

$$\begin{aligned} T_2 &= \mathbb{E} [\hat{\mu}_n \hat{\mu}_n^T] \\ &= \mathbb{E} [(\hat{\mu}_n - \mathbb{E}[\hat{\mu}_n])(\hat{\mu}_n - \mathbb{E}[\hat{\mu}_n])^T] + \mathbb{E}[\hat{\mu}_n]\mathbb{E}[\hat{\mu}_n]^T \\ &= \text{Var}(\hat{\mu}_n) + \mu\mu^T \\ &\stackrel{(a)}{=} \frac{\Sigma}{n} + \mu\mu^T, \end{aligned}$$

where Var means variance.

(a) is true because $\hat{\mu}_n = \frac{\sum_n \mathbf{x}_i}{n}$ and $\text{Var}(\hat{\mu}_n) = \text{Var}(\frac{\sum_n \mathbf{x}_i}{n}) = \frac{\sum_n \text{Var}(\mathbf{x}_i)}{n^2} = \frac{\sum_n \Sigma}{n^2} = \frac{\Sigma}{n}$. Plugging T_1 and T_2 into the first equation, we have

$$\begin{aligned} \mathbb{E} [\Sigma_n] &= \frac{1}{n} \sum_{i=1}^n (\Sigma + \mu\mu^T) - \left(\frac{\Sigma}{n} + \mu\mu^T \right) \\ &= \Sigma - \frac{\Sigma}{n} \\ &= \frac{n-1}{n} \Sigma. \end{aligned}$$

So $\hat{\Sigma}$ is a biased estimate of the true covariance matrix Σ .

1. (25 points) Let $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ with $\mathbf{x}_i \in \mathbb{R}^d$ be n samples drawn independently from a univariate Gaussian distribution with mean $\mu \in \mathbb{R}$ and variance $\sigma^2 \in \mathbb{R}_{++}$, i.e., $\sigma^2 > 0$.

(a) (10 points) Starting from the density of a univariate Gaussian distribution, clearly explain how to compute the maximum likelihood estimates $\hat{\mu}_n$ and $\hat{\sigma}_n^2$ for the mean μ and variance σ^2 .

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | \mu, \sigma^2) = \prod_{i=1}^n p(\mathbf{x}_i | \mu, \sigma^2) \stackrel{(1)}{=} \exp(-\frac{1}{2\sigma^2} \sum_{i=1}^n (\mathbf{x}_i - \mu)^2) \stackrel{(2)}{=} \exp(-\frac{1}{2\sigma^2} \sum_{i=1}^n (\mathbf{x}_i - \mu)^2) \stackrel{(3)}{=} \exp(-\frac{1}{2\sigma^2} \sum_{i=1}^n (\mathbf{x}_i - \mu)^2)$$

$$\max_{(\mu, \sigma^2)} L(\mu, \sigma^2) = \min_{(\mu, \sigma^2)} -\log p(\mathbf{x}_1, \dots, \mathbf{x}_n | \mu, \sigma^2) \stackrel{(1)}{=} \min_{(\mu, \sigma^2)} \sum_{i=1}^n -\log p(\mathbf{x}_i | \mu, \sigma^2)$$

$$\Rightarrow \min_{(\mu, \sigma^2)} \sum_{i=1}^n -\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\mathbf{x}_i - \mu)^2}{2\sigma^2}\right) \right) \stackrel{(1)}{=} \min_{(\mu, \sigma^2)} \left(\frac{1}{2} \sum_{i=1}^n \log 2\pi + n \log \sigma^2 \right)$$

$$\Rightarrow \min_{(\mu, \sigma^2)} \sum_{i=1}^n -\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\mathbf{x}_i - \mu)^2}{2\sigma^2}\right) \right) \stackrel{(1)}{=} \min_{(\mu, \sigma^2)} \left(\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \mu)^2 \right)$$

$$\Rightarrow \frac{\partial L}{\partial \mu} = 0 \Rightarrow \sum_i (\mathbf{x}_i - \mu) = 0 \Rightarrow \hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

$$\frac{\partial L}{\partial \sigma^2} = 0 \Rightarrow \frac{n}{2} - \frac{\sum_i (\mathbf{x}_i - \mu)^2}{2\sigma^2} = 0 \Rightarrow \hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \mu)^2$$

(b) (7 points) Is the estimate $\hat{\mu}_n$ of μ unbiased, i.e., $E[\hat{\mu}_n] = \mu$? Clearly explain your answer.

$$E[\hat{\mu}_n] = E\left[\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i\right] = \frac{1}{n} \sum_{i=1}^n E[\mathbf{x}_i] = \frac{1}{n} n \mu = \mu$$

is unbiased

(c) (8 points) Is the estimate $\hat{\sigma}_n^2$ of σ^2 unbiased, i.e., $E[\hat{\sigma}_n^2] = \sigma^2$? Clearly explain your answer.

$$E[\hat{\sigma}_n^2] = E\left[\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu}_n)^2\right] = \frac{1}{n} \sum_{i=1}^n E[(\mathbf{x}_i - \hat{\mu}_n)^2]$$

$$E[(\mathbf{x}_i - \hat{\mu}_n)^2] \stackrel{(1)}{=} E[(\mathbf{x}_i - \mu + \hat{\mu}_n - \mu)^2] \stackrel{(2)}{=} E[(\mathbf{x}_i - \mu)^2] + 2E[(\mathbf{x}_i - \mu)(\hat{\mu}_n - \mu)] + E[(\hat{\mu}_n - \mu)^2]$$

$$E[(\mathbf{x}_i - \mu)^2] = E[(\mathbf{x}_i - \mu)^2] \stackrel{(1)}{=} \sigma^2$$

$$E[(\mathbf{x}_i - \mu)(\hat{\mu}_n - \mu)] = E\left[\left(\mathbf{x}_i - \mu\right) \sum_{j=1}^n \mathbf{x}_j\right] = \frac{1}{n} \sum_{j=1}^n E[\mathbf{x}_i \mathbf{x}_j] = \frac{1}{n} \left[(n-1) \sigma^2 + \mu^2 \right]$$

$$E[(\hat{\mu}_n - \mu)^2] = E\left[\left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i - \mu\right)^2\right] = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n E[\mathbf{x}_i \mathbf{x}_j] = \frac{1}{n^2} \left[(n^2-1) \sigma^2 + \mu^2 \right]$$

$$\Rightarrow E[\hat{\sigma}_n^2] = \frac{1}{n} \sum_{i=1}^n \left[(\sigma^2 + \mu^2) - 2(\sigma^2 + \frac{1}{n} \mu^2) + (\sigma^2 + \frac{1}{n} \mu^2) \right]$$

$$= \frac{1}{n} \sum_{i=1}^n (6\sigma^2 + \frac{1}{n} \mu^2) = (n-1)\sigma^2 + \frac{1}{n} \mu^2$$

$\hat{\sigma}_n^2$ is biased

3. (20 points) Suppose there are three kinds of bags of candies:

- $\frac{1}{3}$ are type h_1 : 100% cherry candies,
- $\frac{1}{3}$ are type h_2 : 50% cherry candies and 50% lime candies,
- $\frac{1}{3}$ are type h_3 : 100% lime candies.

We have one type of bag, but we don't know which type it is. We draw a candy from the bag, and it turns out to be lime. What are the posterior probabilities $p(h_i | \text{cherry} = \text{lime}, p(h_2 | \text{cherry} = \text{lime}), p(h_3 | \text{cherry} = \text{lime})$ of each type of bag? Clearly explain your answer, e.g., how you are using Bayes rule, and show your calculations.¹

$$p(h_i | \text{lime}) = \frac{p(\text{cherry} | h_i) \cdot p(h_i)}{p(\text{cherry})}$$

$$p(\text{cherry}) = \frac{3}{3} \sum_{i=1}^3 p(h_i) p(\text{cherry} | h_i) = \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{2} + \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{2} + \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{2} = \frac{1}{2}$$

required derivation

$$p(h_1 | \text{lime}) = \frac{0 \times \frac{1}{2}}{\frac{1}{2}} = 0$$

$$p(h_2 | \text{lime}) = \frac{\frac{1}{2} \times \frac{1}{2}}{\frac{1}{2}} = \frac{1}{2}$$

$$p(h_3 | \text{lime}) = \frac{1 \times \frac{1}{2}}{\frac{1}{2}} = \frac{1}{2}$$

4. (25 points) Consider a k -class multivariate parametric classification with independently drawn training data $Z = \{(\mathbf{x}_1, r_1), \dots, (\mathbf{x}_N, r_N)\}$. Each data point \mathbf{x}_i is d -dimensional with binary features, i.e., $x_{ij} \in \{0, 1\}, j = 1, \dots, d$. We will focus on a Bernoulli Naive Bayes model for the dataset. In particular, for any feature vector $\mathbf{x} = [x_1 \dots x_d]$ and any class C_i , we assume

$$p(\mathbf{x}|C_i) = \prod_{j=1}^d p(x_j|C_i).$$

(a) (10 points) For a given class $C_i = 1, \dots, k$ and a given feature $j, j = 1, \dots, d$, clearly explain how to compute the maximum likelihood estimate \hat{p}_{ij} corresponding to the Bernoulli probability $p(x_j = 1|C_i)$ using the training data.

Indices : $i = 1 \dots k \rightarrow \text{class}$
 $t = 1 \dots N \rightarrow \text{data point}$
 $j = 1 \dots d \rightarrow \text{feature}$

* Note for Likelihood, maximizing \hat{p}_{ij} for class i and feature j is independent from maximizing other \hat{p}_{ij} for a different class or feature because \hat{p}_{ij} only affects probabilities with class i and feature j .

Thus : $\hat{p}_{ij}(\mathbf{x}|r_j) = \frac{N}{t} \left[\prod_{i=1}^N p_{ij}^{x_{ij}} (1-p_{ij})^{1-x_{ij}} \right]^{\frac{1}{N}}$, where p_{ij} is $\begin{cases} 1, & \text{if } x_{ij} = 1 \\ 0, & \text{otherwise} \end{cases}$

$\log \hat{p}_{ij} = \frac{N}{t} \sum_{i=1}^N (x_{ij} \log p_{ij} + (1-x_{ij}) \log (1-p_{ij}))$

Taking derivative and setting it to zero :

$$\frac{\partial \log \hat{p}_{ij}}{\partial p_{ij}} = \frac{N}{t} \sum_{i=1}^N \left(\frac{x_{ij}}{p_{ij}} - \frac{1-x_{ij}}{1-p_{ij}} \right) = 0$$

$$\Rightarrow \sum_{i=1}^N x_{ij} = \frac{N}{t} \frac{p_{ij} - p_{ij} x_{ij}}{1-p_{ij}}$$

$$\Rightarrow (1-p_{ij}) \sum_{i=1}^N x_{ij} = p_{ij} \frac{N}{t} (1-p_{ij})$$

$$\Rightarrow \sum_{i=1}^N x_{ij} = p_{ij} \frac{N}{t}$$

$$\Rightarrow \hat{p}_{ij} = \frac{\sum_{i=1}^N x_{ij}}{\frac{N}{t}}$$

Answer with a final solution (e.g. from slide) but without a clear derivation receives no credit.

Name: _____ Page 7 of 8

(b) (7 points) Given the training data, clearly explain how to compute the maximum likelihood estimates $\hat{p}(C_i)$ of the class prior probabilities $p(C_i)$ for each class, $i = 1, \dots, k$.

* Note that for prior there is an implicit constraint : $\sum_{i=1}^k p(C_i) = 1$ and $p(C_i) \geq 0$

so $\hat{p}(C_i)$ can't be infinitely large

* Similar to (a), prior $p(C_i)$ only affects points that are in class i

$$\hat{p}(C_i) = \frac{1}{t} \sum_{i=1}^t p(C_i) \stackrel{(1)}{=} \frac{1}{t} \left[\prod_{i=1}^t (1-p_{Ci})^{1-p_{Ci}} p_{Ci}^{p_{Ci}} \right]$$

$$\log \hat{p}(C_i) = \frac{1}{t} \sum_{i=1}^t \left[\log p_{Ci} + (1-p_{Ci}) \log (1-p_{Ci}) \right]$$

Taking derivative and setting to zero :

$$\sum_{i=1}^t \frac{p_{Ci}}{1-p_{Ci}} = \sum_{i=1}^t \frac{1-p_{Ci}}{p_{Ci}}$$

$$\Rightarrow (1-p_{Ci}) \sum_{i=1}^t p_{Ci} = p_{Ci} \sum_{i=1}^t (1-p_{Ci})$$

$$\Rightarrow \sum_{i=1}^t p_{Ci} = p_{Ci} \frac{t}{\sum_{i=1}^t p_{Ci}}$$

$$\Rightarrow \hat{p}(C_i) = \frac{t}{N}$$

Answer with a final solution (e.g. from slide) but without a clear derivation receives no credit.

(c) (8 points) Given a new test point \mathbf{x}_{test} with binary features, clearly explain how to predict a class label \hat{y}_{test} using the estimated \hat{p}_{ij} and $\hat{p}(C_i)$.

$$p(\mathbf{x} | \mathbf{x}_{\text{test}}) \propto p(\mathbf{x}_{\text{test}} | C_i) \cdot p(C_i)$$

$$\Rightarrow p(C_i | \mathbf{x}_{\text{test}}) \propto \left[\prod_{j=1}^d \hat{p}_{ij}^{x_{ij}} (1-\hat{p}_{ij})^{1-x_{ij}} \right] \cdot \hat{p}(C_i)$$

① For each class i , compute $p(C_i | \mathbf{x}_{\text{test}})$

② class label $\hat{y}_{\text{test}} = \underset{i}{\operatorname{arg\,max}} p(C_i | \mathbf{x}_{\text{test}})$

the derivation and formulas are required

A very high level answer (e.g. step ① & ②)

without mathematic details receives no credit.

using \hat{p}_{ij} and $\hat{p}(C_i)$