

Name: Haoye He
NYU ID: h2537

Midterm Group A

Problem 1 (20 points)

Recall D-dimensional regression problem when your model performs label

prediction for the i-th example x_i in the training data set using linear function:

$$f(x_i; \theta_0, \theta_1, \theta_2, \dots, \theta_D) = \sum_{d=1}^D \theta_d x_i(d) + \theta_0.$$

In this case x_i is D-dimensional. Let y_i denote the true label of the i-th example and let N be the total number of training examples. Parameters of the model $(\theta_0, \theta_1, \theta_2, \dots, \theta_D)$ are obtained by minimizing the empirical risk provided below:

$$R(\theta) = \frac{1}{2N} \|y - X\theta\|_2^2 + \theta^T H \theta + \theta^T \theta + a^T \theta,$$

where a is a vector and H is a matrix that satisfies the condition: $H = H^T$. Both a and H are given. Write what is y, X, and θ in the formula above. Compute the optimal setting of parameters by setting the gradient of the risk to 0. Explain all steps in your derivations.

Problem 2 (15 points)

A kernel is an efficient way to write out an inner product between two feature vectors computed from a pair of input vectors as follows:

$$K(x, y) = \phi(x)^\top \phi(y).$$

Assume that both inputs are 2-dimensional and write out the explicit mapping ϕ that mimics the kernel value for a 3rd-order polynomial kernel as follows:

$$K(x, y) = (x^\top y + 1)^3.$$

Problem 3 (15 points)

The exponential distribution has density given as

$$p_\lambda(x) = \lambda e^{-\lambda x}$$

Find the maximum likelihood estimator for λ . Calculate an estimate using this estimator when $x_1 = 1, x_2 = 2, x_3 = 4, x_4 = 2$.

Problem 4 (15 points)

Consider 2d family of classifiers given by an origin-centered circles $f(x) = \text{sign}(ax^\top x + b)$. What is the VC dimension of this family? Prove it.

Problem 5 (15 points)

Using the principle of Lagrange multipliers, find the maximum and minimum values of $f(x, y) = x^2 - y^2$ subject to the constraint, $x^2 + y^2 = 1$.

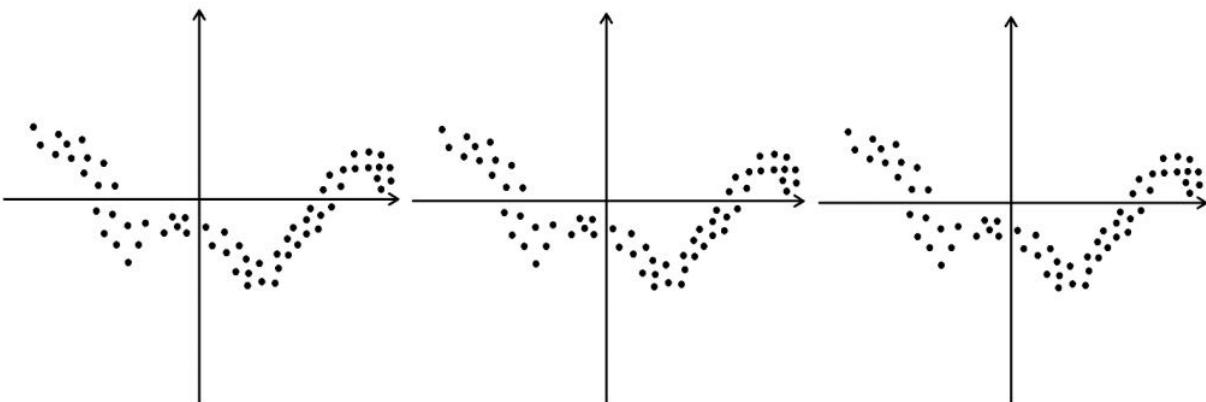
Problem 6 (10 points)

What is the maximum likelihood estimator of the mean μ and covariance matrix Σ of the following 2-dimensional data set \mathcal{X} ? Justify your answer.

$$\mathcal{X} = \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \begin{bmatrix} 5 \\ 6 \end{bmatrix}, \begin{bmatrix} 7 \\ 8 \end{bmatrix}, \begin{bmatrix} 9 \\ 10 \end{bmatrix} \right\}$$

Problem 7 (10 points)

Explain the difference between overfitting and underfitting. For a picture given below, show an example of underfitting on the left plot, proper fit in the middle plot, and overfitting on the right plot.



Q1.

Step 1, explain the parameters:-

$\vec{y} \in \mathbb{R}^{N \times 1}$: matrix contains all the true labels

$\vec{X} \in \mathbb{R}^{N \times D}$: matrix contains all the training data.

$\vec{\theta} \in \mathbb{R}^{D \times 1}$: parameters used to fit the training data and give the result.

Step 2.

$$R(\theta) = \frac{1}{N} \|\vec{y} - \vec{X}\vec{\theta}\|_2^2 + \vec{\theta}^\top H \vec{\theta} + \vec{\theta}^\top \vec{\theta} + \vec{\alpha}^\top \vec{\theta}$$

// Original formula

$$\begin{aligned}
 &= \frac{1}{2N} (\vec{y} - \vec{\theta})^T (\vec{y} - \vec{\theta}) + \vec{\theta}^T H \vec{\theta} + \vec{\theta}^T \vec{\alpha} + \vec{\alpha}^T \vec{\theta} \quad // \text{Append 1st term} \\
 &= \frac{1}{2N} (\vec{y}^T \vec{y} - 2\vec{y}^T \vec{\theta} + \vec{\theta}^T \vec{\theta}) + \vec{\theta}^T H \vec{\theta} + \vec{\theta}^T \vec{\alpha} + \vec{\alpha}^T \vec{\theta} \quad // \text{Append 1st term} \\
 \frac{\partial R(\theta)}{\partial \theta} &= \frac{1}{2N} (2\vec{x}^T \vec{y} + 2\vec{x}^T H \vec{\theta}) + (H + H^T) \vec{\theta} + 2\vec{\theta} + \vec{\alpha}^T = 0 \quad // \text{Start Derivation} \\
 \rightarrow \vec{x}^T \vec{y} + \vec{x}^T H \vec{\theta} + 2N(H^T \vec{\theta} + 2\theta + \alpha) &= 0 \quad // H^T = H \text{ according to Question} \\
 \rightarrow \vec{x}^T \vec{y} + \vec{x}^T H \vec{\theta} + 2NH^T \vec{\theta} + 2N\vec{\theta} + Na &= 0 \quad // \text{simplify} \\
 \vec{x}^T H \vec{\theta} + 2NH^T \vec{\theta} + 2N\vec{\theta} &= \vec{x}^T \vec{y} - Na \quad // \text{simplify} \\
 (\vec{x}^T H + 2NH^T + 2N\vec{\theta}) \vec{\theta} &= \vec{x}^T \vec{y} - Na \\
 \vec{\theta} &= (\vec{x}^T H + 2NH^T + 2N\vec{\theta})^{-1} (\vec{x}^T \vec{y} - Na) \quad // \text{solution}
 \end{aligned}$$

Q2.

$$\begin{aligned}
 k(x, y) &= (x^T y + 1)^3 = (x^T y + 1)(x^T y)^2 + 2x^T y + 1 \\
 &= (x^T y)^3 + 1 + 3(x^T y)^2 + 3x^T y \\
 &= (x_1 y_1 + x_2 y_2)^3 + 1 + 3(x_1 y_1 + x_2 y_2)^2 + 3(x_1 y_1 + x_2 y_2) \\
 &= x_1^3 y_1^3 + x_2^3 y_2^3 + 3x_1^2 y_1^2 x_2 y_2 + 3x_1 y_1 x_2^2 y_2^2 + 1 + \\
 &\quad 3(x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 y_1 x_2 y_2) + 3(x_1 y_1 + x_2 y_2) \\
 &= x_1^3 y_1^3 + x_2^3 y_2^3 + 3x_1 y_1 x_2^2 y_2^2 + 3x_1^2 y_1^2 x_2 y_2 + 3x_1^2 y_1^2 + 3x_2^2 y_2^2 + \\
 &\quad 6x_1 y_1 x_2 y_2 + 3x_1 y_1 + 3x_2 y_2 + 1
 \end{aligned}$$

Hence:

$$\vec{\theta}(x) = [x_1^3 \quad x_2^3 \quad \sqrt{3}x_1 x_2 \quad \sqrt{3}x_1 x_2^2 \quad \sqrt{3}x_1^2 \quad \sqrt{3}x_2^2 \quad \sqrt{6}x_1 x_2 \quad \sqrt{3}x_1 \quad \sqrt{3}x_2 \quad 1]$$

s.t.

$$\vec{\theta}(x) \vec{\theta}(y) = (x^T y + 1)^3$$

Q3.

$$L(\lambda) = \prod_{i=1}^N \lambda \cdot e^{-\lambda x_i}$$

Instead of directly maximizing the likelihood. We instead maximize the log-likelihood.
the log-likelihood.

$$\log L(\lambda) = \log \prod_{i=1}^N \lambda \cdot e^{-\lambda x_i} = \sum_{i=1}^N (\log \lambda + (-\lambda x_i)) = N \log \lambda + \sum_{i=1}^N (-\lambda x_i)$$

$$\frac{\partial \log L(\lambda)}{\partial \lambda} = \frac{N}{\lambda} + \sum_{i=1}^N (-x_i)$$

Set the derivative equal to 0:

$$\lambda = \frac{N}{\sum_{i=1}^n x_i}$$

Using the given data.

$$\lambda = \frac{N}{\sum_{i=1}^n x_i} = \frac{4}{1+2+4+2} = \frac{4}{9}$$

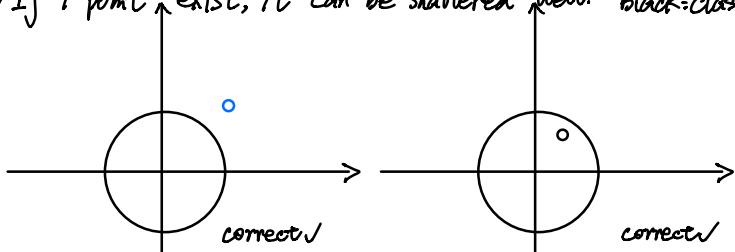
Then, verify it's maximum by taking second derivative

$$\frac{\partial(\log L(\lambda))}{\partial^2(\lambda)} = -\frac{N}{\lambda^2} < 0$$

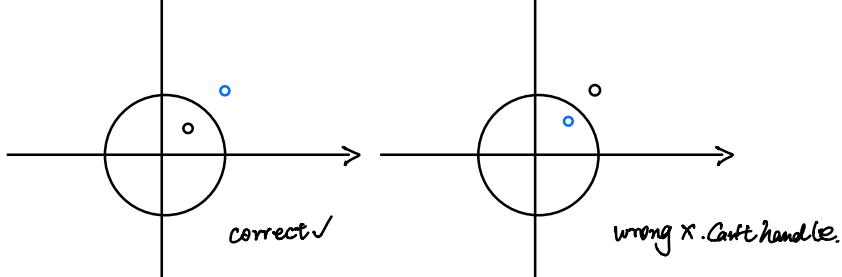
Hence, proved.

Q4. VC Dimension = 1

① If 1 point exist, it can be shattered well. Black: class 1 Blue: class 2



② If 2 points exist, it can't be shattered



Hence, VC Dimension = 1

Q5

$$f(x, y) = x^2 - y^2 \quad \textcircled{1}$$

$$g(x, y) = x^2 + y^2 - 1 = 0 \quad \textcircled{2}$$

$$F(x) = x^2 - y^2 + \lambda(x^2 + y^2 - 1)$$

$$\frac{\partial F(x)}{\partial x} = 2x + 2\lambda x = 0 \Rightarrow 2x(1 + \lambda) = 0 \quad \textcircled{3}$$

$$\frac{\partial F(x)}{\partial y} = -2y + 2\lambda y = 0 \Rightarrow -2y(1 - \lambda) = 0 \quad \textcircled{4} \Rightarrow 4xy = 0 \Rightarrow x=0 \text{ or } y=0$$

$$\frac{\partial F(x)}{\partial \lambda} = x^2 + y^2 - 1 = 0 \quad \textcircled{2}$$

Assume $x=0 \Rightarrow$ According to \textcircled{2}, $y=\pm 1$
 $y=0 \Rightarrow$ According to \textcircled{2}, $x=\pm 1$

$$\text{Hence, } \text{Max} = x^2 + y^2 = 1^2 + 0^2 = 1 \\ \text{Min} = x^2 + y^2 = 0^2 + 1^2 = 1$$

Q6. To get μ and Σ , we try to maximum likelihood:

$$P(\vec{X} | \theta) = \prod_{i=1}^N P(\vec{x}_i | \mu_i, \Sigma_i) = \prod_{i=1}^N P(\vec{x}_i | \vec{\mu}, \Sigma) \text{ i.i.d.}$$

$$\log L(\theta) = \sum_{i=1}^N \log P(\vec{x}_i | \vec{\mu}, \Sigma) = \sum_{i=1}^N \log \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \right\}$$

① Max over $\vec{\mu}$:

$$\frac{\partial \log L(\theta)}{\partial \vec{\mu}} = \sum_{i=1}^N (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} = \vec{0} \Rightarrow \vec{\mu} = \frac{1}{N} \sum_{i=1}^N \vec{x}_i = \begin{bmatrix} 5 \\ 6 \end{bmatrix}$$

② Max over Σ :

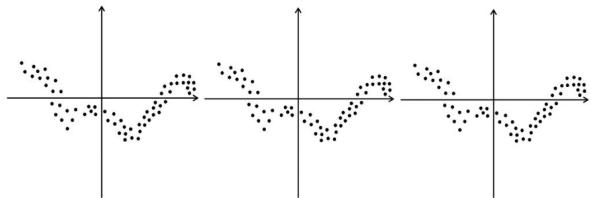
$$\begin{aligned} \log L(\theta) &= \sum_{i=1}^N -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\vec{x}_i - \vec{\mu})^T \Sigma^{-1} (\vec{x}_i - \vec{\mu}) \\ &= -\frac{ND}{2} \log 2\pi + \frac{N}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^N \text{tr} [(\vec{x}_i - \vec{\mu})^T (\vec{x}_i - \vec{\mu}) \Sigma^{-1}] \\ \text{We set } A = \Sigma^{-1} &= -\frac{ND}{2} \log 2\pi + \frac{N}{2} \log |A| - \frac{1}{2} \sum_{i=1}^N \text{tr} [(\vec{x}_i - \vec{\mu})^T (\vec{x}_i - \vec{\mu}) A] \end{aligned}$$

$$\begin{aligned} \frac{\partial \log L(\theta)}{\partial A} &= \frac{N}{2} \sum_{i=1}^N -\frac{1}{2} \sum_{j=1}^D (\vec{x}_{ij} - \mu_j)^2 \cdot (\vec{x}_{ij} - \mu_j) = 0 \\ &\sum = \frac{1}{N} \sum_{i=1}^N (\vec{x}_i - \vec{\mu})^T (\vec{x}_i - \vec{\mu}) \\ &= \frac{1}{5} \cdot [32 + 8 + 0 + 8 + 32] \\ &= 16 \end{aligned}$$

Hence, $\mu = \begin{bmatrix} 5 \\ 6 \end{bmatrix}$, $\Sigma = 16$ maximum the likelihood.

Q7. Problem 7 (10 points)

Explain the difference between overfitting and underfitting. For a picture given below, show an example of underfitting on the left plot, proper fit in the middle plot, and overfitting on the right plot.



Difference - Underfitting means a machine learning model can not properly present the trend and feature of data.

Overfitting means a machine learning model describes the noise and random error of training set instead of the proper trend of the data.

② Overfitting occurs when the Machine Learning Model is over-complex.

underfitting occurs when the model is not complex enough or the model is not converged.

