

SMLDM HT 2014 - MSc Problem Sheet 3

1. In this question we will take a Bayesian approach to learning decision trees. Assume that we have a binary classification problem, with classes $\{0, 1\}$ and we have n data items $(x_i, y_i)_{i=1}^n$.

Recall the greedy tree growing heuristics for decision trees: we start with the root, for each leaf node of the tree we find an optimal feature j and split point v , according to some criteria, to split the node, and recurse on both sides.

- (a) Consider a decision tree T where the data space has been split into R regions $\mathcal{R}_1, \dots, \mathcal{R}_R$. Each region \mathcal{R}_m corresponds to a leaf in T and is a rectangle, with side $[a_j, b_j]$ in the j th dimension, $j = 1, \dots, p$, and has a parameter β_m which gives the probability of class 1 for data vectors in that region.

What is a conjugate prior for the parameters β ? Calculate the marginal probability $p(\mathbf{Y}|\mathbf{X}, T)$ of the responses given the data vectors and tree T , marginalizing out β .

- (b) Consider a greedy model selection procedure for determining the structure of T :
- i. We start with a trivial tree with a single node.
 - ii. At each iteration we consider expanding a leaf node m of the tree by creating a split at feature j , value v . This produces a tree T' with two more nodes than T , both children of node m .
 - iii. We compute the marginal probability of \mathbf{Y} under T' , for each j and v , and find the split producing the highest marginal probability.
 - iv. If the marginal probability of the resulting T' is larger than T , we split the node, otherwise we consider expanding other nodes.
 - v. We stop once all leaf nodes of the current tree T have been considered for expansion, but all lead to trees T' with lower marginal probability than T .

Calculate the marginal probability $p(\mathbf{Y}|\mathbf{X}, T')$ of the responses given the data vectors under tree T' .

- (c) Explain how the ratio of marginal probabilities under T' and T (the so-called **Bayes factor**) simplifies to a function which depends only on the data items under region \mathcal{R}_m .
- (d) For each j , explain why the marginal probability under T' is a piecewise constant function of v .
- (e) Describe an algorithm that can determine the optimal split of node m , with computational cost $p \times N_m$, where N_m is the number of data items in region m , and p the number of features.
2. Recall the definition of a 1 hidden layer neural network for binary classification in the lectures. The objective function is:

$$J = - \sum_{i=1}^n y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i) + \frac{1}{2} \sum_{jk} C |W_{jk}^h|^2 + \frac{1}{2} \sum_k C |W_k^o|^2$$

and the network definition is:

$$\hat{y}_i = s \left(b^o + \sum_{k=1}^m W_k^o h_{ik} \right) \quad h_{ik} = s \left(b_k^h + \sum_{j=1}^p W_{jk}^h x_{ij} \right)$$

(a) Verify that the derivatives needed for gradient descent are:

$$\frac{dJ}{dW_k^o} = CW_k^o + \sum_{i=1}^n (\hat{y}_i - y_i) h_{ik}$$

$$\frac{dJ}{dW_{jk}^h} = CW_{jk}^h + \sum_{i=1}^n (\hat{y}_i - y_i) W_k^o h_{ik} (1 - h_{ik}) x_{ij}$$

(b) Suppose instead that you have an L layer neural network for binary classification, with each hidden layer having m neurons with logistic nonlinearity. Define carefully the network, giving the parameterization of each layer, and derive the backpropagation algorithm to compute the derivatives of the objective with respect to the parameters. You may ignore bias terms for simplicity.

3. A **mixture of experts** is a type model in which a number of experts “compete” to predict a label.

Consider a regression problem with dataset $(x_i, y_i)_{i=1}^n$ and $y_i \in \mathbb{R}$. We have E experts, each being a parametrized function $f_j(x; \theta_j)$, for $j = 1, \dots, E$. For example each expert could be a neural network. Each expert $f_j(x; \theta_j)$ tries to predict the response y corresponding to data vector x .

(a) A simple mixture of experts model uses as it’s objective function

$$J(\pi, \sigma^2, (\theta_j)_{j=1}^E) = \sum_{i=1}^n \log \sum_{j=1}^E \pi_j e^{-\frac{1}{2\sigma^2} \|f_j(x_i; \theta_j) - y_i\|^2}$$

where $\pi = (\pi_1, \dots, \pi_E)$ are mixing proportions and σ^2 is a parameter.

Relate the objective function to the log likelihood of a mixture model where each component is a conditional distribution of Y given $X = x$.

- (b) Differentiate the objective function with respect to θ_j , interpreting the computation of the derivative as a generalized EM algorithm, where in the E step the posterior distribution is computed, and in the M step gradient descent is used to update the expert parameters θ_j .
- (c) A mixture of experts allows each expert to specialize in predicting the response in a certain part of the data space, with the overall model having better predictions than any one of the experts.

However to encourage this specialization, it is useful also for the mixing proportions to depend on the data vectors x , i.e. to model $\pi_j(x; \phi)$ as a function of x with parameters ϕ . The idea is that this **gating network** controls where each expert specializes. To ensure $\sum_{j=1}^E \pi_j(x; \phi) = 1$, we can use the softmax nonlinearity:

$$\pi_j(x; \phi) = \frac{\exp(g_j(x; \phi_j))}{\sum_{\ell=1}^E \exp(g_\ell(x; \phi_\ell))}$$

where $g_j(x; \phi_j)$ are parameterized functions for the gating network.

The previous generalized EM algorithm extends to this scenario easily. Describe the latent variables you will need to introduce into the system, the free energy lower bound on the log likelihood, and derive the E step and generalized M step for ϕ_j from the free energy.