

Module 4F10: DEEP LEARNING AND STRUCTURED DATA

Examples Paper 1

Straightforward questions are marked †

*Tripos standard (but not necessarily Tripos length) questions are marked **

Bayes Risk and Decision Boundaries

1. In many pattern classification problems, there is the option either to assign the pattern to one of the c classes, or to reject it as being unrecognizable. If the cost to reject is not too high, rejection may be a desirable action. For particular pattern classification task the cost of classification be defined as

$$\lambda(\omega_i|\omega_j) = \begin{array}{ll} 0 & \omega_i = \omega_j \quad (\text{i.e. (Correct classification)}) \\ \lambda_r & \omega_i = \omega_0 \quad (\text{i.e. Rejection}) \\ \lambda_s & \text{Otherwise} \quad (\text{i.e. Substitution Error}) \end{array}$$

Show that for minimum risk classification, the decision rule should associate a test vector \mathbf{x}^* with class ω_i , if $P(\omega_i|\mathbf{x}^*) \geq P(\omega_j|\mathbf{x}^*)$ for all j **and** $P(\omega_i|\mathbf{x}^*) \geq 1 - \lambda_r/\lambda_s$, and reject otherwise.

2. A logistic regression discriminative classifier is to be used for a two-class, ω_1 and ω_2 , classification problem. Given a test vector \mathbf{x}^* the posterior probability of class one, ω_1 , is given by

$$P(\omega_1|\mathbf{x}^*; \mathbf{a}, b) = \phi(\mathbf{a}^\top \mathbf{x}^* + b) = \frac{1}{1 + \exp(-(\mathbf{a}^\top \mathbf{x}^* + b))}$$

where \mathbf{a} and b are the model parameters. To train these parameters there is supervised training data, $\mathcal{D} = \{\{\mathbf{x}_1, y_1\}, \dots, \{\mathbf{x}_N, y_N\}\}$, where $y_i \in \{\omega_1, \omega_2\}$.

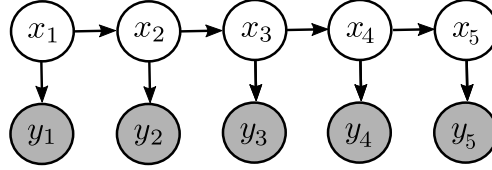
- (a) Give the equation for a point \mathbf{x} that lies on the decision boundary between classes ω_1 and ω_2 .
- (b) Give an expression that can be used to train the model parameters, \mathbf{a} and b , using the training data \mathcal{D} . Briefly discuss how the parameters can be found.
- (c) An alternative generative model is proposed using multivariate Gaussian class-conditional PDFs. Here for class ω_i ($i = 1, 2$)

$$p(\mathbf{x}|\omega_i; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

Discuss how the model parameters for this generative model can be trained. What constraints on the model parameters are required to yield the same form of decision boundary as the logistic regression classifier?

Graphical Models

3. Let $p(x_1, \dots, x_5, y_1, \dots, y_5)$ be specified by the following Bayesian network:



- What is the factorization of $p(x_1, \dots, x_5, y_1, \dots, y_5)$ implied by the graph?
- Write the set of conditional independencies (CIs) implied by the graph.
- Let x_1, \dots, x_5 be discrete random variables taking ten possible values each. What is the cost of computing $p(x_1, \dots, x_5 | y_1, \dots, y_5)$ in terms of sum operations if we ignore the CIs present in $p(x_1, \dots, x_5, y_1, \dots, y_5)$, that is, if we assume full generality for this distribution?
- What is the cost of computing $p(x_1, \dots, x_5 | y_1, \dots, y_5)$ in terms of sum operations if we do make use of the factorization implied by the graph?

Latent Variable Models

4. A hidden Markov model (HMM) is to be used to model a sequence of T , d -dimensional, vectors, $\mathbf{x}_1, \dots, \mathbf{x}_T$. Each element of the vector \mathbf{x}_t has a binary value, $x_{ti} \in \{0, 1\}$. The HMM comprises 2 non-emitting states, \mathbf{s}_1 and \mathbf{s}_N , and $N - 2$ emitting states, $\mathbf{s}_2, \dots, \mathbf{s}_{N-2}$. The transition probabilities (and model structure) are assumed known, but the state output distribution are to be trained using EM. The required auxiliary function (considering only a single sequence) can be expressed as

$$\mathcal{Q}(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}) = C + \sum_{t=1}^T \sum_{j=2}^{N-1} P(q_t = \mathbf{s}_j | \mathbf{x}_1, \dots, \mathbf{x}_T; \boldsymbol{\theta}) \log \left(P(\mathbf{x}_t | \mathbf{s}_j, \hat{\boldsymbol{\theta}}) \right)$$

where C is independent of the values of the parameters to be estimated.

- Briefly describe the meaning of each of the terms in the auxiliary function above.
- Show that the forward, $\alpha_j(t)$, and backward, $\beta_j(t)$, probabilities described in the lectures can be used to find $P(q_t = \mathbf{s}_j | \mathbf{x}_1, \dots, \mathbf{x}_T; \boldsymbol{\theta})$.
- If the output probability distribution has the form

$$P(\mathbf{x}_t | \mathbf{s}_j; \boldsymbol{\theta}) = \prod_{i=1}^d \lambda_{ji}^{x_{ti}} (1 - \lambda_{ji})^{1-x_{ti}}$$

where state \mathbf{s}_j PMF has parameters $\lambda_{j1}, \dots, \lambda_{jd}$. Show that the estimate for the “new” parameters that maximise the auxiliary function $\hat{\lambda}_{ji}$, are given by

$$\hat{\lambda}_{ji} = \frac{\sum_{t=1}^T P(q_t = \mathbf{s}_j | \mathbf{x}_1, \dots, \mathbf{x}_T; \boldsymbol{\theta}) x_{ti}}{\sum_{t=1}^T P(q_t = \mathbf{s}_j | \mathbf{x}_1, \dots, \mathbf{x}_T; \boldsymbol{\theta})}$$

5. * Expectation maximisation can be used to train a factor analysis model. In factor analysis the d -dimensional observations \mathbf{x} is assumed to be generated by a process of the form

$$\mathbf{x} = \mathbf{C}\mathbf{z} + \mathbf{v}$$

where \mathbf{z} is p -dimensional and Gaussian distributed (zero-mean identity covariance matrix), \mathbf{C} is a $d \times p$ matrix ($d > p$), and \mathbf{v} is d -dimensional. The parameters to estimate are the loading matrix \mathbf{C} and the diagonal covariance matrix, Σ_{diag} , where

$$p(\mathbf{v}) = \mathcal{N}(\mathbf{v}; \mathbf{0}, \Sigma_{\text{diag}})$$

A set of n independent training samples $\mathbf{x}_1, \dots, \mathbf{x}_N$ are available. The mean of the data is known to be zero, $\mathbf{0}$.

- (a) The joint distribution of \mathbf{x} and \mathbf{z} is multivariate Gaussian distributed. Show that this distribution can be expressed as

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{C}\mathbf{C}^\top + \Sigma_{\text{diag}} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{I} \end{bmatrix} \right)$$

- (b) Using the previous result, show that the posterior probability of the hidden variables given the current model parameters, \mathbf{C} and Σ_{diag} is

$$p(\mathbf{z}|\mathbf{x}; \mathbf{C}, \Sigma_{\text{diag}}) = \mathcal{N}(\mathbf{z}; \beta\mathbf{x}, \mathbf{I} - \beta\mathbf{C})$$

where

$$\beta = \mathbf{C}^\top (\Sigma_{\text{diag}} + \mathbf{C}\mathbf{C}^\top)^{-1}$$

You can make use of the following standard equalities

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

$$\mathbf{x}_1|\mathbf{x}_2 \sim \mathcal{N}(\boldsymbol{\mu}_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})$$

- (c) Derive the re-estimation formulae for the new model parameters $\hat{\boldsymbol{\theta}} = \{\hat{\mathbf{C}}, \hat{\Sigma}_{\text{diag}}\}$ based on the current model parameters, $\boldsymbol{\theta} = \{\mathbf{C}, \Sigma_{\text{diag}}\}$, using the following auxiliary function

$$\mathcal{Q}(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}) = \sum_{i=1}^N \int p(\mathbf{z}|\mathbf{x}_i; \boldsymbol{\theta}) \log(p(\mathbf{x}_i|\mathbf{z}; \hat{\boldsymbol{\theta}})) d\mathbf{z}$$

- (d) Describe why factor analysis is a useful tool for covariance modelling and data visualisation.

Deep Learning

6. † A multi-layer perceptron (feed-forward, fully connected, neural network) consists of d inputs, L hidden layers with M hidden units in each hidden layer, and K output nodes. Write down an expression for the total number of weights (including biases) in the network. Describe the factors that influence the number of hidden layers, the activation functions on the output layer, and the number of hidden units.
7. † For the logistic regression function, $\phi(z)$, show that

$$\frac{\partial}{\partial z}\phi(z) = \phi(z)(1 - \phi(z))$$

How does the nature of the activation function affect the computational cost of the error-back propagation algorithm?

8. A *leaky ReLU* activation function is to be used in a multi-layer perceptron. This activation function has the form

$$\phi(z_i) = \begin{cases} z_i; & z_i \geq 0; \\ \alpha z_i & z_i < 0 \end{cases}$$

A large number of samples, generated from a Gaussian distribution with zero mean and a variance of σ^2 , are passed through this activation function. What is the variance of the data at the output of the activation function?

How could this information be used when initialising the network with N nodes per layer?

9. * The Hessian is a useful matrix for use in the optimisation of the weights of multi-layer perceptrons.

(a) Describe how the Hessian may be used for optimising the weights of a multi-layer perceptron. Discuss the limitations for the practical implementation of such schemes.

(b) For the least squares error function

$$E = \sum_{p=1}^n E^{(p)} = \frac{1}{2} \sum_{p=1}^n (y(x_p) - t(x_p))^2$$

show that the elements of the Hessian matrix can be expressed as

$$\frac{\partial^2 E}{\partial w_{ij} \partial w_{lk}} = \sum_{p=1}^n \frac{\partial y(x_p)}{\partial w_{ij}} \frac{\partial y(x_p)}{\partial w_{lk}} + \sum_{p=1}^n (y(x_p) - t(x_p)) \frac{\partial^2 y(x_p)}{\partial w_{ij} \partial w_{lk}}$$

For the case of well trained, sufficiently powerful, network, with an infinitely large training set, show that at the minimum the second term may be ignored. This is called the *outer-product* approximation.

(c) The Hessian after the N^{th} data point is approximated by

$$\mathbf{H}_N = \sum_{p=1}^N \mathbf{g}^{(p)} (\mathbf{g}^{(p)})^\top$$

where

$$\mathbf{g}^{(p)} = \nabla y(x_p)|_{\mathbf{w}_{[\tau]}}$$

By using the equality

$$(\mathbf{A} + \mathbf{BC})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{I} + \mathbf{CA}^{-1} \mathbf{B})^{-1} \mathbf{CA}^{-1}$$

where \mathbf{I} is the identity matrix, show that

$$\mathbf{H}_{N+1}^{-1} = \mathbf{H}_N^{-1} - \frac{\mathbf{H}_N^{-1} \mathbf{g}^{(N+1)} (\mathbf{g}^{(N+1)})^\top \mathbf{H}_N^{-1}}{1 + (\mathbf{g}^{(N+1)})^\top \mathbf{H}_N^{-1} \mathbf{g}^{(N+1)}}$$

Why is this a useful approximation to estimate the inverse Hessian during multi-layer perceptron training.