

Degrees of Freedom in Regression Ensembles

Henry WJ Reeve

Gavin Brown *

University of Manchester - School of Computer Science
Kilburn Building, University of Manchester, Oxford Rd, Manchester M13 9PL

Abstract. Negative correlation learning is an effective approach to ensemble learning in which model diversity is encouraged through a correlation penalty term. The level of emphasis placed upon the correlation penalty term is controlled by the diversity parameter. We shall provide a degrees of freedom analysis of negative correlation learning. Our contributions are as follows: we give an exact formula for the effective degrees of freedom in a negative correlation ensemble with fixed basis functions; we show that the effective degrees of freedom is a continuous, convex and monotonically increasing function of the diversity parameter; finally, we show that the degrees of freedom formula gives rise to an efficient way to tune the diversity parameter on large data sets.

1 Introduction

Ensemble methods are a cornerstone of modern machine learning. Numerous applications have shown that by combining a multiplicity of models we are able to train powerful estimators from large data sets in a tractable way. Successful ensemble performance emanates from a fruitful trade-off between the individual accuracy of the models and their diversity [1]. Typically diversity is introduced implicitly, by sub-sampling the data or varying the architecture of the models. Negative correlation learning takes an alternative approach in which diversity is explicitly encouraged by appending a covariance penalty term to the cost function [3]. The level of emphasis placed upon the covariance penalty term is controlled by the diversity parameter.

We shall provide a degrees of freedom analysis of negative correlation learning. Our contributions are as follows: we give an exact formula for the effective degrees of freedom in a negative correlation ensemble with fixed basis functions; we use this formula to address the following two questions:

- How does the complexity of the ensemble estimator vary as function of the diversity parameter?
- How can we efficiently optimise the diversity parameter on large data sets?

We begin with the necessary background before addressing these two questions in sections 3 and 4, respectively.

*The authors gratefully acknowledge the support of the EPSRC for the LAMBDA project (EP/N035127/1) and the Manchester Centre for Doctoral Training (EP/1038099/1). We would also like to thank Kit Elliot and the anonymous reviewers for useful feedback.

2 Background

In this section we shall introduce the relevant background on negative correlation learning and degrees of freedom. We begin by setting the scene. Throughout this paper we consider a regression problem in which we are given a data set $\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N$ with $(x_n, y_n) \in \mathcal{X} \times \mathbb{R}$ generated independently from a distribution \mathbb{P} over $\mathcal{X} \times \mathbb{R}$. Our goal is to provide an estimator $\hat{\mu} : \mathcal{X} \rightarrow \mathbb{R}$ of the conditional expectation $\mu(x) = \mathbb{E}[y|x]$, where expectation is taken with respect to \mathbb{P} .

2.1 Negative correlation learning and diversity

We briefly review negative correlation learning in the regression context. For a detailed exposition see [1]. Suppose we have an ensemble $\mathcal{F} = \{f_m\}_{m=1}^M$ consisting of M functions $f_m : \mathcal{X} \rightarrow \mathbb{R}$ each parametrised by θ_m . Our ensemble estimator $\hat{\mu}$ is given by the average $F = (1/M) \cdot \sum_{m=1}^M f_m$. The negative correlation learning rule, introduced by Liu and Yao [3] proceeds as follows. First each parameter vector θ_m is randomly initialised. Then, for each training example $(x_n, y_n) \in \mathcal{D}$, we update each θ_m in parallel according to

$$\theta_m \leftarrow \theta_m - \alpha \cdot \frac{\partial f_m}{\partial \theta_m} \cdot \left(\overbrace{(f_m(x_n) - y)}^{\text{accuracy}} - \lambda \cdot \overbrace{(f_m(x_n) - F(x_n))}^{\text{diversity}} \right),$$

where α is a learning rate. Brown et al. [1] observed that the negative correlation learning rule is equivalent to stochastic gradient descent with respect to the following loss function (with a scaled learning rate),

$$L_\lambda(\mathcal{F}, x, y) := \overbrace{\frac{1}{M} \sum_{m=1}^M (f_m(x) - y)^2}^{\text{accuracy}} - \lambda \cdot \overbrace{\frac{1}{M} \sum_{m=1}^M (f_m(x) - F(x))^2}^{\text{diversity}}. \quad (1)$$

We shall refer to L_λ as the negative correlation loss. The diversity parameter λ explicitly manages a trade-off between two competing factors: accuracy and diversity. The key focus of this paper will be understanding the behaviour of the ensemble as a function of the diversity parameter λ . When $\lambda = 0$ each function f_m is trained individually. On the other hand, it follows from the ambiguity decomposition that when $\lambda = 1$, L_λ is the squared error for the average F . Hence, negative correlation learning scales smoothly between training each of the functions f_m individually and training as a single combined estimator F . Brown et al. conducted a detailed analysis of negative correlation learning, relating the behaviour of the ensemble to the bias-variance-covariance decomposition [1]. In addition, Brown et al. gave an upper bound on the diversity parameter λ , showing that for $\lambda > M/(M-1) > 1$ the Hessian matrix of the weights is non-positive semi-definite. It was subsequently shown that for any $\lambda > 1$, minimising L_λ causes the weights to diverge [4, Theorem 3]. Thus, we should restrict the diversity parameter λ to the region $\lambda \in [0, 1]$.

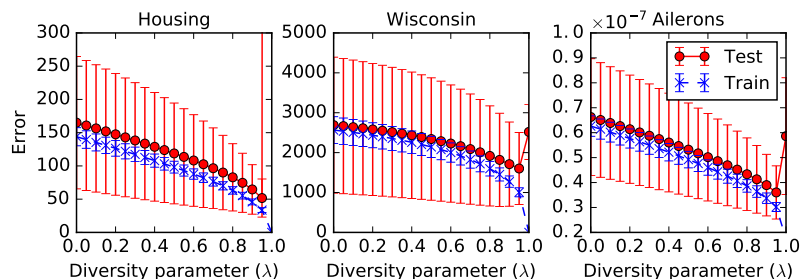


Fig. 1: The mean square error for negative correlation ensembles on three benchmark datasets. For experimental details see Section 4.

Figure 1 displays the typical behaviour of the mean squared error as a function of the diversity parameter λ . As the diversity parameter moves from zero to one, the training error declines. This is to be expected given that L_λ with $\lambda = 1$ corresponds to the square error for the ensemble F . The test error typically declines initially before rising sharply as λ approaches one. It is this phenomenon that we intend to explain. Note that this increase in the test error cannot be explained by the upper bounds in [1] and [4, Theorem 3], since these relate to behaviour on the training set \mathcal{D} for $\lambda > 1$. Taking $\lambda < 1$ appears to act as a regulariser, reducing the discrepancy between test and train error. This is intuitively plausible. When $\lambda = 0$ we are independently training a collection of M simple models and aggregating the result, an approach which is likely to under-fit. When $\lambda = 1$ we are minimising the training error for a single complex model F , which is likely to overfit. Choosing λ between zero and one blends between these extremes, providing an effective balance between underfitting and overfitting. Nonetheless, the underlying hypothesis class does not change with λ . Hence, the usual formalisms of VC dimension or Rademacher complexity do not apply. We shall show that the regularising behaviour of negative correlation learning may be explained through the notion of effective degrees of freedom.

2.2 Degrees of freedom

In this section we shall introduce the relevant background on degrees of freedom.

Definition 1 (Degrees of freedom [5]). *Suppose that we have estimation procedure \mathcal{M} which, given a data set \mathcal{D} outputs an estimator $\hat{\mu} : \mathcal{X} \rightarrow \mathbb{R}$. The effective degrees of freedom of \mathcal{M} is defined by $df(\mathcal{M}, \mathcal{D}) := \sum_{n=1}^N \partial \hat{\mu}(x_n) / \partial y_n$.*

In the ordinary least squares setting $df(\mathcal{M}, \mathcal{D})$ is simply the number of non-trivial parameters in the model, hence the nomenclature. The degrees of freedom quantifies the complexity of an estimation procedure via the sensitivity of model outputs to the data. Stein's unbiased risk estimate shows that under the assumption of homoscedastic Gaussian noise in the fixed design setting, the expected difference between the test and train error is twice the degrees of freedom multiplied by the noise variance [5]. Hence, the degrees of freedom is directly related to the estimator's potential for over-fitting.

One of the primary motivations for understanding model complexity and degrees of freedom is to estimate the out of sample error for the purpose of model selection. Stein's unbiased risk estimate may be used in this way [5]. However, this requires an estimate for the noise variance, which is rarely known in practice. Non-parametric techniques such as cross validation do not have this limitation. The generalised cross validation procedure introduced by Cravan and Wahba gives an efficient rotation invariant form of the leave-one-out cross validation error, based upon the degrees of freedom [2].

$$\text{GCV}(\mathcal{M}, \mathcal{D}) := \left(\frac{1}{N} \sum_{n=1}^N (\hat{\mu}(x_n) - y_n)^2 \right) / (1 - \text{df}(\mathcal{M}, \mathcal{D})/N)^2. \quad (2)$$

In Section 4 we use the GCV formula for tuning the diversity parameter.

3 Degrees of freedom for negative correlation learning

In this section we analyse the degrees of freedom for a special class of negative correlation ensembles \mathcal{F} in which each function f_m consists of a linear map applied to a fixed basis function. We consider ensembles \mathcal{F} where for each $m = 1, \dots, M$ there exists a fixed function $\phi_m : \mathcal{X} \rightarrow \mathbb{R}^H$ with $f_m(x) = \mathbf{w}_m \phi_m(x)$, where $\mathbf{w}_m \in \mathbb{R}^{1 \times H}$ is a trainable weight vector. In this situation the ensemble F which minimises L_λ averaged over the data has a simple closed form solution. We first introduce some notation. Let $Q = H \cdot M$ and let $\phi : \mathcal{X} \rightarrow \mathbb{R}^Q$ be the function defined by $\phi(x) = [\phi_1(x)^T, \dots, \phi_M(x)^T]^T$. We let

$$\begin{aligned} \langle \phi, \phi \rangle_{\mathcal{D}} &:= \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \in \mathbb{R}^{Q \times Q} \\ \langle \mathbf{y}, \phi \rangle_{\mathcal{D}} &:= \frac{1}{N} \sum_{n=1}^N y_n \phi(\mathbf{x}_n)^T \in \mathbb{R}^{1 \times Q} \end{aligned}$$

Similarly, for each $m = 1, \dots, M$, we define

$$\begin{aligned} \langle \phi_m, \phi_m \rangle_{\mathcal{D}} &:= \frac{1}{N} \sum_{n=1}^N \phi_m(\mathbf{x}_n) \phi_m(\mathbf{x}_n)^T \in \mathbb{R}^{H \times H} \\ \langle \phi, \phi \rangle_{\mathcal{D}}^{\text{diag}} &:= \text{diag}(\langle \phi_1, \phi_1 \rangle_{\mathcal{D}}, \dots, \langle \phi_M, \phi_M \rangle_{\mathcal{D}}) \in \mathbb{R}^{Q \times Q}. \end{aligned}$$

We shall assume that $\langle \phi_m, \phi_m \rangle_{\mathcal{D}}$ is of full rank H . Under these assumptions the negative correlation estimator has a closed form solution.

Theorem 2. *Suppose that each function f_m is a linear map of a fixed basis function. For each $\lambda \in [0, 1]$, the minimiser F_λ of the negative correlation loss averaged over a data set \mathcal{D} is given by $F_\lambda(\mathbf{x}) = \beta_\lambda \phi(\mathbf{x})$ where*

$$\beta_\lambda = \langle \mathbf{y}, \phi \rangle_{\mathcal{D}} \left(M(1 - \lambda) \cdot \langle \phi, \phi \rangle_{\mathcal{D}}^{\text{diag}} + \lambda \langle \phi, \phi \rangle_{\mathcal{D}} \right)^+.$$

The closed form solution allows us to efficiently locate the minimiser of the negative correlation loss. The proof of Theorems 2 & 3 is given in the supplementary material. ¹

¹Available at <http://www.cs.man.ac.uk/~reeveh/>

Theorem 3. Given any $\lambda \in [0, 1]$, the degrees of freedom of the minimiser F_λ of the negative correlation loss L_λ averaged over the data data set \mathcal{D} is given by

$$df(F_\lambda) = \text{trace} \left(\langle \phi, \phi \rangle_{\mathcal{D}} \left(M(1 - \lambda) \cdot \langle \phi, \phi \rangle_{\mathcal{D}}^{\text{diag}} + \lambda \langle \phi, \phi \rangle_{\mathcal{D}} \right)^+ \right).$$

Moreover, $df(F_\lambda)$ is a continuous non-decreasing and convex function of λ with $df(F_0) = H$ and $\lim_{\lambda \rightarrow 1} df(F_\lambda) = \text{rank}(\langle \phi, \phi \rangle_{\mathcal{D}})$. In addition, if $H < \text{rank}(\langle \phi, \phi \rangle_{\mathcal{D}})$ then $df(F_\lambda)$ is strictly increasing and strictly convex as a function of λ .

Theorem 3 is the central result of the paper. Theorem 3 explains the regularisation behaviour of negative correlation ensembles and shows that the the complexity of the estimator increases continuously and monotonically as a function of the diversity parameter. This behaviour is demonstrated on three benchmark data sets in Figure 2.

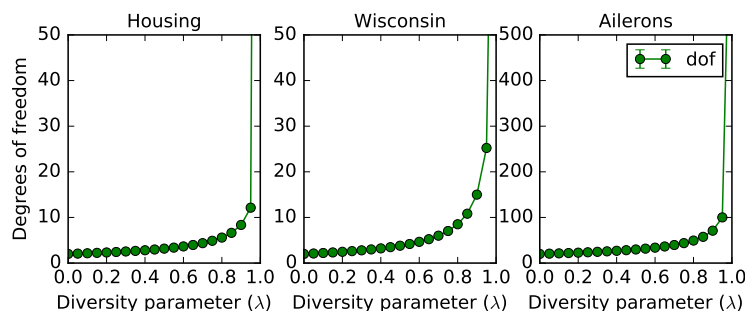


Fig. 2: The degrees of freedom for negative correlation ensembles on three benchmark data sets. For experimental details see Section 4.

4 Efficient optimisation of the diversity parameter

A persistent challenge in machine learning is hyper-parameter optimisation. We shall address this challenge for the diversity parameter. We would like to choose the diversity parameter which minimises the test error. Of course, the generalisation error is not accessible, so we must choose an appropriate proxy measure. The standard choice is the k -fold cross-validation error. However, this requires re-fitting the model k -times, which is prohibitively expensive for large data sets with limited time and computational resources. The generalised cross validation formula gives an efficient proxy for the test error without the need for multiple re-fits [2]. We shall apply Theorem 3 to compute the generalised cross validation error (Equation (2)).

We compared the performance of the generalised cross validation error with the five-fold-cross validation error as a metric for choosing the diversity parameter, across six benchmark data sets. The basis functions are generated with random Fourier features with $M = 500$ and $H = 2$ for data sets of less than a thousand examples (Housing, Machine and Wisconsin) and $M = 50$ and $H = 20$ for the remainder. In each case the diversity parameter is chosen by applying the Brent minimisation routine to the corresponding metric. We evaluate the test mean square error with the selected diversity

parameter along with the total optimisation time. Table 1 shows the results. Whilst there is no statistically significant difference in generalisation performance, using the generalised cross validation formula leads to a significant speed up.

Data set	Error (gcv)	Error (5-fold)	Gcv time (s)	5-fold time (s)
Housing	$3.41e+01 \pm 2e+01$	$3.41e+01 \pm 2e+01$	25.3 ± 0.6	126.2 ± 1.2
Machine	$1.94e+04 \pm 2e+04$	$1.94e+04 \pm 2e+04$	25.0 ± 0.7	123.8 ± 1.0
Wisconsin	$1.56e+03 \pm 9e+02$	$1.57e+03 \pm 9e+02$	25.7 ± 0.6	124.2 ± 1.1
Ailerons	$3.23e-08 \pm 1e-08$	$3.23e-08 \pm 1e-08$	25.1 ± 0.7	139.9 ± 0.3
Elevators	$6.04e-06 \pm 3e-06$	$6.04e-06 \pm 3e-06$	24.7 ± 0.6	146.5 ± 0.5
Kinematic	$2.29e-02 \pm 1e-03$	$2.29e-02 \pm 1e-03$	24.7 ± 0.7	141.5 ± 1.1

Table 1: A comparison between the generalised cross validation method, exploiting the degrees of freedom formula and standard five fold cross validation for setting the diversity parameter. For experimental details see Section 4.

5 Discussion

In this paper we have given a degrees of freedom analysis for negative correlation ensembles with fixed basis functions. Our analysis has enabled us to characterise the complexity of the estimator as a function of the diversity parameter and obtain an efficient approach to tuning the diversity parameter via the generalised cross validation error. In future work we intend to extend this analysis to ensembles of networks with multiple trainable hidden layers by applying the Monte-Carlo estimate for generalised degrees of freedom as suggested by Ye [6].

References

- [1] Gavin Brown, Jeremy L Wyatt, and Peter Tiño. Managing diversity in regression ensembles. *The Journal of Machine Learning Research*, 6:1621–1650, 2005.
- [2] Peter Craven and Grace Wahba. Smoothing noisy data with spline functions. *Numerische Mathematik*, 31(4):377–403, 1978.
- [3] Yong Liu and Xin Yao. Ensemble learning via negative correlation. *Neural Networks*, 12(10):1399–1404, 1999.
- [4] Henry WJ Reeve and Gavin Brown. Modular autoencoders for ensemble feature extraction. *Feature Extraction: Modern Questions and Challenges*, 1:242–259, 2015.
- [5] Charles M Stein. Estimation of the mean of a multivariate normal distribution. *The annals of Statistics*, pages 1135–1151, 1981.
- [6] Jianming Ye. On measuring and correcting the effects of data mining and model selection. *Journal of the American Statistical Association*, 93(441):120–131, 1998.