

# Advanced Data Analytics Combining Models

Ian T. Nabney

- Bias and variance
- Bayesian model averaging
- Committees
- Boosting

# Bias-Variance Decomposition

- The phenomenon of over-fitting is really an unfortunate property of maximum likelihood and does not arise when we marginalize over parameters in a Bayesian setting.
- It is instructive to consider a frequentist viewpoint of the model complexity issue, known as the **bias-variance** trade-off.
- We use the the squared loss function, for which the optimal prediction is given by the conditional expectation, which we denote by  $h(\mathbf{x})$  and which is given by

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt. \quad (1)$$

- The expected squared loss can be written in the form

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt. \quad (2)$$

# Model uncertainty

In a frequentist treatment, we make a point estimate of  $\mathbf{w}$  based on the data set  $\mathcal{D}$ , and interpret the uncertainty of this estimate as follows:

- Suppose we had a large number of data sets each of size  $N$  and each drawn independently from the distribution  $p(t, \mathbf{x})$ .
- For any given data set  $\mathcal{D}$ , we can run our learning algorithm and obtain a prediction function  $y(\mathbf{x}; \mathcal{D})$ .
- Different data sets from the ensemble will give different functions and consequently different values of the squared loss.
- The performance of a particular learning algorithm is then assessed by taking the average over this ensemble of data sets.
- The key variable is the first term in (1).

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} \left[ \{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 \right] \\ = \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}} \left[ \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 \right]}_{\text{variance}}. \end{aligned} \quad (3)$$

We see that the expected squared difference between  $y(\mathbf{x}; \mathcal{D})$  and the regression function  $h(\mathbf{x})$  can be expressed as the sum of two terms.

- 1 The squared **bias**, represents the extent to which the average prediction over all data sets differs from the desired regression function.
- 2 The **variance**, measures the extent to which the solutions for individual data sets vary around their average, and hence this measures the extent to which the function  $y(\mathbf{x}; \mathcal{D})$  is sensitive to the particular choice of data set.

# Full loss function decomposition

We substitute this expansion back into (1), we obtain the following decomposition of the expected squared loss

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise} \quad (4)$$

where

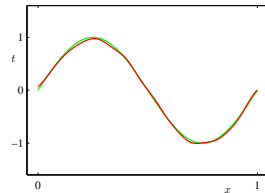
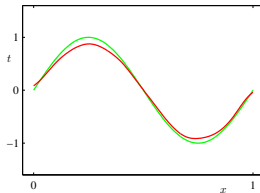
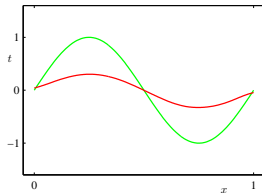
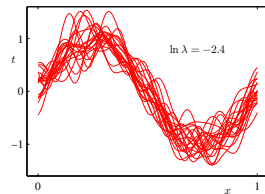
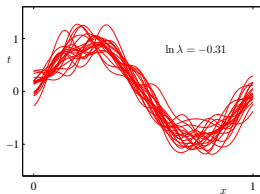
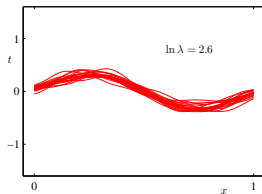
$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} \quad (5)$$

$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} \left[ \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 \right] p(\mathbf{x}) d\mathbf{x} \quad (6)$$

$$\text{noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt \quad (7)$$

and the bias and variance terms now refer to integrated quantities.

# Dependence of bias and variance on model complexity



# Bias-variance trade-off

The average prediction is estimated from

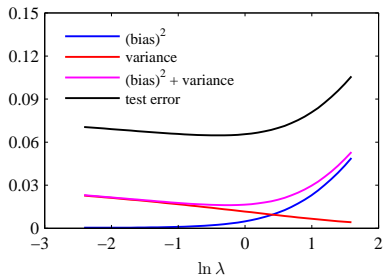
$$\bar{y}(x) = \frac{1}{L} \sum_{l=1}^L y^{(l)}(x) \quad (8)$$

and the integrated squared bias and integrated variance are then given by

$$(\text{bias})^2 = \frac{1}{N} \sum_{n=1}^N \{\bar{y}(x_n) - h(x_n)\}^2 \quad (9)$$

$$\text{variance} = \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L \{y^{(l)}(x_n) - \bar{y}(x_n)\}^2 \quad (10)$$

where the integral over  $x$  weighted by the distribution  $p(x)$  is approximated by a finite sum over data points drawn from that distribution.





# Mixture models

- It is important to distinguish between model combination methods and Bayesian model averaging. Consider the example of density estimation using a mixture of Gaussians in which several Gaussian components are combined probabilistically.
- The model contains a binary latent variable  $\mathbf{z}$  that indicates which component of the mixture is responsible for generating the corresponding data point.
- density over the observed variable  $\mathbf{x}$  is obtained by marginalizing over the latent variable

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}). \quad (11)$$

In the case of our Gaussian mixture example, this leads to a distribution of the form

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (12)$$

and a marginal distribution of a dataset

$$p(\mathbf{X}) = \prod_{n=1}^N p(\mathbf{x}_n) = \prod_{n=1}^N \left[ \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n) \right] \quad (13)$$

# Bayesian model averaging

- Now suppose we have several different models indexed by  $h = 1, \dots, H$  with prior probabilities  $p(h)$ . For instance one model might be a mixture of Gaussians and another model might be a mixture of Cauchy distributions.
- The marginal distribution over the data set is given by

$$p(\mathbf{X}) = \sum_{h=1}^H p(\mathbf{X}|h)p(h). \quad (14)$$

- The interpretation of this summation over  $h$  is that just one model is responsible for generating the whole data set, and the probability distribution over  $h$  simply reflects our uncertainty as to which model that is.
- This highlights the key difference between Bayesian model averaging and model combination, because in Bayesian model averaging the whole data set is generated by a single model. By contrast, when we combine multiple models, as in (13), we see that different data points within the data set can potentially be generated from different values of the latent variable  $\mathbf{z}$  and hence by different components.

- One lesson that we can take from the bias-variance dilemma is that there can be value in averaging multiple models.
- When we averaged a set of low-bias models (corresponding to higher order polynomials), we obtained accurate predictions for the underlying sinusoidal function from which the data were generated.
- The simplest way to construct a committee is to average the predictions of a set of individual models.
- In practice, of course, we have only a single data set, and so we have to find a way to introduce variability between the different models within the committee.

# Bootstrap revisited

- Earlier we saw the use of the bootstrap (sampling with replacement) as a means to obtain a better estimate of generalisation performance.
- Consider a regression problem in which we are trying to predict the value of a single continuous variable, and suppose we generate  $M$  bootstrap data sets and then use each to train a separate copy  $y_m(\mathbf{x})$  of a predictive model where  $m = 1, \dots, M$ .
- The committee prediction is given by

$$y_{\text{COM}}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M y_m(\mathbf{x}). \quad (15)$$

This procedure is known as bootstrap aggregation or **bagging**.

# Bagging performance

- Suppose the true regression function that we are trying to predict is given by  $h(\mathbf{x})$ , so that the output of each of the models can be written as the true value plus an error in the form

$$y_m(\mathbf{x}) = h(\mathbf{x}) + \epsilon_m(\mathbf{x}). \quad (16)$$

- The average sum-of-squares error then takes the form

$$\mathbb{E}_{\mathbf{x}} \left[ \{y_m(\mathbf{x}) - h(\mathbf{x})\}^2 \right] = \mathbb{E}_{\mathbf{x}} [\epsilon_m(\mathbf{x})^2] \quad (17)$$

where  $\mathbb{E}_{\mathbf{x}}[\cdot]$  denotes a frequentist expectation with respect to the distribution of the input vector  $\mathbf{x}$ .

- If we assume that the errors have zero mean and are uncorrelated, then we obtain

$$E_{\text{COM}} = \frac{1}{M} E_{\text{AV}}. \quad (18)$$

# Bagging performance revisited

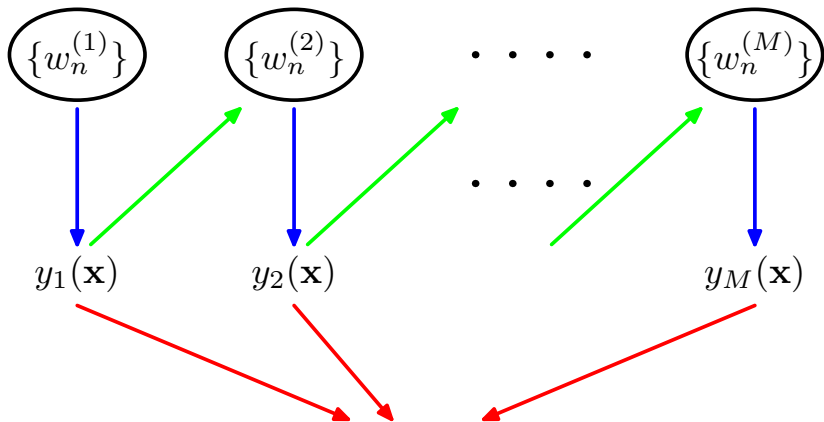
- This apparently dramatic result suggests that the average error of a model can be reduced by a factor of  $M$  simply by averaging  $M$  versions of the model. Unfortunately, it depends on the key assumption that the errors due to the individual models are uncorrelated.
- In practice, the errors are typically highly correlated, and the reduction in overall error is generally small.
- It can, however, be shown that the expected committee error will not exceed the expected error of the constituent models, so that  $E_{\text{COM}} \leq E_{\text{AV}}$ .
- Note that this is a very simple way of combining trained models that does not depend on their internal structure: often beneficial just to try it out.

- **Boosting** is a powerful technique for combining multiple 'base' classifiers to produce a form of committee whose performance can be significantly better than that of any of the base classifiers.
- Here we describe the most widely used form of boosting algorithm called **AdaBoost**, short for 'adaptive boosting'.
- Boosting can give good results even if the base classifiers have a performance that is only slightly better than random, and hence sometimes the base classifiers are known as **weak learners**.
- Originally designed for solving classification problems, boosting can also be extended to regression.

- The base classifiers are trained in sequence, and each base classifier is trained using a weighted form of the data set in which the weighting coefficient associated with each data point depends on the performance of the previous classifiers.
- In particular, points that are misclassified by one of the base classifiers are given greater weight when used to train the next classifier in the sequence.
- Once all the classifiers have been trained, their predictions are then combined through a weighted majority voting scheme.

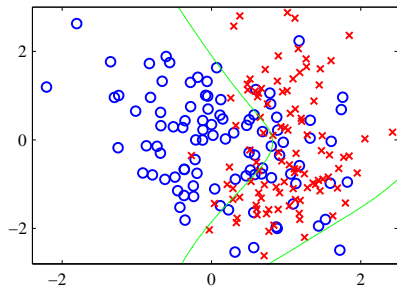


# Boosting schematic

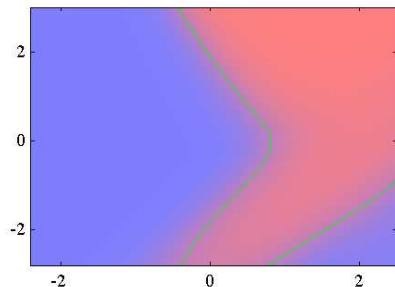


$$Y_M(\mathbf{x}) = \text{sign} \left( \sum_m^M \alpha_m y_m(\mathbf{x}) \right)$$

# Synthetic dataset

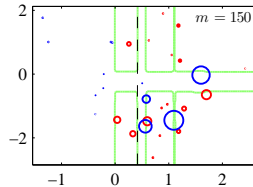
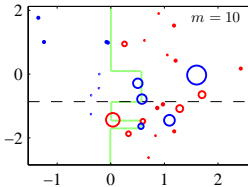
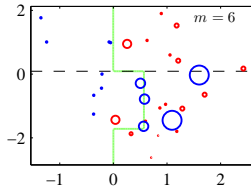
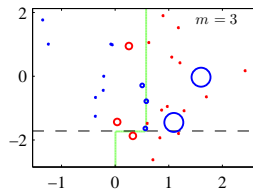
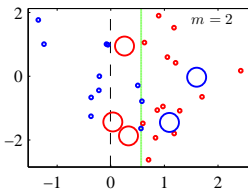
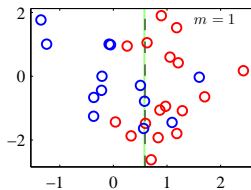


Use 30 points.



True posterior probabilities  
(red/blue scale).

# Boosting in action



- Bias and variance
- Bayesian model averaging
- Committees
- Boosting