# 1   Advanced Evaluation

## 1.1   Slide 5

Schematic illustration of the joint probabilities $p(x, \mathcal{C}_k)$ for each of two classes plotted against $x$, together with the decision boundary $x = \widehat{x}$. Values of $x \geqslant \widehat{x}$ are classified as class $\mathcal{C}_2$ and hence belong to decision region $\mathcal{R}_2$, whereas points $x < \widehat{x}$ are classified as $\mathcal{C}_1$ and belong to $\mathcal{R}_1$. Errors arise from the blue, green, and red regions, so that for $x < \widehat{x}$ the errors are due to points from class $\mathcal{C}_2$ being misclassified as $\mathcal{C}_1$ (represented by the sum of the red and green regions), and conversely for points in the region $x \geqslant \widehat{x}$ the errors are due to points from class $\mathcal{C}_1$ being misclassified as $\mathcal{C}_2$ (represented by the blue region). As we vary the location $\widehat{x}$ of the decision boundary, the combined areas of the blue and green regions remains constant, whereas the size of the red region varies. The optimal choice for $\widehat{x}$ is where the curves for $p(x, \mathcal{C}_1)$ and $p(x, \mathcal{C}_2)$ cross, corresponding to $\widehat{x} = x_0$, because in this case the red region disappears. This is equivalent to the minimum misclassification rate decision rule, which assigns each value of $x$ to the class having the higher posterior probability $p(\mathcal{C}_k|x)$.

## 1.2   Slide 6

Which cost is greater in each case?

- Medical diagnostic tests: does X have leukaemia? Misclassifying someone with cancer.

- Loan decisions: approve mortgage for X? Approving a bad credit risk.

- Web mining: will X click on this link? Neutral?

- Promotional mailing: will X buy the product? Not predicting someone will buy the product.

## 1.3   Slide 10

Illustration of the reject option. Inputs $x$ such that the larger of the two posterior probabilities is less than or equal to some threshold $\theta$ will be rejected.

## 1.4   Slide 13

The class-conditional densities may contain a lot of structure that has little effect on the posterior probabilities Example of the class-conditional densities for two classes having a single input variable $x$ (left plot) together with the corresponding posterior probabilities (right plot). Note that the left-hand mode of the class-conditional density $p(\mathbf{x}|\mathcal{C}_1)$, shown in blue on the left plot, has no effect on the posterior probabilities. The vertical green line in the right plot shows the decision boundary in $x$ that gives the minimum misclassification rate.

## 2   Combining Models

### 2.1   Slide 3

The second term, which is independent of $y(\mathbf{x})$, arises from the intrinsic noise on the data and represents the minimum achievable value of the expected loss.

The first term depends on our choice for the function $y(\mathbf{x})$, and we will seek a solution for $y(\mathbf{x})$ which makes this term a minimum. Because it is nonnegative, the smallest that we can hope to make this term is zero.

If we had an unlimited supply of data (and unlimited computational resources), we could in principle find the regression function $h(\mathbf{x})$ to any desired degree of accuracy, and this would represent the optimal choice for $y(\mathbf{x})$. However, in practice we have a data set $\mathcal{D}$ containing only a finite number $N$ of data points, and consequently we do not know the regression function $h(\mathbf{x})$ exactly.

### 2.2   Slide 5

Consider the integrand of the first term in the squared loss equation which for a particular data set $\mathcal{D}$ takes the form

$$\{y(\mathbf{x};\mathcal{D}) - h(\mathbf{x})\}^2. \tag{2.1}$$

Because this quantity will be dependent on the particular data set $\mathcal{D}$, we take its average over the ensemble of data sets. If we add and subtract the quantity $\mathbb{E}_{\mathcal{D}}[y(\mathbf{x};\mathcal{D})]$ inside the braces, and then expand, we obtain

$$
\begin{aligned}
\{y(\mathbf{x};\mathcal{D}) &- \mathbb{E}_{\mathcal{D}}[y(\mathbf{x};\mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x};\mathcal{D})] - h(\mathbf{x})\}^2 \\
= \ & \{y(\mathbf{x};\mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x};\mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x};\mathcal{D})] - h(\mathbf{x})\}^2 \\
& + 2\{y(\mathbf{x};\mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x};\mathcal{D})]\}\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x};\mathcal{D})] - h(\mathbf{x})\}.
\end{aligned}
\tag{2.2}
$$

We now take the expectation of this expression with respect to $\mathcal{D}$ and note that the final term will vanish, giving the equation on the slide.

### 2.3   Slide 7

Here we generate 100 data sets, each containing $N = 25$ data points, independently from the sinusoidal curve $h(x) = \sin(2\pi x)$. The data sets are indexed by $l = 1, \ldots, L$, where $L = 100$, and for each data set $\mathcal{D}^{(l)}$ we fit a model with 24 Gaussian basis functions by minimizing the regularized error function to give a prediction function $y^{(l)}(x)$.

There are 24 Gaussian basis functions in the model so that the total number of parameters is $M = 25$ including the bias parameter. The first row shows the result of fitting the model to the data sets for various values of $\ln \lambda$ (for clarity, only 20 of the 100 fits are shown). The second row shows the corresponding average of the 100 fits (red) along with the sinusoidal function from which the data sets were generated (green).

The top row corresponds to a large value of the regularization coefficient $\lambda$ that gives low variance (because the red curves in the left plot look similar) but high bias (because the two curves in the right plot are very different). Conversely on the bottom row, for which $\lambda$ is small, there is large variance (shown by the high variability between the red curves in the left plot) but low bias (shown by the good fit between the average model fit and the original sinusoidal function). Note that the result of averaging many solutions for the complex model with $M = 25$ is a very good fit to the regression function, which suggests that averaging may be a beneficial procedure. Indeed, a weighted averaging of multiple solutions lies at the heart of a Bayesian approach, although the averaging is with respect to the posterior distribution of parameters, not with respect to multiple data sets.

## 2.4 Slide 8

Plot of squared bias and variance, together with their sum, corresponding to the results shown in the Figure. Also shown is the average test set error for a test data set size of 1000 points. The minimum value of $(\text{bias})^2 + \text{variance}$ occurs around $\ln \lambda = -0.31$, which is close to the value that gives the minimum error on the test data.

## 2.5 Slide 10

As the size of the data set increases, this uncertainty reduces, and the posterior probabilities $p(h|\mathbf{X})$ become increasingly focussed on just one of the models.

## 2.6 Slide 12

You can read more about bagging in [Breiman, 1996].

## 2.7 Slide 13

The average error made by the models acting individually is therefore

$$E_{\text{AV}} = \frac{1}{M} \sum_{m=1}^{M} \mathbb{E}_{\mathbf{x}} \left[ \epsilon_m(\mathbf{x})^2 \right]. \tag{2.3}$$

Similarly, the expected error from the committee is given by

$$
\begin{aligned}
E_{\text{COM}} &= \mathbb{E}_{\mathbf{x}} \left[ \left\{ \frac{1}{M} \sum_{m=1}^{M} y_m(\mathbf{x}) - h(\mathbf{x}) \right\}^2 \right] \\
&= \mathbb{E}_{\mathbf{x}} \left[ \left\{ \frac{1}{M} \sum_{m=1}^{M} \epsilon_m(\mathbf{x}) \right\}^2 \right]
\end{aligned}
\tag{2.4}
$$

## 2.8   Slide 15

Original paper on AdaBoost is [Freund, 1996]. Extension to regression models can be found in [Friedman, 2001].

## 2.9   Slide 17

Each base classifier $y_m(\mathbf{x})$ is trained on a weighted form of the training set (blue arrows) in which the weights $w_n^{(m)}$ depend on the performance of the previous base classifier $y_{m-1}(\mathbf{x})$ (green arrows). Once all base classifiers have been trained, they are combined to give the final classifier $Y_M(\mathbf{x})$ (red arrows).

Details of the AdaBoost algorithm.

1. Initialize the data weighting coefficients $\{w_n\}$ by setting $w_n^{(1)} = 1/N$ for $n = 1, \ldots, N$.

2. For $m = 1, \ldots, M$:

    (a) Fit a classifier $y_m(\mathbf{x})$ to the training data by minimizing the weighted error function

    $$J_m = \sum_{n=1}^{N} w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n) \tag{2.5}$$

    where $I(y_m(\mathbf{x}_n) \neq t_n)$ is the indicator function and equals 1 when $y_m(\mathbf{x}_n) \neq t_n$ and 0 otherwise.

    (b) Evaluate the quantities

    $$\epsilon_m = \frac{\sum_{n=1}^{N} w_n^{(m)} I(y_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^{N} w_n^{(m)}} \tag{2.6}$$

    and then use these to evaluate

    $$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}. \tag{2.7}$$

    (c) Update the data weighting coefficients

    $$w_n^{(m+1)} = w_n^{(m)} \exp\{\alpha_m I(y_m(\mathbf{x}_n) \neq t_n)\} \tag{2.8}$$

3. Make predictions using the final model, which is given by

$$Y_M(\mathbf{x}) = \text{sign}\left(\sum_{m=1}^{M} \alpha_m y_m(\mathbf{x})\right). \tag{2.9}$$

### 2.10    Slide 19

In this example, each base learner consists of a threshold on one of the input variables. This simple classifier corresponds to a form of decision tree known as a 'decision stumps', i.e., a decision tree with a single node. Thus each base learner classifies an input according to whether one of the input features exceeds some threshold and therefore simply partitions the space into two regions separated by a linear decision surface that is parallel to one of the axes.

The figure illustrates boosting in which the base learners consist of simple thresholds applied to one or other of the axes. Each figure shows the number $m$ of base learners trained so far, along with the decision boundary of the most recent base learner (dashed black line) and the combined decision boundary of the ensemble (solid green line). Each data point is depicted by a circle whose radius indicates the weight assigned to that data point when training the most recently added base learner. Thus, for instance, we see that points that are misclassified by the $m = 1$ base learner are given greater weight when training the $m = 2$ base learner.

# References

[Breiman, 1996]  Breiman, L. 1996. Bagging predictors. *Machine learning* **24** (2), 123–140.

[Freund, 1996]  Freund, Y. 1996. Experiments with a new boosting algoritm. In *13th International Conference on Machine Learning, 1996.*

[Friedman, 2001]  Friedman, J. H. 2001. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 1189–1232.