

Machine Learning Foundations

Maheesan Niranjana

School of Electronics and Computer Science
University of Southampton

Summer School 2020

Machine Learning as Data-driven Modelling

Single-slide overview of the subject and challenging questions

Data	$\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N$	$\{\mathbf{x}_n\}_{n=1}^N$
Function Approximator	$\mathbf{y} = f(\mathbf{x}, \boldsymbol{\theta}) + v$	
Parameter Estimation	$E_0 = \sum_{n=1}^N \{\ \mathbf{y}_n - f(\mathbf{x}_n; \boldsymbol{\theta})\ \}^2$	
Prediction	$\hat{\mathbf{y}}_{N+1} = f(\mathbf{x}_{N+1}, \hat{\boldsymbol{\theta}})$	
Regularization	$E_1 = \sum_{n=1}^N \{\ \mathbf{y}_n - f(\mathbf{x}_n)\ \}^2 + g(\ \boldsymbol{\theta}\)$	
Modelling Uncertainty	$p(\boldsymbol{\theta} \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N)$	
Probabilistic Inference	$\mathbf{E}[g(\boldsymbol{\theta})] = \int g(\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} = \frac{1}{N_s} \sum_{n=1}^{N_s} g(\boldsymbol{\theta}^{(n)})$	
Sequential Estimation	$\boldsymbol{\theta}(n-1 n-1) \longrightarrow \boldsymbol{\theta}(n n-1) \longrightarrow \boldsymbol{\theta}(n n)$ Kalman & Particle Filters; Reinforcement Learning	

Bayesian Decision Theory

- Classes: $\omega_i, i = 1, \dots, K$
- Prior Probabilities: $P[\omega_1], \dots, P[\omega_K];$
 $P[\omega_i] \geq 0, \sum_{i=1}^K P[\omega_i] = 1$
- Likelihoods (class conditional probabilities): $p(\mathbf{x}|\omega_i), i = 1, \dots, K$
- Posterior Probability: $P[\omega_j | \mathbf{x}]$

$$P[\omega_j | \mathbf{x}] = \frac{p(\mathbf{x} | \omega_j) P[\omega_j]}{\sum_{i=1}^K p(\mathbf{x} | \omega_i) P[\omega_i]}$$

- From prior knowledge: $P[\omega_i]$; From training data: $p(\mathbf{x}|\omega_i)$
- Decision rule: Assign \mathbf{x} to the class that maximizes posterior probability.
- The denominator is a constant; i.e. does not depend on ω_j
- Hence the decision rule becomes:

$$\mathbf{x} \in \max_j p(\mathbf{x} | \omega_j) P[\omega_j]$$

Bayes' Classifier for Gaussian Densities

Make assumptions, cancel common terms when making comparisons...

- Decision rule from: $p(\mathbf{x} | \omega_j) P[\omega_j]$
- Assume the two classes are Gaussian distributed with distinct means and identical covariance matrices
 $p(\mathbf{x} | \omega_j) = \mathcal{N}(\mathbf{m}_j, \mathbf{C})$
- Substitute into Bayes' classifier decision rule

$$\begin{aligned} P[\omega_1 | \mathbf{x}] &\leq P[\omega_2 | \mathbf{x}] \\ p(\mathbf{x} | \omega_1) P[\omega_1] &\leq p(\mathbf{x} | \omega_2) P[\omega_2] \end{aligned}$$

$$\begin{aligned} \frac{1}{(2\pi)^{p/2}(\det(\mathbf{C}))^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mathbf{m}_1)^t \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m}_1) \right\} P[\omega_1] &\leq \\ \frac{1}{(2\pi)^{p/2}(\det(\mathbf{C}))^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mathbf{m}_2)^t \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m}_2) \right\} P[\omega_2] \end{aligned}$$

Bayes' classifier for simple densities (cont'd)

Distinct Means; Equal, isotropic covariance matrix

- Suppose the densities are isotropic and priors are equal
i.e. $\mathbf{C} = \sigma^2 \mathbf{I}$ and $P[\omega_1] = P[\omega_2]$
- The comparison simplifies to (see algebra on board):

$$\begin{aligned}(\mathbf{x} - \mathbf{m}_1)^t (\mathbf{x} - \mathbf{m}_1) &\leq (\mathbf{x} - \mathbf{m}_2)^t (\mathbf{x} - \mathbf{m}_2) \\ |\mathbf{x} - \mathbf{m}_1| &\leq |\mathbf{x} - \mathbf{m}_2|\end{aligned}$$

- The above is a simple *distance to mean* classifier
- Under the above simplistic assumptions, we only need to store one template per class (the means)!

Bayes' classifier for simple densities (cont'd)

Distinct Means; Common covariance matrix (but not isotropic)

- Cancel common terms and take log

$$(\mathbf{x} - \mathbf{m}_1)^t \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}_1) \leq (\mathbf{x} - \mathbf{m}_2)^t \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}_2) - \log \left\{ \frac{P[\omega_1]}{P[\omega_2]} \right\}$$

- Also simplifies to a linear classifier!

$$\mathbf{w}^t \mathbf{x} + b \leq 0$$

$$\mathbf{w} = 2\mathbf{C}^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

$$b = (\mathbf{m}_1^t \mathbf{C}^{-1} \mathbf{m}_1 - \mathbf{m}_2^t \mathbf{C}^{-1} \mathbf{m}_2) - \log \left\{ \frac{P[\omega_1]}{P[\omega_2]} \right\}$$

- Also a distance to template classifier, where the distance is

$$(\mathbf{x} - \mathbf{m}_1)^t \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m}_1)$$

Known as Mahalanobis distance

Posterior probabilities for simple Gaussian cases

- Bayes classifier:

$$P[\omega_1|\mathbf{x}] = \frac{p(\mathbf{x}|\omega_1) P[\omega_1]}{p(\mathbf{x}|\omega_1) P[\omega_1] + p(\mathbf{x}|\omega_2) P[\omega_2]}$$

- Restrictive assumptions:
 - Gaussian $p(\mathbf{x}|\omega_j) = \mathcal{N}(\mathbf{m}_j, \mathbf{C}_j)$
 - Equal covariance matrices: $\mathbf{C}_1 = \mathbf{C}_2 = \mathbf{C}$
- Substitute, divide through by numerator term and cancel common terms to get

$$P[\omega_1|\mathbf{x}] = \frac{1}{1 + \exp\{-\mathbf{w}^t \mathbf{x} + w_0\}}$$

- The functional form $1/(1 + \exp(-\alpha))$ is known as sigmoid / logistic (See lab class for W3)

Bayesian decision theory is fundamental to machine learning; we started from a probabilistic setting in which we described the posterior probability of membership and derived several results starting from this premise. You need to be able to:

- Derive and demonstrate that under certain assumptions...
 - The class boundary is linear
 - The posterior probability has a sigmoidal shape
 - The optimal classifier reduces to a distance to template classifier
- ...and under certain other assumptions...
 - The best classifier is still a distance to template classifier, but instead of Euclidean distance we need to use Mahalanobis distance.
- ...and under certain other assumptions...
 - The optimal classifier is a quadratic classifier