\mathcal{L}}^{\max} = \delta_{\mathcal{L}}$, $d^* = d$, $s^* = s$, $\lambda_L^* = \lambda_L$, $\lambda_R^* = \lambda_R$;

end

end

end

end

if $\delta_{\mathcal{L}}^{\max} > 0$ **then**

 Randomly select M_L and M_R from Pool_{NumHiddenNodes};

$I_L^* = \{n | x_{n,d^*} \leq s^*, n \in I\}$;

$I_R^* = \{n | x_{n,d^*} > s^*, n \in I\}$;

$f_L^* = \text{FitModelELM}(\{\mathbf{x}_n, y_n\}_{n \in I_L^*}, F_I, \lambda_L^*, M_L)$;

$f_R^* = \text{FitModelELM}(\{\mathbf{x}_n, y_n\}_{n \in I_R^*}, F_I, \lambda_R^*, M_R)$;

 Calculate $\delta_{\mathcal{L}}$ in (9);

if $\delta_{\mathcal{L}} > 0$ **then**

 node.leftChild = {data = $\{\mathbf{x}_n, y_n\}_{n \in I_L^*}$, model = f_L^* , leftChild = None, rightChild = None};

 node.rightChild = {data = $\{\mathbf{x}_n, y_n\}_{n \in I_R^*}$, model = f_R^* , leftChild = None, rightChild = None};

 NumLeaf = NumLeaf + 1;

end

if MaxNumLeaf is not supplied, or NumLeaf \leq MaxNumLeaf **then**

 For each leaf, the cross entropy loss (classification) or RMSE (regression) multiplied by the number of samples at that leaf is used as the impurity score;

 Identify node*, the leaf node with the highest impurity score;

 split(node*);

end

end

Supplementary Algorithm 8: FitModelELM for regression.

Input: $\{\mathbf{x}_n, y_n\}_{n \in I}$, sample set of the current node;

F_I , ensemble of the models along the path from the root node to the parent node of the current node;

λ , the regularization parameter;

M , the number of hidden nodes.

Output: ELM model f_m for the current node.

$\tilde{y}_n = y_n - F_I(\mathbf{x}_n)$, $n \in I$;

Fit $f_m = \text{ELM}(\{\mathbf{x}_n, \tilde{y}_n\}_{n \in I}, \lambda, M)$ using ELM with regularization parameter λ and M hidden nodes.

Supplementary Algorithm 9: BoostForest-ELM for regression.

Input: Data = $\{\mathbf{x}_n, y_n\}_{n=1}^N$, N training samples, where $\mathbf{x}_n \in \mathbb{R}^D$;

NumEstimators, the number of trees;

PoolMinSamplesLeaf, candidate value pool of the minimum number of samples at a leaf;

PoolNumHiddenNodes, candidate value pool of the number of hidden nodes;

Pool $_{\lambda}$, candidate value pool of the regularization parameter λ .

Output: A BoostForest-ELM.

BoostForest-ELM = {};

for $i = 1 : \text{NumEstimators}$ **do**

 Bootstrap Data' from Data;

 Train BoostTree-ELM $_i$ on Data' using Supplementary Algorithm 7;

 Add BoostTree-ELM $_i$ to BoostForest-ELM;

end

Supplementary Algorithm 10: BoostTree-SVR for regression.

Input: Data = $\{\mathbf{x}_n, y_n\}_{n=1}^N$, N training samples, where $\mathbf{x}_n \in \mathbb{R}^D$;
 Pool_{MinSamplesLeaf}, candidate value pool of the minimum number of samples at a leaf;
 Pool_C, candidate value pool of the regularization parameter;
 Pool _{ϵ} , candidate value pool of the slack variable;
 (optional) MaxNumLeaf, the maximum number of leaves.

Output: A BoostTree-SVR

NumLeaf = 1;

$f(\mathbf{x}) \equiv 0$;

root \leftarrow {data = Data, model = f , leftChild = None, rightChild = None};

BoostTree-SVR \leftarrow split(root);

split(node){

Let t be the index of the current node;

$\{\mathbf{x}_n, y_n\}_{n \in I} \leftarrow$ node.data;

$F_I(\mathbf{x}) \leftarrow \sum_{m \in \text{Path}(t)} \text{node}_m.\text{model}(\mathbf{x})$;

Randomly select N_{\min}^L and N_{\min}^R from Pool_{MinSamplesLeaf};

$\delta_{\mathcal{L}}^{\max} = 0$;

Randomly select \sqrt{D} features to create a feature subset \mathcal{F} ;

for $d \in \mathcal{F}$ **do**

$S_d = \{x_{n,d} | n \in I\}$;

for s in S_d **do**

$I_L = \{n | x_{n,d} \leq s, n \in I\}$;

$I_R = \{n | x_{n,d} > s, n \in I\}$;

if $|I_L| \geq N_{\min}^L$ and $|I_R| \geq N_{\min}^R$ **then**

 Randomly select C_L and C_R from Pool_C;

 Calculate $\delta_{\mathcal{L}}$ in (12);

if $\delta_{\mathcal{L}} > \delta_{\mathcal{L}}^{\max}$ **then**

$\delta_{\mathcal{L}}^{\max} = \delta_{\mathcal{L}}$, $d^* = d$, $s^* = s$, $C_L^* = C_L$, $C_R^* = C_R$;

end

end

end

end

if $\delta_{\mathcal{L}}^{\max} > 0$ **then**

 Randomly select ϵ_L and ϵ_R from Pool _{ϵ} ;

$I_L^* = \{n | x_{n,d^*} \leq s^*, n \in I\}$;

$I_R^* = \{n | x_{n,d^*} > s^*, n \in I\}$;

$f_L^* = \text{FitModelSVR}(\{\mathbf{x}_n, y_n\}_{n \in I_L^*}, F_I, C_L^*, \epsilon_L)$;

$f_R^* = \text{FitModelSVR}(\{\mathbf{x}_n, y_n\}_{n \in I_R^*}, F_I, C_R^*, \epsilon_R)$;

 Calculate $\delta_{\mathcal{L}}$ in (23);

if $\delta_{\mathcal{L}} > 0$ **then**

 node.leftChild = {data = $\{\mathbf{x}_n, y_n\}_{n \in I_L^*}$, model = f_L^* , leftChild = None, rightChild = None};

 node.rightChild = {data = $\{\mathbf{x}_n, y_n\}_{n \in I_R^*}$, model = f_R^* , leftChild = None, rightChild = None};

 NumLeaf = NumLeaf + 1;

end

if MaxNumLeaf is not supplied, or NumLeaf \leq MaxNumLeaf **then**

 For each leaf, the cross entropy loss (classification) or RMSE (regression) multiplied by the number of samples at that leaf is used as the impurity score;

 Identify node*, the leaf node with the highest impurity score;

 split(node*);

end

end

Supplementary Algorithm 11: FitModelSVR for regression.

Input: $\{\mathbf{x}_n, y_n\}_{n \in I}$, sample set of the current node;
 F_I , ensemble of the models along the path from the root node to the parent node of the current node;

C , the regularization parameter of SVR;

ϵ , the slack variable of SVR.

Output: The SVR model f_m for the current node.

$\tilde{y}_n = y_n - F_I(\mathbf{x}_n)$, $n \in I$;

Fit $f_m = \text{SVR}(\{\mathbf{x}_n, \tilde{y}_n\}_{n \in I}, C, \epsilon)$ using SVR.

Supplementary Algorithm 12: BoostForest-SVR for regression.

Input: Data = $\{\mathbf{x}_n, y_n\}_{n=1}^N$, N training samples, where $\mathbf{x}_n \in \mathbb{R}^D$;

NumEstimators, the number of trees;

Pool_{MinSamplesLeaf}, candidate value pool of the minimum number of samples at a leaf;

Pool _{ϵ} , candidate value pool of the slack variable;

Pool _{C} , candidate value pool of the regularization parameter.

Output: A BoostForest-SVR

BoostForest-SVR = {};

for $i = 1 : \text{NumEstimators}$ **do**

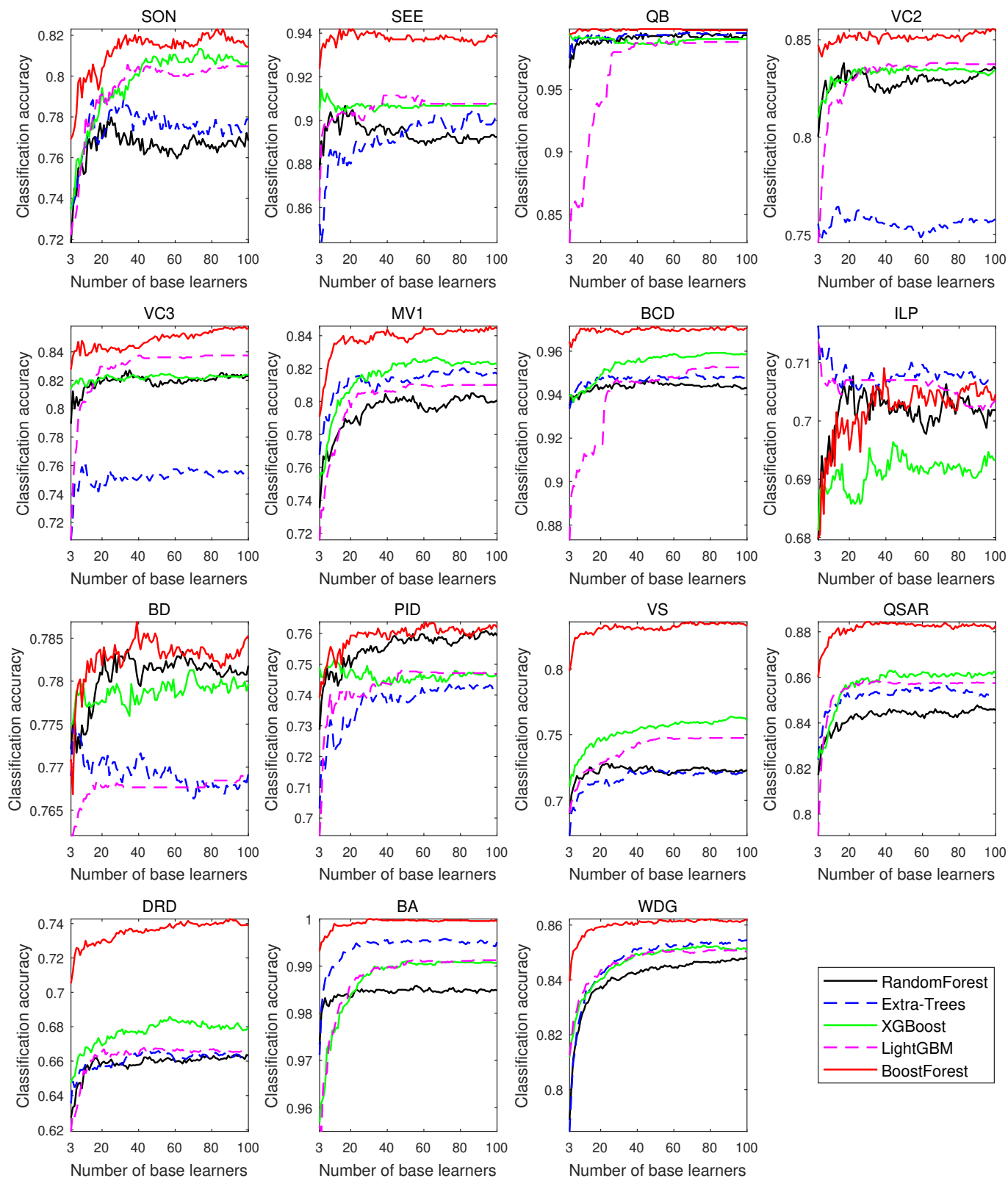
 Bootstrap Data' from Data;

 Train BoostTree-SVR _{i} on Data' using Supplementary Algorithm 10;

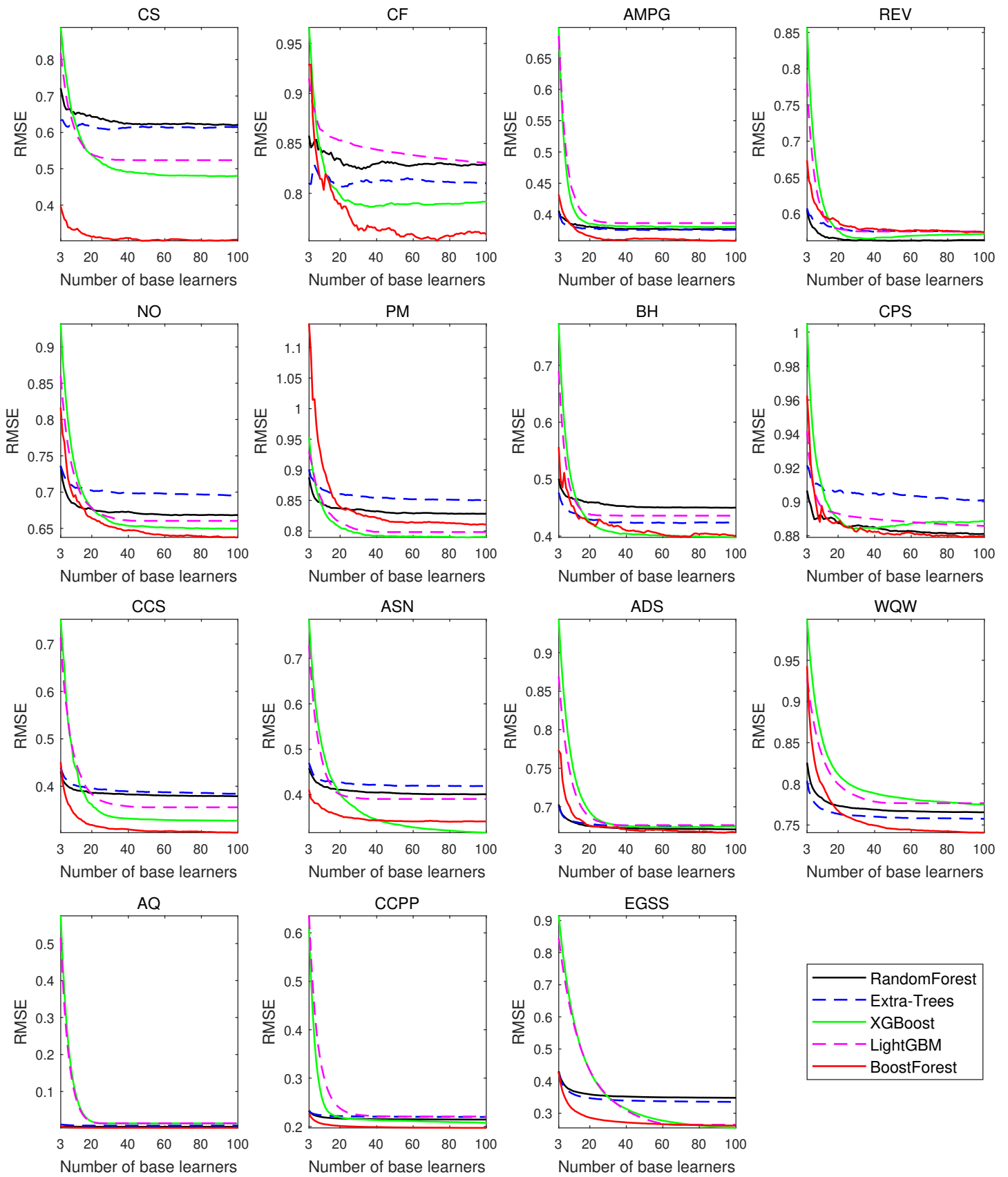
 Add BoostTree-SVR _{i} to BoostForest-SVR;

end

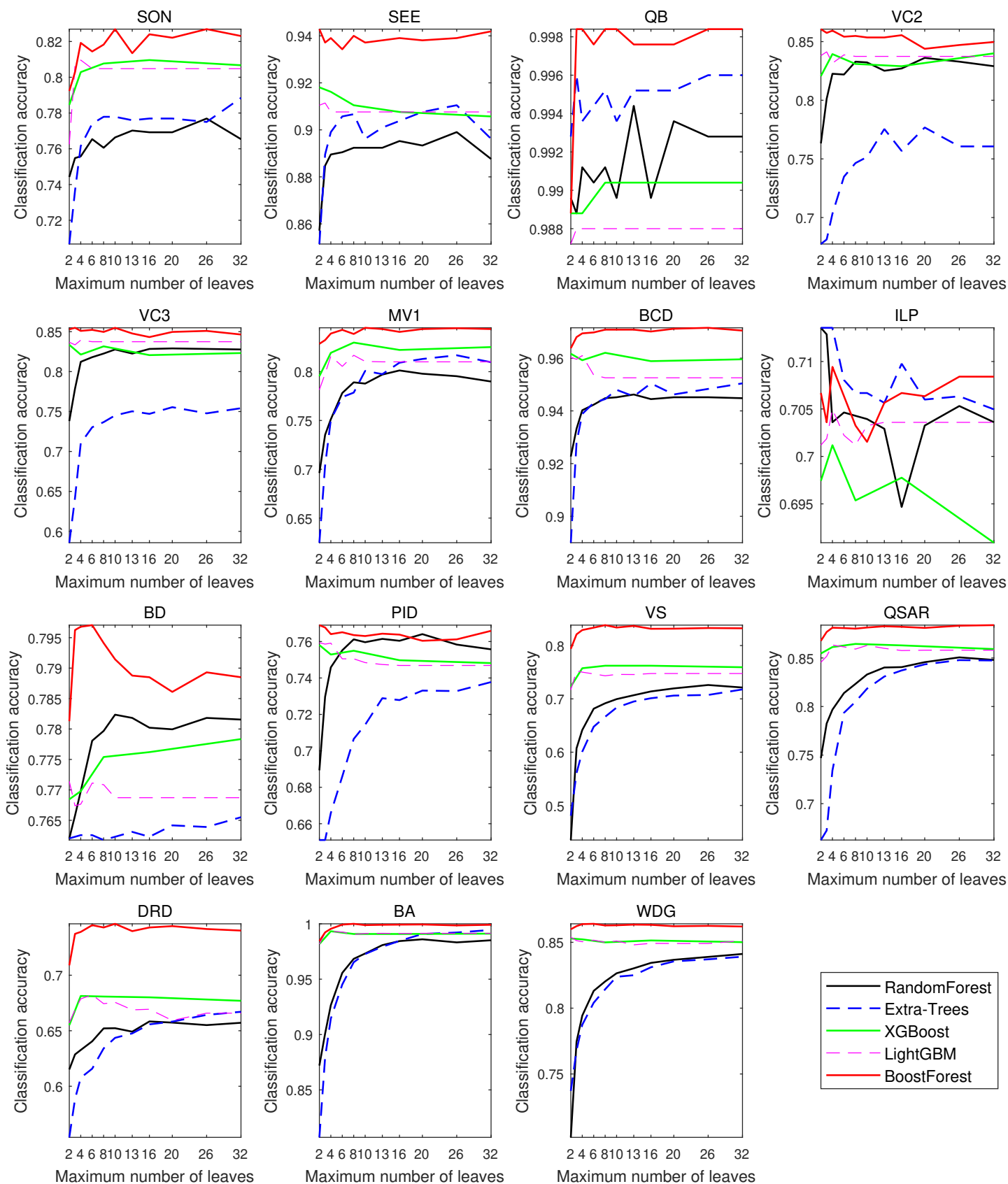
2 Supplementary Figures



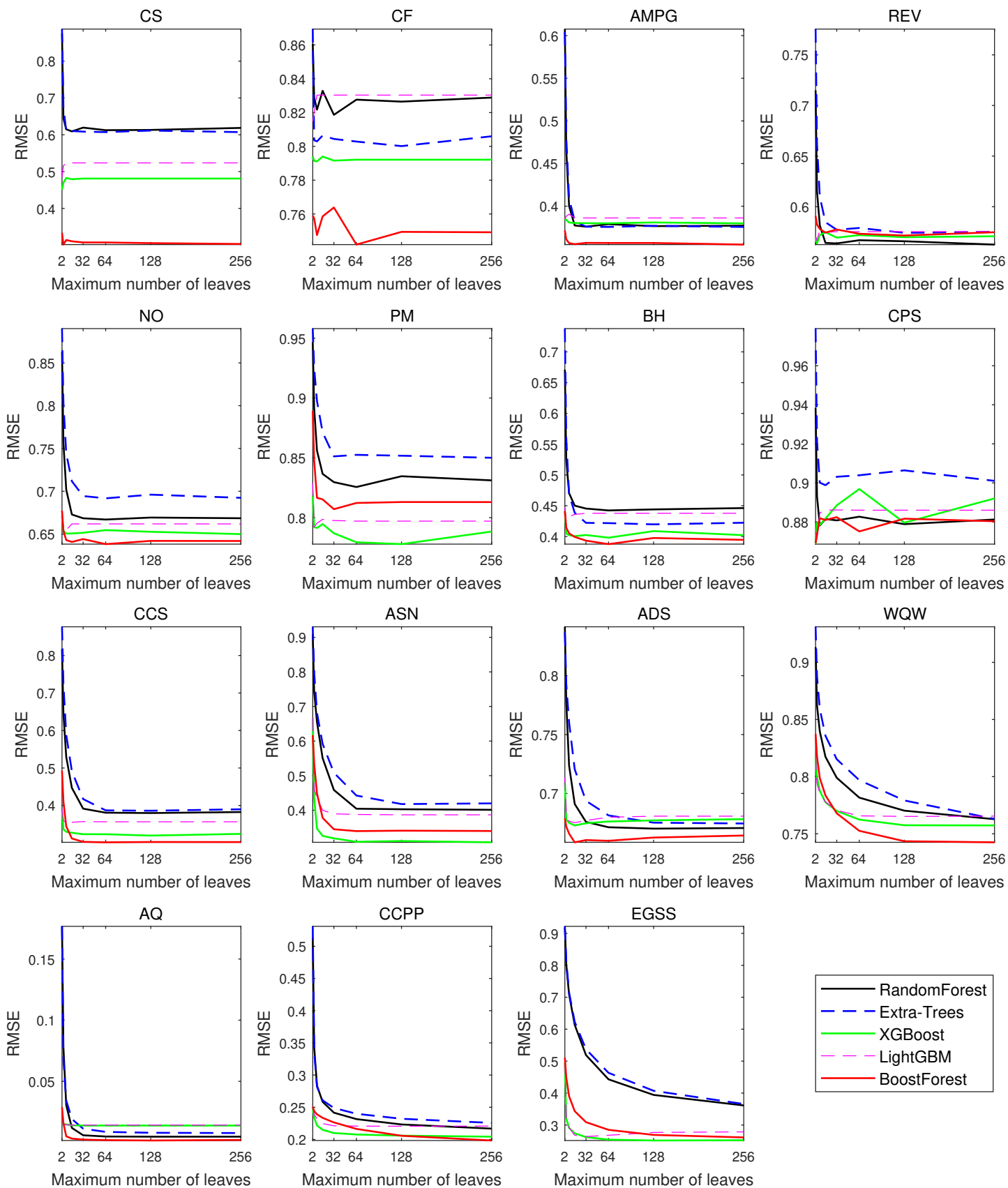
Supplementary Figure 1: Average classification accuracies on the 15 classification datasets, with different number of base learners.



Supplementary Figure 2: Average RMSEs on the 15 regression datasets, with different number of base learners.



Supplementary Figure 3: Average classification accuracies on the 15 classification datasets, with different maximum number of leaves.



Supplementary Figure 4: Average RMSEs on the 15 regression datasets, with different maximum number of leaves.

3 Supplementary Tables

This section introduces the statistics of the 34 datasets used in our experiments, and parameters of the ensemble algorithms.

Supplementary Table 1: The 34 datasets used in our experiments.

Task	Dataset	Abbreviation	#Samples	#Features	#Classes
Classification	Sonar	SON	208	60	2
	Seeds	SEE	210	7	3
	Qualitative Bankruptcy	QB	250	6	2
	Vertebral Column-2	VC2	310	6	2
	Vertebral Column-3	VC3	310	6	3
	Musk Version 1	MV1	476	166	2
	Breast Cancer Diagnosis	BCD	569	30	2
	Indian Liver Patient	ILP	583	11	2
	Blood Donations	BD	748	4	2
	Pima Indians Diabetes	PID	768	8	2
	Vehicle Silhouettes	VS	846	18	4
	QSAR Biodegradation	QSAR	1,055	41	2
	Diabetic Retinopathy Debrecen	DRD	1,151	19	2
	Banknote Authentication	BA	1,372	4	2
	Waveform Database Generator	WDG	5,000	21	3
	Letter Recognition	LR	20,000	16	26
MNIST	MNIST	70,000	784	10	
Regression	Concrete Slump	CS	103	7	
	Concrete Flow	CF	103	7	
	autoMPG	AMPG	392	7	
	Real Estate Valuation	REV	414	6	
	NO2	NO	500	7	
	PM10	PM	500	7	
	Boston Housing	BH	506	13	
	CPS	CPS	534	11	
	Concrete Compressive Strength	CCS	1,030	8	
	Airfoil Self-Noise	ASN	1,503	5	
	Abalone Data Set	ADS	4,177	10	
	Wine Quality White	WQW	4,898	11	
	Air Quality	AQ	9,357	8	
	Combined Cycle Power Plant	CCPP	9,568	4	
	Electrical Grid Stability Simulated	EGSS	10,000	12	
	Protein Tertiary Structure	PTS	45,730	9	
	Relative Location of CT	RLCT	53,500	384	

Supplementary Table 2: Parameter settings for RandomForest, Extra-Trees, XGBoost and LightGBM. More details about the parameters for RandomForest and Extra-Trees can be found at <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>. More details about the parameters for XGBoost and LightGBM can be found at <https://xgboost.readthedocs.io/en/latest/parameter.html>.

Parameter	RandomForest Extra-Trees	XGBoost LightGBM
If the number of samples in the dataset is < 10,000		
<i>nEstimators</i> (number of trees)	{80, 100, 120}	{80, 100, 120}
<i>minSamplesLeaf</i> (minimum number of samples at a leaf node)	[5, 15]	
η (learning rate)		{0.01, 0.1, 0.2}
γ (minimum loss reduction required to further split a leaf node)		{0.25, 0.5, 1}
α (L1 regularization coefficient on weights)		{0.1, 0.5, 1}
λ (L2 regularization coefficient on weights)		{0.1, 0.5, 1}
<i>subsample</i> (subsample ratio of the training samples)		{0.8, 1}
<i>colsample</i> (subsample ratio of columns)		{0.8, 1}
<i>maxDepth</i> (maximum depth of a tree)		{4, 5, 6}
<i>minChildWeight</i> (minimum sum of sample weight needed in a child)		{0.5, 1, 3} (classification) {5, 10, 15} (regression)
If the number of samples in the dataset is \geq 10,000		
<i>nEstimators</i> (number of trees)	{80, 100, 120}	{1000}
<i>minSamplesLeaf</i> (minimum number of samples at a leaf node)	{5, 10, 15}	
η (learning rate)		{0.01}
γ (minimum loss reduction required to further split a leaf node)		default
α (L1 regularization coefficient on weights)		{0.1, 1}
λ (L2 regularization coefficient on weights)		{0.1, 1}
<i>subsample</i> (subsample ratio of the training samples)		default
<i>colsample</i> (subsample ratio of columns)		default
<i>maxDepth</i> (maximum depth of a tree)		{5, 6}
<i>minChildWeight</i> (minimum sum of sample weight needed in a child)		{5}
<i>early_stopping_rounds</i>		{50}

4 Supplementary Experiments

This section presents additional experimental results on using other base learners in BoostForest, and comparisons with MultiBoosting².

Use other base learners in BoostForest. We studied if the strategy that BoostForest uses to combine multiple BoostTrees (data replica by bootstrapping, and random parameters selection from a parameter pool) can also be extended to other tree models, i.e., whether we can still achieve good ensemble learning performance when BoostTree is replaced by another base learner, e.g., LMT or Model Tree. The resulting forests are denoted as LMForest and ModelForest, respectively.

LMT³ extends model tree from regression to classification by integrating logistic regression into the tree model. It uses Stepwise Model Tree Induction⁴ to construct the tree. The final logistic regression model at a leaf consists of all linear models at the nodes in the path from the root to that leaf. SimpleLogistic³ (a variant of LogitBoost) is used to incrementally refine the linear logistic model. In each iteration, instead of using all features to perform linear regression, SimpleLogistic uses only one feature to train the model. In this way, only relevant features are selected, and the risk of over-fitting is reduced.

LMT¹ and M5P² implementations in Weka were used in our experiments. We stopped SimpleLogistic training if the minimum error on the validation set had not changed for 20 iterations. For each M5P model tree in ModelForest, the number of samples at a leaf node should not be too small, so the minimum number of samples at its leaf was randomly sampled from the parameter pool {10, 20, 30}. To improve its stability, RR models with regularization coefficient $\lambda = 0.001$ were trained at every leaf.

Supplementary Table 3 compares the performances of LMT with LMForest on the 15 classification datasets. Supplementary Table 4 compares the performances of M5P with ModelForest on the 15 regression datasets. All results were averaged over five repeats of 2-fold cross-validations. LMForest outperformed LMT on 10 of the 15 datasets, among which three were statistically significant. ModelForest outperformed M5P on 14 of the 15 datasets, among which 12 were statistically significant.

Comparison with MultiBoosting. Considering that BoostForest is an ensemble method integrating Bagging and Boosting, it is necessary to compare its performance with other algorithms using similar strategies.

MultiBoosting² is a representative method that combines waggings and Adaboost⁵ to reduce both the bias and the variance. It adopts AdaBoost as the sub-committee classifiers' producer and decision tree (CART in our experiment) as the base learner to decide committees. By wagging a set of sub-committees, MultiBoosting achieves a lower error than AdaBoost. Its model complexity can be controlled by the number of estimators, the maximum depth of each tree, and MinSampleLeaf. We set their candidate values to {80, 100, 120}, {5, 7, 10} and {5, 10, 15}, respectively. The best combination of parameters was determined by grid-search from its candidate values using inner 5-fold cross-validation.

There are three main differences between MultiBoosting and BoostForest:

1. MultiBoosting uses AdaBoost to train multiple tree models for boosting, and can only be used for classification, whereas BoostForest uses gradient boosting to train the base models, which can handle both classification and regression. BoostForest optimizes the loss function using gradient descent, and only one tree model is generated for boosting.
2. MultiBoosting first uses AdaBoost to train tree models to reduce the bias, and then wagging to train tree models to reduce the variance, when the training error is relatively small. BoostForest generates only one

¹<http://weka.sourceforge.net/doc.dev/weka/classifiers/trees/LMT.html>

²<http://weka.sourceforge.net/doc.dev/weka/classifiers/trees/M5P.html>

Supplementary Table 3: Mean and standard deviation (in parenthesis) of the classification accuracy, when LMT is used to replace BoostTree in BoostForest. The best performance is marked in bold. • indicates statistically significant win for LMForest.

Dataset	LMT	LMForest
SON	0.762 (0.043)	0.796 (0.030)
SEE	0.941 (0.020)	0.933 (0.020)
QB	0.986 (0.013)	0.994 (0.012)
VC2	0.847 (0.027)	0.846 (0.021)
VC3	0.857 (0.023)	0.850 (0.026)
MV1	0.794 (0.026)•	0.822 (0.026)
BCD	0.967 (0.007)	0.971 (0.007)
ILP	0.706 (0.016)	0.700 (0.021)
BD	0.778 (0.016)	0.787 (0.009)
PID	0.761 (0.019)	0.762 (0.010)
VS	0.792 (0.017)	0.800 (0.012)
QSAR	0.862 (0.010)•	0.871 (0.009)
DRD	0.702 (0.016)	0.707 (0.013)
BA	0.999 (0.002)•	0.999 (0.002)
WDG	0.866 (0.005)	0.865 (0.007)
Average ACC	0.841	0.847
Average rank	1.667	1.333

Supplementary Table 4: Mean and standard deviation (in parenthesis) of the regression RMSE, when M5P is used to replace BoostTree in BoostForest. The best performance is marked in bold. • indicates statistically significant win for ModelForest.

Dataset	M5P	ModelForest
CS	0.384 (0.039)•	0.330 (0.030)
CF	0.782 (0.084)	0.763 (0.076)
AMPG	0.379 (0.020)•	0.365 (0.015)
REV	0.594 (0.067)•	0.568 (0.069)
NO	0.697 (0.030)•	0.652 (0.038)
PM	0.908 (0.075)•	0.833 (0.082)
BH	0.459 (0.077)•	0.389 (0.071)
CPS	0.871 (0.132)	0.913 (0.122)
CCS	0.398 (0.020)•	0.336 (0.017)
ASN	0.495 (0.060)•	0.342 (0.016)
ADS	0.672 (0.029)	0.667 (0.024)
WQW	0.832 (0.021)•	0.748 (0.020)
AQ	0.089 (0.075)•	0.004 (0.001)
CCPP	0.239 (0.003)•	0.210 (0.003)
EGSS	0.373 (0.006)•	0.281 (0.005)
Average RMSE	0.545	0.493
Average rank	1.933	1.067

Supplementary Table 5: Mean and standard deviation (in parenthesis) of the classification accuracy, when MultiBoosting is compared with BoostForest. The best performance is marked in bold. • indicates statistically significant win for BoostForest.

Dataset	MultiBoosting	BoostForest
SON	0.804 (0.026)•	0.822 (0.026)
SEE	0.919 (0.029)•	0.940 (0.022)
QB	0.991 (0.009)•	0.998 (0.005)
VC2	0.839 (0.025)•	0.853 (0.026)
VC3	0.841 (0.039)	0.850 (0.027)
MV1	0.837 (0.030)	0.843 (0.025)
BCD	0.966 (0.008)	0.972 (0.009)
ILP	0.697 (0.020)	0.709 (0.019)
BD	0.766 (0.018)•	0.785 (0.007)
PID	0.748 (0.008)•	0.762 (0.011)
VS	0.755 (0.020)•	0.830 (0.016)
QSAR	0.868 (0.007)•	0.882 (0.006)
DRD	0.679 (0.014)•	0.740 (0.008)
BA	0.997 (0.003)	0.998 (0.003)
WDG	0.845 (0.007)•	0.861 (0.005)
Average accuracy	0.837	0.856
Average rank	2.000	1.000

tree for boosting, and uses Random Forest to reduce the variance, which is less likely to be trapped in a local minimum.

3. MultiBoosting uses traditional decision trees as the base learners, whereas BoostForest uses BoostTrees.

Table 5 compares the performances of BoostForest and MultiBoosting on the 15 classification datasets. BoostForest achieved higher accuracy on all datasets, among which 10 were statistically significant.

References

1. Friedman, J., Hastie, T. & Tibshirani, R. Additive logistic regression: A statistical view of boosting. *Annals of Statistics* **28**, 337–407 (2000).
2. Webb, G. I. Multiboosting: A technique for combining boosting and wagging. *Machine Learning* **40**, 159–196 (2000).
3. Landwehr, N., Hall, M. & Frank, E. Logistic model trees. *Machine Learning* **59**, 161–205 (2005).

4. Malerba, D., Appice, A., Ceci, M. & Monopoli, M. Trading-off local versus global effects of regression nodes in model trees. In *Proc. 13th Int'l Symposium on Foundations of Intelligent Systems*, 393–402 (Lyon, France, 2002).
5. Freund, Y. & Schapire, R. E. Experiments with a new boosting algorithm. In *Proc. 13th Int'l Conf. on Machine Learning*, 148–156 (Bari, Italy, 1996).