**Instructions:**

- This midterm is designed to be finished **within 50 minutes**. The additional 20 minutes are designed for scanning and uploading your submission and any potential technical difficulty.

- Follow directions and answer questions with requested supporting work.

- Clearly indicate your answer in the allotted space or by putting a box around it.

- The midterm exam will be posted on May 4th 3:00pm PST. You will have 24 hours to finish and upload your solution to CCLE by May 5th 2:59pm PST. You can use the textbook, any course material posted on CCLE, and your hand-written notes; you are not allowed to use calculators nor the Internet, and you cannot work with anyone else (classmate, family member, private tutor, etc.). You can scan or take high-resolution photos of your hand-written solutions, but the uploaded submission must be a single PDF file.

**Problem 1**

Given a data set
$$\mathcal{D} = \{(\mathbf{x}^{(i)}, t^{(i)})\}_{i=1}^{N}, \quad \mathbf{x}^{(i)} \in \mathbb{R}^{D}, t^{(i)} \in \mathbb{R}.$$

Fixing an integer $M \in \mathbb{N}$ and a basis function $\boldsymbol{\phi} : \mathbb{R}^{D} \to \mathbb{R}^{M}$. The regression problem is about finding the best parameter $\mathbf{w} \in \mathbb{R}^{M}$ so

$$t^{(i)} = \mathbf{w}^{T} \boldsymbol{\phi}(\mathbf{x}^{(i)}) + \boldsymbol{\epsilon}^{(i)}$$

where $\epsilon^{(i)} \sim \mathcal{N}(0, \beta^{-1})$ are independent identical (unbiased) Gaussian noise.

(a) Write down a formula for the likelihood function $p(\mathcal{D}|\mathbf{w}, \beta)$. (10 points)

(b) Show that the maximum likelihood solution

$$\mathbf{w}_{\beta}^{*} = \arg\max_{\mathbf{w}} p(\mathcal{D}|\mathbf{w}, \beta)$$

for any value of $\beta > 0$ is the same as the least square solution

$$\overline{\mathbf{w}} = \arg\min_{\mathbf{w}} \frac{1}{2} \left| t_{i} - y(x^{(i)}, \mathbf{w}) \right|^{2}. \qquad \text{(10 points)}$$

(c) Fixing the model complexity $M \in \mathbb{N}$, give three examples of the basis function $\boldsymbol{\phi}(\mathbf{x})$. (10 points)

**Problem 2**

Suppose a data set $\mathcal{D} = \{(\mathbf{x}^{(i)}, t^{(i)})\}_{i=1}^{N}$ is given. $\mathbf{x}^{(i)} \in \mathbb{R}^D, t^{(i)} \in \mathbb{R}$ for $i = 1, \cdots, N$.

(a) Show that the optimal solution $\mathbf{w}^* = \arg\min J(\mathbf{w})$ for a regularized sum-of-squares error function

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^{N} \left(t^{(i)} - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)})\right)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2,$$

where $\lambda > 0$, is a linear combination of the vectors $\{\boldsymbol{\phi}(\mathbf{x}^{(i)})\}_{i=1}^{N}$. In other words, show that

$$\mathbf{w}^* = \sum_{i=1}^{N} a^{(i)} \boldsymbol{\phi}(\mathbf{x}^{(i)})$$

for some scalars $a^{(i)} \in \mathbb{R}, i = 1, \cdots, N$. (10 points)

(b) We define the Gram matrix

$$\mathbf{K} = [\, K_{ij} \,] = [\, \boldsymbol{\phi}(\mathbf{x}^{(i)})^T \boldsymbol{\phi}(\mathbf{x}^{(j)}) \,] \in \mathbb{R}^{N \times N}.$$

Show that $\mathbf{K}$ is symmetric semi-positive definite. (10 points)

(c) Show that the coefficients from part (a) satisfy

$$(\mathbf{K} + \lambda \mathbf{I}) \begin{bmatrix} a^{(1)} \\ a^{(2)} \\ \vdots \\ a^{(N)} \end{bmatrix} = \begin{bmatrix} t^{(1)} \\ t^{(2)} \\ \vdots \\ t^{(N)} \end{bmatrix}. \qquad (10 \text{ points})$$

## Problem 3

Consider the two-class classification problem. Denote the data set $\mathcal{D} = \{(\mathbf{x}^{(i)}, t^{(i)})\}_{i=1}^{N}$ where

$$t^{(i)} = \begin{cases} 1, & i \in C_{[1]} \\ 0, & i \in C_{[2]} \end{cases}$$

is the target variable encoding the class membership.

(a) Suppose $p(\mathbf{x}|C_{[1]}) \sim \mathcal{N}(\boldsymbol{\mu}_{[1]}, \boldsymbol{\Sigma})$ and $p(\mathbf{x}|C_{[2]}) \sim \mathcal{N}(\boldsymbol{\mu}_{[2]}, \boldsymbol{\Sigma})$, that is, data from two classes scatter around different class-specific mean but share the same covariance matrix. Denote $p(C_{[1]}) = \pi$, hence $p(C_{[2]}) = 1 - \pi$. The likelihood function is given by

$$p(\mathcal{D}|\pi, \boldsymbol{\mu}_{[1]}, \boldsymbol{\mu}_{[2]}, \boldsymbol{\Sigma}) = \left( \prod_{i \in C_{[1]}} \pi \mathcal{N}(\mathbf{x}^{(i)}|\boldsymbol{\mu}_{[1]}, \boldsymbol{\Sigma}) \right) \cdot \left( \prod_{i \in C_{[2]}} (1 - \pi) \mathcal{N}(\mathbf{x}^{(i)}|\boldsymbol{\mu}_{[2]}, \boldsymbol{\Sigma}) \right)$$

Show that the maximum likelihood estimate of the class probability $\pi$ is given by the fraction of data points in $C_{[1]}$, i.e.

$$\arg \max_{\pi} p(\mathcal{D}|\pi, \boldsymbol{\mu}_{[1]}, \boldsymbol{\mu}_{[2]}, \boldsymbol{\Sigma}) = \frac{\#\{i : i \in C_{[1]}\}}{N}. \qquad \text{(10 points)}$$

(b) The *logistic sigmoid* function is defined by

$$\sigma(b) = \frac{1}{1 + e^{-b}}.$$

Show that (i) $\sigma(-b) = 1 - \sigma(b)$, (ii) $\sigma$ is a monotonically increasing function, and (iii) $\sigma$ maps all of $\mathbb{R}$ onto the interval $(0, 1)$. (15 points)

(c) In an approach different from part (a), we suppose that $p(C_{[1]}|\mathbf{x}) = \sigma(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}))$ where $\mathbf{w} \in \mathbb{R}^M$ is a coefficient vector to be trained. According to part (b), $p(C_{[2]}|\mathbf{x}) = \sigma(-\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}))$. The likelihood function is then given by this different formula,

$$p(\mathcal{D}|\mathbf{w}) = \left( \prod_{i \in C_{[1]}} \sigma(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)})) \right) \cdot \left( \prod_{i \in C_{[2]}} \sigma(-\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}^{(i)})) \right).$$

Define the error function $E(\mathbf{w}) = -\log p(\mathcal{D}|\mathbf{w})$ to be the negative logarithm of the likelihood. Compute $\nabla E(\mathbf{w})$ and explain why the maximum likelihood estimate $\nabla E(\mathbf{w}^*) = 0$ doesn't have an analytical solution. (15 points)

(d) **(Bonus)** Provide a strategy to compute the optimal solution $\mathbf{w}^*$ numerically. (up to 10 points)