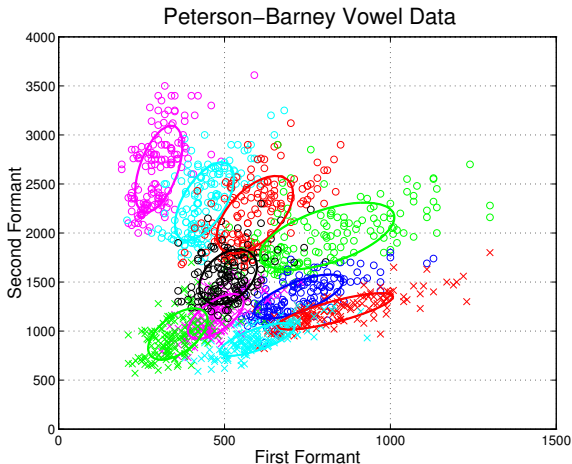


4F10: Probability of Error & Decision Boundaries

Mark Gales

Michaelmas 2021

Vowel Classification Using Formant Features



- No separation of classes given observation \mathbf{x}^* (uncertainty)
 - there is no “perfect” decision

Supervised Training and Evaluation Data

- Interested in supervised training data for classification

$$\mathcal{D} = \{\{\mathbf{x}_1, y_1\}, \dots, \{\mathbf{x}_N, y_N\}\}$$

- \mathbf{x}_i : the observation, feature vector
- $y_i \in \{\omega_1, \dots, \omega_K\}$: class label for observation \mathbf{x}_i
- Samples are “draws” from joint distribution $p(\omega, \mathbf{x})$
- “Standard” process for obtaining data for a task
 - obtain set of observations (features) \mathbf{x}_i and labels y_i

$$\mathbf{x}_i \sim p(\mathbf{x}); \quad y_i \sim P(\omega|\mathbf{x}_i)$$

- BUT care about performance on unseen data: split data
 - training data: used to train the model parameters
 - evaluation data: used to evaluate model on unseen data

Expected Loss

- Given any decision, there will be an associated “Loss”
 - useful for any system if this value is known! Usually

$$\mathcal{L}_{\text{act}} = \int \left[\sum_{i=1}^K \mathcal{L}(f(\mathbf{x}, \boldsymbol{\theta}), \omega_i) P(\omega_i | \mathbf{x}) \right] p(\mathbf{x}) d\mathbf{x} \geq \mathcal{L}_{\text{emp}} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i, \boldsymbol{\theta}), y_i)$$

- $f(\mathbf{x}, \boldsymbol{\theta})$ is the prediction given model parameters $\boldsymbol{\theta}$
 - as $N \rightarrow \infty$ empirical and actual losses the same $\mathcal{L}_{\text{emp}} \rightarrow \mathcal{L}_{\text{act}}$
- Simple approach is a **held-out** evaluation set of N^* samples
 - observation \mathbf{x}_i^* , and label, $y_i^* \in \{\omega_1, \dots, \omega_K\}$, pairs
 - draws from $\mathbf{x}_i^* \sim p(\mathbf{x})$ followed by $y_i^* \sim P(\omega | \mathbf{x}_i^*)$
 - use this data to compute held-out empirical loss $\mathcal{L}_{\text{eval}}$

$$\mathcal{L}_{\text{eval}} = \frac{1}{N^*} \sum_{i=1}^{N^*} \mathcal{L}(f(\mathbf{x}_i^*, \boldsymbol{\theta}), y_i^*) \approx \mathcal{L}_{\text{act}}$$

Bayes' Decision Rule

- Consider making a decision with classifier $f(\mathbf{x}^*, \boldsymbol{\theta})$
 - let $f(\mathbf{x}^*, \boldsymbol{\theta}) = \hat{\omega}$, and assume following loss function

$$\mathcal{L}(\hat{\omega}, \omega_i) = \begin{cases} 0, & \hat{\omega} = \omega_i \\ 1, & \text{otherwise} \end{cases}$$

- Apply Bayes decision rule
 - make decision that minimises loss (probability of error)

$$\begin{aligned} \hat{\omega} &= f(\mathbf{x}^*, \boldsymbol{\theta}) \\ &= \arg \min_{\omega} \left\{ \sum_{i=1}^K \mathcal{L}(\omega, \omega_i) P(\omega_i | \mathbf{x}^*; \boldsymbol{\theta}) \right\} \\ &= \arg \max_{\omega} \{ P(\omega | \mathbf{x}^*; \boldsymbol{\theta}) \} \end{aligned}$$

- select the class with the highest posterior
- but need to train $\boldsymbol{\theta}$ to obtain $P(\omega | \mathbf{x}^*; \boldsymbol{\theta})$

Two-Class (Binary) Classification

- Consider a simple **binary classification** task
 - class ω_1 has label 1, class ω_2 has label -1
 - prediction from $f(\mathbf{x}_i^*, \boldsymbol{\theta})$, $\hat{\omega}$, is either 1 or -1
 - 1/0 loss function

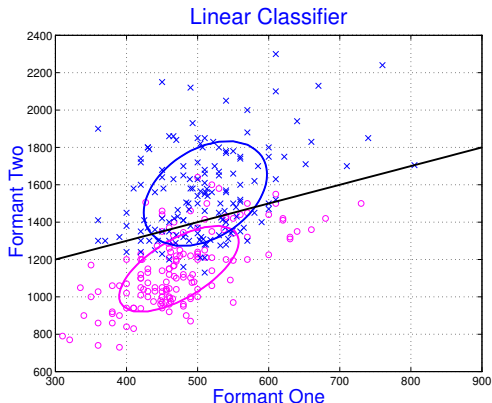
$$\mathcal{L}(\hat{\omega}, \omega_j) = \begin{cases} 0, & \hat{\omega} = \omega_j \\ 1, & \text{otherwise} \end{cases}$$

- Estimate of probability of error (equal “loss”) becomes

$$\begin{aligned} P(\text{error}) \approx \mathcal{L}_{\text{eval}} &= \frac{1}{N^*} \sum_{i=1}^{N^*} \mathcal{L}(f(\mathbf{x}_i^*, \boldsymbol{\theta}), y_i^*) \\ &= \frac{1}{2N^*} \sum_{i=1}^{N^*} |f(\mathbf{x}_i^*, \boldsymbol{\theta}) - y_i^*| \end{aligned}$$

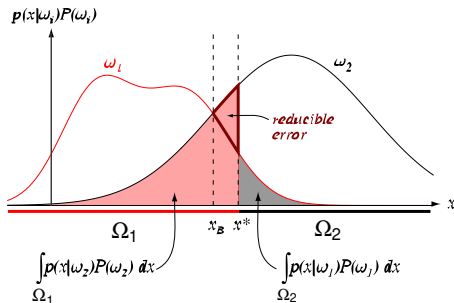
- count the number of errors, divides by total number of samples

Binary Classification



- Classifier partitions feature space into regions
 - for each region there is an associated label
 - simple linear decision boundary illustrated above

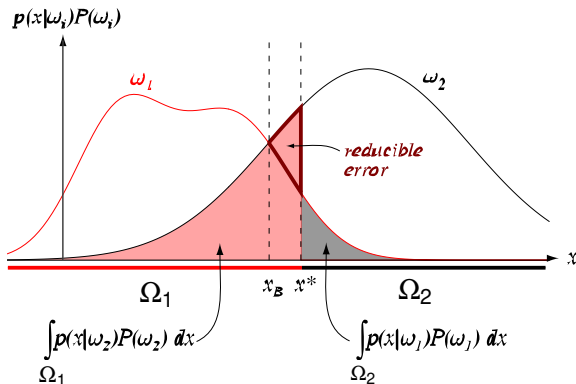
True Probability of Error



- $f(\mathbf{x}, \theta)$ yields two regions $\Omega_1 \rightarrow \omega_1$, and $\Omega_2 \rightarrow \omega_2$,

$$\begin{aligned} P(\text{error}) &= P(\mathbf{x} \in \Omega_2, \omega_1) + P(\mathbf{x} \in \Omega_1, \omega_2) \\ &= P(\mathbf{x} \in \Omega_2 | \omega_1)P(\omega_1) + P(\mathbf{x} \in \Omega_1 | \omega_2)P(\omega_2) \\ &= \int_{\Omega_2} p(\mathbf{x} | \omega_1)P(\omega_1)d\mathbf{x} + \int_{\Omega_1} p(\mathbf{x} | \omega_2)P(\omega_2)d\mathbf{x} \end{aligned}$$

Probability of Error Visualisation (from DHS)



- Optimal **decision boundary** at x_B (suboptimal at x^*)
 - no lower probability of error can be obtained
 - same **decision boundary** as Bayes' decision rule

- Need to decide the form of classifier (with parameters θ)
- Classifiers for making these decisions can be broadly split as:
 - **Generative models:** model the joint distribution of observations and classes is trained, $p(\mathbf{x}, \omega; \theta)$.
The posterior of class ω_i is then obtained from **Bayes' rule**

$$P(\omega_i | \mathbf{x}^*; \theta) = \frac{p(\mathbf{x}^*, \omega_i; \theta)}{\sum_{j=1}^K p(\mathbf{x}^*, \omega_j; \theta)}$$

- **Discriminative models:** a model of the posterior distribution of the class given the observation is trained, $P(\omega | \mathbf{x}^*; \theta)$.
- **Discriminant functions:** a mapping from an observation \mathbf{x}^* to a class ω is directly trained. No posterior probability, $P(\omega | \mathbf{x})$, generated just class label.

See Bishop for a discussion of the merits of these.

- Probability of error marginalises over joint distribution

$$P(\text{error}) = \int_{\Omega_2} p(\mathbf{x}, \omega_1) d\mathbf{x} + \int_{\Omega_1} p(\mathbf{x}, \omega_2) d\mathbf{x}$$

- can be expressed as (generative model)

$$P(\text{error}) = \int_{\Omega_2} p(\mathbf{x}|\omega_1)P(\omega_1)d\mathbf{x} + \int_{\Omega_1} p(\mathbf{x}|\omega_2)P(\omega_2)d\mathbf{x}$$

- can also be expressed as (discriminative model)

$$P(\text{error}) = \int_{\Omega_2} P(\omega_1|\mathbf{x})p(\mathbf{x})d\mathbf{x} + \int_{\Omega_1} P(\omega_2|\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

Unequal loss function

- Loss function may be unequal for particular applications
 - consider speaker verification for bank account access
- Consider the following loss for a binary task
 - here correct classification has zero loss

$$\mathcal{L}(f(\mathbf{x}^*, \boldsymbol{\theta}), \omega_1) = \begin{cases} 0, & f(\mathbf{x}^*, \boldsymbol{\theta}) = \omega_1 \\ C_{21}, & f(\mathbf{x}^*, \boldsymbol{\theta}) = \omega_2 \end{cases}$$

$$\mathcal{L}(f(\mathbf{x}^*, \boldsymbol{\theta}), \omega_2) = \begin{cases} C_{12}, & f(\mathbf{x}^*, \boldsymbol{\theta}) = \omega_1 \\ 0, & f(\mathbf{x}^*, \boldsymbol{\theta}) = \omega_2 \end{cases}$$

- Apply Bayes' decision rule to classification minimising loss

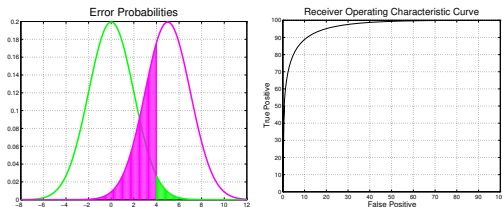
$$\hat{\omega} = \arg \min_{\omega} \left\{ \sum_{i=1}^2 \mathcal{L}(\omega, \omega_i) P(\omega_i | \mathbf{x}^*; \boldsymbol{\theta}) \right\}$$

Receiver Operating Characteristics Curve

- This decision rule can be expressed as classifying \mathbf{x}^* using

$$\frac{C_{21}P(\omega_1|\mathbf{x}^*; \boldsymbol{\theta})}{C_{12}P(\omega_2|\mathbf{x}^*; \boldsymbol{\theta})} \begin{matrix} \omega_1 \\ > \\ \omega_2 \end{matrix} 1, \quad \frac{P(\omega_1|\mathbf{x}^*; \boldsymbol{\theta})}{P(\omega_2|\mathbf{x}^*; \boldsymbol{\theta})} \begin{matrix} \omega_1 \\ > \\ \omega_2 \end{matrix} \frac{C_{12}}{C_{21}}$$

- the “cost” ratio C_{12}/C_{21} is effectively an operating threshold
- Possible to produce curves by changing this threshold
 - ω_1 = positive, ω_2 = negative, equal priors $P(\omega_1) = P(\omega_2)$



Decision Boundaries

- From Bayes decision rule for point \mathbf{x}^*
 - assuming equal loss

$$\hat{\omega} = \arg \max_{\omega} \{P(\omega|\mathbf{x}^*; \boldsymbol{\theta})\}$$

- we can select the form of distribution
 - **but** we don't know the model parameters $\boldsymbol{\theta}$...
- Need to find $\boldsymbol{\theta}$
 - this is the model that we train - use supervised training

$$\mathcal{D} = \{\{\mathbf{x}_1, y_1\}, \dots, \{\mathbf{x}_N, y_N\}\}$$

- draws from $\mathbf{x}_i \sim p(\mathbf{x})$ followed by $y_i \sim P(\omega|\mathbf{x}_i)$
 - use this data to estimate $\boldsymbol{\theta}$

- Need to get the “best” set of model parameters
 - find the parameters **most-likely** to generate the training labels

(Conditional) Maximum Likelihood Estimation (MLE)

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} \{ \log (P(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N; \theta)) \} \\ &= \arg \max_{\theta} \left\{ \sum_{i=1}^N \log (P(y_i | \mathbf{x}_i; \theta)) \right\}\end{aligned}$$

- assumes that all training samples are independent
- Can also be applied to distribution estimation
 - unsupervised training case

$$\hat{\theta} = \arg \max_{\theta} \left\{ \sum_{i=1}^N \log (p(\mathbf{x}_i; \theta)) \right\}$$

- Use **generative model** to approximate $P(\omega|\mathbf{x}^*)$
 - generative model, parameters θ , **joint distribution** $p(\mathbf{x}, \omega)$
 - posterior distribution of class ω_i for unseen observation \mathbf{x}^*

$$P(\omega_i|\mathbf{x}^*) \approx \frac{p(\mathbf{x}^*, \omega_i; \theta)}{\sum_{j=1}^K p(\mathbf{x}^*, \omega_j; \theta)} = \frac{p(\mathbf{x}^*|\omega_i; \theta)P(\omega_i)}{\sum_{j=1}^K p(\mathbf{x}^*|\omega_j; \theta)P(\omega_j)}$$

- $P(\omega_i)$ is the **prior** for class ω_i
 - $p(\mathbf{x}^*|\omega_i; \theta)$ is the **likelihood** of the observation given class ω_i
- Use maximum likelihood (ML) training to estimate θ
 - separate model trained for each class ω_i , θ_i

$$\hat{\theta}_i = \arg \max_{\theta} \left\{ \sum_{j: y_j = \omega_i} \log (p(\mathbf{x}_j|\omega_i; \theta)) \right\}$$

Multivariate Gaussian Class Conditional PDFs

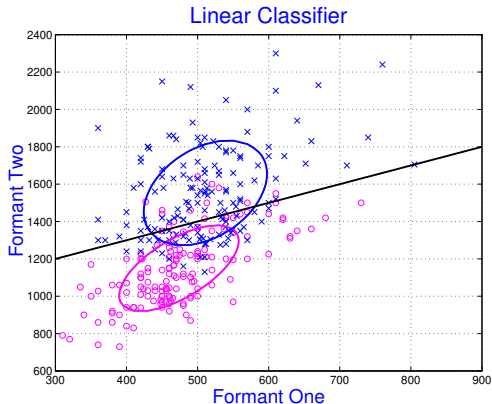
- Need to decide on form and parameters of $p(\mathbf{x}|\omega_i; \boldsymbol{\theta})$
 - use multivariate Gaussian distribution as class-conditional PDF

$$p(\mathbf{x}|\omega_i; \boldsymbol{\theta}_i) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\boldsymbol{\Sigma}_i|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i)\right)$$

- \mathbf{x} : the d -dimensional observation (input) vector
 - $\boldsymbol{\mu}_i$: the d -dimensional mean vector for class ω_i
 - $\boldsymbol{\Sigma}_i$: the $d \times d$ covariance matrix for class ω_i
- Maximum likelihood estimates of the parameters given by

$$\begin{aligned}\hat{\boldsymbol{\mu}}_i &= \frac{\sum_{j:y_j=\omega_i} \mathbf{x}_j}{\sum_{j:y_j=\omega_i} 1} \\ \hat{\boldsymbol{\Sigma}}_i &= \frac{\sum_{j:y_j=\omega_i} (\mathbf{x}_j - \hat{\boldsymbol{\mu}}_i)(\mathbf{x}_j - \hat{\boldsymbol{\mu}}_i)^T}{\sum_{j:y_j=\omega_i} 1}\end{aligned}$$

Two-Class (Binary) Decision Boundaries



- Decision boundary is (hyper-)plane between class labels
 - for two-class task when the class probabilities the same

Decision Boundaries Equation

- Boundary occurs when (log) class posteriors are the same
 - a point \mathbf{x} on the decision boundary (2 class) satisfies

$$\log(P(\omega_1|\mathbf{x}; \boldsymbol{\theta})) = \log(P(\omega_2|\mathbf{x}; \boldsymbol{\theta}))$$

- for a generative classifier (denominator cancels) this yields

$$\log(P(\omega_1)p(\mathbf{x}|\omega_1; \boldsymbol{\theta}_1)) = \log(P(\omega_2)p(\mathbf{x}|\omega_2; \boldsymbol{\theta}_2))$$

- Consider multivariate Gaussian class-conditional PDFs
 - Gaussian parameters: Class $\omega_1 : \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1$, Class $\omega_2 : \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2$
 - prior parameters: Class $\omega_1 : P(\omega_1)$, Class $\omega_2 : P(\omega_2)$

What form does the decision boundary take?

Multivariate Gaussian Decision Boundaries (cont)

- Substituting for class-conditional PDFs and simplifying yields

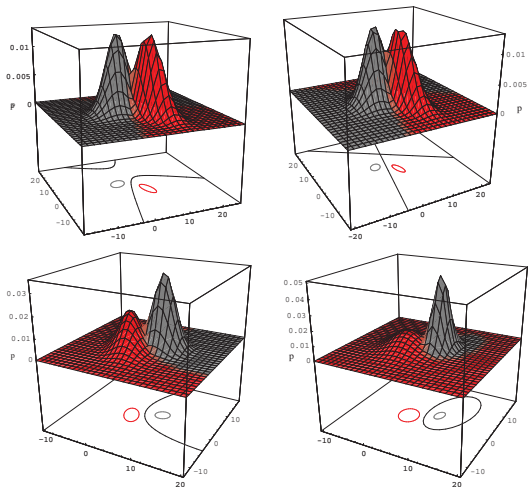
$$\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c = 0$$

where

- $\mathbf{A} = \boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1}$
- $\mathbf{b} = 2 \left(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 \right)$
- $c = \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2 - \log \left(\frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} \right) - 2 \log \left(\frac{P(\omega_1)}{P(\omega_2)} \right)$
- General form yields **hyper-quadratic** decision boundaries
- Constrain covariance matrices to be the same $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}$
 - results in a linear decision boundary

$$\mathbf{b}^T \mathbf{x} + c = 0$$

Example Binary Decision Boundaries (from DHS)



How good is the Classifier?

- Bayes' decision rule minimises the loss (probability of error)
 - **but** relies on the estimates of class posteriors $P(\omega|\mathbf{x}^*; \boldsymbol{\theta})$
 - how close to optimal depends on the accuracy of this estimate
- Class posterior estimates depend on:
 - **training data** (\mathcal{D}): quantity (infinite!), “matched” to test data
 - **form of model**: is the model “correct”?
 - **optimisation**: have the global optimal parameters been found
- In practice none of these are usually true ...
 - many design decisions - good engineering/maths required