# Eliminating Multicollinearity Issues in Neural Network Ensembles: Incremental, Negatively Correlated, Optimal Convex Blending

**Pola-Lydia Lagari**
School of Nuclear Engineering, Purdue University
West Lafayette, IN 47907
lydialagaris@gmail.com

**Lefteri H. Tsoukalas**
School of Nuclear Engineering, Purdue University
West Lafayette, IN 47907, USA
tsoukala@purdue.edu

**Salar Safarkhani**
Krannert School of Management, Purdue University
West Lafayette, In 47907, USA
salar.safarkhani@gmail.com

**Isaac E. Lagaris**
Department of Computer Science and Engineering, University of Ioannina
Ioannina 45500, Greece
lagaris@cs.uoi.gr

## Abstract

Given a {features, target} dataset, we introduce an incremental algorithm that constructs an aggregate regressor, using an ensemble of neural networks. It is well known that ensemble methods suffer from the multicollinearity issue, which is the manifestation of redundancy arising mainly due to the common training-dataset. In the present incremental approach, at each stage we optimally blend the aggregate regressor with a newly trained neural network under a convexity constraint which, if necessary, induces negative correlations. Under this framework, collinearity issues do not arise at all, rendering so the method both accurate and robust.

## 1   Introduction

The combination of estimators has been studied by many research teams. Wolpert (1992)[†] was among the first who studied the effect of estimator combination on the approximation accuracy. His work was followed up by Breiman (1996), where an explanation for imposing convexity constraints on the linear combination coefficients was offered, based on grounds of generalization quality. Perrone and Cooper (1993) started a systematic study, developing the "Basic Ensemble Method" and the "Generalized Ensemble Method", that concluded in Perrone (1993) PhD Thesis. Independently, and around the same time, Hashem and Schmeiser (1993, 1995) developed an optimal linear combination method that led to Hashem (1993) PhD Thesis, and additional publications, Hashem (1996, 1997).

---

[†]This work was available since 1990 as LA-UR-90-3460 technical report

The potential of ensemble methods to offer improved accuracy was also stressed by Krogh and Vedelsby (1994), and by Meir (1995) whose work focused particularly on small and noisy datasets. Early work on combining classification neural networks, was published by Hansen and Salamon (1990) and later by Leblanc and Tibshirani (1996) and Zhou et al. (2002). A problem that emerges when combining many estimators is that of collinearity. This is due to the linear dependence of estimators that are trained on the same data. The relevant correlation matrix involved, becomes near singular or even singular and one has to discard a number of estimators in order to circumvent the problem, for example by using techniques such as principal component analysis, Merz and Pazzani (1999). A different approach was followed by Liu and Yao (1999) and Chen and Yao (2009), termed "negative correlation", a technique that encourages the formation of diversity among the estimators during the training process, by adding a proper term in the loss function. Diversity management through negative correlation has been examined and reviewed by Brown et al. (2005), Brown et al. (2005), Chan and Kasabov (2005) and by Reeve and Brown (2018) among others. For ensembles of classifiers, Tumer and Ghosh (2002), have presented a fusion approach based on order statistics.

The structure of this paper is as follows. In section 2, we first provide a brief background on the collinearity issue, and we proceed with the description and analysis of the proposed method. In section 2.3, we lay out the algorithmic procedure aiming to aid the implementation. In section 3, we present the results of performed numerical experiments on a number of test functions and network architectures. A summary with conclusions and thoughts for future research are given in section 4. There are also two short appendices describing technical matters.

## 2  Analysis and description of the approach

Let us begin with a few definitions. The problem is to approximate an unknown function $y(x)$ by a parametric model $f(x)$, given a training set $(T_r)$ and a test set $(T_s)$:

$$T_r = \{x_i, y_i = y(x_i)\}_{i=1,M}, \tag{1}$$
$$T_s = \{\hat{x}_i, \hat{y}_i = y(\hat{x}_i)\}_{i=1,L}. \tag{2}$$

The expected values of a function $f(x)$ over $T_r$ and over $T_s$, are denoted as $\langle f \rangle$ and $\langle f \rangle_s$ correspondingly, and are given by:

$$\langle f \rangle \equiv \frac{1}{M} \sum_{i=1}^{M} f(x_i), \tag{3}$$

$$\langle f \rangle_s \equiv \frac{1}{L} \sum_{i=1}^{L} f(\hat{x}_i). \tag{4}$$

The "misfit" of a function $f(x)$ is defined as:

$$m_f(x) \equiv f(x) - y(x), \tag{5}$$

In this article we will adopt for the model function $f(x)$, neural networks of various architectures and activation functions, all trained to fit $y(x)$. Let these networks be denoted as:

$$N_j(x), \ \forall j = 1, 2, \cdots, K_e, \tag{6}$$

with corresponding Mean Squared Error (MSE) values:

$$\langle m_j^2 \rangle \equiv \langle (N_j - y)^2 \rangle. \tag{7}$$

We will assume that each network $N_j(x)$ has zero bias over $T_r$, namely:

$$\langle m_j \rangle = \langle N_j - y \rangle = 0, \tag{8}$$

a property that can be easily arranged to hold (see Appendix A).
Let the ensemble estimator be a convex linear combination of ensemble members $N_i(x)$,

$$N_e(x) = \sum_{i=1}^{K_e} a_i N_i(x). \tag{9}$$

2

Convexity imposes the constraints:

$$\sum_{i=1}^{K_e} a_i = 1 \text{ and } a_i \geq 0, \forall i = 1, 2, \cdots, K_e. \tag{10}$$

The ensemble misfit is:

$$M_e(x) = N_e(x) - y(x) = \sum_{i=1}^{K_e} a_i m_i(x), \tag{11}$$

and its mean squared error is given by:

$$\langle M_e^2 \rangle = \sum_{i,j} a_i a_j \langle m_i m_j \rangle. \tag{12}$$

To minimize $\langle M_e^2 \rangle$, subject to the convexity constraints (10), it is necessary to construct the relevant Lagrangian which is written using multipliers $\lambda \in R$ and $\mu \in R^{K_e}$ as:

$$\mathcal{L}(a, \lambda, \mu) = \frac{1}{2} a^T C a - \lambda(e^T a - 1) - \mu^T a, \tag{13}$$

where $e^T = (1, 1, \cdots, 1)$, and $C_{ij} = \langle m_i m_j \rangle$.

The first order optimality conditions and the equality constraint $e^T a = 1$, yield:

$$\lambda = \frac{1 - e^T C^{-1} \mu}{e^T C^{-1} e}, \tag{14}$$

$$a = \frac{C^{-1} e}{e^T C^{-1} e} + \left[ C^{-1} - \frac{C^{-1} e e^T C^{-1}}{e^T C^{-1} e} \right] \mu. \tag{15}$$

Eq. (15) together with $a_i \geq 0$, $\mu_i \geq 0$ and the complimentarity conditions $a_i \mu_i = 0$, determine the sought solution. Note that if the correlation matrix $C$ is rank deficient, which is often the case due to collinearity, the inverse $C^{-1}$ is not defined and the above procedure is not applicable. Some methods, in order to avoid the multicollinearity issue, use a simple ensemble average (i.e. $a_i = \frac{1}{K_e}, \forall i = 1, \cdots, K_e$). These methods lead to non-optimal results since they neglect the relative importance of the various components, the remedy being the consideration of very large ensembles ($K_e \gg 1$), a tactic that renders the procedure excessive and inefficient as well.

In the present article we propose a novel method which determines the coefficients $a_i$ without having to deal with a potentially singular system, avoiding so the above mentioned issues due to multicollinearity. Instead, an aggregate network is being built incrementally from ground up, by optimally blending it with a single new network at every stage, via a convex linear combination. The new network is trained so that its misfit is negatively correlated to the misfit of the current aggregate, a constraint that is sufficient (although not necessary) to guarantee the convexity requirement. The method takes in account the relative importance of the different networks, satisfies the convexity requirements, and eliminates multicollinearity complications and associated numerical side effects.

## 2.1 The General Framework

Consider two zero-bias networks $N(x)$ and $\hat{N}(x)$, with corresponding misfits:

$$M(x) = N(x) - y(x), \text{ and } \hat{M}(x) = \hat{N}(x) - y(x). \tag{16}$$

Then their convex blend may be defined as:

$$\tilde{N}(x) = \beta N(x) + (1 - \beta)\hat{N}(x), \beta \in [0, 1], \tag{17}$$

with misfit $\tilde{M}(x) = \beta M(x) + (1 - \beta)\hat{M}(x)$, and MSE given by:

$$\langle \tilde{M}^2 \rangle = \langle (M - \hat{M})^2 \rangle \beta^2 + 2\langle \hat{M}(M - \hat{M}) \rangle \beta + \langle \hat{M}^2 \rangle. \tag{18}$$

Note that $\langle \tilde{M}^2 \rangle$ being quadratic in $\beta$, has a minimum at:

$$\beta^* = \max\{\min\{1, \beta_u\}, 0\}, \text{ with } \beta_u = \frac{\langle (\hat{M} - M)\hat{M} \rangle}{\langle (M - \hat{M})^2 \rangle}, \tag{19}$$

3

where $\beta_u$ is the unconstrained minimizer of $\langle \tilde{M}^2 \rangle$. If $\beta_u \notin [0, 1]$, then either $\beta^* = 1$ or $\beta^* = 0$ and there is no blending, as can be readily deduced by inspecting eq. (17).

Note from (19) that $\beta_u \in (0, 1)$, implies:

$$\langle M\hat{M} \rangle < \min\{\langle M^2 \rangle, \langle \hat{M}^2 \rangle\}, \tag{20}$$

and in this case $\beta^* = \beta_u$ leading to:

$$\langle \tilde{M}^2 \rangle^* = \frac{\langle M^2 \rangle \langle \hat{M}^2 \rangle - \langle M\hat{M} \rangle^2}{\langle M^2 \rangle + \langle \hat{M}^2 \rangle - 2\langle M\hat{M} \rangle}. \tag{21}$$

From eq. (21), one may deduce, after some manipulation, that $\langle \tilde{M}^2 \rangle^* < \min\{\langle M^2 \rangle, \langle \hat{M}^2 \rangle\}$. The important conclusion is that a convex combination of two zero-bias networks may be arranged so as to yield a zero-bias network with a reduced MSE.

## 2.2 Recursive Convex Blending

The idea is to construct recursively an aggregate network, by combining each time the current aggregate with a newly trained network. Namely, if by $N_k^{(a)}(x)$ we denote the $k^{th}$ aggregate network, then the $(k+1)^{th}$ aggregate will be given by the convex blend:

$$N_{k+1}^{(a)}(x) = \beta_{k+1} N_k^{(a)}(x) + (1 - \beta_{k+1}) N_{k+1}(x). \tag{22}$$

Let $M_k(x) \equiv N_k^{(a)}(x) - y(x)$ be the aggregate misfit. Then $N_{k+1}(x)$ is trained so that its misfit $m_{k+1}(x)$ satisfies:

$$\langle M_k m_{k+1} \rangle < \min\{\langle M_k^2 \rangle, \langle m_{k+1}^2 \rangle\}, \tag{23}$$

in order to guarantee convex blending (i.e. $\beta_{k+1} \in [0, 1]$). This is a constrained optimization problem stated formally as:

$$\min\langle m_{k+1}^2 \rangle, \text{ subject to:} \tag{24}$$

$$\langle M_k m_{k+1} \rangle \leq \min\{\langle M_k^2 \rangle, \langle m_{k+1}^2 \rangle\}. \tag{25}$$

Note that since the positive quantity: $\min\{\langle M_k^2 \rangle, \langle m_{k+1}^2 \rangle\}$, is expected to be small, condition (25) may be satisfied by imposing that the aggregate and the new network misfits are negatively correlated or uncorrelated, i.e. $\langle M_k m_{k+1} \rangle \leq 0$, which is a stronger requirement than (25), and which in addition prevents $\beta_{k+1}$ from assuming the limiting no-blend values of 0 and 1[‡].

The next aggregate misfit is then given by:

$$M_{k+1}(x) = \beta_{k+1} M_k(x) + (1 - \beta_{k+1}) m_{k+1}(x), \tag{26}$$

$$\text{with } \beta_{k+1} = \frac{\langle m_{k+1}^2 \rangle - \langle M_k m_{k+1} \rangle}{\langle (M_k - m_{k+1})^2 \rangle}. \tag{27}$$

Setting $M_1(x) = m_1(x)$ (and hence $\beta_1 = 0$), and noting that $M_k$ is a linear combination of $m_1(x)$, $m_2(x), \cdots, m_k(x)$, one may express $M_k$ as:

$$M_k(x) = \sum_{i=1}^{k} a_i^{(k)} m_i(x), \tag{28}$$

with the coefficients $a_i^{(k)}$ given by:

$$a_i^{(k)} = \begin{cases} 1 - \beta_k, & i = k \\ (1 - \beta_i) \prod_{l=i+1}^{k} \beta_l, & \forall i = 1, 2, \cdots, k-1 \end{cases} . \tag{29}$$

One may verify that $\sum_{i=1}^{k} a_i^{(k)} = 1$ and $a_i^{(k)} \geq 0$, and the aggregate network is then given by:

$$N_k^{(a)}(x) = \sum_{i=1}^{k} a_i^{(k)} N_i(x). \tag{30}$$

---

[‡]For different blending restrictions see Appendix B

Note that the aggregate network $N_k^{(a)}(x)$ is not the simple average of the member networks as is the case in Ahmad and Zhang (2009) and in the negative correlation approaches Liu and Yao (1999); Chan and Kasabov (2005); Chen and Yao (2009); Brown et al. (2005). The combination coefficients $a_i^{(k)}$ given by eq. (29 ), do take in account the relative importance of the different network contributions, and are being built iteratively by optimal pairwise convex blending, avoiding so problems due to multicollinearity. This is in contrast to the approach described by eq. (15) which had been adopted by several authors in the past.

## 2.3 Algorithmic Procedure and Implementation

The above results and ideas are employed to design an algorithm for practical use.

### Initialization

- Create the empty list, `MList`, to store the aggregate misfit in each step.
- Create the empty list, `BList`, to store $\beta$'s.
- Create the empty list, `ModelList` to store the trained neural network model.
- Set $b_L$, $b_U$ both $\in (0, 1)$, the lower and upper bound for the $\beta$'s.
- Set $K_e$, the number of networks to be contained in the ensemble.

### Main Part

1. Train the first zero bias network, $N_1(x)$ and push it to the `ModelList`.
2. Set $M_1(x_i) = N_1(x_i) - y_i, \forall\{x_i, y_i\} \in T_r$ and push it to `MList`.
3. Set $\beta_1 = 0$ and push it to the `BList`.
4. $k \leftarrow 1$
5. Read $M_k(x), \forall x \in T_r$ from `MList`.
6. Train the zero bias network $N_{k+1}(x)$ s.t. $\langle M_k m_{k+1}\rangle \leq \min\{\langle M_k^2\rangle, \langle m_{k+1}^2\rangle\}$. This can be facilitated via a penalty method, i.e. by optimizing:

$$\langle m_{k+1}^2\rangle + \lambda \max\{\langle m_{k+1}M_k\rangle, 0\}$$

   for a sequence of increasing $\lambda$ values, until $\beta_{k+1} \in [b_L, b_U]^*$.
7. Calculate $M_{k+1}$ from eq. (26) and $\langle M_{k+1}^2\rangle$.
8. Push $M_{k+1}$ to `MList`, $\beta_{k+1}$ to `BList`, and $N_{k+1}$ to `ModelList`.
9. $k \leftarrow k + 1$
10. If $k < K_e$ repeat from $\boxed{5}$.
11. Training Task Completed.

### Final Network Assembly

From `BList` and `ModelList`, recover the $\beta$'s and the trained networks respectively.

(a) From $\beta_1, \cdots, \beta_{K_e}$, calculate the member coefficients $a_i^{(K_e)}$ using eq. (29).

(b) Construct the final zero-bias aggregate network $N_{K_e}^{(a)}(x)$, according to eq. (30).

## 3 Numerical Experiments

We consider two model functions in our numerical study,

$$
\begin{aligned}
f_1(x) &= x\sin(x^2), \text{ and} & (31)\\
f_2(\mathbf{x}) &= 10d + \sum_{i=1}^{d}[x_i^2 - 10cos(2\pi x_i)], \text{ with } \mathbf{x} \in [-1.5, 1.5]^d. & (32)
\end{aligned}
$$

We have performed four different experiments to assess the effectiveness of the method. Common features and practices shared across all experiments are listed below.

---

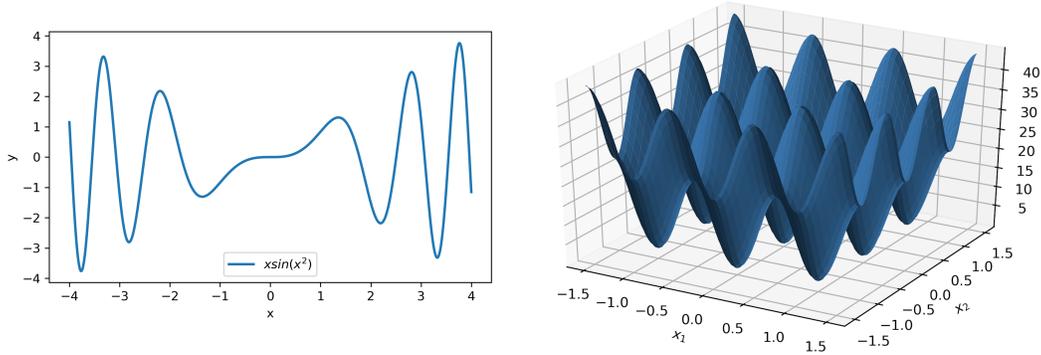$^*\beta_{k+1}$ is calculated according to eq. (27)

Figure 1: Plots for $f_1(x)$ of eq.(31) with $x \in [-4, 4]$, and $f_2(x)$ of eq. (32) with $d = 2$

### Datasets

For every experiment we define a dataset $D$, and we split it in two disjoint sets, the training set $T_r$, and the test set $T_s$, with $D = T_r \cup T_s$ and $T_r \cap T_s = \emptyset$.

### Training

Training was performed using the BFGS Quasi-Newton optimization algorithm, contained in "TensorFlow Probability" Abadi et al. (2015), with the maximum number of iterations set to 20,000.

### Regularization

We have used $L_2$-mean weight decay, with regularization factor $\nu_{reg}$, chosen so as to inhibit excessive growing of the network weights.

### Blending range

We set the lower and upper bounds for $\beta$'s to $b_L = 0$ and $b_U = 0.99$.

### Penalty strategy

In step No. 6 of the section 2.3 algorithm, we set the initial penalty to $\lambda = 4$. At each subsequent iteration, until condition $b_L < \beta_{k+1} < b_U$ is satisfied, $\lambda$ is doubled. If the above condition is not satisfied within 10 iterations, we discard the model.

### ANN architecture

For all the test cases we use a single hidden layer neural network of the form:

$$N(x) = \sum_{i=1}^{Nodes} \gamma_i h(\delta_i^T x + \phi_i), \text{ where } h(x) \text{ is the network's activation function.}$$

We use three types of activation functions:

1. Sigmoid: $\qquad\qquad h(z) = (1 + exp(-z))^{-1}$.
2. Softplus: $\qquad\qquad h(z) = \ln(1 + exp(z))$.
3. Hyperbolic tangent: $h(z) = \tanh(z)$.

### Notation

In all tables, "Nodes" denotes the number of neurons in the hidden layer, "Type" stands for the activation type, "MSE" for the mean squared error over the training set, "$\beta$" for the blending coefficient given by eq. (27), "AG. MSE" for the mean squared error of the aggregate over the training set, "AG. MSE/TE" for the mean squared error of the aggregate over the test set, and "$a$" for the member participation coefficient, given by eq. (29).

### Randomness

To ease the reproduction of the numerical results, we set all random seeds to 12345 using Python's Numpy.

### Code

We made the code publicly available via the Github repository with URL: `https://github.com/salarsk1/Ensemble-NN`, under an MIT license.

## 3.1 Test Case 1

Model function: $f_1(x), x \in [-4, 4]$. We construct a set $D$ containing $M_D = 1000$ equidistant points, $x_i = -4 + (i-1)\frac{8}{999}$. The training set $T_r$ contains $M_r = 38$ equidistant points, $z_j = -4 + (j-1)\frac{216}{999} = x_{27j-26}$. Note that $T_r \subset D$ and the test set is $T_s = D - T_r$. The ensemble contained six networks. Details of this experiment are listed in table 1. Among the 6 models, model No. 5 has the minimum training MSE, which is 0.38498. The training MSE of the final aggregate is 0.14035, which is 63% lower than that of model No. 5. The plots of the six member-networks $N_i(x)$ together with the target function $f_1(x)$ are depicted in fig. 2. Likewise, plots of the corresponding aggregate networks are depicted in fig. 3.

Table 1: Experiment with $f_1(x)$ with $x \in [-4, 4]$ and a training set of 38 equidistant points. Regularization factor: $\nu_{reg} = 0.002$.

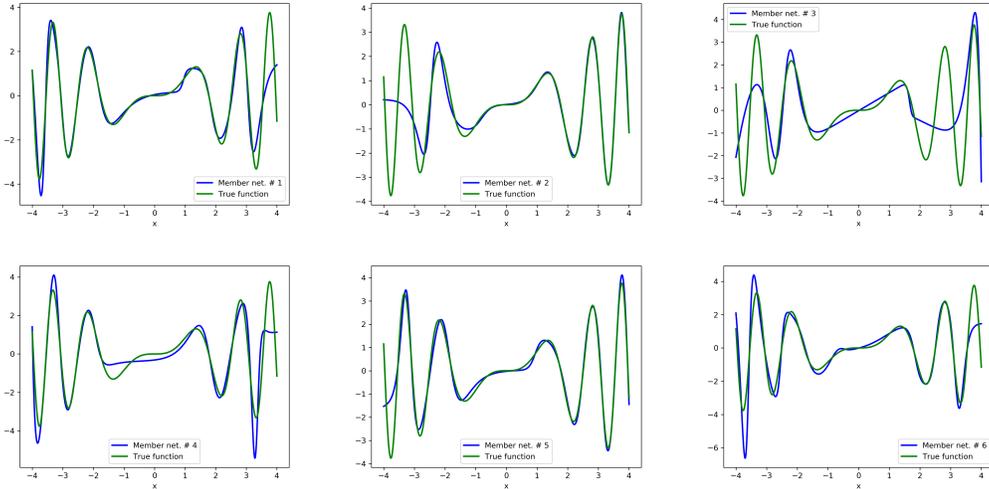| ANN | Nodes | Type | MSE | $\beta$ | AG. MSE | AG. MSE/TE | $a$ |
|-----|-------|------|------|---------|---------|-----------|------|
| 1 | 9 | tanh | 0.45326 | 0.00000 | 0.45326 | 0.32955 | 0.20175 |
| 2 | 11 | sigmoid | 0.85077 | 0.65221 | 0.29531 | 0.21960 | 0.10759 |
| 3 | 11 | softplus | 1.92921 | 0.94300 | 0.28932 | 0.22312 | 0.01870 |
| 4 | 9 | tanh | 0.43278 | 0.65750 | 0.23590 | 0.16807 | 0.17088 |
| 5 | 11 | sigmoid | 0.38498 | 0.63707 | 0.16427 | 0.10779 | 0.28423 |
| 6 | 12 | sigmoid | 0.45224 | 0.78315 | 0.14035 | 0.08970 | 0.21685 |



Figure 2: Plot of the ANN members related to table 1, along with the plot of $f_1(x)$.

## 3.2 Test Case 2

To escalate the difficulty of the task, in this case we have experimented with $f_1(x)$ in a wider range $x \in [-6, 6]$, which introduces intense oscillatory behavior illustrated in fig. 4. We construct a set $D$ containing $M_D = 1201$ equidistant points, $x_i = -6 + \frac{i-1}{100}$. The training set $T_r$ contains $M_r = 121$ equidistant points, $z_j = -6 + \frac{j-1}{10} = x_{10j-9}$. Note that $T_r \subset D$ and the test set is $T_s = D - T_r$. The generalization performance is quite satisfactory but not of the quality of the first example, as can be seen by inspecting table 2. This was expected since the rapid oscillations as $|x|$ grows, render the task harder. We have noticed that if one of the member networks $N_i(x)$ overfits the data, then a number of additional networks are necessary in the ensemble in order to eliminate the overfit effect. So it seems that a preferred tactic would be to use small networks, that are less prone to over-fitting, rather than large ones. This also serves the need of the so called "*few-shot learning*", that refers to problems with limited training examples Wang et al. (2020).
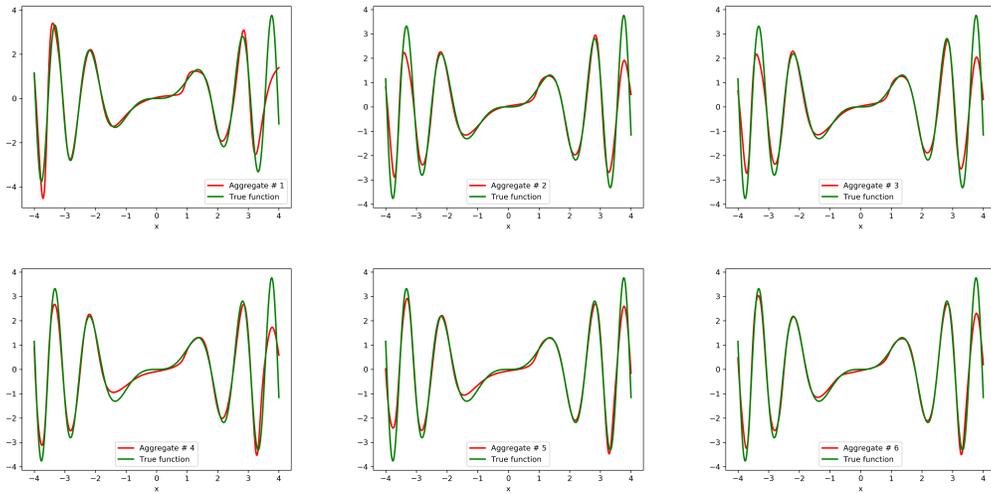
Figure 3: Plot of the aggregate ANNs related to table 1, along with the plot of $f_1(x)$.
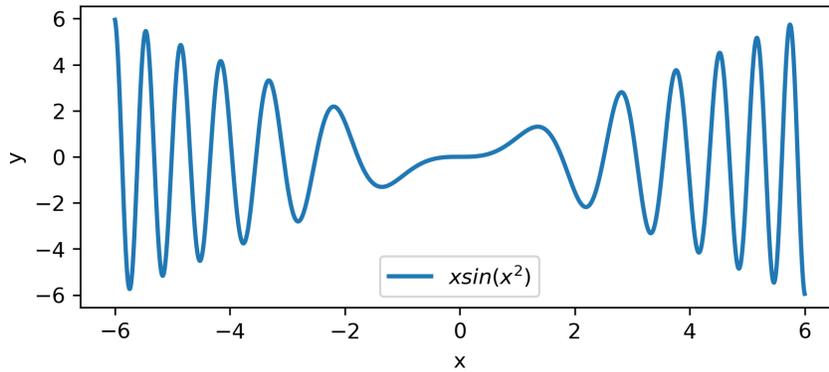


Figure 4: Plot of $f_1(x)$ with $x \in [-6, 6]$, illustrating the rapidly oscillating pattern as $|x|$ grows

Table 2: Experiment with $f_1(x)$ with $x \in [-6, 6]$ and a training set of 121 equidistant points. Regularization factor: $\nu_{reg} = 0.003$.

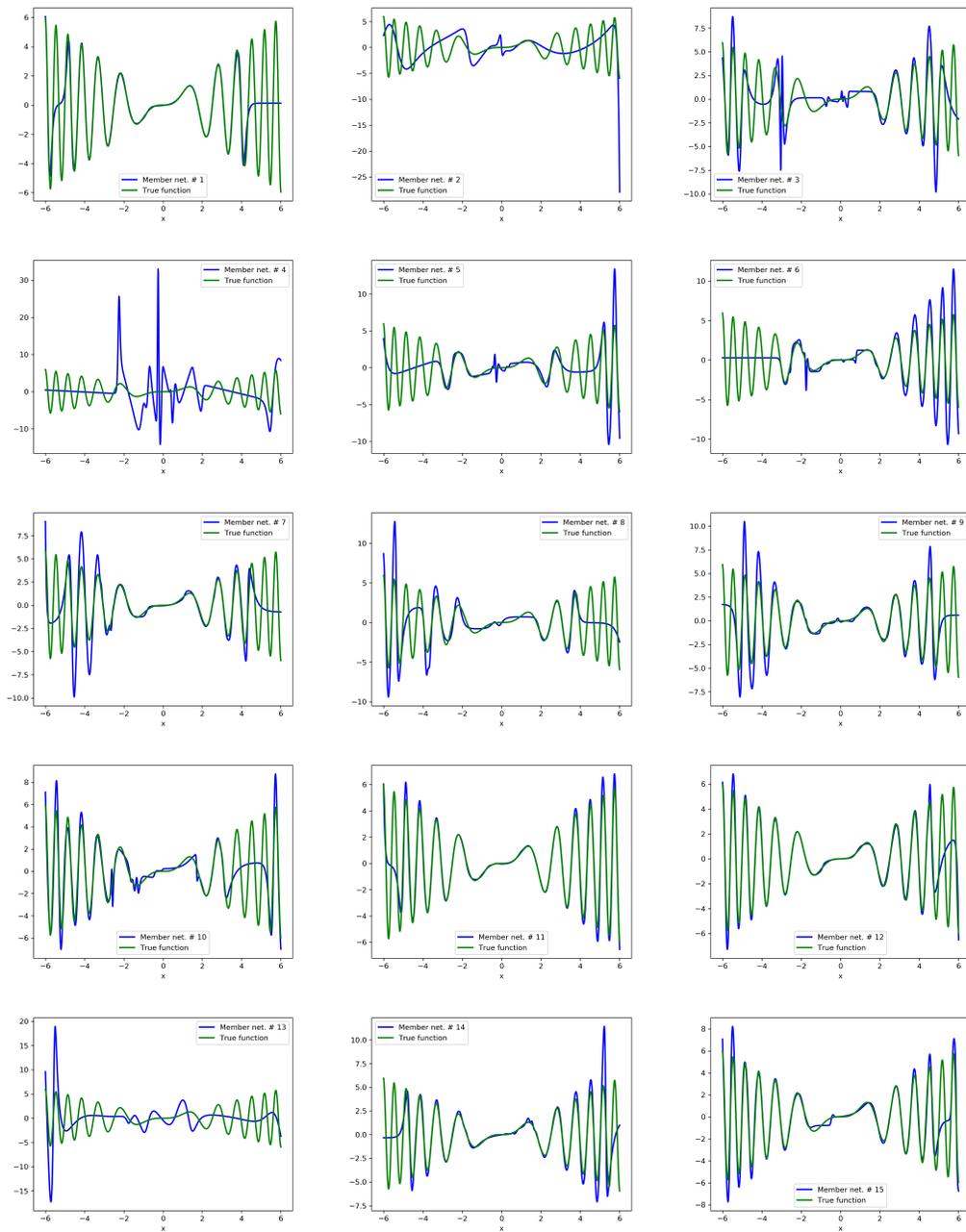| ANN | Nodes | Type | MSE | $\beta$ | AG. MSE | AG. MSE/TE | $a$ |
|---|---|---|---|---|---|---|---|
| 1 | 23 | tanh | 2.87200 | 0.00000 | 2.87200 | 2.73192 | 0.04995 |
| 2 | 25 | sigmoid | 12.93168 | 0.86162 | 2.60566 | 2.62689 | 0.00802 |
| 3 | 27 | sigmoid | 3.19976 | 0.55508 | 1.53797 | 1.60549 | 0.04646 |
| 4 | 23 | softplus | 18.90312 | 0.95861 | 1.50554 | 1.56396 | 0.00451 |
| 5 | 24 | sigmoid | 5.05116 | 0.79406 | 1.24986 | 1.30413 | 0.02825 |
| 6 | 29 | tanh | 4.18786 | 0.81902 | 1.09904 | 1.13224 | 0.03032 |
| 7 | 26 | tanh | 4.07782 | 0.83087 | 0.97027 | 1.00334 | 0.03410 |
| 8 | 23 | sigmoid | 4.48691 | 0.88752 | 0.91286 | 0.93800 | 0.02555 |
| 9 | 24 | tanh | 3.70554 | 0.86650 | 0.84497 | 0.84990 | 0.03500 |
| 10 | 25 | tanh | 2.00931 | 0.73126 | 0.66317 | 0.68001 | 0.09634 |
| 11 | 28 | sigmoid | 0.98583 | 0.61088 | 0.44286 | 0.45987 | 0.22836 |
| 12 | 27 | tanh | 1.27605 | 0.78000 | 0.37085 | 0.38586 | 0.16552 |
| 13 | 26 | softplus | 10.20572 | 0.97406 | 0.36387 | 0.37846 | 0.02004 |
| 14 | 26 | tanh | 2.98190 | 0.95094 | 0.35688 | 0.37098 | 0.03985 |
| 15 | 26 | tanh | 1.12186 | 0.81227 | 0.31371 | 0.32395 | 0.18773 |

8

Figure 5: Plot of the ANN members related to table 2, along with the plot of $f_1(x)$.
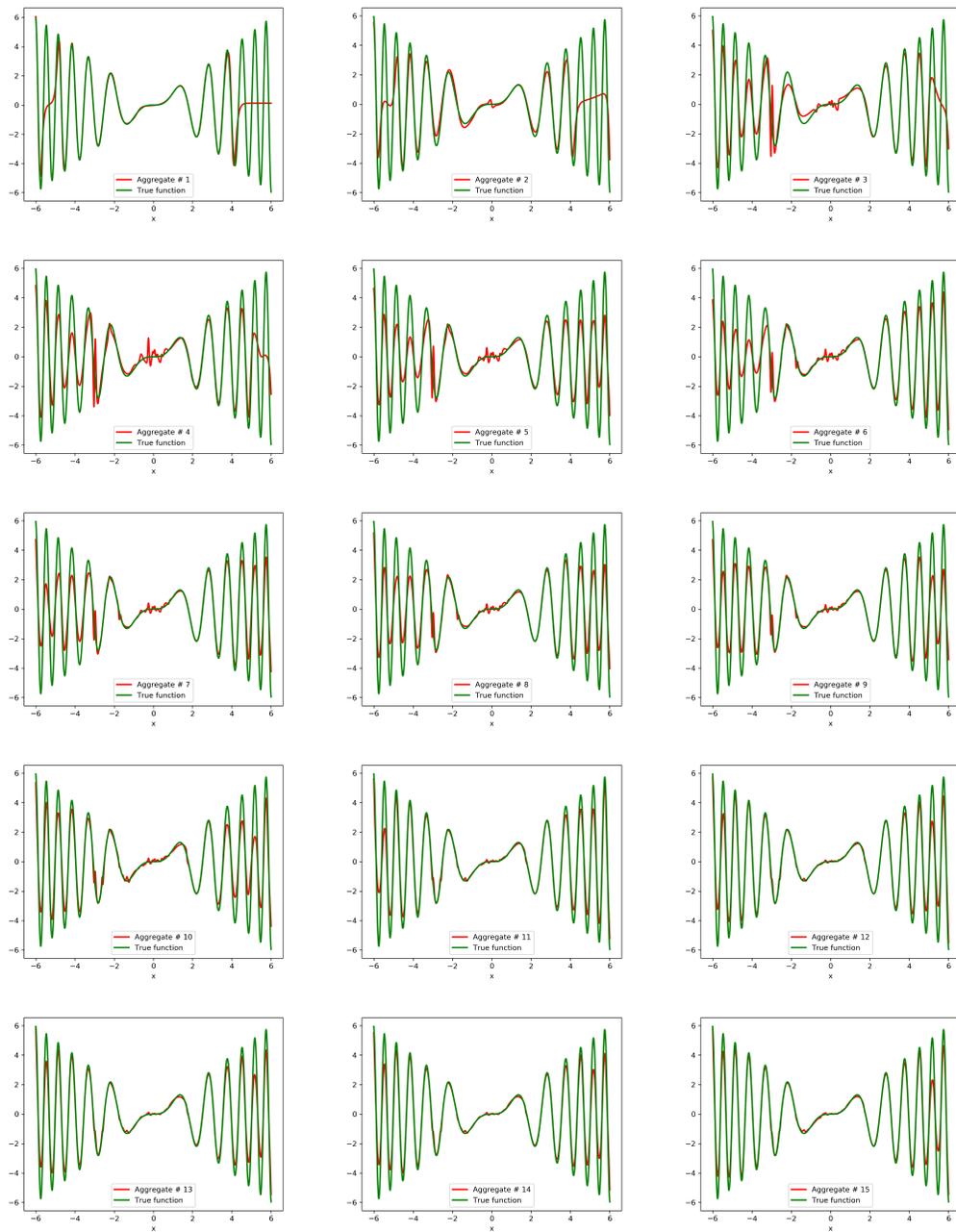
Figure 6: Plot of the aggregate ANNs related to table 2, along with the plot of $f_1(x)$.

In fig. 5, the member networks are plotted. Note that networks $N_2(x)$, $N_4(x)$, $N_5(x)$ and $N_{13}(x)$, clearly overfit $f_1(x)$. However the aggregate networks, plotted in figure 6, eliminate this effect gradually, by appropriate weighting. Note that the corresponding coefficients $a_2, a_4, a_5$ and $a_{13}$ are among the lowest ones in the linear combination composing the final aggregate network. Among the 15 models, model No. 11 has the minimum training MSE, which is 0.98583. The training MSE of final aggregate model is 0.31371, which is $68\%$ lower than that of model No. 11.

### 3.3 Test Case 3

We have also experimented with the multidimensional function, $f_2(\mathbf{x})$, where we have set $d = 4$. Again, we use MLP's with one hidden layer. For this function, we have used a grid of 15 equidistant points along each dimension, producing so a set $D$ containing a total of 50,625 points. We have chosen 2,000 of these points at random (seed is set to 12345) for the training set $T_r$, and the rest were used for the test set $T_s$. Detailed results of the experiments with $f_2(\mathbf{x})$, are listed in Table 3. Model No. 6 has a training MSE of 0.04537, which is the lowest among the 10 trained models. The training MSE of the final aggregate is 0.02547, which is $43\%$ lower than that of model No. 6.

Table 3: Experiment with $f_2(\mathbf{x})$ with $\mathbf{x} \in [-1.5, 1.5]^4$ and a training set of 2,000 points chosen randomly from a set of 50,625 equidistant points. Regularization factor: $\nu_{reg} = 0.05$.

| ANN | Nodes | Type | MSE | $\beta$ | AG. MSE | AG. MSE/TE | $a$ |
|-----|-------|------|------|---------|---------|------------|-----|
| 1 | 38 | sigmoid | 0.25331 | 0.00000 | 0.25331 | 0.29718 | 0.06849 |
| 2 | 38 | tanh | 0.14980 | 0.37494 | 0.09163 | 0.10700 | 0.11418 |
| 3 | 37 | sigmoid | 0.35718 | 0.80263 | 0.07454 | 0.08540 | 0.04492 |
| 4 | 37 | sigmoid | 0.31343 | 0.82541 | 0.06335 | 0.07445 | 0.04814 |
| 5 | 39 | sigmoid | 0.54320 | 0.91340 | 0.05900 | 0.06983 | 0.02614 |
| 6 | 39 | tanh | 0.04537 | 0.40159 | 0.03420 | 0.04083 | 0.44982 |
| 7 | 40 | sigmoid | 0.17129 | 0.86819 | 0.03096 | 0.03691 | 0.11413 |
| 8 | 40 | tanh | 1.09517 | 0.95779 | 0.02889 | 0.03429 | 0.03815 |
| 9 | 41 | sigmoid | 0.33329 | 0.93684 | 0.02750 | 0.03331 | 0.06094 |
| 10 | 41 | tanh | 1.56689 | 0.96492 | 0.02547 | 0.03061 | 0.03508 |

### 3.4 Test Case 4

In this case we have experimented with noisy data. Model function: $f_1(x)$ with $x \in [-5, 5]$. We construct a set $D$ containing $M_D = 1001$ equidistant points, $x_i = -5 + \frac{i-1}{100}$. The training set $T_r$ contains $M_r = 101$ equidistant points, $z_j = -5 + \frac{j-1}{10} = x_{10j-9}$ and the target values are contaminated with white normal noise $(0.3N(0, 1))$. The test set is left intact, i.e. without any added noise. The ensemble comprises 15 different networks. In fig. 8, the plots of $f_1(x)$ and of the trained member networks are displayed. The aggregate ANN's are plotted in fig. 9. The noisy data are depicted in fig. 7, where the dots correspond to the noisy training observations, while the solid line is a plot of $f_1(x)$.
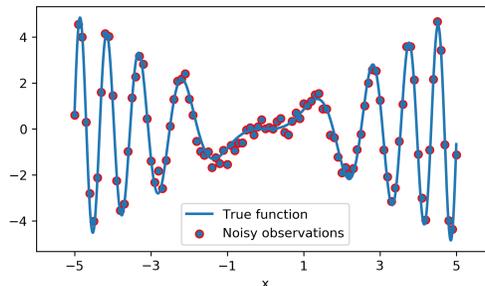


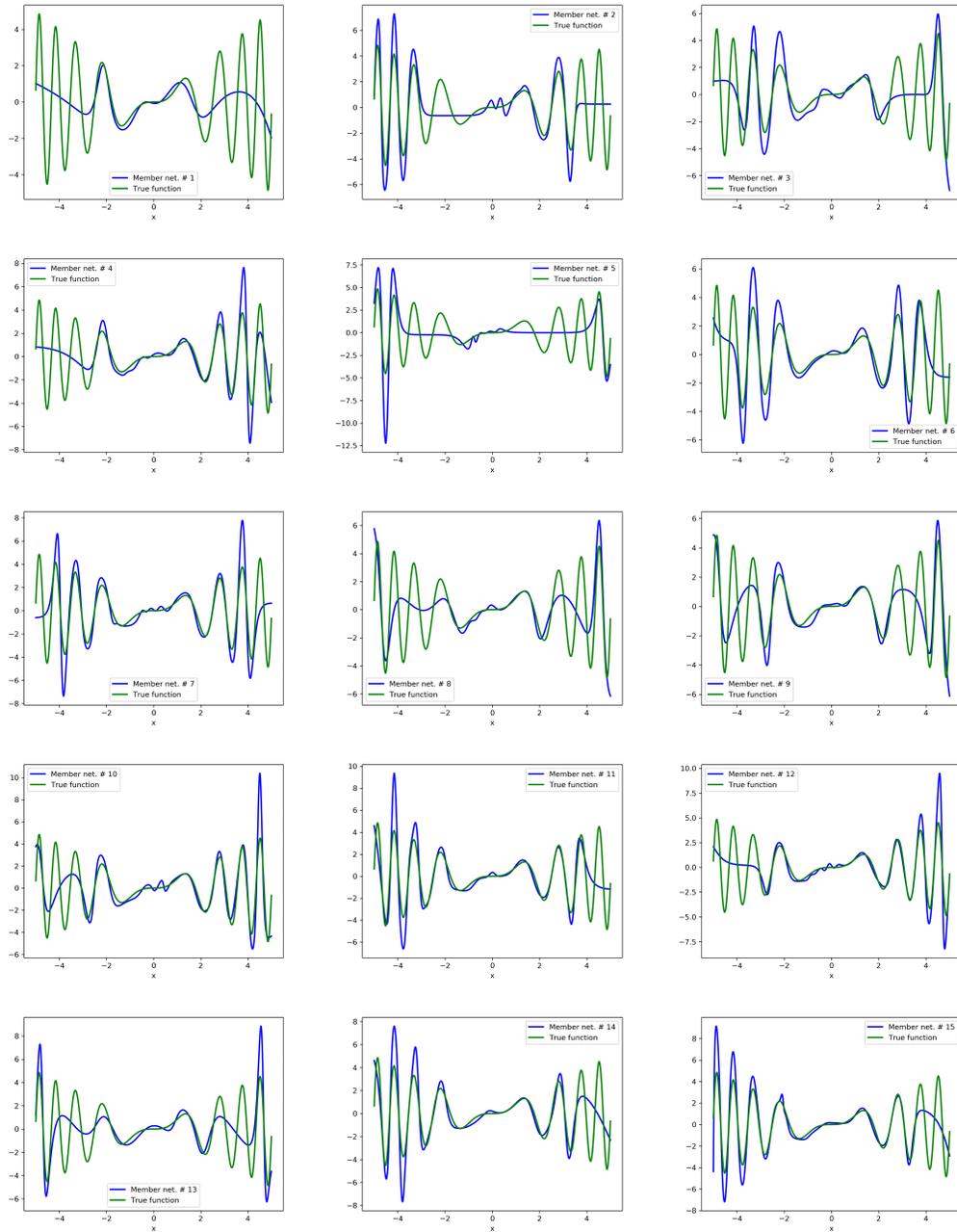Figure 7: Plot of the noise-corrupted observations used for training, alongside $f_1(x)$

Figure 8: Plot of the ANN members related to table 4, along with the plot of $f_1(x)$.

Detailed results are listed in table 4. As expected, the aggregate training error decreases by adding more networks and similar is the behavior of the test error. From column "AG. MSE/TE" of table 4, we note that the aggregate network generalizes extremely well. Note, model No. 12 has the lowest training MSE among all models which is 1.43211. The training MSE of the final aggregate is 0.45040, which is 68% lower than that of model No. 12.
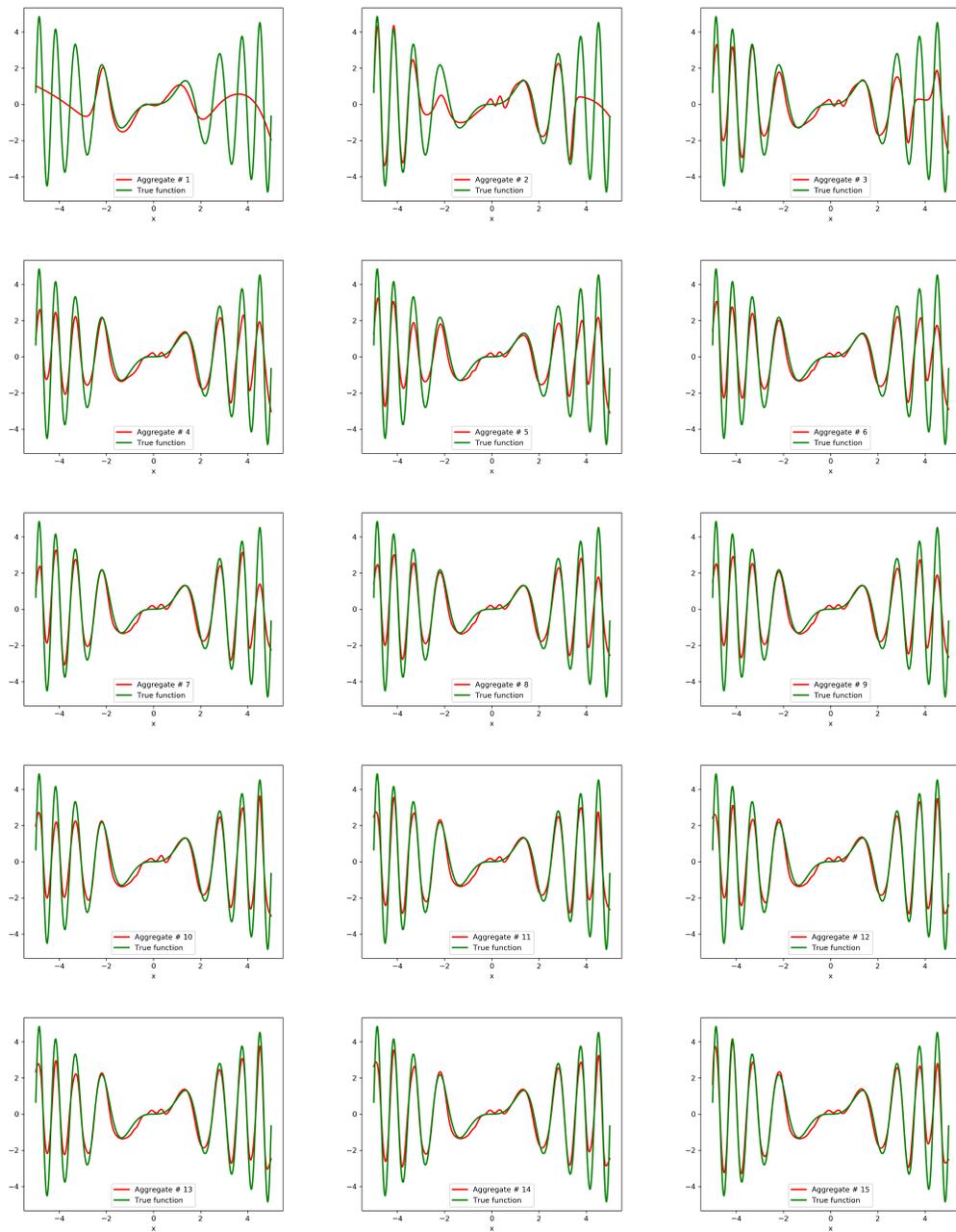
Figure 9: Plot of the aggregate ANNs related to table 4, along with the plot of $f_1(x)$.

Table 4: Experiment with $f_1(x)$ with $x \in [-5, 5]$ and a training set of 101 equidistant points with added normal white noise. Regularization factor: $\nu_{reg} = 0.1$.

| ANN | Nodes | Type | MSE | $\beta$ | AG. MSE | AG. MSE/TE | $a$ |
|---|---|---|---|---|---|---|---|
| 1 | 24 | tanh | 3.42492 | 0.00000 | 3.42492 | 3.60286 | 0.06919 |
| 2 | 25 | sigmoid | 12.89758 | 0.81382 | 2.90180 | 3.04288 | 0.01583 |
| 3 | 27 | sigmoid | 4.71427 | 0.62760 | 1.91684 | 1.92987 | 0.05045 |
| 4 | 23 | softplus | 5.21665 | 0.98091 | 1.91559 | 1.93207 | 0.00264 |
| 5 | 25 | sigmoid | 10.43508 | 0.92285 | 1.85563 | 1.88489 | 0.01155 |
| 6 | 29 | tanh | 2.28159 | 0.55624 | 1.10990 | 1.12316 | 0.11939 |
| 7 | 26 | softplus | 7.02311 | 0.96502 | 1.10213 | 1.11369 | 0.00975 |
| 8 | 24 | sigmoid | 5.98359 | 0.95641 | 1.09196 | 1.10721 | 0.01271 |
| 9 | 25 | tanh | 3.65851 | 0.86828 | 1.03151 | 1.05286 | 0.04422 |
| 10 | 25 | tanh | 2.72174 | 0.74489 | 0.80691 | 0.79448 | 0.11498 |
| 11 | 25 | softplus | 12.74256 | 0.96355 | 0.78981 | 0.77719 | 0.01705 |
| 12 | 27 | tanh | 1.43211 | 0.66732 | 0.57738 | 0.58926 | 0.23319 |
| 13 | 26 | softplus | 9.88673 | 0.97901 | 0.57310 | 0.57708 | 0.01503 |
| 14 | 25 | tanh | 2.31302 | 0.84593 | 0.51341 | 0.51266 | 0.13040 |
| 15 | 24 | tanh | 2.36287 | 0.84637 | 0.45040 | 0.45868 | 0.15363 |

# 4    Summary and Conclusions

In this article we have presented an ensemble method that is free of the multicollinearity issue, and without having to resort to the ad-hoc choice of uniform weights. Instead, we build an aggregate network, that does take in account the relative importance of the member networks. The aggregate network is incrementally built by blending it at every stage with a newly trained network under a negative correlation constraint. The blending is a convex linear combination as suggested in Breiman (1996), in order to maintain high generalization performance. Indeed, this may be confirmed by inspecting the "AG. MSE" and "AG.MES/TE" columns in tables 1, 2, 3 and 4. The experiment with noisy data, summarized in table 4, indicates that the method is robust and keeps on delivering quality generalization, as well as partial noise filtering.

Ensemble methods are useful in many applications. They have been used with success in econometrics and statistics for quite some time, Granger (1989), Wallis (2011). Since ensembles may contain many small networks, which are capable of being trained over datasets limited in size without over-fitting, they seem to be suitable for the interesting "Few-Shot Learning" case, Wang et al. (2020).
The present method has been developed with regression problems in mind. Further work is necessary to extend it properly for handling classification tasks as well.

# Acknowledgements

# References

Wolpert, D. H. Stacked Generalization. *Neural Networks* **1992**, *5*, 241–259.

Breiman, L. Stacked regressions. *Machine Learning* **1996**, *24*, 49–64.

Perrone, M. P.; Cooper, L. N. When Networks Disagree: Ensemble Methods for Hybrid Neural Networks. Artificial Neural Networks for Speech and Vision. 1993; pp 126–142.

Perrone, M. P. Improving Regression Estimation: Averaging Methods for Variance Reduction with Extensions to General Convex Measure Optimization. Ph.D. thesis, Physics Department, Brown University, Providence, RI, USA, 1993.

Hashem, S.; Schmeiser, B. Approximating a Function and its Derivatives Using MSE-Optimal Linear Combinations of Trained Feedforward Neural Networks. In Proceedings of the Joint Conference on Neural Networks. 1993; pp 617–620.

Hashem, S.; Schmeiser, B. Improving Model Accuracy Using Optimal Linear Combinations of Trained Neural Networks. *IEEE Transactions on Neural Networks* **1995**, *6*, 792–794.

Hashem, S. Optimal Linear Combinations of Neural Networks. Ph.D. thesis, School of Industrial Engineering, Purdue University, West Lafayette, IN, USA, 1993.

Hashem, S. Effects of Collinearity on Combining Neural Networks. *Connection Science* **1996**, *8*, 315–336.

Hashem, S. Optimal Linear Combinations of Neural Networks. *Neural Networks* **1997**, *10*, 599–614.

Krogh, A.; Vedelsby, J. Neural Network Ensembles, Cross Validation and Active Learning. Proceedings of the 7th International Conference on Neural Information Processing Systems. Cambridge, MA, USA, 1994; pp 231–238.

Meir, R. Bias, Variance and the Combination of Least Squares Estimators. NIPS 7. Cambridge MA, 1995.

Hansen, L. K.; Salamon, P. Neural network ensembles. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **1990**, *12*, 993–1001.

Leblanc, M.; Tibshirani, R. Combining Estimates in Regression and Classification. *Journal of the American Statistical Association* **1996**, *91*, 1641–1650.

Zhou, Z.-H.; Wu, J.; Tang, W. Ensembling neural networks: Many could be better than all. *Artificial Intelligence* **2002**, *137*, 239 – 263.

Merz, C. J.; Pazzani, M. J. A Principal Components Approach to Combining Regression Estimates. *Machine Learning* **1999**, *36*, 9–32.

Liu, Y.; Yao, X. Ensemble learning via negative correlation. *Neural Networks* **1999**, *12*, 1399–1404.

Chen, H.; Yao, X. Regularized Negative Correlation Learning for Neural Network Ensembles. *IEEE Transactions on Neural Networks* **2009**, *20*, 1962–1979.

Brown, G.; Wyatt, J. L.; Tiňo, P. Managing Diversity in Regression Ensembles. *J. Mach. Learn. Res.* **2005**, *6*, 1621–1650.

Brown, G.; Wyatt, J.; Harris, R.; Yao, X. Diversity creation methods: a survey and categorization. *Information Fusion* **2005**, *6*, 5–20.

Chan, Z. S. H.; Kasabov, N. Fast Neural Network Ensemble Learning via Negative-Correlation Data Correction. *IEEE Transactions on Neural Networks* **2005**, *16*, 1707–1710.

Reeve, H. W.; Brown, G. Diversity and degrees of freedom in regression ensembles. *Neurocomputing* **2018**, *298*, 55–68.

Tumer, K.; Ghosh, J. Robust Combining of Disparate Classifiers through Order Statistics. *Pattern Analysis & Applications* **2002**, *8*, 189–200.

Ahmad, Z.; Zhang, J. Selective combination of multiple neural networks for improving model prediction in nonlinear systems modelling through forward selection and backward elimination. *Neurocomputing* **2009**, *72*, 1198 – 1204.

Abadi, M. et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. 2015; https://www.tensorflow.org/, Software available from tensorflow.org.

Wang, Y.; Yao, Q.; Kwok, J.; Ni, L. M. Generalizing from a Few Examples: A Survey on Few-Shot Learning. 2020.

Granger, C. W. J. Invited review combining forecasts—twenty years later. *Journal of Forecasting* **1989**, *8*, 167–173.

Wallis, K. F. Combining forecasts – forty years later. *Applied Financial Economics* **2011**, *21*, 33–41.

## A Zero-Bias Networks

Let a network $\underline{\mathbf{N}}(x)$ and a training set $T_r = \{x_i, y_i = y(x_i)\}_{i=1,M}$. Then the network

$$N(x) \equiv \underline{\mathbf{N}}(x) - \langle \underline{\mathbf{N}} \rangle + \langle y \rangle, \tag{33}$$

is by construction a zero-bias network on $\{T_r\}$. Clearly the misfit is

$$m(x) = N(x) - y(x) = \underline{\mathbf{N}}(x) - \langle \underline{\mathbf{N}} \rangle + \langle y \rangle - y(x), \text{ satisfying } \langle m \rangle = 0. \tag{34}$$

At this point one may rewrite the misfit as:

$$m(x) = (\underline{\mathbf{N}}(x) - y(x)) - (\langle \underline{\mathbf{N}} - y \rangle) \equiv \underline{\mathbf{m}}(x) - \langle \underline{\mathbf{m}} \rangle, \tag{35}$$

and the MSE becomes:

$$\langle m^2 \rangle = \langle (\underline{\mathbf{m}} - \langle \underline{\mathbf{m}} \rangle)^2 \rangle = \langle \underline{\mathbf{m}}^2 \rangle - \langle \underline{\mathbf{m}} \rangle^2. \tag{36}$$

Thus minimizing $\langle m^2 \rangle$ subject to $\langle m \rangle = 0$, is equivalent to minimizing $\langle \underline{\mathbf{m}}^2 \rangle - \langle \underline{\mathbf{m}} \rangle^2$, subject to no constraints. The zero-bias network $N(x)$, is then recovered from eq. (33).

## B Blending Restriction

The blending coefficient may be further restricted to lay in a shorter interval, e.g.

$$\epsilon \leq \beta_{k+1} \leq 1 - \epsilon, \quad \text{with } \epsilon \in [0, 0.5) \tag{37}$$

This imposes a slightly different constraint on the correlation $\langle M_k m_{k+1} \rangle$, instead of (25), i.e.:

$$\langle M_k m_{k+1} \rangle \leq \frac{1}{1 - 2\epsilon} \left[ \min\{\langle M_k^2 \rangle, \langle m_{k+1}^2 \rangle\} - \epsilon \left( \langle M_k^2 \rangle + \langle m_{k+1}^2 \rangle \right) \right]. \tag{38}$$

A typical value, $\epsilon = 0.1$ is restricting $\beta_{k+1} \in [0.1, 0.9]$.