

Generalization Error of Ensemble Estimators

Naonori UEDA and Ryohei NAKANO

NTT Communication Science Laboratories
Hikaridai Seika-cho Soraku-gun Kyoto 619-02 Japan
Tel: +81 774 95 1823 Fax: +81 774 95 1839
e-mail: ueda@cslab.kecl.ntt.jp

ABSTRACT

It has been empirically shown that a better estimate with less generalization error can be obtained by averaging outputs of multiple estimators. This paper presents an analytical result for the generalization error of ensemble estimators. First, we derive a general expression of the *ensemble generalization error* by using factors of interest (bias, variance, covariance, and noise variance) and show how the generalization error is affected by each of them. Then, some special cases are investigated. Next, the result of a simulation is shown to verify our analytical result. A practically important problem of the ensemble approach, *ensemble dilemma*, is also discussed.

1. Introduction

One of the most important issues to learning systems is to construct estimators with high generalization performance. The most popular way is to find the best model for a target task. In general, however, it is difficult to estimate the best model with given finite samples. On the other hand, as compared to that obtained by a single estimator [1]-[4], combining multiple regression estimators has been shown to improve generalization error. This combining approach has recently attracted major concern in the neural network community because of its simplicity and theoretical interest. This approach has also been successfully applied to pattern classification tasks [5][6].

The output of the combined regression estimator for some input is defined as the linear combination of outputs of multiple estimators. Recent work has provided optimal combination weights under some assumptions [3]. The most popular way to combine multiple estimators is the simple averaging of outputs of multiple estimators. Several theoretical studies on generalization error of these ensemble estimators have recently been made. Perrone was the first to show that the simple averaging of M estimators reduces the generalization error by a factor of M under uncorrelated and unbiased assumptions [2].

In this paper, we provide more general and explicit expressions and show how the generalization error of the ensemble estimator is affected by each of factors of interest (variance, covariance, bias, and noise variance). Some special cases where the models of estimators are the same and/or the outputs of estimators are uncorrelated are also investigated. Based on these expressions, we will provide some important observations for improving the generalization error due to the averaging. We verify our analytical result by showing the result of a simulation, using standard feedforward neural networks as regression estimators. An important open problem of the ensemble approach is discussed at the end of the paper.

2. Generalization Error of Single Estimator

The purpose of regression, particularly for nonlinear regressions such as feedforward neural networks, is to construct an estimator $f(\mathbf{x}; \Theta)$ that approximates an unknown target function $g(\mathbf{x})$. This is done by using a given set of N training samples $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$, where $\mathbf{x}_i \in \mathcal{R}^d, y_i \in \mathcal{R}$. Θ is a parameter vector and, in the case of a feedforward neural network, Θ corresponds to the weight vector. We assume that there is a functional relationship between the training pair \mathbf{x}_i and y_i : $y_i = g(\mathbf{x}_i) + \epsilon$, where ϵ is the additive noise with zero mean ($E\{\epsilon\} = 0$) and the finite variance ($\text{Var}\{\epsilon\} = \sigma^2 < \infty$). For convenience, let $z^N = \{z_1, \dots, z_N\}$ denote the training set, where $z_i = (\mathbf{x}_i, y_i)$. z^N is a realization of a random sequence

$Z^N = \{Z_1, \dots, Z_N\}$, whose i th component consists of a random vector $Z_i = (X_i, Y_i)$. In other words, each z_i is an independent and identically distributed (i.i.d.) sample from an unknown joint probability distribution $p(\mathbf{x}, y)$. Usually, given z^N , the parameter is estimated as

$$\hat{\Theta}(z^N) = \operatorname{argmin}_{\Theta} \sum_{i=1}^N (y_i - f(\mathbf{x}_i; \Theta))^2 / N. \quad (1)$$

Since the estimate $\hat{\Theta}$ depends on given z^N , we write $\hat{\Theta}(z^N)$ to clarify the dependency of z^N . Therefore, an output value of the estimator f for an input \mathbf{x} should be written as $f(\mathbf{x}; \hat{\Theta}(z^N))$; however for simplicity we write the output as $f(\mathbf{x}; z^N)$. Note that $f(\mathbf{x}; z^N)$ is also a realization of random variable $f(\mathbf{x}; Z^N)$.

Now, we introduce a new random vector $Z_0 = (X_0, Y_0) \in \mathcal{R}^{d+1}$, which has a distribution identical to that of Z_i , but is independent of Z_i for all i . Then, the *generalization error* (GErr) of the estimator f can be defined as the following mean squared error (MSE) averaged over all possible realizations of Z^N and Z_0 :

$$\text{GErr}(f) = \mathbb{E}_{Z^N} \{ \mathbb{E}_{Z_0} \{ [Y_0 - f(X_0; Z^N)]^2 \} \}. \quad (2)$$

Here $\mathbb{E}_{Z_0} \{ \}$ and $\mathbb{E}_{Z^N} \{ \}$ represent expectation¹ with respect to the distribution of Z_0 and Z^N , respectively. From the above definition, it is clear that $\text{GErr}(f)$ depends on neither particular training set z^N nor an unknown sample z_0 , and only depends on the sample size N and the model of the estimator f . With simple algebra, (2) can also be expressed by the following familiar “bias/variance” decomposition [7].

$$\text{GErr}(f) = \mathbb{E}_{X_0} \{ \text{Var}\{f|X_0\} + \text{Bias}\{f|X_0\}^2 \} + \sigma^2, \quad (3)$$

where $\text{Var}\{f|X_0 = \mathbf{x}_0\}$ and $\text{Bias}\{f|X_0 = \mathbf{x}_0\}$ are conditional variance and conditional bias given $X_0 = \mathbf{x}_0$:

$$\begin{aligned} \text{Var}\{f|X_0\} &= \mathbb{E}_{Z^N} \{ (f(X_0; Z^N) - \mathbb{E}_{Z^N} \{ f(X_0; Z^N) \})^2 \}, \\ \text{Bias}\{f|X_0\} &= \mathbb{E}_{Z^N} \{ f(X_0; Z^N) \} - g(X_0). \end{aligned} \quad (4)$$

Note that $\text{Var}\{f|X_0\}$ and $\text{Bias}\{f|X_0\}$ are random variables dependent on the random variable² X_0 .

3. Generalization Error of Ensemble Estimator

Let f_1, \dots, f_M denote M estimators, where the m th estimator is separately trained on $z_{(m)}^N$, $m = 1, \dots, M$. For simplicity, we assume that the sample size of each training set is uniformly N . Note that training set $z_{(m)}^N$ is a realization of a random sequence $Z_{(m)}^N$ and that $Z_{(m)}^N$, $m = 1, \dots, M$, have the same distribution $p(\mathbf{x}, y)$; however, they cannot always be assumed to be mutually independent. The independent case will be discussed later as a special case.

The output of the ensemble estimator for some input \mathbf{x} is defined as the simple average of outputs of M estimators for \mathbf{x} after they have been separately trained. Specifically,

$$f_{\text{ens}}^{(M)}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M f_m(\mathbf{x}; z_{(m)}^N). \quad (5)$$

Then, we have the following theorem about the generalization error of the ensemble estimator.

Theorem 1 (Generalization error in general case)

Let $\text{GErr}(f_{\text{ens}}^{(M)})$ denote the generalization error of the ensemble estimator given by (5). We have

$$\text{GErr}(f_{\text{ens}}^{(M)}) = \mathbb{E}_{X_0} \left\{ \frac{1}{M} \overline{\text{Var}}(X_0) + \left(1 - \frac{1}{M} \right) \overline{\text{Cov}}(X_0) + \overline{\text{Bias}}(X_0)^2 \right\} + \sigma^2, \quad (6)$$

¹In order to clarify the random variable, we will use notation such as $\mathbb{E}_{Z^N} \{ \}$ instead of simply $\mathbb{E} \{ \}$.

²Strictly speaking, they should be defined by using a conditional expectation like $\mathbb{E}_{Z^N|X_0} \{ \varphi(X_0, Z^N) | X_0 \}$. However, since Z^N is independent of X_0 , $\mathbb{E}_{Z^N|X_0} \{ \varphi(X_0, Z^N) | X_0 \}$ becomes $\mathbb{E}_{Z^N} \{ \varphi(X_0, Z^N) \}$.

where $\overline{\text{Var}}(X_0)$, $\overline{\text{Cov}}(X_0)$, and $\overline{\text{Bias}}(X_0)$ are average conditional variance, conditional covariance, and conditional bias, averaged over M estimators, respectively. That is,

$$\begin{cases} \overline{\text{Var}}(X_0) = \frac{1}{M} \sum_{m=1}^M \text{Var}\{f_m|X_0\}, \\ \overline{\text{Cov}}(X_0) = \frac{1}{M(M-1)} \sum_m \sum_{m' \neq m} \text{Cov}\{f_m, f_{m'}|X_0\}, \\ \overline{\text{Bias}}(X_0) = \frac{1}{M} \sum_{m=1}^M \text{Bias}\{f_m|X_0\}. \end{cases} \quad (7)$$

Here f_m denotes $f_m(X_0; Z_{(m)}^N)$. Using the following average generalization error averaged over M estimators,

$$\overline{\text{GErr}} = \frac{1}{M} \sum_{m=1}^M (\mathbb{E}_{X_0} \{\text{Var}\{f_m|X_0\} + \text{Bias}\{f_m|X_0\}^2\} + \sigma^2),$$

we can also present the relationship between $\text{GErr}(f_{\text{ens}}^{(M)})$ and $\overline{\text{GErr}}$ as follows.

$$\text{GErr}(f_{\text{ens}}^{(M)}) = \frac{1}{M} \overline{\text{GErr}} + \left(1 - \frac{1}{M}\right) \sigma^2 + \mathbb{E}_{X_0} \left\{ \left(1 - \frac{1}{M}\right) \overline{\text{Cov}}(X_0) + \frac{1}{M^2} \sum_m \sum_{m' \neq m} \text{Bias}\{f_m|X_0\} \text{Bias}\{f_{m'}|X_0\} \right\}. \quad (8)$$

Proof Using $f_{\text{ens}}^{(M)}$ instead of f in (3), we have

$$\text{GErr}(f_{\text{ens}}^{(M)}) = \mathbb{E}_{X_0} \left\{ \text{Var}\{f_{\text{ens}}^{(M)}|X_0\} + \text{Bias}\{f_{\text{ens}}^{(M)}|X_0\}^2 \right\} + \sigma^2, \quad (9)$$

Since $f_{\text{ens}}^{(M)}$ depends on $z_{(1)}^N, \dots, z_{(M)}^N$, $\text{Var}\{f_{\text{ens}}^{(M)}|X_0\}$ and $\text{Bias}\{f_{\text{ens}}^{(M)}|X_0\}$ can be given by extending (4) as

$$\begin{aligned} \text{Var}\{f_{\text{ens}}^{(M)}|X_0\} &= \mathbb{E}_{Z_{(1)}^N, \dots, Z_{(M)}^N} \left\{ \left[\frac{1}{M} \sum_{m=1}^M f_m - \mathbb{E}_{Z_{(1)}^N, \dots, Z_{(M)}^N} \left\{ \frac{1}{M} \sum_{m=1}^M f_m \right\} \right]^2 \right\} \\ &= \frac{1}{M^2} \sum_{m=1}^M \mathbb{E}_{Z_{(m)}^N} \left\{ [f_m - \mathbb{E}_{Z_{(m)}^N} \{f_m\}]^2 \right\} + \\ &\quad \frac{1}{M^2} \sum_m \sum_{m' \neq m} \mathbb{E}_{Z_{(m)}^N, Z_{(m')}^N} \left\{ [f_m - \mathbb{E}_{Z_{(m)}^N} \{f_m\}] [f_{m'} - \mathbb{E}_{Z_{(m')}^N} \{f_{m'}\}] \right\} \\ &= \frac{1}{M} \overline{\text{Var}}(X_0) + \left(1 - \frac{1}{M}\right) \overline{\text{Cov}}(X_0) \end{aligned} \quad (10)$$

$$\begin{aligned} \text{Bias}\{f_{\text{ens}}^{(M)}|X_0\} &= \mathbb{E}_{Z_{(1)}^N, \dots, Z_{(M)}^N} \left\{ \frac{1}{M} \sum_{m=1}^M f_m \right\} - g \\ &= \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{Z_{(m)}^N} \{f_m - g\} = \overline{\text{Bias}}(X_0). \end{aligned} \quad (11)$$

Substitution of (10) and (11) into (9) gives (6). Once (6) is proven, the derivation of (8) is trivial. ■

One can see that when $M = 1$, (6) reduces to (3). Theorem 1 shows that the generalization error of the ensemble estimator, unlike that of a single estimator, is affected not only by variance, bias and noise variance, but also by covariance. This may be intuitively reasonable. Moreover, these explicit expressions yield much information about how the generalization error of the ensemble estimator is affected by each of factors of interest (Var, Cov, Bias, and σ^2), and they also quantify the improvement due to the ensemble.

Corollary 1 (Equal model case)

If models of M estimators are the same,³ then we have

$$\text{GErr}(f_{\text{ens}}^{(M)}) = \mathbb{E}_{X_0} \left\{ \frac{1}{M} \text{Var}\{f|X_0\} + \left(1 - \frac{1}{M}\right) \overline{\text{Cov}}(X_0) + \text{Bias}\{f|X_0\}^2 \right\} + \sigma^2, \quad (12)$$

³In feedforward neural networks, the same model means the same network structure.

An equivalent expression is

$$\text{GErr}(f_{\text{ens}}^{(M)}) = \frac{1}{M} \text{GErr}(f) + \left(1 - \frac{1}{M}\right) \mathbb{E}_{X_0} \{ \text{Bias}\{f|X_0\}^2 + \overline{\text{Cov}}(X_0) + \sigma^2 \}. \quad (13)$$

where $\text{GErr}(f)$ is the generalization error of a single estimator given in (3).

Proof Since $Z_{(m)}^N, m = 1, \dots, M$ have the same distribution, if models of M estimators are the same, then clearly variances (biases) of all estimators become identical. That is,

$$\text{Var}\{f_1|X_0\} = \dots = \text{Var}\{f_M|X_0\} \equiv \text{Var}\{f|X_0\}, \quad \text{Bias}\{f_1|X_0\} = \dots = \text{Bias}\{f_M|X_0\} \equiv \text{Bias}\{f|X_0\}.$$

Then, (6) and (8) reduce to (12) and (13), respectively. \blacksquare

Corollary 2 (Equal model & uncorrelated case)

If models of M estimators are the same and outputs of M estimators are mutually uncorrelated,⁴ then we have

$$\text{GErr}(f_{\text{ens}}^{(M)}) = \mathbb{E}_{X_0} \left\{ \frac{1}{M} \text{Var}\{f|X_0\} + \text{Bias}\{f|X_0\}^2 \right\} + \sigma^2. \quad (14)$$

An equivalent expression is

$$\text{GErr}(f_{\text{ens}}^{(M)}) = \frac{1}{M} \text{GErr}(f) + \left(1 - \frac{1}{M}\right) \mathbb{E}_{X_0} \{ \text{Bias}\{f|X_0\}^2 + \sigma^2 \}. \quad (15)$$

Proof Since f_m and $f_{m'}$ are mutually independent, $\text{Cov}\{f_m, f_{m'}|X_0\} = 0$. Hence, setting $\overline{\text{Cov}} \equiv 0$ in (12) and (13) lead to (14) and (15), respectively. \blacksquare

These results can easily be extended to the *general ensemble estimators* defined by

$$f_{\text{gens}}^{(M)}(\mathbf{x}) = \sum_{m=1}^M \alpha_m f_m(\mathbf{x}; z_{(m)}^N), \quad \text{where } \alpha_m > 0, \sum_{m=1}^M \alpha_m = 1. \quad (16)$$

Due to lack of space, only the result will be provided in the appendix.

Remarks At this point, the following important observations can be made :

- If f_m and $f_{m'}$ are *negatively* correlated, then the correlation contributes to a decrease in the generalization error. Conversely, if f_m and $f_{m'}$ are *positively* correlated, then the correlation increases the generalization error.
- Under the uncorrelated assumption, from (10), $\text{Var}\{f_{\text{ens}}^{(M)}|X_0\}$ reduces to $\overline{\text{Var}}(X_0)/M$ due to the averaging, and this reduction decreases the generalization error⁵. On the other hand, since $\text{Bias}\{f_{\text{ens}}^{(M)}|X_0\}$ is a simple average of M biases ($\text{Bias}\{f_m|X_0\}, m = 1, \dots, M$), the bias term basically does not contribute to a decrease in the generalization error. Therefore, combining estimators that each has a higher ratio of variance to bias is more useful for decreasing generalization error than combining estimators that each has a lower ratio of variance to bias.
- Even if we increase the number of ensembles M , the generalization error cannot be less than $\mathbb{E}_{X_0} \{ \text{Bias}\{f|X_0\}^2 \} + \sigma^2$. That is, the amount gives the lower bound of $\text{GErr}(f_{\text{ens}}^{(M)})$.

4. Simulation

In order to verify the theoretical results presented so far, we performed a simulation for the case of Corollary 2 that is easily simulated. We used standard feedforward neural networks with one input, ten hidden units and one output (1-10-1) as regression estimators. Figure 1 shows a target function

⁴Uncorrelated estimators can be obtained by training each estimator on a different (independent) training set.

⁵This variance reduction by a factor of M has been already mentioned by Perrone [2]. Perrone's result for basic ensemble agrees with (15) when $\text{Bias} \equiv 0$ and $\sigma^2 \equiv 0$.

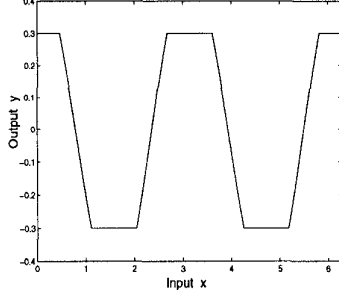


Fig. 1: Target function.

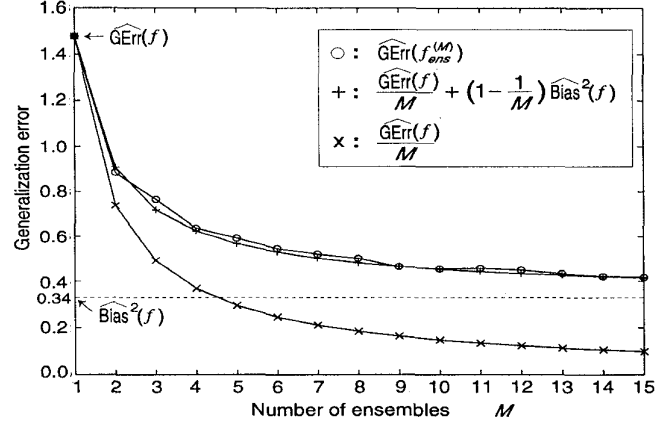


Fig. 2: Number of ensembles vs. generalization error.

($y = g(x), x, y \in \mathcal{R}$). We set $N = 40$. That is, $z_{(m)}^{40} = \{(x_1, y_1), \dots, (x_{40}, y_{40})\}$, $m = 1, \dots, M$ are artificially and independently generated so that each (x_i, y_i) pair satisfies $y_i = g(x_i)$. This means $\sigma^2 \equiv 0$. A test set (z_0^{1000}) containing 1000 samples is also generated in the same way. For $M = 1, \dots, 15$, the estimate of $\text{GErr}(f_{\text{ens}}^{(M)})$, denoted as $\widehat{\text{GErr}}(f_{\text{ens}}^{(M)})$, is calculated by the following approximation: $\widehat{\text{GErr}}(f_{\text{ens}}^{(M)}) = \frac{1}{1000} \sum_{i=1}^{1000} \frac{1}{100} \sum_{j=1}^{100} (y_i^{(0)} - f_{\text{ens}}^{(M)}(x_i^{(0)}; z_{(j)}^{40}))^2$. That is, $E_{Z^{40}}$ is approximated by 100 sample sets. Similarly, the estimate of $E_{X_0}\{\text{Bias}\{f|X_0\}^2\}$, denoted as $\widehat{\text{Bias}}^2(f)$, is calculated as: $\widehat{\text{Bias}}^2(f) = \frac{1}{1000} \sum_{i=1}^{1000} (\frac{1}{100} \sum_{j=1}^{100} f(x_i^{(0)}; z_{(j)}^{40}) - y_i^{(0)})^2$. Clearly, $\widehat{\text{GErr}}(f)$ can be obtained as $\widehat{\text{GErr}}(f_{\text{ens}}^{(1)})$.

We summarize the above simulation result in Fig. 2. In Fig. 2, 'o' indicates $\widehat{\text{GErr}}(f_{\text{ens}}^{(M)})$ and '+' indicates the theoretical result obtained by substituting values of $\widehat{\text{GErr}}(f)$ and $\widehat{\text{Bias}}^2(f)$ and $\sigma^2 = 0$ into (15). For comparison, $\widehat{\text{GErr}}(f)/M, M = 1, \dots, 15$ are also plotted in Fig. 2 as 'x'. One can see that the experimental curve ('o') fits into the theoretical one (far from $\widehat{\text{GErr}}(f)/M$) and that both curves approach $\widehat{\text{Bias}}^2(f) = 0.34$ as M increases.

5. Discussion

We should notice that these results do not always support the usefulness of the ensemble approach. In the case of Corollary 2, one can see that an increase in M necessitates an increase in the total number of training samples. In other words, if fixed N samples are available,⁶ we must divide the data set into M disjoint sets. Each set contains N/M samples such that M sets are mutually uncorrelated. In this case, it is natural to ask, "Why don't you train a single estimator by using all N training samples at once? Does combining estimators, each of which is trained on divided N/M samples, still outperform the single estimator trained on N samples?" In general, both bias and variance depend on the sample size and the dimensionality of the parameter to be estimated. This means that bias (variance) obtained by N training samples is definitely less than that obtained by N/M training samples. Hence, recalling that the generalization error of a single estimator is expressed by the sum of the variance and the squared bias, it is quite possible to develop the situation where a single estimator outperforms an ensemble estimator. An interesting result is presented in Table 1, which shows that a single estimator really outperforms an ensemble estimator.

In Table 1, $\widehat{\text{GErr}}(f)$ and $\widehat{\text{Bias}}^2(f)$ were obtained from a single estimator trained on N samples. For example, $N = 400$ corresponds to $M = 10$, because each of M estimators is trained on 40 samples in this

⁶In many practical situations, only one small sample set is available, because collecting samples often involves a large cost and is sometimes impossible.

Table 1: Sample size effect on generalization error

N	40	80	120	160	200	280	400
$\widehat{\text{GE}}\text{rr}(f)$	1.48	0.40	0.35	0.16	0.14	0.09	0.08
$\widehat{\text{Bias}}^2(f)$	0.34	0.17	0.11	0.10	0.08	0.06	0.05

simulation. As previously mentioned, bias monotonically decreases as N increases. Comparing $\widehat{\text{GE}}\text{rr}(f_{\text{ens}}^{(M)})$ in Fig. 2 with corresponding $\widehat{\text{GE}}\text{rr}(f)$ in Table 1, one can see that the generalization error of the single estimator trained on larger samples is much smaller than that of the ensemble estimator. For example, when $M = 10$, $\widehat{\text{GE}}\text{rr}(f_{\text{ens}}^{(10)}) = 0.46$, while the corresponding $\widehat{\text{GE}}\text{rr}(f)$ of a single estimator is 0.08 ($N = 400$). In this experiment, the single estimator was superior to the ensemble estimator for all M . These results suggest that when we train each estimator, we should utilize as many samples as possible. The straightforward way is to divide samples into overlapped ones. However, this will produce positive correlations among these sample sets, and increase the generalization error as mentioned before. This problem may be called “ensemble dilemma.” Future work will focus on solving this problem.

Appendix

Let $\text{GErr}(f_{\text{gens}}^{(M)})$ denote the generalization error of the general ensemble estimator given by (16). Then, we have

$$\text{GErr}(f_{\text{gens}}^{(M)}) = \mathbb{E}_{X_0} \left\{ \sum_{m=1}^M \sum_{l=1}^M \alpha_m^* \alpha_l^* R_{ml} \right\}, \quad \text{where } \alpha_i^* = \sum_j R_{ij}^{-1} / \sum_k \sum_j R_{kj}^{-1}. \quad (17)$$

α_i^* denotes the optimal weight that minimizes the generalization error of the general ensemble estimator given in (16). R_{ij}^{-1} indicates the i, j the component of the inverse matrix of \mathbf{R} . The i, j th component of matrix \mathbf{R} is given by

$$R_{ij} = \begin{cases} \text{Var}\{f_i|X_0\} + \text{Bias}\{f_i|X_0\}^2, & \text{if } i = j \\ \text{Cov}\{f_i, f_j|X_0\} + \text{Bias}\{f_i|X_0\}\text{Bias}\{f_j|X_0\}, & \text{otherwise.} \end{cases} \quad (18)$$

Note that if M estimators have the same model, then α_i^* reduces to $1/M$.

References

- [1] Wolpert D. H.: “Stacked generalization,” *Neural Networks*, **5**, 2, pp. 241–259, 1992.
- [2] Perrone M. P.: “Improving regression estimates: Averaging methods for variance reduction with extensions to general convex measure optimization,” *PhD Thesis*, Brown University, 1993.
- [3] Tresp V. and Taniguchi M.: “Combining estimators using non-constant weighting functions,” In Tesauro G. *et al.* eds., *Advances in Neural Information Processing Systems 7*, MIT Press, Cambridge MA, 1995.
- [4] Krogh A. and Vedelsby J.: “Neural network ensembles, cross validation, and active learning,” In Tesauro G. *et al.* eds., *Advances in Neural Information Processing Systems 7*, MIT Press, Cambridge MA, 1995.
- [5] Hansen L. K. and Salamon P.: “Neural network ensembles,” *IEEE Trans. Pattern Anal. Machine Intell.*, **PAMI-12**, 10, pp. 993–1001, 1990.
- [6] Hampshire J. B. and Waibel A.: “The meta-pi network: Building distributed knowledge representations for robust multisource pattern recognition,” *IEEE Trans. Pattern Anal. Machine Intell.*, **PAMI-14**, 7, pp. 751–769, 1992.
- [7] Geman S., Bienenstock E. and Cooper L. N.: “Neural networks and the bias/variance dilemma,” *Neural Computation*, **4**, 1, pp. 1–58, 1992.