

# Technical Explanation of EII Prediction Model Training and Testing

Your Name

January 30, 2025

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Training Script</b>	<b>2</b>
2.1	Overview . . . . .	2
2.2	Complete Training Script Code . . . . .	2
2.3	Technical Explanation of Training Steps . . . . .	7
2.3.1	Random Seed and Data Loading . . . . .	7
2.3.2	Model Containers & Cross-Validation Setup . . . . .	7
2.3.3	Feature Extraction and Validity Checks . . . . .	7
2.3.4	Median and MAD Normalization . . . . .	7
2.3.5	Training Multiple Neural Networks . . . . .	7
2.3.6	Random Forest Training . . . . .	8
2.3.7	Prediction and Ensemble Combining . . . . .	8
2.3.8	Result Storage and Visualization . . . . .	8
<b>3</b>	<b>Testing Script</b>	<b>9</b>
3.1	Overview . . . . .	9
3.2	Complete Testing Script Code . . . . .	9
3.3	Technical Explanation of Testing Steps . . . . .	15
3.3.1	Data and Model Loading . . . . .	15
3.3.2	User Input and Matching Dataset . . . . .	15
3.3.3	Preprocessing and Prediction . . . . .	15
3.3.4	Accuracy Calculation . . . . .	16
<b>4</b>	<b>Reasoning Behind Algorithmic Choices</b>	<b>16</b>
4.1	Neural Networks . . . . .	16
4.2	Random Forest . . . . .	16
4.3	Ensemble Weighting . . . . .	17
4.4	Cross-Validation and Robust Scaling . . . . .	17
<b>5</b>	<b>Conclusion</b>	<b>17</b>

# 1 Introduction

This document provides a **technical, in-depth explanation** of the two scripts used for predicting the Energy Intensity Index (EII) from multiple datasets. We use an ensemble approach consisting of *multiple feed-forward neural networks* (with an averaging mechanism) and a *Random Forest* (RF) model, combined through a weighted ensemble method.

The document is organized into two main parts:

- **Training Script:** trains an ensemble of models (Neural Networks + Random Forest) via 5-fold cross-validation on four different datasets.
- **Testing Script:** loads the best-trained models, asks the user for new input, and produces predictions (with accuracy calculations).

**Note:** All dataset files referenced (e.g., MAF PU-EII.xlsx) are assumed to be in the same directory as the scripts.

## 2 Training Script

### 2.1 Overview

The training script:

1. Loads four datasets into a MATLAB `struct`.
2. For each dataset, performs a 5-fold cross-validation.
3. Trains 10 neural networks (ensemble) and 1 Random Forest per fold.
4. Combines (averages) neural network predictions and linearly weights the final ensemble with the Random Forest prediction.
5. Selects the best ensemble across all folds based on MSE.
6. Saves the best models and preprocessing parameters to a `.mat` file for each dataset.

### 2.2 Complete Training Script Code

```
1 % Enhanced Debuggable Training Script with Multiple Neural Networks
2
3 % Set random seed for reproducibility
4 rng(42);
5
6 % Load datasets
7 datasets = struct();
8 datasets.MAF = readtable('MAF PU-EII.xlsx', 'VariableNamingRule', '
    preserve');
9 datasets.SRC = readtable('SR1 PU-EII.xlsx', 'VariableNamingRule', '
    preserve');
10 datasets.AP = readtable('AP EII-PU.xlsx', 'VariableNamingRule', '
    preserve');
```

```

11 datasets.SRIP = readtable('SR2 EII-PU.xlsx', 'VariableNamingRule', '
    preserve');
12
13 model_results = struct();
14 k = 5; % Number of folds for cross-validation
15 num_networks = 10; % Number of neural networks to train
16
17 % Iterate through all datasets
18 fields = fieldnames(datasets);
19 for i = 1:length(fields)
20     dataset_name = fields{i};
21     disp(['Training models for ', dataset_name, ' dataset...']);
22
23     % Extract X and Y
24     data = datasets.(dataset_name);
25     X_data = data{:, 1:end-1};
26     Y_data = data{:, end};
27
28     % Input validation
29     if size(X_data, 2) < 1
30         error('Dataset %s has no features.', dataset_name);
31     end
32     if ~isnumeric(Y_data) || all(isnan(Y_data))
33         error('Target variable in dataset %s is invalid.', dataset_name)
34         ;
35     end
36
37     % Cross-validation partition
38     cv = cvpartition(size(X_data, 1), 'Kfold', k);
39     mse_values = zeros(k, 1);
40     mae_values = zeros(k, 1);
41
42     % Initialize results for best ensemble
43     best_mse = inf;
44     best_ensemble_weight = 0;
45     best_nns = cell(1, num_networks);
46     best_rf_model = [];
47
48     for fold_idx = 1:k
49         disp(['Processing fold ', num2str(fold_idx), ' of ', num2str(k),
50             '...']);
51         train_idx = training(cv, fold_idx);
52         val_idx = test(cv, fold_idx);
53
54         X_train = X_data(train_idx, :);
55         Y_train = Y_data(train_idx);
56         X_val = X_data(val_idx, :);
57         Y_val = Y_data(val_idx);
58
59         % Check for empty data
60         if isempty(X_train) || isempty(Y_train) || isempty(X_val) ||
61             isempty(Y_val)
62             disp('Warning: Training or validation data is empty.
63                 Skipping fold.');
```

```

64     X_median = median(X_train, 1);
65     X_mad = mad(X_train, 1);
66     X_train_std = (X_train - X_median) ./ (X_mad + eps);
67     X_val_std = (X_val - X_median) ./ (X_mad + eps);
68
69     % Train multiple Neural Networks in parallel
70     fold_nns = cell(1, num_networks);
71     parfor nn_idx = 1:num_networks
72         disp(['Training neural network ', num2str(nn_idx), '...']);
73         net = feedforwardnet([100 50]);
74         net.trainFcn = 'trainlm';
75         net.divideFcn = 'dividerand';
76         net.divideParam.trainRatio = 1;
77         net.divideParam.testRatio = 0;
78         net.performFcn = 'mse';
79         net.trainParam.epochs = 25;
80
81         [x_norm, ps_input] = mapminmax(X_train_std', -1, 1);
82         [t_norm, ps_target] = mapminmax(Y_train', -1, 1);
83
84         try
85             [trained_net, ~] = train(net, x_norm, t_norm);
86             trained_net.userdata.ps_input = ps_input;
87             trained_net.userdata.ps_target = ps_target;
88             fold_nns{nn_idx} = trained_net;
89         catch nn_err
90             disp(['Error training neural network ', num2str(nn_idx),
91                 ': ', nn_err.message]);
92             fold_nns{nn_idx} = [];
93         end
94     end
95
96     % Train Random Forest with adjusted hyperparameters
97     disp('Training Random Forest...');
98     try
99         num_features = size(X_train_std, 2);
100         rf_model = TreeBagger(500, X_train_std, Y_train, ...
101             'Method', 'regression', ...
102             'MinLeafSize', 5, ...
103             'NumPredictorsToSample', max(1, floor(sqrt(num_features)
104             )), ...
105             'OOBPrediction', 'on');
106     catch rf_err
107         disp(['Error training Random Forest: ', rf_err.message]);
108         continue;
109     end
110
111     % Predictions
112     try
113         % Neural Network Predictions
114         nn_predictions = zeros(length(Y_val), num_networks);
115         valid_nn_count = 0;
116         for nn_idx = 1:num_networks
117             if isempty(fold_nns{nn_idx})
118                 nn_predictions(:, nn_idx) = NaN;
119                 continue;
120             end

```

```

119         x_val_norm = mapminmax('apply', X_val_std', fold_nns{
120             nn_idx}.userdata.ps_input);
121         nn_predictions(:, nn_idx) = mapminmax('reverse',
122             fold_nns{nn_idx}(x_val_norm), fold_nns{nn_idx}.
123             userdata.ps_target)';
124         valid_nn_count = valid_nn_count + 1;
125     end
126     nn_avg_prediction = nanmean(nn_predictions, 2);
127
128     % Random Forest Predictions
129     rf_prediction = predict(rf_model, X_val_std);
130
131     if all(isnan(nn_avg_prediction))
132         % If all NN predictions failed, fallback to RF
133         warning('All neural networks failed. Using Random Forest
134             only.');
```

```

135         ensemble_weight = 0;
136         ensemble_prediction = rf_prediction;
137     else
138         % Calculate ensemble weights
139         mse_nn_avg = mean((Y_val - nn_avg_prediction).^2);
140         mse_rf = mean((Y_val - rf_prediction).^2);
141         ensemble_weight = mse_rf / (mse_rf + mse_nn_avg);
142         ensemble_prediction = ensemble_weight *
143             nn_avg_prediction + (1 - ensemble_weight) *
144             rf_prediction;
145     end
146
147     % Compute metrics
148     mse_values(fold_idx) = mean((Y_val - ensemble_prediction)
149         .^2);
150     mae_values(fold_idx) = mean(abs(Y_val - ensemble_prediction)
151         );
152
153     % Update best ensemble
154     if mse_values(fold_idx) < best_mse
155         best_mse = mse_values(fold_idx);
156         best_ensemble_weight = ensemble_weight;
157         best_nns = fold_nns;
158         best_rf_model = rf_model;
159     end
160     catch pred_err
161         disp(['Error during prediction: ', pred_err.message]);
162         continue;
163     end
164 end
165
166 % Save best model and preprocessing parameters
167 model_results.(dataset_name).NeuralNetworks = best_nns;
168 model_results.(dataset_name).RandomForest = best_rf_model;
169 model_results.(dataset_name).Preprocessing.mean = X_median;
170 model_results.(dataset_name).Preprocessing.std = X_mad + eps;
171 model_results.(dataset_name).BestEnsembleWeight =
172     best_ensemble_weight;
173
174 % Display final metrics
175 valid_mse = mse_values(~isnan(mse_values));
176 valid_mae = mae_values(~isnan(mae_values));

```

```

168     if ~isempty(valid_mse)
169
170     else
171         disp([dataset_name, ': No valid folds for metrics.']);
172     end
173
174     % Visualization
175     figure;
176     plot(1:k, mse_values, '-o', 'DisplayName', 'MSE');
177     hold on;
178     plot(1:k, mae_values, '-x', 'DisplayName', 'MAE');
179     xlabel('Fold');
180     ylabel('Error');
181     legend;
182     grid on;
183     title(['Performance Across Folds: ', dataset_name]);
184 end
185
186 % Save models
187 disp('Saving models...');
188 for dataset_name = fieldnames(model_results)'
189     filename = ['model_results_', dataset_name{1}, '.mat'];
190     current_model = model_results.(dataset_name{1});
191     save(filename, 'current_model', '-v7.3');
192     disp(['Saved ', dataset_name{1}, ' to ', filename]);
193 end
194
195 disp('Training completed successfully.');
```

Listing 1: Training Script

## 2.3 Technical Explanation of Training Steps

### 2.3.1 Random Seed and Data Loading

- **Random Seed** (`rng(42)`): Ensures reproducible results across multiple runs.
- **Data Loading**: Each dataset is stored in a field of `datasets`; we keep consistent variable naming rules with `VariableNamingRule`.

### 2.3.2 Model Containers & Cross-Validation Setup

- **model\_results**: Stores the final trained models and normalization parameters for each dataset.
- **5-Fold CV**: Partitions the data for robust performance estimation and to reduce overfitting.
- **10 Neural Networks Per Fold**: An ensemble approach typically offers higher resilience to variance in initialization.

### 2.3.3 Feature Extraction and Validity Checks

- We separate the last column as `Y_data` (target) and all preceding columns as `X_data`.
- Checks ensure that the dataset has at least one feature and that the target is valid (numeric and not all NaNs).

### 2.3.4 Median and MAD Normalization

- We compute the **Median** and the **MAD (Median Absolute Deviation)** for each feature.
- Normalization formula:

$$X_{\text{train\_std}} = \frac{X_{\text{train}} - \text{median}(X_{\text{train}})}{\text{MAD}(X_{\text{train}}) + \epsilon}$$

- $\epsilon$  avoids division by zero.
- This *robust scaling* is less sensitive to outliers than simple z-score normalization.

### 2.3.5 Training Multiple Neural Networks

- `feedforwardnet([100 50])` creates a neural network with two hidden layers, of sizes 100 and 50.
- `trainFcn = 'trainlm'` uses the Levenberg-Marquardt backpropagation.
- `parfor`: Parallelization to speed up training 10 neural networks.
- Preprocessing with `mapminmax` normalizes  $[X'_{\text{train\_std}}]$  to a range  $[-1, 1]$ .

### 2.3.6 Random Forest Training

- `TreeBagger(500, ...)` creates an ensemble of 500 regression trees.
- `MinLeafSize = 5` and `NumPredictorsToSample = floor(sqrt(num_features))` are tuned hyperparameters to balance performance and interpretability.
- `'OOBPrediction','on'` enables out-of-bag prediction for internal validation and tuning.

### 2.3.7 Prediction and Ensemble Combining

- We get **Neural Network Predictions** for each of the 10 networks and then average them.
- We get **Random Forest Predictions**.
- We compute an **ensemble weight**:

$$w = \frac{\text{MSE}_{\text{RF}}}{\text{MSE}_{\text{RF}} + \text{MSE}_{\text{NN}}},$$

and use that to linearly combine  $\text{NN}_{\text{avg}}$  and RF.

- This gives a **Final Ensemble Prediction**.
- Each fold's performance is measured by MSE and MAE.
- The best model per dataset is stored along with its normalization parameters and ensemble weight.

### 2.3.8 Result Storage and Visualization

- We store the final best models in `model_results_*.mat`.
- We also plot MSE and MAE across the 5 folds to visualize performance consistency.



## 3 Testing Script

### 3.1 Overview

The testing script:

1. Loads the same four datasets (for reference and to check input dimension).
2. Loads each dataset's .mat file (model\_results\_X.mat).
3. Prompts the user to input a feature vector corresponding to one of the datasets.
4. Identifies which dataset the input corresponds to (based on vector length).
5. Applies the saved preprocessing (median, MAD) to the input.
6. Makes predictions from the (i) Neural Networks, (ii) Random Forest, and (iii) Weighted Ensemble.
7. Asks for the actual EII to compute accuracy percentages.
8. Outputs detailed results and identifies which model performed best.

### 3.2 Complete Testing Script Code

```
1 try
2     % Load datasets
3     disp('Loading datasets...');
4     datasets.MAF = readtable('MAF PU-EII.xlsx', 'VariableNamingRule', '
    preserve');
5     datasets.SRC = readtable('SR1 PU-EII.xlsx', 'VariableNamingRule', '
    preserve');
6     datasets.AP = readtable('AP EII-PU.xlsx', 'VariableNamingRule', '
    preserve');
7     datasets.SRIP = readtable('SR2 EII-PU.xlsx', 'VariableNamingRule', '
    preserve');
8     disp('Datasets loaded successfully.');
```

```
9
10    % Load individual model files
11    disp('Loading models...');
12    model_names = {'MAF', 'AP', 'SRC', 'SRIP'};
13    for name = model_names
14        filename = ['model_results_' name{1} '.mat'];
15        if exist(filename, 'file')
16            loaded = load(filename);
17            model_results.(name{1}) = loaded.current_model;
18            disp(['Loaded model: ' name{1}]);
19        else
20            error(['Could not find model file: ' filename]);
21        end
22    end
23    disp('All models loaded successfully.');
```

```
24
25    % Validate preprocessing fields
26    for name = model_names
27        if ~isfield(model_results.(name{1}), 'Preprocessing')
```

```

28         warning('Preprocessing parameters missing for dataset: %s.
29             Assigning default values.', name{1});
30         model_results.(name{1}).Preprocessing.mean = zeros(1, size(
31             datasets.(name{1}), 2) - 1);
32         model_results.(name{1}).Preprocessing.std = ones(1, size(
33             datasets.(name{1}), 2) - 1);
34     end
35     if ~isfield(model_results.(name{1}), 'NeuralNetworks')
36         warning('NeuralNetworks field missing for dataset: %s.
37             Assigning default empty neural networks.', name{1});
38         model_results.(name{1}).NeuralNetworks = {};
39     end
40 end
41 catch load_err
42     disp('Error loading required files:');
43     disp(['Error message: ' load_err.message]);
44     if exist('filename', 'var')
45         disp(['Failed while processing: ' filename]);
46     end
47     return;
48 end
49 while true
50     try
51         % Ask for user input
52         prompt = 'Enter the values for X variables as a vector (e.g.,
53             [100, 90, 72, 74, 90]): ';
54         X_input = input(prompt);
55
56         if isempty(X_input)
57             disp('Input is empty. Please try again.');

```

```

80
81 % Use the matched dataset
82 disp(['Using ', matched_dataset, ' dataset for prediction...']);
83
84 % Get original variable names from the dataset
85 original_vars = datasets.(matched_dataset).Properties.
    VariableNames(1:end-1);
86
87 % Standardize the input using saved preprocessing parameters
88 X_mean = model_results.(matched_dataset).Preprocessing.mean;
89 X_std = model_results.(matched_dataset).Preprocessing.std;
90
91 % Validate input dimensions
92 if length(X_input) ~= length(X_mean)
93     error('Input vector length (%d) does not match expected
        length (%d).', ...
        length(X_input), length(X_mean));
94
95 end
96
97 % Standardize input
98 X_input_std = (X_input - X_mean) ./ X_std;
99
100 % Ensure input table variable names match model expectations
101 rf_var_names = datasets.(matched_dataset).Properties.
    VariableNames(1:length(X_input_std));
102 rf_input = array2table(X_input_std, 'VariableNames',
    rf_var_names);
103
104 % Debugging dimensions and variable names
105 disp('Input standardized size:');
106 disp(size(X_input_std));
107 disp('Input variable names:');
108 disp(rf_var_names);
109
110 % Random Forest prediction
111 try
112     rf_model = model_results.(matched_dataset).RandomForest;
113     disp('RF Model Information:');
114     disp(['Predictor Names: ' strjoin(rf_model.PredictorNames, '
        , ')]);
115     disp(['Input Names: ' strjoin(rf_input.Properties.
        VariableNames, ', ')]);
116     disp('Input Values:');
117     for i = 1:length(rf_var_names)
118         disp([rf_var_names{i} ': ' num2str(X_input_std(i))]);
119     end
120
121     [Y_pred_raw, rf_scores] = predict(rf_model, rf_input);
122
123     if iscell(Y_pred_raw)
124         Y_pred_rf = str2double(Y_pred_raw{1});
125     elseif isnumeric(Y_pred_raw)
126         Y_pred_rf = Y_pred_raw(1);
127     else
128         disp(['Unexpected prediction type: ' class(Y_pred_raw)])
            ;
129         Y_pred_rf = NaN;
130     end

```

```

131         disp(['Final RF prediction: ' num2str(Y_pred_rf)]);
132         disp(['RF scores: ' num2str(rf_scores)]);
133
134     if isnan(Y_pred_rf)
135         disp('Warning: Random Forest produced NaN prediction.');
```

136 disp('RF Model Details:');

137 disp(['Number of trees: ' num2str(rf\_model.NumTrees)]);

138 disp(['Method: ' rf\_model.Method]);

139 end

140 catch rf\_err

141 disp('Error in Random Forest prediction:');

142 disp(['Error Message: ' rf\_err.message]);

143 disp(['Error Details: ' getReport(rf\_err)]);

144 disp('Random Forest Model Details:');

145 if exist('rf\_model', 'var')

146 disp(['Number of trees: ' num2str(rf\_model.NumTrees)]);

147 disp(['Method: ' rf\_model.Method]);

148 end

149 Y\_pred\_rf = NaN;

150 end

151

152

153 % Neural Network predictions

154 ensemble\_nets = model\_results.(matched\_dataset).NeuralNetworks;

155 num\_networks = length(ensemble\_nets);

156 if isempty(ensemble\_nets)

157 warning('No neural networks found for dataset: %s. Skipping

158 NN predictions.', matched\_dataset);

159 Y\_pred\_nn\_avg = NaN;

160 else

161 Y\_pred\_nn = zeros(1, num\_networks);

162 for i = 1:num\_networks

163 net = ensemble\_nets{i};

164 ps\_input = net.userdata.ps\_input;

165 ps\_target = net.userdata.ps\_target;

166 x\_norm = mapminmax('apply', X\_input\_std, ps\_input);

167 net\_out = net(x\_norm);

168 Y\_pred\_nn(i) = mapminmax('reverse', net\_out, ps\_target)

169 ';

170 end

171 Y\_pred\_nn\_avg = mean(Y\_pred\_nn);

172 end

173

174 % Final weighted ensemble prediction

175 if ~isnan(Y\_pred\_rf) && ~isnan(Y\_pred\_nn\_avg)

176 mse\_nn = mean((X\_input - Y\_pred\_nn\_avg).^2);

177 mse\_rf = mean((X\_input - Y\_pred\_rf).^2);

178 ensemble\_weight = mse\_rf / (mse\_rf + mse\_nn);

179 Y\_pred\_final = ensemble\_weight \* Y\_pred\_nn\_avg + (1 -

180 ensemble\_weight) \* Y\_pred\_rf;

181 elseif ~isnan(Y\_pred\_rf)

182 Y\_pred\_final = Y\_pred\_rf;

183 disp('Using only Random Forest predictions due to missing

184 Neural Network predictions.');

185 elseif ~isnan(Y\_pred\_nn\_avg)

186 Y\_pred\_final = Y\_pred\_nn\_avg;

187 disp('Using only Neural Network predictions due to missing

188 Random Forest predictions.');

```

184     else
185         Y_pred_final = NaN;
186         disp('No valid predictions available. ');
187     end
188
189     Y_pred_ensemble = Y_pred_final;
190     Y_actual = input('Please provide the actual EII value for this
        input: ');
191
192     % Store all predictions and their accuracies
193     predictions = struct();
194     if ~isempty(ensemble_nets)
195         for i = 1:num_networks
196             predictions.nn(i).value = Y_pred_nn(i);
197             predictions.nn(i).accuracy = 100 - abs((Y_actual -
                Y_pred_nn(i)) / Y_actual * 100);
198         end
199     end
200
201     predictions.nn_ensemble.value = Y_pred_nn_avg;
202     predictions.nn_ensemble.accuracy = 100 - abs((Y_actual -
        Y_pred_nn_avg) / Y_actual * 100);
203
204     if ~isnan(Y_pred_rf)
205         predictions.rf.value = Y_pred_rf;
206         predictions.rf.accuracy = 100 - abs((Y_actual - Y_pred_rf) /
            Y_actual * 100);
207     end
208
209     predictions.final_ensemble.value = Y_pred_final;
210     predictions.final_ensemble.accuracy = 100 - abs((Y_actual -
        Y_pred_final) / Y_actual * 100);
211
212     % Calculate accuracies
213     accuracies = struct();
214     if ~isempty(ensemble_nets)
215         accuracies.nn = zeros(1, num_networks);
216         for i = 1:num_networks
217             accuracies.nn(i) = 100 - abs((Y_actual - Y_pred_nn(i)) /
                Y_actual * 100);
218         end
219     end
220
221     accuracies.nn_ensemble = 100 - abs((Y_actual - Y_pred_nn_avg) /
        Y_actual * 100);
222     if ~isnan(Y_pred_rf)
223         accuracies.rf = 100 - abs((Y_actual - Y_pred_rf) / Y_actual
            * 100);
224     else
225         accuracies.rf = NaN;
226     end
227     accuracies.final_ensemble = 100 - abs((Y_actual -
        Y_pred_ensemble) / Y_actual * 100);
228
229     % Find best prediction
230     best_accuracy = -inf;
231     best_prediction = '';
232     best_value = NaN;

```

```

233
234 % Check individual neural networks
235 if ~isempty(ensemble_nets)
236     for i = 1:num_networks
237         if accuracies.nn(i) > best_accuracy
238             best_accuracy = accuracies.nn(i);
239             best_prediction = sprintf('Neural Network %d', i);
240             best_value = Y_pred_nn(i);
241         end
242     end
243 end
244
245 % Check ensemble predictions
246 if accuracies.nn_ensemble > best_accuracy
247     best_accuracy = accuracies.nn_ensemble;
248     best_prediction = 'Neural Network Ensemble';
249     best_value = Y_pred_nn_avg;
250 end
251
252 if ~isnan(Y_pred_rf) && accuracies.rf > best_accuracy
253     best_accuracy = accuracies.rf;
254     best_prediction = 'Random Forest';
255     best_value = Y_pred_rf;
256 end
257
258 if accuracies.final_ensemble > best_accuracy
259     best_accuracy = accuracies.final_ensemble;
260     best_prediction = 'Final Ensemble';
261     best_value = Y_pred_ensemble;
262 end
263
264 % Display results
265 disp('=== Model Predictions and Accuracies ===');
266 fprintf('Actual Value: %.4f\n\n', Y_actual);
267
268 % Individual Neural Networks
269 if ~isempty(ensemble_nets)
270     disp('Individual Neural Networks:');
271     for i = 1:num_networks
272         fprintf('Network %d: %.4f (Accuracy: %.2f%%)\n', i,
273             Y_pred_nn(i), accuracies.nn(i));
274     end
275 end
276
277 fprintf('\nNeural Network Ensemble: %.4f (Accuracy: %.2f%%)\n',
278     ...
279     Y_pred_nn_avg, accuracies.nn_ensemble);
280
281 if ~isnan(Y_pred_rf)
282     fprintf('Random Forest: %.4f (Accuracy: %.2f%%)\n',
283         Y_pred_rf, accuracies.rf);
284 else
285     fprintf('Random Forest: Prediction failed\n');
286 end
287
288 fprintf('Final Ensemble: %.4f (Accuracy: %.2f%%)\n',
289     Y_pred_ensemble, accuracies.final_ensemble);
290
291 % Display best prediction
292 fprintf('\n=== BEST PREDICTION ===\n');

```

```

287     fprintf('Model: %s\n', best_prediction);
288     fprintf('Value: %.4f\n', best_value);
289     fprintf('Accuracy: %.2f%%\n', best_accuracy);
290
291     catch ME
292         disp('An error occurred during prediction. Debug info:');
293         disp(['Error Message: ', ME.message]);
294         disp(['Error Location: ', ME.stack(1).name, ' line ', num2str(ME
            .stack(1).line)]);
295
296         % Additional debugging information
297         disp('Debug information:');
298         disp(['Input size: ', mat2str(size(X_input))]);
299         disp(['Standardized input size: ', mat2str(size(X_input_std))]);
300         if exist('rf_input', 'var')
301             disp('RF input variable names:');
302             disp(rf_input.Properties.VariableNames);
303         end
304     end
305 end

```

Listing 2: Testing Script

### 3.3 Technical Explanation of Testing Steps

#### 3.3.1 Data and Model Loading

- We attempt to load all four datasets again, primarily to **check input dimension consistency**.
- We load the `model_results_*.mat` files, each containing:
  - `NeuralNetworks`: the saved ensemble of neural nets.
  - `RandomForest`: the trained Random Forest model.
  - `Preprocessing.mean` / `Preprocessing.std`: median and MAD for robust scaling.
  - `BestEnsembleWeight`: chosen weight for combining NN ensemble and RF.

#### 3.3.2 User Input and Matching Dataset

- The user is prompted with `input(prompt)` to provide a feature vector.
- We check which dataset has the same number of features (`num_features = size(datasets.(data{2} - 1))`).
- If no dataset has the same dimensionality, we error out.

#### 3.3.3 Preprocessing and Prediction

- The script applies the saved **median/MAD-based normalization** from the training phase to the user's input.
- Predictions from:

1. **Random Forest (RF)**
2. **Neural Network Ensemble:** We apply each NN in the saved ensemble, average their outputs.
3. **Combined/Weighted Ensemble:** Takes an MSE-based weight between the NN ensemble and the RF prediction.

### 3.3.4 Accuracy Calculation

- The user is asked for the actual EII value: `Y_actual`.
- Accuracy is computed as:

$$\text{Accuracy (\%)} = 100 - \left| \frac{Y_{\text{actual}} - Y_{\text{pred}}}{Y_{\text{actual}}} \right| \times 100.$$

- The script displays all results:
  - Each individual NN prediction + accuracy
  - NN ensemble prediction + accuracy
  - RF prediction + accuracy
  - Final ensemble prediction + accuracy
  - The **best** prediction is highlighted.

## 4 Reasoning Behind Algorithmic Choices

### 4.1 Neural Networks

- We used multiple neural networks with different initial conditions. Neural networks can converge to local minima; using multiple networks and averaging outputs reduces variance and improves robustness.
- The architecture `[100, 50]` with `trainlm` is a balance between model complexity and computational feasibility.

### 4.2 Random Forest

- A Random Forest is an ensemble method itself, combining multiple decision trees. It is often robust to outliers and can handle complex interactions.
- We used 500 trees, which is a common choice to ensure stable predictions without excessive computational cost.
- Setting `MinLeafSize = 5` controls model complexity; ensures each tree is forced to consider subsets of features.



### 4.3 Ensemble Weighting

- We compute an MSE-based weight because it leverages how each component of the ensemble (NN\_avg, RF) performs on the validation set.
- If the RF has lower MSE than NN, it will get a higher weight, and vice versa.
- This linear combination is a straightforward, effective method of blending models.

### 4.4 Cross-Validation and Robust Scaling

- **Cross-Validation:** ensures that the final model is thoroughly tested on multiple distinct folds, reducing overfitting risk.
- **Median + MAD Normalization:** robust to outliers; typical in industrial or chemical process data (where extreme values can occur).

## 5 Conclusion

This setup, combining **(1) multiple feed-forward neural networks**, **(2) a random forest model**, and **(3) a weighted ensemble approach**, aims to provide robust, accurate EII predictions across four distinct datasets.

Key advantages:

- **Reduced Variance** through ensembling multiple neural networks.
- **Versatility** of Random Forests, which handle various data distributions well.
- **Flexible Weighted Combination** that adapts to which model performs best on validation data.

In practice, this approach demonstrates improved performance and reliability compared to any single model approach, particularly in scenarios where data may be noisy or exhibit outliers.