# Coopetitive Soft Gating Ensemble

Stephan Deist, Maarten Bieshaar, Jens Schreiber, André Gensler, and Bernhard Sick
Intelligent Embedded Systems Group, University of Kassel (Germany)
email: {stephan.deist, mbieshaar, jens.schreiber, gensler, bsick}@uni-kassel.de

*Abstract*—In this article, we propose the *Coopetititve Soft Gating Ensemble* or *CSGE* for general machine learning tasks and interwoven systems. The goal of machine learning is to create models that generalize well for unknown datasets. Often, however, the problems are too complex to be solved with a single model, so several models are combined. Similar, Autonomic Computing requires the integration of different systems. Here, especially, the local, temporal online evaluation and the resulting (re-)weighting scheme of the CSGE makes the approach highly applicable for self-improving system integrations. To achieve the best potential performance the *CSGE* can be optimized according to arbitrary loss functions making it accessible for a broader range of problems. We introduce a novel training procedure including a hyper-parameter initialisation at its heart. We show that the *CSGE* approach reaches state-of-the-art performance for both classification and regression tasks. Further on, the *CSGE* provides a human-readable quantification on the influence of all base estimators employing the three weighting aspects. Moreover, we provide a scikit-learn compatible implementation.

## I. Introduction

The primary goal of *machine learning (ML)* is to create models from training data, which have a high generalization capability for unseen data. Often the problems are so complex that one single *estimator* cannot handle the whole scope. These problems can, e.g., be tackled with a combination of multiple estimators instead of an individual estimator. This attempt of combining multiple estimators is called *ensemble*. In many fields, ensembles can achieve state-of-the-art performance. Popular ensemble methods are *Boosting* [1], *Bagging*, [2] or *Stacking* [3]. [4] shows that ensembles often lead to better results than using a single estimator. When considering the *Biasvariance tradeoff* [5], ensembles can reduce both variance and bias and therefore result in stronger models.

The combination of different models in an ensemble can be compared with the integration of different systems. Due to the common weighting of different models, the prediction is decisively determined by the mutual influence of individual models, as with interwoven systems [6]. In addition to this similarity, models are usually heterogeneous, e.g., linear and non-linear models. Ultimately, similar to interwoven system predictions are linked to uncertainty [6]. Recently, an ensemble method called *Coopetitive Soft Gating Ensemble* or *CSGE* was proposed for wind power forecasts. However, CSGE also offers a basic technology to make the functionality of *self-improving system integration (SISSY)* more robust against interference and to better avoid uncertainty by its weighting and optimization scheme [7].

In [8], [9], and [10] it is statistically shown that the CSGE can achieve state-of-the-art performance in the area of power forecasting. In this article, we aim to extend the original

approach to general *ML* problems to show its potential for a wide range of problems including SISSY applications. The idea of the CSGE is to gradually weight individual ensemble members according to their historically observed performance of different aspects. In particular, there are three aspects which take influence on the weight: First, the overall performance of the estimator. Second, the local performance of the estimator in similar historical situations. Third, time-dependent effects modeling the autocorrelation in the estimator's outcome. Aspect two and three are optimized w.r.t (current) online data, and the first aspect is determined based on the training data beforehand.

Hence, the ability to assess the individual performance of each base estimator for those three aspects can be interpreted, regarding Organic Computing (OC) [11] and Autonomic Computing (AC) [12], as self-awareness and self-improving capabilities for SISSY.

## II. Main Contribution

The main contribution of this article is an extended coopetitive soft gating ensemble approach. It generalizes the original CSGE method proposed in [8], [9] for wind power forecasting to other *ML* tasks including regression, classification, and time series forecasting. The main contributions of this article are:

- The loss function of the *CSGE* can be chosen by the user with only minimal constraints allowing optimization of arbitrary loss functions to make it available for SISSY and *ML* applications.
- A novel heuristic to choose the hyper-parameters of the *CSGE* training algorithm is reducing the required number of adjustable parameters.
- An extensive evaluation of our approach on common real-world reference datasets, in which we show that our *CSGE* approach reaches state-of-the-art performance compared to other ensembles methods. Additionally, the *CSGE* allows quantifying the influence of all base estimators utilizing the three weighting aspects in a human-readable way.
- A scikit-learn compatible implementation of the *CSGE*[1].

The remainder of this article is structured as follows. In Section III, we review the related work in the field of ensemble method for *ML*. Afterward, in Section IV, we introduce our *CSGE* approach. In Section V, we present the evaluation of our *CSGE* on three synthetic datasets, four reference classification, and real-world regression datasets. Therefore, showing its applicability to a wide range of problems. Finally, in Section VI, the conclusion and open issues for future work are discussed.

---

[1] https://git.ies.uni-kassel.de/csge/csge

## III. RELATED WORK

The following section limits the discussion of related work to ensemble methods; this allows better comparability of the CSGE compared to self-improving systems. In *ML* the term *ensemble* describes the combination of multiple models. The ensemble comprises a finite set of estimators, whose predictions are aggregated forming the ensemble prediction. The theoretical justification of why ensembles can increase the overall predictive performance is given by the bias-variance decomposition [5]. The key to ensemble methods is model diversification, i.e., how to create sufficiently different models from sample data. A comprehensive review of ensembles is given in [13]. The most important design principles for ensembles are: Data, parameter, and structural diversity. Data diversity comprises ensembles trained on different subsets of the data. Well known representatives of this type are bagging [2], boosting [1], and random forest [14]. The idea of parameter diversity is to induce diversity into the ensemble by varying the parameters of the ensemble members.

A representative of this type is the multiple kernel learning algorithm [15] in which multiple kernels are combined. Lastly, structural diversity comprises the combination of different models, e.g., obtained by applying different learning algorithms or variable model types. These ensembles are also referred to as heterogeneous ensembles [16]. A well-known representative of this type is the stacking algorithm [3]. Another ensemble technique is Bayesian model averaging (BMA) [17] accounts for this model uncertainty when deriving parameter estimates. Hence, the ensemble estimate comprises the weighted estimate of the various model hypothesis. Another method not to be confused with BMA is Bayesian model combination [18]. It overcomes the shortcoming of BMA to converge to a single model. Recently, a mixture of expert models, which comprise a gating model weighting the outputs of different submodels, gained much attention, as they determine state-of-the-art performance in language modelling [19] and multi-source machine translation [20]. These approaches are based on deep neural networks. Hence, they require many training samples and their weightings be barely interpretable.

In [8], [9], the CSGE was presented in the context of renewable energy power forecasting. It comprises a hierarchical two-stage ensemble prediction system and weights the ensemble member's predictions based on three aspects, namely global, local, and time-dependent performance. In [10], the system was extended to handle probabilistic forecasts. The approach presented in this article is a generalization of the approach to other *ML* tasks.

## IV. METHOD

In this section the novel *Coopetitive Soft Gating Ensemble* method or short *CSGE*, as proposed in [10], is introduced. After a brief general overview, we detail the different characteristics of the ensemble method namely soft gating, global-, local- and time-dependent-weighting. In the final sections, we give details on the (self-) optimization process and recommendations for training.
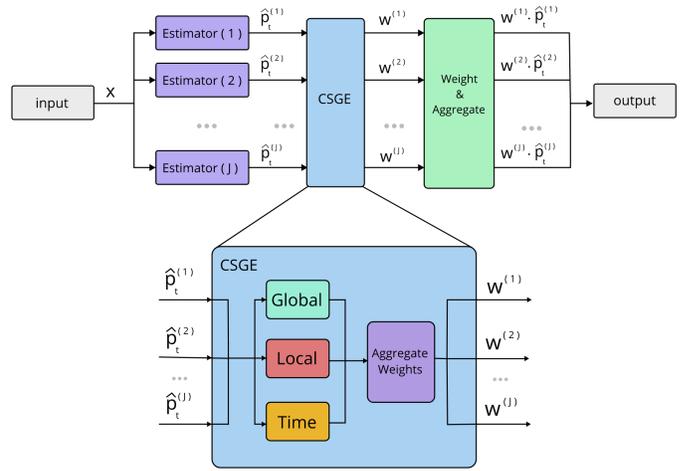


Fig. 1. The architecture of the *CSGE*. The predictions $\hat{p}_t^{(j)}$ of the input $x$ are passed to the CSGE module. Weights are calculated regarding *global-*, *local-* and *time-dependent weighting*. In the next step, the predictions are weighted and aggregated.

### A. Coopetitive Soft Gating Ensemble

The architecture of the *CSGE*, as depicted in Fig. 1, highlights the three weighting aspects: global-, local- and time-dependent-weighting. For each of the weighting methods the novel *coopetetive soft gating* principle is applied. *Coopetetive soft gating* is a conglomerate of cooperation and competetion. The ensemble combines two well known principles in ensemble methods, *weighting* and *gating*. *Weighting* combines all ensemble members in a linear combination, while *gating* selects only one of all ensemble members. The idea of the CSGE is to have the possibility to have a mixture of both *weighting* and *gating* and let the ensemble optimize which concept to use for the combination of different predictions.

Each of the three weighting aspects is calculate by the the predictions from *J*-ensemble members. Each ensemble member provides estimations $\hat{p}_t^{(j)}$ for the input $x$. $t$ denotes the timestamp $t$, also called *leadtime*, when operating on timeseries for the $j$-th ensemble member. For each prediction and estimator the *CSGE* calculates the local, global and time-dependent *weighting* and aggregates their results. After normalization of $\omega^{(j)}$ each prediction $\hat{p}_t^{(j)}$ is weighted to obtain the final prediction as follows:

$$\bar{\bar{p}}_t = \sum_{j=1}^{J} \omega^{(j)} \cdot \hat{p}_t^{(j)} \tag{1}$$

To ensure that the prediction is not distorted weights have the following constraint:

$$\sum_{j=1}^{J} \omega^{(j)} = 1 \tag{2}$$

The optimal weights $w^{(j)}$ with $j \in \{1, \ldots, J\}$ are obtained by the *CSGE* w.r.t. an arbitrary loss function, e.g. *mean squared error*, *cross-entropy* etc. Each weighting aspect has different characteristics related to the loss function summarised as follows

- **Global weights** are determined by observed training performance for each ensemble member and is a fixed weighting after training. Thereby, overall strong models have more influence than weaker models.
- **Local weighting** considers the fact that different ensemble members have various prediction quality over the complete feature space. As an example, when considering the problem of renewable energy prediction, an ensemble member could perform well on rainy weather inputs but has worse quality when using sunny weather inputs. Therefore, the *local weighting* rewards ensemble members with a higher weighting, which performed well on similar input data. These weights are adjusted online for each prediction during runtime.
- The **time-dependent weight** aspect is used when performing predictions on time series. Ensemble members may perform differently for different lead times. E.g., one method might achieve superior results on short time horizons, while losing quality for larger lead times. Other methods may perform worse on short time horizons, but have greater stability on larger lead times. Again, these weights are calculated online for each prediction during runtime.

To combine these three weighting aspects for an individual ensemble member we use the the following equation:

$$\omega^{(j)} = \omega_g^{(j)} \cdot \omega_l^{(j)} \cdot \omega_k^{(j)} \tag{3}$$

where $\omega_g^{(j)}$ is the *global weighting*, $\omega_l^{(j)}$ is the *local weighting* and $\omega_k^{(j)}$ is the *time-dependent weighting*. To calculate the final weighting the values are normalized for the $j$-th ensemble member $\omega^{(j)}$ as follows:

$$\omega^{(j)} = \frac{\omega^{(j)}}{\sum_{\tilde{j}=1}^{J} \omega^{(\tilde{j})}}. \tag{4}$$

This equation ensures that constraint of Eq. 2 is fulfilled.

### B. Soft Gating Principle

The primary goal of the *CSGE* is to increase the quality of the prediction by weighting robust predictors greater than predictors with worse quality results. Traditionally in ensemble methods, one of the two paradigms *weighting* or *gating* are used to combine individual ensemble members. The *soft gating* approach of the CSGE introduces a novel method, which allows the mixture of both *weighting* and *gating* and a (self-) optimization process to select the optimal combination of different predictions. Moreover, the soft gating approach applies to all three weighting aspects.

To evaluate the quality of an individual ensemble member, we need to relate the error of the prediction to its respective weighting. This mapping is achieved by the function $\varsigma'_\eta(\Omega, \rho)$ to determine the weights of the estimator $j$ as follows:

$$\varsigma'_\eta(\Omega, \rho) = \frac{\sum_{j=1}^{J} \Omega_j}{\rho^\eta + \epsilon}, \eta \in \mathbb{R}_0^+. \tag{5}$$

$\Omega$ contains reference errors of all $J$ estimators, while $\rho$ is the individual error of the estimator $j$; the user chooses parameter
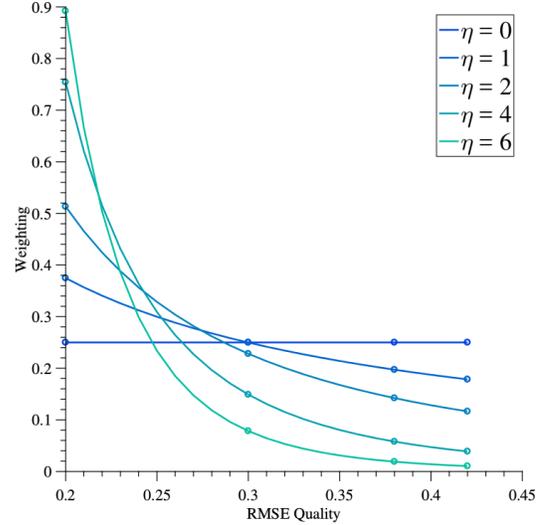


Fig. 2. The error (RMSE) of a predictor is drawn on the x-axis, while the y-axis contains the corresponding weights computed by $\varsigma'_\eta$. For greater $\eta$ a higher error gets more regulated with less weighting, than for smaller $\eta$.

$\eta$. It controls the linearity of the weighting. For greater $\eta$ the *CSGE* tends to work as *gating*, while smaller $\eta$ results in a *weighting* approach.

In Fig. 2 we observe the following characteristics of $\varsigma'_\eta$:

- $\varsigma'_\eta$ is falling monotonously.
- $\varsigma_\eta$ returns smaller weightings for an ensemble with larger errors $\rho$.
- For $\eta = 0$ every ensemble members are weighted with $\frac{1}{J}$, due to the later explained normalization. Respectively, disrespecting the individual errors.
- $\epsilon$ is a small constant to prevent a division by zero.

To ensure that $\sum_{j=1}^{J} \omega^{(j)} = 1$, $\varsigma'_\eta(\Omega, \rho)$ is adjusted in the following way

$$\varsigma_\eta(\Omega, \rho) = \frac{\varsigma'_\eta(\Omega, \rho)}{\sum_{j=1}^{J} \varsigma'_\eta(\Omega, \rho_j)} \tag{6}$$

Besides the advantage on having only one parameter ($\eta$) to tune, the soft gating offers a direct relation between the weighting and the errors of the ensemble members providing a strong correlation to the actual data.

### C. Global Weighting

The *global weighting* is calculated during ensemble training and then remains constant. Ensemble members that perform well on the training data get larger weights compared to those who showed a worse performance. Therefore, the difference between estimation and ground truth is calculated with

$$e_n^{(j)} = \theta(\hat{y}_n^{(j)}, y_n) \tag{7}$$

$\hat{y}_n^{(j)}$ is the prediction of the $j$-th ensemble member, while $y_n$ is the corresponding ground-truth. $\theta$ is an a arbitrary scoring function, which could for example be the *root mean-squared error (RMSE)* for regression or the *accuracy score (ACC)* for classification. The only condition is that the loss needs to be

falling monotonously with increasing errors to work correctly with the *soft gating principle*, see Eq. 6. The error score $R^{(j)}$ of the $j$-th ensemble member is calculated by:

$$R^{(j)} = \frac{1}{N} \cdot \sum_{n=1}^{N} (e_n^{(j)}) \tag{8}$$

$$R = (R^{(1)}, \ldots, R^{(j)}, \ldots, R^{(J)}) \tag{9}$$

By applying the *soft gating* principle to the vector $R$ of all error scores of the $J$ ensemble members we obtain the final global weighting with

$$w_g^{(j)} = \varsigma_\eta(R, R^{(j)}) \tag{10}$$

### D. Local Weighting

The *local weighting* considers the quality difference between the predictors for distinct situations over the whole *feature space*. Therefore, the *local weighting* rewards ensemble members with a higher weighting, which performed well on similar input data. In contrast to the *global weighting* the *local weighting* is calculated online for each estimation during runtime.

For similar situations, we consider the distances in the input feature space. Therefore, we assume situations with low distance have more in common compared to situations with a more significant distance. $X_H$ contains all data that is used during ensemble training. Often the features of $X_H$ vary in their ranges and information value. Since we use the distances of features to determine situations which are similar, it can be useful to apply a *principal component analysis* (*PCA*) on the training data $X_H$.

$$X_{H_{PCA}} = PCA(X_H, N_{dim}) \tag{11}$$

$X_{H_{PCA}}$ is the transformed training dataset, which has a dimension of $N_{dim}$. The parameter is chosen by the user and is in the range of $1, \ldots, N_f$, where $N_f$ is the number of features of $X_H$. To calculate the local weight of a new prediction, we have to transform the input data $x$ into the transformed feature space $\hat{x}$ by applying the PCA:

$$\hat{x} = applyPCA(x) \tag{12}$$

By using, e.g., *k-nearest neighbor* we determine $c$ similar situations in the input data.

$$\alpha = knn(\hat{x}, X_{H_{PCA}}, c) \tag{13}$$

The vector $\alpha$ of similar situation in the input data is used to derive the errors for each situation $a$ with $e_a^j = \theta(\hat{y}_a^{(j)}, y_n)$ to obtain the average local error with:

$$q^{(j)} = \frac{1}{c} \cdot \sum_{a \in \alpha} |e_a^{(j)}| \tag{14}$$

This equation is applied to each ensemble member to obtain all local error scores $Q$ for all $J$ ensemble members.

$$Q = (q^{(1)}, \ldots, q^{(j)}, \ldots, q^{(J)}) \tag{15}$$

Finally, the local weight $\omega_l^{(j)}$ is calculated by using the *soft gating* principle to derive the best possible local weighting:

$$\omega_l^{(j)} = \varsigma_\eta(Q, q^{(j)}) \tag{16}$$

### E. Time-Dependent Weighting

The *time-dependent weighting* considers the fact that the quality of an ensemble member varies over leadtime. Similar to *local weighting*, *time-dependent weighting* is calculated for each estimation. $\hat{Y}^{(T,j)}$ contains all predictions of estimator $j$ starting at time $t = 0$ to time $t = T$.

$$\hat{Y}^{(T,j)} = (\hat{y}_0^{(j)}, \hat{y}_1^{(j)}, \ldots, \hat{y}_T^{(j)}) \tag{17}$$

The error for a specific time $t \in \{0, 1, \ldots, T\}$ is calculated by the average error over all training samples with:

$$R_t^{(j)} = \frac{1}{N} \cdot \sum_{n=1}^{N} e_n^{(t,j)} \text{ and} \tag{18}$$

$$e_n^{(t,j)} = \theta(\hat{y}_n^{(t,j)}, y_n^{(t)}) \tag{19}$$

With $\hat{y}_n^{(t,j)} \in \hat{Y}^{(T,j)}$ and $y_n^{(t)}$ as ground truth for time $t = t$. To estimate the error score for time $t$ of estimator $j$, we use the following equation:

$$r_t^{(j)} = \frac{R_t^{(j)}}{\frac{1}{t+1} \cdot \sum_{t*=0}^{T} R_{t*}^{(j)}} \tag{20}$$

$r_t^{(j)}$ is a measure that compares the error of the prediction with $t = t$ to the average error in the time interval $t \in \{0, 1, \ldots, T\}$. The weight $\omega_k^{(j)}$ is calculated analogous to *global-* and *local weighting* using the soft gating principle with

$$P_t = (r_t^{(1)}, \ldots, r_t^{(j)}, \ldots, r_t^{(J)}), \tag{21}$$

$$\omega_k^{(j)} = \varsigma_\eta(P_t, r_t^{(j)}), \tag{22}$$

to derive the potentially best time-dependent weighting.

### F. Model Fusion and Ensemble Training

To find the optimal set of parameters for the predictions (including all weighting aspects) we aim to optimize the prediction of Eq. 1. Since there are three aspects, *global-*, *local-* and *time-dependent weighting*, it follows that there are also three $\eta = (\eta_0, \eta_1, \eta_2)$ to be chosen. As mentioned previously the parameter $\eta$ is chosen by the user and controls the non-linearity of the system. Therefore, the following minimization problem solves the task to adjust $\eta$ with:

$$\sum_{n=1}^{N} [y_n - f_{CSGE}(x_n, \eta)]^2 + c \cdot \sum_{s=1}^{3} \eta_s, \tag{23}$$

where $f_{CSGE}(x_n, \eta) = \bar{\hat{p}}_t$ is the prediction from Eq. 1 given its current weights. $\sum_{n=1}^{N} [y_n - f_{CSGE}(x_n, \eta)]^2$ are the summed errors of the training data, while $c \cdot \sum_{s=1}^{3} \eta_s$ is a regularisation term to control overfitting.

Fig. 3. The two training sets must be distinct in order to get information about the quality of each ensemble member.

However, to optimize Eq. 23, adjust $\eta = (\eta_0, \eta_1, \eta_2)$ and calculate the *global weighting*, we need training data $E_T$. In general, the ensemble members are trained on a training dataset and validated on a validation dataset. By using the same training dataset to train the *CSGE* it will often become overfitted and not generalize well. Therefore, we need training data for the *CSGE* that is not used to train the $J$ ensemble members. A simple Method is shown in Fig. 3. The training data gets split into two sets of data. One to train the ensemble members and one to train the *CSGE* itself. Even though the setup is straightforward, it has a disadvantage. The training data is wasted because the training data for the ensemble members and the one for the *CSGE* need to be distinct.

A more advanced approach shown in Fig. 4, allows using the training data more efficiently. Since the *CSGE* uses the output data of the predictors we need those data for training. Therefore, a *cross validation* with K-folds is used to generate this data. The training data in the $k$-th step of this k-fold is split in a set $E_T^{(k)}$ for training and a set $E_P^{(k)}$ for prediction. Then a copy of the $j$-th ensemble member is trained by using the set $E_T^{(k)}$. This temporary predictor is denoted with $f_j^{(k)}$, where $k$ is the $k$-th step and $j$ the indices for the $j$-th ensemble member. The temporary predictor $f_j^{(k)}$ is used to predict $E_{P_j}^{(k)}$, to concatenate all predictions in $k$-iterations.

Afterward, all $J$ ensemble members are trained on the whole training set. The training data now consists of the output data $E_{P_j}$ of the estimators; this requires us to adjust the calculation of the *CSGE*. Therefore, we have to store the predictions in an $N \times J$ dimensional matrix, where $N$ is the number of samples and $J$ the number of estimators. $t$ is the timestamp when operating on time series.

$$C_t = \begin{bmatrix} C_{(t,0)}^{(0)} & \cdots & C_{(t,0)}^{(J)} \\ \vdots & \ddots & \vdots \\ C_{(t,N)}^{(0)} & \cdots & C_{(t,N)}^{(J)} \end{bmatrix} \qquad (24)$$

Now, we have to adjust the Eq. 7, in which the difference between prediction and ground truth is calculated. We can use $t = 0$ since *global-* and *local weighting* do not consider the time aspect.

$$e_n^{(j)} = \theta(C_{(0,n)}^{(j)}, y_n) \qquad (25)$$

Eq. 17, where the set $\hat{Y}_n^{(t,j)}$ is defined, which contains all predictions of the training point $n$ of the ensemble member $j$ over the timerange 0 to $t$.

$$\hat{Y}_n^{(T,j)} = (C_{(0,n)}^{(j)}, C_{(1,n)}^{(j)}, \ldots, C_{(T,n)}^{(j)}) \qquad (26)$$

### G. Regularisation Heuristic

The ensemble learning tends to choose high $\eta$ for one single aspect and therefore $\eta = 0$ for other aspects. As an



Fig. 4. In the $k$-th step we divide the the training data in a distinct set $E_T^{(k)}$ and $E_P^{(k)}$. $E_T^{(k)}$ is used to train a copy of the ensemble members, while $E_P^{(k)}$ is used to make predictions. After the $k$-th iteration every element of our training data is predicted and we can use these predictions for ensemble training.
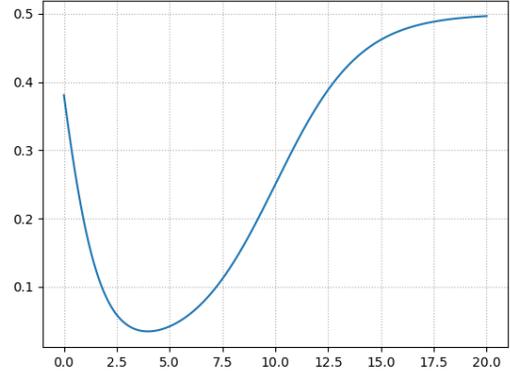


Fig. 5. $a(x)$ weights the chosen $\eta$'s to avoid choosing too high or too low values for $\eta$.

example, the $\eta$ for *local weighting* often are chosen very high. This example means that the local aspect of the *CSGE* works as a selecting ensemble, which chooses one of the $J$ ensemble members. In order to minimise the regularisation term $c \cdot \sum_{s=1}^{S} \eta_s$, the $\eta$ of *global-* and *time-dependent weighting* is chosen very low. This regularization leads to an averaging ensemble for the global and time-dependent aspect, that weights all $j$ ensemble member equally with $\frac{1}{j}$, which disables these two aspects. Even though it can be necessary to disable some aspects, it is often better to distribute the values for the $\eta$'s more evenly, to get a more generalized ensemble model. We propose the function $a(x)$ to prevent this problem. $a(x)$ weights the $\eta$ and penalises when choosing $\eta = 0$ or $\eta$ too high. Typically the parameter $\eta$ lies in the range of $1.0 \leq \eta \leq 6.0$ [8].

$$a(x) = \frac{1}{1 + e^{-\frac{1}{2} \cdot (x - 10)}} + \frac{1}{2 \cdot (1 + e^{x^{\frac{1}{2}}})} \qquad (27)$$

The minimization Eq. 23 must be adjusted in the following way

$$\sum_{n=1}^{N} [y_n - f_{CSGE}(x_n, \eta)]^2 + c \cdot \sum_{s=1}^{3} a(\eta_s) \qquad (28)$$

## V. EXPERIMENTAL EVALUATION

In this Section, we present the evaluation of the *CSGE*. We split the evaluation into two steps. First, we show the proper functionality of each of the three different weighting aspects with a distinct synthetic dataset to show its interpretability
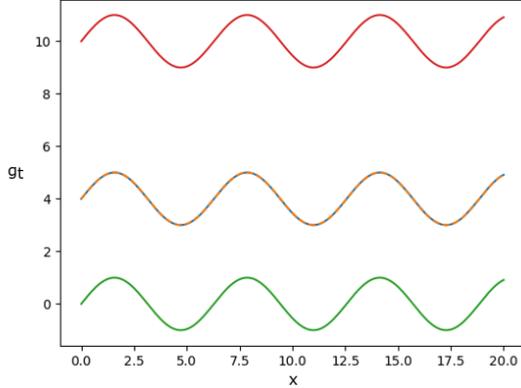
Fig. 6. $g_t$ is the orange function, $f_1$ is the green function and $f_2$ is drawn in red. The predicted values by the *CSGE* are drawn blue dotted (partly showing the function $f_1$, $f_2$ and $g_t$ below), which exactly match the yellow $g_t$.
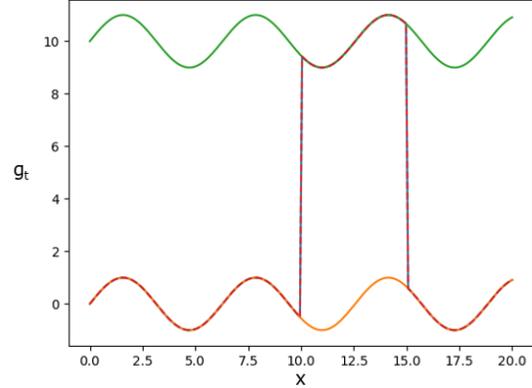


Fig. 7. $g_t$ is the red function, $f_1$ is the orange function and $f_2$ is drawn in green. The predicted values by the *CSGE* are drawn blue dotted (partly showing the function $f_1$, $f_2$ and $g_t$ below), which exactly match the red $g_t$.

aspects. Second, we evaluate the *CSGE* on four real-world reference datasets for both regression and classification. The evaluation includes a comparison with other state-of-the-art ensemble methods for a wide range of problems and shows the potential as a basic technique for SISSY systems.

*A. Synthetic Datasets*

We created synthetic datasets in order to evaluate each aspect of the *CSGE*, i.e., *global-*, *local-* and *time dependent weighting*, separately. For each synthetic dataset, we created a data generating function $g_t$. Since we are interested in the general interpretation and functionality, we do not consider any additional noise. Furthermore, we defined mathematical estimators $f_j$ which have to be combined by the *CSGE* properly to match the function $g_t$.

*1) Global Weighting:* For evaluation of the *global Weighting*, we created $g_t$ in the following way:

$$g_t(x) = sin(x) + 4 \tag{29}$$

We use two estimators as ensemble members who are defined as follows:

$$f_1(x) = sin(x) \text{ and } f_2(x) = sin(x) + 10$$

The result after training the *CSGE* is depicted in Fig. 6. Since there are neither local nor time-dependent aspects, the learning procedure chooses $\eta_{time} = 0$ and $\eta_{local} = 0$. When interpreting the chosen $\omega_1$ and $\omega_2$ from a mathematical point of view we observe that the chosen weights are correct, i.e., $0.6 \cdot sin(x) + 0.4 \cdot (sin(x) + 10)$. The *CSGE* perfectly matches $g_t(x)$.

*2) Local Weighting:* For evaluation of the *Local Weighting* we created $g_t$ in the following way:

$$g_t(x) = \begin{cases} sin(x) & x < 10 \\ sin(x) + 10 & 10 \le x \le 15 \\ sin(x) & x > 15 \end{cases} \tag{30}$$

We use two estimators as ensemble members who are defined as follows:

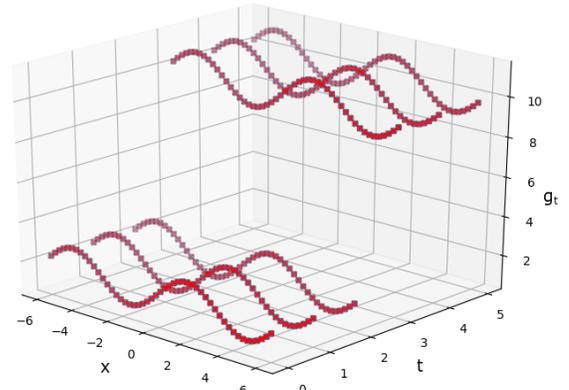$$f_1(x) = sin(x) \text{ and } f_2(x) = sin(x) + 10$$



Fig. 8. $g_t$ is drawn with blue points and the predictions using the *CSGE* are drawn with red crosses which exactly match the $g_t$.

This experiment has no global and time dependent aspect, therefore the learning algorithm chooses $\eta_{global} = 0$ and $\eta_{time} = 0$. Since $g_t$ can only be approximated by picking either $f_1$ or $f_2$ depending on the feature space, the chosen $\eta_{local}$ should be larger than zero. In Fig. 7, we see the results of the experiment. We observe, that the *CSGE* is able to perfectly reconstruct the reference model $g_t(x)$.

*3) Time-dependent Weighting:* For evaluation of the *time-dependent Weighting* we created $g_t$ in the following way:

$$g_t(x,t) = \begin{cases} sin(x) & t < 3 \\ sin(x) + 10 & t \ge 3 \end{cases} \quad t \in \mathbb{N} \tag{31}$$

We use two estimators as ensemble members who are defined as follows:

$$f_1(x,t) = sin(x) \text{ and } f_2(x,t) = sin(x) + 10$$

Fig. 8 shows the results of the experiment. Since there is no global and local aspect, the learning algorithm picks $\eta_{global} = 0$ and $\eta_{local} = 0$, $\eta_{time} > 0$. We observe that after training, *CSGE* perfectly matches the reference function. These evaluations on synthetic data show that the *CSGE* works properly.

#### TABLE I
#### RMSE ON BOSTON HOUSING DATASET
**Ensemble Members**

|  | Linear Regression | SVM | Decision Tree |
|---|---|---|---|
| Mean | 24.6673 | 82.3656 | **21.7313** |
| Standard Deviation | **5.8063** | 14.7024 | 9.9349 |
| Minimum | 17.2139 | 66.0964 | **13.6710** |
| Maximum | **33.9569** | 107.5636 | 46.7434 |

**Ensemble Methods**

|  | CSGE | Linear Stacking | ANN Stacking | Averaging |
|---|---|---|---|---|
| Mean | 18.9079 | **15.9885** | 18.1849 | 23.2753 |
| Standard Deviation | 8.7271 | **5.5606** | 5.6026 | 8.3834 |
| Minimum | **9.8018** | 10.9614 | 13.7156 | 15.7049 |
| Maximum | 34.9028 | **27.6670** | 32.0342 | 39.4708 |

#### TABLE II
#### RMSE ON DIABETES DATASET
**Ensemble Members**

|  | Linear Regression | SVM | Decision Tree |
|---|---|---|---|
| Mean | **3083.1198** | 6356.0135 | 6518.2421 |
| Standard Deviation | **322.2800** | 405.2114 | 728.0636 |
| Minimum | **2641.9339** | 5620.8063 | 5487.6391 |
| Maximum | **3419.9466** | 6837.5063 | 7819.4436 |

**Ensemble Methods**

|  | CSGE | Linear Stacking | ANN Stacking | Averaging |
|---|---|---|---|---|
| Mean | 3333.9250 | **3099.9531** | 3465.7875 | 3916.1540 |
| Standard Deviation | 453.4425 | **323.7526** | 553.6934 | 364.5056 |
| Minimum | 2738.1390 | **2664.3987** | 2795.0731 | 3380.1302 |
| Maximum | 4273.5943 | **3459.9381** | 4878.9214 | 4392.5727 |

#### TABLE III
#### ACCURACY ON IRIS DATASET
**Ensemble Members**

|  | Linear Classifier | SVM | Decision Tree |
|---|---|---|---|
| Mean | 0.6000 | **0.9711** | 0.9333 |
| Standard Deviation | 0.2071 | 0.0183 | **0.0181** |
| Minimum | 0.2889 | **0.9333** | 0.9111 |
| Maximum | 0.9778 | **1.0000** | 0.9556 |

**Ensemble Methods**

|  | CSGE | Linear Stacking | ANN Stacking | Voting |
|---|---|---|---|---|
| Mean | **0.9578** | 0.6756 | 0.9378 | 0.9511 |
| Standard Deviation | 0.0221 | 0.1582 | 0.0888 | **0.0204** |
| Minimum | **0.9333** | 0.4000 | 0.6889 | **0.9333** |
| Maximum | **0.9778** | 0.9333 | 0.9778 | **0.9778** |

### B. Real-world Regression Datasets

In order to evaluate the *CSGE* on regression problems, we chose Boston Housing[2] and Diabetes datasets[2]. As ensemble members we used a *Support Vector Regression (SVR)* with radial basis function (RBF) kernel, a *Neural Network Regressor* and a *Decision Tree Regressor*. As composition proceeding to the *CSGE* we chose *Stacking* and *Voting* (i.e., *Averaging*). For *Stacking* we used a *Neural Network* (i.e., referred to as *ANN Stacking*) and a *Linear Regression* (i.e., referred to as *Linear Stacking*) as meta learner. We chose the RMSE loss to optimize the *CSGE*. For each dataset, we performed ten-fold cross-validation with ten different random seeds.

We used default model parameters for the ensemble members as supplied by the scikit-learn library, the parameters for the ensemble methods are optimized for each experiment. To adjust the regularisation parameter $C$ and the number of neighbors $k$ of the *CSGE*, we used a grid search. Since the layer size of the *ANN Stacking* also needs to be optimized, we applied a grid search, too. The *Linear Regression Model* has no hyper-parameters to be optimized. As reference to *CSGE* and *Stacking*, we used a simple *Averaging* approach.

*1) Boston Housing:* The overall result, i.e., RMSE, can be seen in Tbl. I. We observe, that both *CSGE* and *Stacking* achieve better results than each ensemble member. The *Stacking* approach with a *Linear Regression* meta learner achieves best results. Even though the *CSGE* has sligthly worse results compared to *Linear Stacking*, it has similar performance to *ANN Stacking*.

*2) Diabetes:* The overall result can be seen in Tbl. II. We can see, that every ensemble method achieved worse results compared to the best ensemble member (*Linear Regression*). The *Linear Stacking* accomplished the best results of all ensemble methods. Nevertheless, the *CSGE* performed better than *Averaging* and *ANN Stacking*.

### C. Real-world Classification Datasets

For each dataset, we performed ten-fold cross-validation with ten different random seeds.

In order to evaluate the *CSGE* on classification tasks, we chose Iris[3] and Wine[3] datasets. As ensemble members we used a *Support Vector Classification (SVC)* with linear and RBF

[2]http://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_[boston|diabetes].html (last accessed: 2018/06/25)

[3]http://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_[iris|wine].html (last accessed: 2018/06/25)

kernel and a *Decision Tree Classifier*. The SVC with linear kernel is referred to as *Linear Classifier*, while the SVC with RBF is referred to as *SVC*. As composition proceeding to the *CSGE* we chose *Stacking* and majority *Voting*. For *Stacking*, we used a *Neuronal Network* (i.e., referred to as *ANN Stacking*) and *SVC* with linear kernel (i.e., referred to as *Linear Stacking*) as meta learner. We chose the accuracy loss to optimize the *CSGE*.

As before with the regression, we used default model parameters for the ensemble members and only optimized the ensemble's parameters using a grid search. As a reference to *CSGE* and *Stacking*, we used the majority *Voting* ensemble.

### D. Iris

The overall results, i.e., classification accuracies, are depicted in Tbl. III. We observe that the *CSGE* is superior to both *Stacking* ensembles and *Voting*. All ensemble methods results are worse than the single ensemble member, i.e., *SVM*.

Fig. 9 shows the ROC curve of the iris dataset, we can see that the *CSGE* achieves the best results compared to *Stacking* and *Voting*.

### E. Wine

The resulting accuracies are depicted in Tbl. IV. We can see, that both *CSGE* achieved the best results compared to *Stacking* and *Voting*. Since the *Decision Tree* is by far best ensemble member, the *CSGE* worked as a gating ensemble by selecting the predictions of the *Decision Tree*, only.

Fig. 10 shows the ROC curve of the classifiers on the wine dataset. We observe that the *CSGE* achieves the best results compared to *Stacking* and *Voting*.
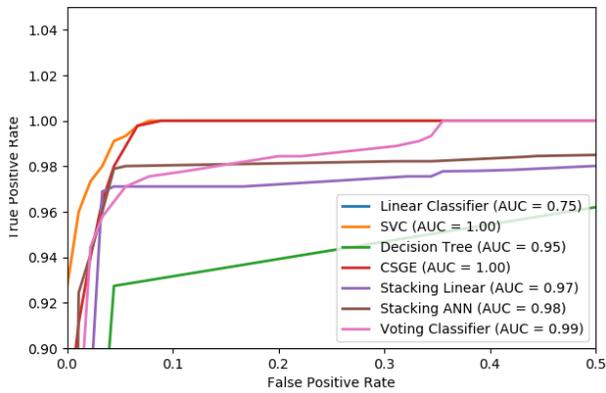
Fig. 9. ROC Curve of the different ensemble approaches and their ensemble members on the iris dataset. The Linear classifier not visible since it is below the clipping.

TABLE IV
ACCURACY ON WINE DATASET
**Ensemble Members**

|                    | Linear Classifier | SVM    | Decision Tree |
|--------------------|-------------------|--------|---------------|
| Mean               | 0.4944            | 0.4204 | **0.9148**    |
| Standard Deviation | 0.1197            | 0.0800 | **0.0265**    |
| Minimum            | 0.3704            | 0.3148 | **0.8704**    |
| Maximum            | 0.6852            | 0.5185 | **0.9444**    |

**Ensemble Methods**

|                    | CSGE       | Linear Stacking | ANN Stacking | Voting |
|--------------------|------------|-----------------|--------------|--------|
| Mean               | **0.9148** | 0.6907          | 0.8759       | 0.6704 |
| Standard Deviation | **0.0265** | 0.1615          | 0.0509       | 0.1776 |
| Minimum            | **0.8704** | 0.4444          | 0.7593       | 0.3704 |
| Maximum            | **0.9444** | **0.9444**      | **0.9444**   | 0.9259 |

## VI. CONCLUSION AND FUTURE WORK

In this article, we proposed the *CSGE* for general machine learning tasks and interwoven systems. The CSGE is an ensemble method which comprises human-understandable weightings based on the three basic aspects as there are *global-*, *local-* and *time-dependent* weights.

The CSGE can be optimized according to arbitrary loss functions making it accessible for a broader range of problems and provides a self-improving scheme based on previously seen data. This self-improving scheme can be applied to the self-integration problem and consequently constitutes a possible basic technique for SISSY systems as outlined in [6]. Moreover, we introduced a novel hyper-parameter initialization heuristics, enhancing the training process. We showed

the applicability and easy interpretability of the approach for synthetic datasets as well as real-world data sets. For the real-world datasets, we showed that our *CSGE* approach reaches state-of-the-art performance compared to other ensembles methods for both classification and regression tasks.

For future work, we intend to apply the approach to more real-world problems in various domains, such as trajectory forecasting of vulnerable road users, and further investigate its applicability in other domains of *AC*.

## VII. ACKNOWLEDGMENT

## REFERENCES

[1] R. E. Schapire, "The strength of weak learnability," *Machine Learning*, vol. 5, no. 2, pp. 197–227, 1990.
[2] L. Breiman, "Bagging predictors," *Machine Learning*, vol. 24, no. 2, pp. 123–140, 1996.
[3] P. Smyth and D. Wolpert, "Linearly combining density estimators via stacking," *Machine Learning*, vol. 36, no. 1, pp. 59–83, 1999.
[4] L. K. Hansen and P. Salamon, "Neural network ensembles," *TPAMI*, vol. 12, no. 10, pp. 993–1001, 1990.
[5] R. Kohavi and D. Wolpert, "Bias plus variance decomposition for zero-one loss functions," in *ICML*, vol. 13, Bari, Italy, 1996, pp. 275–283.
[6] S. Tomforde, S. Rudolph, K. L. Bellman, and R. P. Würtz, "An organic computing perspective on self-improving system interweaving at runtime," in *ICAC*, 2016, pp. 276–284.
[7] K. L. Bellman, S. Tomforde, and R. P. Würtz, "Interwoven systems: Self-improving systems integration," in *SASOW*, 2014, pp. 123–127. [Online]. Available: https://doi.org/10.1109/SASOW.2014.21
[8] A. Gensler and B. Sick, "Forecasting wind power - an ensemble technique with gradual coopetitive weighting based on weather situation," in *IJCNN*, Vancouver, BC, 2016, pp. 4976–4984.
[9] ——, "A multi-scheme ensemble using coopetitive soft gating with application to power forecasting for renewable energy generation," *CoRR*, vol. arXiv:1803.06344, 2018.
[10] ——, "Probabilistic wind power forecasting: A multi-scheme ensemble technique with gradual coopetitive soft gating," in *SSCI*, Honolulu, HI, 2017, pp. 1–10.
[11] C. Müller-Schloer, H. Schmeck, and T. Ungerer, Eds., *Organic Computing – A Paradigm Shift for Complex Systems*, ser. Autonomic Systems. Basel, Switzerland: Birkhäuser Verlag, 2011.
[12] J. O. Kephart and D. M. Chess, "The vision of autonomic computing," *Computer*, no. 1, pp. 41–50, 2003.
[13] Y. Ren, L. Zhang, and P. N. Suganthan, "Ensemble classification and regression-recent developments, applications and future directions," *CIM*, vol. 11, no. 1, pp. 41–53, 2016.
[14] L. Breimann, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
[15] M. Gönen and E. Alpaydn, "Multiple kernel learning algorithms," *J. Mach. Learn. Res.*, vol. 12, pp. 2211–2268, 2011.
[16] J. a. Mendes-Moreira, C. Soares, A. M. Jorge, and J. F. D. Sousa, "Ensemble approaches for regression: A survey," *ACM Comput. Surv.*, vol. 45, no. 1, pp. 10:1–10:40, 2012.
[17] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*, ser. Information Science and Statistics, M. Jordan, J. Kleinberg, and B. Schökopf, Eds. Secaucus, NJ: Springer-Verlag New York, 2006, vol. 1.
[18] K. Monteith, J. L. Carroll, K. Seppi, and T. Martinez, "Turning bayesian model averaging into bayesian model combination," in *IJCNN*, San Jose, CA, 2011, pp. 2657–2663.
[19] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean, "Outrageously large neural networks: The sparsely-gated mixture-of-experts layer," in *ICLR*, Toulon, France, 2017.
[20] E. Garmash and C. Monz, "Ensemble learning for multi-source neural machine translation," in *COLING: Technical Pap*, Osaka, Japan, 2016, pp. 1409–1418.
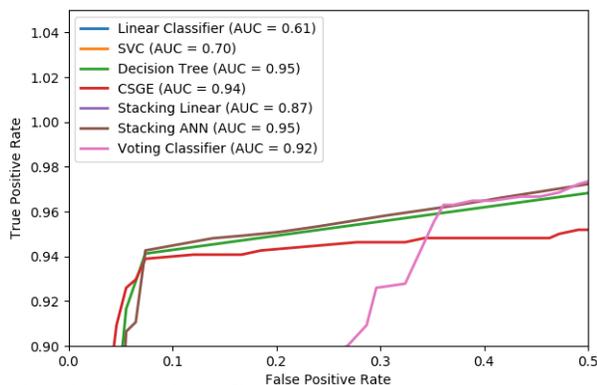


Fig. 10. ROC Curve of the different ensemble approaches and their ensemble members on the wine dataset. The Linear classifier is not visible since it is below the clipping.