



Course Name: Machine Learning Dept.: Computer Science and Engineering
Exam Duration: 48 hours

Question No.	1	2	3	4	5	6	7	8
Score	10	10	10	10	10	20	20	10

This exam paper contains 8 questions and the score is 100 in total (Please hand in your answer sheet in the digital form).

Problem I. Least Square (10 points)

- Consider $Y = AX + V$ and $V \sim \mathcal{N}(\mathbf{v}|\mathbf{0}, Q)$, what is the least square solution of X ?
- If there is a constraint of $b^T X = c$, what is the optimal solution of X ?
- If there is an *additional* constraint of $X^T X = d$, in addition to the constraint in b), what is the optimal solution of X ?

Problem II. Linear Gaussian System (10 points)

Consider $Y = AX + V$, where X and V are Gaussian, $X \sim \mathcal{N}(\mathbf{x}|\mathbf{m}_0, \Sigma_0)$, $V \sim \mathcal{N}(\mathbf{v}|\mathbf{0}, \beta^{-1}\mathbf{I})$.

What are the conditional distribution, $p(Y|X)$, the joint distribution $p(Y, X)$, the marginal distribution, $p(Y)$, the posterior distribution, $p(X|Y = \mathbf{y}, \beta, \mathbf{m}_0, \Sigma_0)$, the posterior predictive distribution, $p(\hat{Y}|Y = \mathbf{y}, \beta, \mathbf{m}_0, \Sigma_0)$, and the prior predictive distribution, $p(Y|\beta, \mathbf{m}_0, \Sigma_0)$, respectively?

Problem III. Linear Regression (10 points)

Consider $y = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + v$, where v is Gaussian, i.e., $v \sim \mathcal{N}(v|0, \beta^{-1})$, and \mathbf{w} has a Gaussian *priori*, i.e., $\mathbf{w} \sim \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \alpha^{-1}\mathbf{I})$. Assume that $\boldsymbol{\phi}(\mathbf{x})$ is known, please derive the posterior distribution, $p(\mathbf{w}|D, \beta, \mathbf{m}_0, \alpha)$, the posterior predictive distribution, $p(\hat{y}|\hat{x}, D, \beta, \mathbf{m}_0, \alpha)$, and the prior predictive distribution, $p(D|\beta, \mathbf{m}_0, \alpha)$, respectively, where $D = \{\phi_n, y_n\}$, $n = 1, \dots, N$, is the training data set and $\phi_n = \phi(\mathbf{x}_n)$.

Problem IV. Logistics Regression (10 points)

Consider a two-class classification problem with the logistic sigmoid function, $y = \sigma(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}))$, for a given data set $D = \{\phi_n, t_n\}$, where $t_n \in \{0, 1\}$, $\phi_n = \phi(\mathbf{x}_n)$, $n = 1, \dots, N$, and the likelihood function is given by

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}$$

where \mathbf{w} has a Gaussian *priori*, i.e., $\mathbf{w} \sim \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \alpha^{-1}\mathbf{I})$. Please derive the posterior distribution, $p(\mathbf{w}|D, \mathbf{m}_0, \alpha)$, the posterior predictive distribution, $p(t|x, D, \mathbf{m}_0, \alpha)$, and the prior predictive distribution, and $p(D|\mathbf{m}_0, \alpha)$, respectively. (*Hint*: using Delta approximation and Laplace approximation properly).

Problem V. Neural Network (10 points)

Consider a two-layer neural network described by following equations:

$$a_1 = \mathbf{w}^{(1)} \mathbf{x}, \quad a_2 = \mathbf{w}^{(2)} \mathbf{z}, \quad z = h(a_1), \quad y = \sigma(a_2)$$

where \mathbf{x} and y are the input and output, respectively, of the neural network, $h(\bullet)$ is a nonlinear function, and $\sigma(\bullet)$ is the sigmoid function.

(1) Please derive the following gradients: $\frac{\partial y}{\partial \mathbf{w}^{(1)}}$, $\frac{\partial y}{\partial \mathbf{w}^{(2)}}$, $\frac{\partial y}{\partial a_1}$, $\frac{\partial y}{\partial a_2}$, and $\frac{\partial y}{\partial \mathbf{x}}$.

- (2) Please derive the updating rules for $\mathbf{w}^{(1)}$ and $\mathbf{w}^{(2)}$ given the classification errors between y and t , where t is the ground truth of the output y .

Problem VI. Bayesian Neural Network (20 points)

- a) Consider a neural network for regression, $t = y(\mathbf{w}, \mathbf{x}) + v$, where v is Gaussian, i.e., $v \sim \mathcal{N}(v|0, \beta^{-1})$, and \mathbf{w} has a Gaussian *priori*, i.e., $\mathbf{w} \sim \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \alpha^{-1}\mathbf{I})$. Assume that $y(\mathbf{w}, \mathbf{x})$ is the neural network output please derive the posterior distribution, $p(\mathbf{w}|D, \beta, \mathbf{m}_0, \alpha)$, the posterior predictive distribution, $p(t|x, D, \beta, \mathbf{m}_0, \alpha)$, and the prior predictive distribution, $p(D|\beta, \mathbf{m}_0, \alpha)$, where $D = \{x_n, t_n\}$, $n = 1, \dots, N$, is the training data set.
- b) Consider a neural network for two-class classification, $y = \sigma(a(\mathbf{w}, \mathbf{x}))$ and a data set $D = \{x_n, t_n\}$, where $t_n \in \{0, 1\}$, \mathbf{w} has a Gaussian *priori*, i.e., $\mathbf{w} \sim \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$, and $a(\mathbf{w}, \mathbf{x})$ is the neural network model. Please derive the posterior distribution, $p(\mathbf{w}|D, \alpha)$, posterior predictive distribution, $p(t|x, D, \alpha)$, and the prior predictive distribution, $p(D|\alpha)$, respectively.

Problem VII. Critical Analyses (20 Points)

- a) Please explain, in terms of cost functions, constraints and predictions, **i)** what are the differences between SVM classification and logistic regression; **ii)** what are the differences between v-SVM regression and least square regression.
- b) Please explain why neural network (NN) based machine learning algorithms use *logistic* activation functions ?
- c) Please explain **i)** what are the differences between the *logistic* activation function and other activation functions (e.g., *relu*, *tanh*); and **ii)** when these activation functions should be used.

- d) Please explain why Jacobian and Hessian matrices are useful for machine learning algorithms.
- e) Please explain why exponential family distributions are so common in engineering practice.
Please give some examples which are **NOT** exponential family distributions.
- f) Please explain why KL divergence is useful for machine learning? Please provide two examples of using KL divergence in machine learning.
- g) Please explain why data augmentation techniques are a kind of regularization skills for NNs.
- h) Please explain why Gaussian distributions are preferred over other distributions for many machine learning models?
- i) Please explain why the MAP model is usually more preferred than the ML model?

Problem VIII. Discussion (10 Points)

What are the generative and discriminative approaches to machine learning, respectively? Can you explain the advantages and disadvantages of these two approaches and provide a detailed example to illustrate your points?