

浙江大学 2012 - 2013 学年 春季 学期

《Artificial Intelligence》课程期末考试试卷

课程号： 21190770 ， 开课学院： 计算机科学与技术学院

考试试卷： A 卷、B 卷（请在选定项上打√）

考试形式： 闭、开卷（请在选定项上打√）， 允许带_____入场

考试日期： 2013 年 4 月 29 日, 考试时间： 120 分钟

诚信考试，沉着应考，杜绝违纪。

考生姓名： _____ 学号： _____ 所属院系： _____

题序	一	二	三	四	五	六	七	八	总 分
得分									
评卷人									

1. Fill in the blanks (30 points)

- 1) AI attempts not just to understand but also to build intelligent entities. Historically, there are many definitions of AI that can be organized into four categories: Thinking Humanly, _____, _____ and Acting Rationally.
- 2) In 1950, Alan Turing proposed an approach, named _____, to test whether a computer would really be intelligent if a human interrogator, after posing some written questions, cannot tell whether the written responses come from a person or from a computer.
- 3) The assertion that machines could act as if they were intelligent is called the _____ hypothesis by philosophers, and assertion that machines that do so are actually thinking (not just simulating thinking) is called the _____ hypothesis.
- 4) Machine learning is usually divided into two main types. Applications in which the training data comprises examples of the input vectors along with their corresponding target vectors are known as _____ learning problem. Applications in which the training data consists of a set of input vectors without any corresponding target values are known as _____ learning problem.

- 5) Assume the sum-of-squares error function is defined by $E(w)$, then the root-mean-square error function can be defined by _____.
- 6) Data points that are drawn independently from the same distribution are said to be _____, which is often abbreviated to _____.
- 7) Given multivariate Gaussian distribution $N(x|\mu, \Sigma)$, then Mahalanobis distance Δ is defined by _____.
- 8) For a given probability distribution $p(x|w)$, if we choose a prior $p(w)$, then the posterior distribution $p(w|x)$ will have the same functional form as the prior. This property is called _____.
- 9) In generalized linear models, if activation function is logistic sigmoid function $f(\cdot)$, then the corresponding link function is _____.
- 10) The generalized linear model based on a probit activation function is known as _____. The inverse probit function can be constructed by the _____ function of a zero mean, unit variance Gaussian $N(x|0,1)$.

2. Multiple Choice (20 points)

- 1) Which of the following statements about ML problems is false? _____
 - A. The regression is one of unsupervised learning problems.
 - B. The classification is one of supervised learning problems.
 - C. The clustering is one of unsupervised learning problems.
 - D. The density estimation is one of unsupervised learning problems.
- 2) In regression problems, we often need to minimize an error function that measures the misfit between the function output and the training set data points. For a given model, assume we have evaluated the sum-of-squares error $E_1(w)$ and $E_2(w)$ for two test data sets D_1 and D_2 with different size. We also computed the corresponding root-mean-square error $E_{rms1}(w)$ and $E_{rms2}(w)$. Which of the following discriminant condition can lead to the conclusion that this model has better fitting performance on test set D_1 ? _____
 - A. $E_1(w) < E_{rms1}(w)$ B. $E_1(w) < E_2(w)$
 - C. $E_{rms1}(w) < E_{rms2}(w)$ D. $E_{rms1}(w) < E_1(w)$

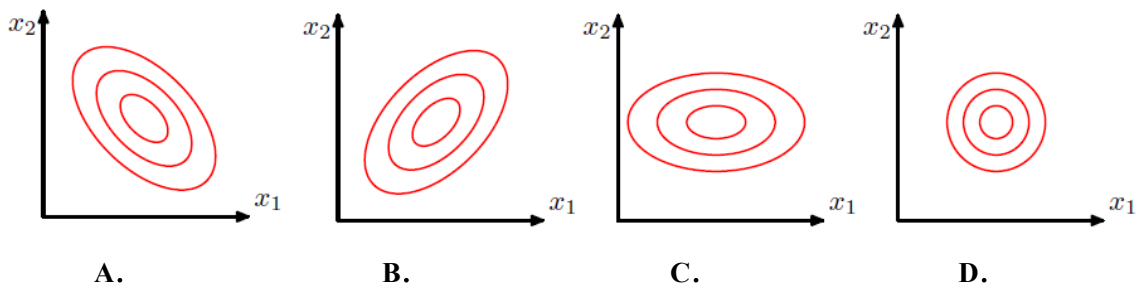
- 3) Consider a polynomial curve fitting problem. If the fitted curve oscillates wildly through each point and achieve bad generalization by making accurate predictions for new data, we say this behavior is over-fitting. Which of the following methods cannot be used to control over-fitting? _____
- A. Use fewer training data
 - B. Add validation set, use Cross-validation
 - C. Add a regularization term to an error function
 - D. Use Bayesian approach with suitable prior
- 4) Assume D is the observed data set and w is model parameter. Which of the following statements about likelihood function $p(D|w)$ is false? _____
- A. It expresses how probable the observed data set is for different settings of the parameter vector w .
 - B. The likelihood is not a probability distribution over w .
 - C. Its integral with respect to w must be equal to one.
 - D. Maximizing the likelihood function is equivalent to minimizing the error.
- 5) Which of the following statements about the Fisher's criterion is correct? _____
- A. It maximizes the separation between the projected class means as well as the total within-class variance.
 - B. It minimizes the separation between the projected class means as well as the total within-class variance.
 - C. It maximizes the separation between the projected class means as well as the inverse of the total within-class variance.
 - D. It minimizes the separation between the projected class means as well as the inverse of the total within-class variance.
- 6) Given two Gaussian distribution $N(x|0,1)$ and $N(x|1,1)$, which of the following formula is correct? _____
- A. $N(0.5|0,1) > N(0.5|1,1)$
 - B. $N(1|0,1) = N(0|1,1)$
 - C. $N(0.5|0,1) < N(0.5|1,1)$
 - D. $N(0|0,1) = N(0|1,1)$
- 7) Which of the following statements about the kernel function is false? _____
- A. The kernel function is a symmetric function.
 - B. The simplest example of a kernel function is $k(x, x') = x^T x'$.
 - C. The feature vector that corresponds to the Gaussian kernel has infinite dimensionality.

D. We cannot construct new kernels by using simpler kernels.

- 8) Assume the precision matrix is given by $\mathbf{R} = \begin{pmatrix} \mathbf{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A} & -\mathbf{A}^T \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{pmatrix}$, then the corresponding covariance matrix $\mathbf{\Sigma} = \mathbf{R}^{-1} = \begin{pmatrix} \mathbf{\Sigma}_{11} & \mathbf{\Lambda}^{-1} \mathbf{A}^T \\ \mathbf{A} \mathbf{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^T \end{pmatrix}$, find the expression of $\mathbf{\Sigma}_{11} =$ _____.

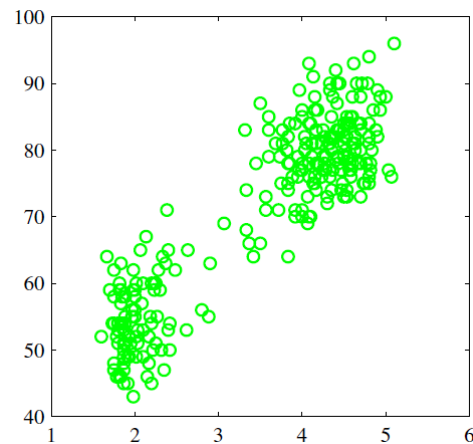
A. \mathbf{L}^{-1} B. $\mathbf{\Lambda}^{-1}$ C. \mathbf{A}^T D. \mathbf{A}^{-1}

- 9) For a 2D multivariate Gaussian distribution $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{\Sigma})$, if $\mathbf{\Sigma} = \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix}$, $\sigma_1^2 > \sigma_2^2$, which of the following figures is the contour of constant probability density of this Gaussian distribution? _____



- 10) For the 'Old Faithful' data shown on the following figure, which probabilistic model can represent it more accurately? _____

A. Gaussian distribution
B. Mixture of Gaussians
C. Dirichlet distribution
D. Wishart distribution



3. Calculus, Analysis and Proof (50 points)

- 1) Consider the multivariate Gaussian distribution given by Appendix 1.(b). By writing the precision matrix (inverse covariance matrix) $\Lambda = \Sigma^{-1}$ as the sum of symmetric matrix $S = (\Lambda + \Lambda^T) / 2$ and anti-symmetric matrix $A = (\Lambda - \Lambda^T) / 2$, show that:

(a) the inverse matrix S^{-1} is symmetric. (3 points)

(b) the anti-symmetric term A does not appear in the exponent of the Gaussian for

$$\Lambda = S + A, \text{ such that } (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^T S (\mathbf{x} - \boldsymbol{\mu}). (4 \text{ points})$$

- 2) Assume an eigenvalue decomposition of the covariance matrix Σ is given by $\Sigma = U \Lambda U^{-1}$, where $U = (\mathbf{u}_1, \dots, \mathbf{u}_D)$, $U^T U = I$, $\Lambda = \text{diag}(\lambda_i), i = 1, \dots, D$, show that:

(a) $\Sigma^{-1} = U \Lambda^{-1} U^T$ (4 points)

$$(b) |\Sigma| = \prod_{i=1}^D \lambda_i \quad (4 \text{ points})$$

- 3) Given a data set $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ in which the observations $\{\mathbf{x}_n\}$ are assumed to be drawn independently from a multivariate Gaussian distribution given by Appendix 1.(b).

(a) Find the likelihood function $p(\mathbf{X} | \boldsymbol{\mu}, \Sigma)$ (2 points)

(b) Find the log likelihood function $\ln p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$ (4 points)

(c) Find the maximum likelihood solution of the mean $\boldsymbol{\mu}_{ML}$ (4 points)

(d) Find the maximum likelihood solution of the covariance $\boldsymbol{\Sigma}_{ML}$ (5 points)

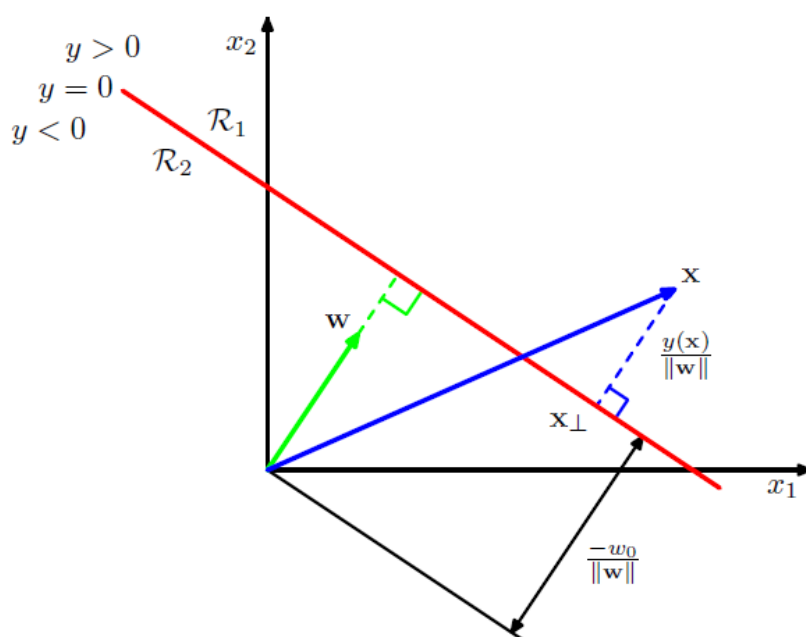
4) Assume the error function with a regularization term in regression has given by

$E(\mathbf{w}) = \frac{1}{2} \|\mathbf{t} - \boldsymbol{\Phi} \mathbf{w}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$, where \mathbf{t} is the target vector, λ is the regularization coefficient and $\boldsymbol{\Phi}$ is the design matrix. Find the solution of \mathbf{w} by minimizing $E(\mathbf{w})$. (10 points)

5) Consider a linear discriminant function $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ for two classes in 2-dimensional input space, the geometry is shown in the following figure. An input vector \mathbf{x} is assigned to class C_1 if $y(\mathbf{x}) \geq 0$ and to class C_2 otherwise. The corresponding decision boundary is therefore defined by the relation $y(\mathbf{x}) = 0$, which corresponds to the line Ω in figure. Prove that:

(a) The vector \mathbf{w} is orthogonal to every vector lying within the decision surface Ω . (5 points)

(b) The perpendicular distance r of arbitrary input vector \mathbf{x} from the decision surface is $r = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}$. (5 points)



Appendix:

1. Probability distributions:

(a) **Single variable Gaussian:** $N(x | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$

(b) **D-dimensional multivariate Gaussian:**

$$N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\},$$

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \text{const}$$

(c) **Beta:** $Beta(\mu | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$

(d) **Dirichlet:**

$$Dir(\boldsymbol{\mu} | \boldsymbol{\alpha}) = C(\boldsymbol{\alpha}) \prod_{k=1}^K \mu_k^{\alpha_k-1}, \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_K \end{pmatrix}, \boldsymbol{\alpha} = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_K \end{pmatrix}, 0 \leq \mu_k \leq 1, \sum_{k=1}^K \mu_k = 1$$

(e) **Gamma:** $Gam(\tau | a, b) = \frac{1}{\Gamma(a)} b^a \tau^{a-1} e^{-b\tau}, a > 0, b > 0, \tau > 0$

2. Matrix calculus

(a) $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T, \quad \mathbf{A}^{-1} \mathbf{A} = \mathbf{A} \mathbf{A}^{-1} = \mathbf{I}, \quad (\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$

(b) $Tr(\mathbf{AB}) = Tr(\mathbf{BA}), \quad Tr(\mathbf{ABC}) = Tr(\mathbf{CAB}) = Tr(\mathbf{BCA}), \quad |\mathbf{AB}| = |\mathbf{A}||\mathbf{B}|$

(c) $\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{MBD}^{-1} \\ -\mathbf{D}^{-1}\mathbf{CM} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{CMBD}^{-1} \end{pmatrix}, \mathbf{M} = (\mathbf{A} - \mathbf{BD}^{-1}\mathbf{C})^{-1}$

(d) $\frac{\partial(\mathbf{x}^T \mathbf{a})}{\partial \mathbf{x}} = \frac{\partial(\mathbf{a}^T \mathbf{x})}{\partial \mathbf{x}} = \mathbf{a} \quad \frac{\partial \mathbf{a}^T \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -\mathbf{X}^{-T} \mathbf{a} \mathbf{b}^T \mathbf{X}^{-T} \quad \frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$

(e) Assume $\boldsymbol{\Lambda}$ is symmetric matrix, then $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}) = 2\boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu})$

(f) $\frac{\partial}{\partial \mathbf{w}} (\mathbf{t} - \boldsymbol{\Phi} \mathbf{w})^T (\mathbf{t} - \boldsymbol{\Phi} \mathbf{w}) = -2\boldsymbol{\Phi}^T (\mathbf{t} - \boldsymbol{\Phi} \mathbf{w}), \quad \frac{\partial}{\partial \mathbf{w}} \|\mathbf{w}\|^2 = \frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T \mathbf{w} = 2\mathbf{w}$

(g) $\frac{\partial}{\partial \mathbf{W}} Tr[(\mathbf{T} - \boldsymbol{\Phi} \mathbf{W})(\mathbf{T} - \boldsymbol{\Phi} \mathbf{W})^T] = -2\boldsymbol{\Phi}^T (\mathbf{T} - \boldsymbol{\Phi} \mathbf{W})$

(h) $\frac{\partial}{\partial \mathbf{A}} \ln |\mathbf{A}| = (\mathbf{A}^{-1})^T, \quad \text{if } \mathbf{A} = \text{diag}(\lambda_i), i = 1, \dots, D, \text{ then } |\mathbf{A}| = \prod_{i=1}^D \lambda_i$

Final Examination Answer Sheet

Section	1	2	3	Total
Score				
Reviewer				

- 1). _____, _____
- 2). _____
- 3). _____, _____
- 4). _____, _____
- 5). _____
- 6). _____, _____
- 7). _____
- 8). _____
- 9). _____
- 10). _____, _____

[illegible]

3. Calculus, Analysis and Proof (50 points)

1) (7 points)

(a) (3 points)

(b) (4 points)

2) (8 points)

(a) (4 points)

(b) (4 points)

3) (15 points)

(a) (2 points)

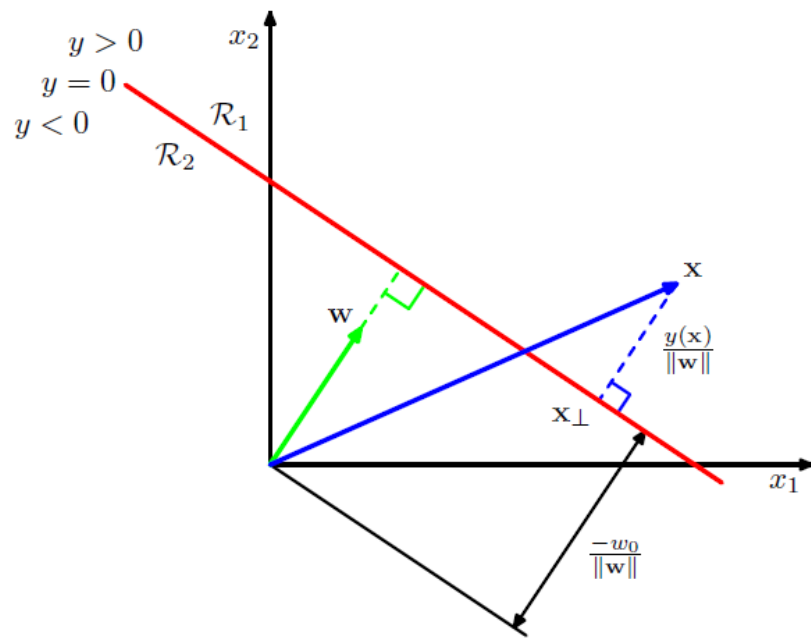
(b) (4 points)

(c) (4 points)

(d) (5 points)

4) (10 points)

5) (10 points)



(a) (5 points)

(b) (5 points)