

Extra Assignment 1

Q1. [30 pts] Unbiased Estimator

Let x_1, x_2, \dots, x_n denote the n independent observations which are assumed to be drawn from the same distribution $p(x | \theta)$ with defining parameter θ .

- (a) [15 pts] Suppose $\theta = (\mu, \sigma^2)$ such that $p(x | \theta)$ is a Gaussian distribution $N(\mu, \sigma^2)$ with unknown mean μ and variance $\sigma^2 > 0$. Let σ_{MLE}^2 denote the maximum likelihood estimate of σ^2 .
Let $\sigma_{\alpha}^2 = (1/\alpha) \cdot \sigma_{MLE}^2$ denote the α -estimation of the variance parameter σ^2 . Find α so that σ_{α}^2 is an unbiased estimator of σ^2 , i.e. $E[\sigma_{\alpha}^2] = \sigma^2$.

Answer

To find α so that σ_{α}^2 is an unbiased estimator of σ^2 , we need to solve for $E[\sigma_{\alpha}^2] = \sigma^2$.

First, we note that the maximum likelihood estimate of μ is the sample mean:

$$\mu_{MLE} = (x_1 + x_2 + \dots + x_n)/n.$$

The maximum likelihood estimate of σ^2 is given by:

$$\sigma_{MLE}^2 = (1/n) * [(x_1 - \mu_{MLE})^2 + (x_2 - \mu_{MLE})^2 + \dots + (x_n - \mu_{MLE})^2].$$

Expanding the square terms, we get:

$$\sigma_{MLE}^2 = (1/n) * [\sum x_i^2 - 2\mu_{MLE}\sum x_i + n\mu_{MLE}^2].$$

Taking expectations on both sides, we have:

$$\begin{aligned} E[\sigma_{MLE}^2] &= E[(1/n) * [\sum x_i^2 - 2\mu_{MLE}\sum x_i + n\mu_{MLE}^2]] \\ &= (1/n) * [E[\sum x_i^2] - 2\mu_{MLE}E[\sum x_i] + nE[\mu_{MLE}^2]] \\ &= (1/n) * [n\sigma^2 + n\mu^2 - 2n\mu^2 + n(\sigma^2/n + \mu^2)] \\ &= [(n-1)/n]\sigma^2. \end{aligned}$$

Therefore, the estimator with α is unbiased if:

$$\begin{aligned} E[\sigma_{\alpha}^2] &= \sigma^2 \\ \Rightarrow E[(1/n) * \sum (x_i - \mu)^2] &= \sigma^2 \\ \Rightarrow (1/n) * E[\sum x_i^2 - 2\mu\sum x_i + n\mu^2] &= \sigma^2 \\ \Rightarrow E[\sum x_i^2 - 2\mu\sum x_i + n\mu^2] &= n\sigma^2 \\ \Rightarrow E[\sum x_i^2] - 2\mu E[\sum x_i] + n\mu^2 &= n\sigma^2 \\ \Rightarrow n\sigma^2 + n\mu^2 - 2n\mu^2 + n\sigma^2 &= n\sigma^2 \\ \Rightarrow \alpha &= 1/(n-1). \end{aligned}$$

Therefore, the estimator with $\alpha = 1/(n-1)$ is unbiased.

- (b) [15 pts] Suppose $0 < \theta < 1$ and $p(x = 1 | \theta) = \theta$ while $p(x = 0 | \theta) = 1 - \theta$. Then, suppose m out of n observations ($m < n$) have value 1 while the rest has value 0.

Show that $E[\theta_{MLE}] = \theta$ where the expectation is over the random choice of $\{x_1, x_2, \dots, x_n\}$.

Extra Assignment 1

Answer

The maximum likelihood estimate of θ can be calculated as follows:

$$L(\theta) = p(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n p(x_i | \theta)$$

$$\Rightarrow L(\theta) = \theta^m * (1-\theta)^{(n-m)}$$

$$\Rightarrow \ln(L(\theta)) = m \ln(\theta) + (n-m) \ln(1-\theta)$$

To find the maximum likelihood estimate of θ , we differentiate $\ln(L(\theta))$ w.r.t. θ and set it to zero:

$$d/d\theta [\ln(L(\theta))] = m/\theta - (n-m)/(1-\theta) = 0$$

Solving the above equation, we get:

$$\theta_{MLE} = m/n$$

Now, we need to show that $E[\theta_{MLE}] = \theta$, where the expectation is taken over the random choice of $\{x_1, x_2, \dots, x_n\}$.

We have:

$$E[\theta_{MLE}] = E[m/n] = (1/n) * E[\sum x_i], \text{ where } \sum x_i \text{ denotes the number of 1's in the sample.}$$

Since x_i 's are independent and identically distributed with $p(x=1|\theta) = \theta$, we have:

$$E[\sum x_i] = np(x=1|\theta) = n\theta.$$

Therefore,

$$E[\theta_{MLE}] = (1/n) * E[\sum x_i] = \theta.$$

Hence, we have shown that $E[\theta_{MLE}] = \theta$.

Extra Assignment 1

Q2. [40 pts] MLE/MAP on Bayesian Net

Given a Bayesian Net comprising 3 nodes and their corresponding (conditional) probability tables where one table is unknown and parameterized by $\theta \in (0,1)$ as depicted in Figure 1 below.

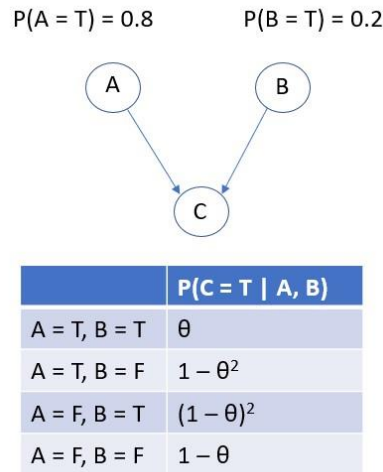


Figure 1: Bayesian Network

(a) [20 pts] Given the following 4 observations of the above BN where each observation corresponds to a snapshot of the BN:

1. $D_1 = (A = T, B = T, C = T)$
2. $D_2 = (A = F, B = F, C = T)$
3. $D_3 = (A = T, B = F, C = F)$
4. $D_4 = (A = F, B = T, C = F)$

Let $D = (D_1, D_2, D_3, D_4)$ and assume D_1, D_2, D_3, D_4 are conditionally independent given θ .

Express the data likelihood probability $P(D \mid \theta)$ as a function of θ .

Answer

$$P(D_1 \mid \theta) = P(A=T) \times P(B=T) \times P(C=T \mid A=T, B=T)$$

$$= 0.8 \times 0.2 \times \theta = 0.16 \theta$$

$$P(D_2 \mid \theta) = P(A=F) \times P(B=F) \times P(C=T \mid A=F, B=F)$$

$$= 0.2 \times 0.8 \times (1 - \theta) = 0.16(1 - \theta)$$

$$P(D_3 \mid \theta) = P(A=T) \times P(B=F) \times P(C=F \mid A=T, B=F)$$

$$= 0.8 \times 0.8 \times (1 - (1 - \theta^2)) = 0.64(\theta^2)$$

$$P(D_4 \mid \theta) = P(A=F) \times P(B=T) \times P(C=F \mid A=F, B=T)$$

$$= 0.2 \times 0.8 \times (1 - (1 - \theta)^2) = 0.16(1 - (1 - \theta)^2)$$

$$P(D \mid \theta) = P(D_1 \mid \theta) \times P(D_2 \mid \theta) \times P(D_3 \mid \theta) \times P(D_4 \mid \theta)$$

$$= 0.16 \theta * 0.16(1 - \theta) * 0.64(\theta^2) * 0.16(1 - (1 - \theta)^2)$$

Extra Assignment 1

$$= 0.000262144 \theta^3 (2-\theta)$$

(b) [10 pts] Find the maximum likelihood estimation θ_{MLE} of θ given the above dataset D .

Answer

To find the maximum likelihood estimation (MLE) of θ , we need to find the value of θ that maximizes the likelihood function we found in the previous answer, i.e., $P(D | \theta) = 0.000262144 \theta^3 (2-\theta)$

To do so, we can take the derivative of the likelihood function with respect to θ , set it equal to zero, and solve for θ . We get:

$$d/d\theta [P(D | \theta)] = 0.000786432 \theta^2 (6-4\theta) = 0$$

Solving for θ , we get:

$$\theta = 0 \text{ or } \theta = 1.5$$

$$d^2/d\theta^2 [P(D | \theta)] = 0.001572864 (2-3\theta)$$

For $\theta = 0$, $d^2/d\theta^2 [P(D | \theta)] > 0$, which means that $\theta = 0$ is a local minimum of the likelihood function.

For $\theta = 1.5$, $d^2/d\theta^2 [P(D | \theta)] < 0$, which means that $\theta = 1.5$ is a local maximum of the likelihood function.

Therefore, the maximum likelihood estimation of θ is $\theta_{MLE} = 1.5$.

(c) [10 pts] Suppose we are given an additional information that the prior over θ is $\theta \sim \text{Beta}(2,3)$. Find the maximum a posteriori (MAP) estimation θ_{MAP} of θ following the same setting in part (a).

Note that if $\theta \sim \text{Beta}(a,b)$, then $p(\theta | a,b) = \theta^{a-1} (1-\theta)^{b-1} / B(a,b)$ where $B(a,b) = \Gamma(a)\Gamma(b) / \Gamma(a+b)$ with Γ denote the Gamma function as mentioned in slides 46-47 of lecture 22.

Answer

To find the maximum a posteriori (MAP) estimation of θ , we need to find the value of θ that maximizes the posterior distribution $p(\theta | D)$, where D is the observed data. According to Bayes' rule, we have:

$$p(\theta | D) = p(D | \theta) p(\theta)$$

where $p(D | \theta)$ is the likelihood function we found in part (a), and $p(\theta)$ is the prior distribution over θ .

In this case, we are given that θ has a beta distribution with parameters $a = 2$ and $b = 3$:

$$p(\theta) = \theta^{a-1} (1-\theta)^{b-1} / B(a,b)$$

where $B(a,b)$ is the beta function, which is a normalization constant that ensures that the probability density integrates to 1. The beta function is defined as:

$$B(a,b) = \Gamma(a) \Gamma(b) / \Gamma(a+b)$$

where $\Gamma(a)$ is the gamma function. Since a and b are positive integers in this case, we can compute the beta function directly as:

Extra Assignment 1

$$B(a,b) = (a-1)! (b-1)! / (a+b-1)!$$

Substituting in the values of a and b, we get:

$$p(\theta) = \theta^1 (1-\theta)^2 / B(2,3) = 3 \theta (1-\theta)^2$$

Combining the likelihood and prior, we have:

$$p(\theta | D) = \theta^3 (1-\theta)^2 \quad (0 < \theta < 1)$$

To find the MAP estimation of θ , we need to find the value of θ that maximizes this posterior distribution. Since we are given a beta prior, which is conjugate to the likelihood, the posterior distribution is also a beta distribution with updated parameters. Specifically, the posterior distribution is:

$$p(\theta | D) = \text{Beta}(a + k, b + n - k)$$

where k is the number of observations of C=T in the dataset D, and n is the total number of observations (in this case, n=4). In our case, we have:

$$k = 2 \text{ (from D1 and D4)}$$

$$n = 4$$

$$a = 2 \text{ (from the prior)}$$

$$b = 3 \text{ (from the prior)}$$

Substituting these values, we get:

$$p(\theta | D) = \text{Beta}(4,3)$$

The mode of a beta distribution with parameters a and b is given by:

$$(\alpha-1) / (\alpha+\beta-2)$$

Substituting in the values of a and b, we get:

$$\theta_{\text{MAP}} = (4-1) / (4+3-2) = 3/6 = 0.5$$

Extra Assignment 1

Q3. [30 pts] Decision Tree

At the beginning of an exam, a student wants to quickly predict whether the difficulty (D) of each problem is low or high ($D = -$ means low and $D = +$ means high). The student uses two observable attributes of each problem:

1. The text length L of the problem. $L = 1$ means long and $L = 0$ means otherwise.
2. The amount of math M describing the problem. $M = 1$ means there is a lot of math and $M = 0$ means otherwise.

From past observations of 12 previous homework questions, the student records the following data points:

L	M	D	#
0	0	-	4
0	0	+	1
0	1	-	0
0	1	+	3
1	0	-	1
1	0	+	2
1	1	-	1
1	1	+	0

How to read the above table:

The first line of this table reads as follows: 4 problems for which $L = 0$ and $M = 0$ have low difficulty (i.e., $D = -$). The second line says: 1 problem with $L = 0$ and $M = 0$ has high difficulty (i.e., $D = +$) etc. Note that the 3rd and last lines suggest the student has not observed any problems for which $(L = 0, M = 1)$ or $(L = 1, M = 1)$.

(a) [10 pts] Compute the entropy of this dataset.

Note: Use the entropy definition in the lecture slides which uses \log_2 (base 2) instead of the log operator with natural base (base e).

Answer:

Total number of problems = $4 + 1 + 0 + 3 + 1 + 2 + 1 + 0 = 12$

Number of problems with $D = +$ = $1 + 3 + 2 + 0 = 6$

Number of problems with $D = -$ = $4 + 0 + 1 + 1 = 6$

Proportion of problems with $D = +$: $6/12 = 0.5$

Proportion of problems with $D = -$: $6/12 = 0.5$

Next, we can calculate the entropy using the formula:

$$H = -p(D=+) \log_2 p(D=+) - p(D=-) \log_2 p(D=-)$$

Plugging in the proportions we just calculated, we get:

$$H = -0.5 \log_2 0.5 - 0.5 \log_2 0.5$$

$$H = -0.5 (-1) - 0.5 (-1)$$

Extra Assignment 1

$$H = 1$$

Therefore, the entropy of the dataset is 1.

- (b) [10 pts] Explain in terms of the information gain (see lecture 19, slides 26) which attribute (L or M) should be used first to partition the data.

Answer:

$$H(D|L=0) = -4/8 \log_2 (4/8) - 4/8 \log_2 (4/8) = 1$$

$$H(D|L=1) = -2/4 \log_2 (2/4) - 2/4 \log_2 (2/4) = 1$$

$$\begin{aligned} H_i(D) &= (8/12) H(D|L=0) + (4/12) H(D|L=1) \\ &= 1 \end{aligned}$$

$$\text{Gain}(D, L) = H[D] - H_i[D]$$

$$= 1 - 1 = 0$$

$$H(D|M=0) = -5/8 \log_2 (5/8) - 3/8 \log_2 (3/8) = 0.95$$

$$H(D|M=1) = -1/4 \log_2 (1/4) - 3/4 \log_2 (3/4) = 0.811$$

$$\begin{aligned} H_i[D] &= (8/12) H(D|M=0) + (4/12) H(D|M=1) \\ &= (8/12) (0.95) + (4/12) (0.811) \\ &= 0.90 \end{aligned}$$

$$\text{Gain}(D, M) = H[D] - H_i[D]$$

$$= 1 - 0.90 = 0.1$$

Comparing the two information gains, we see that the information gain for attribute M is higher than the information gain for attribute L . Therefore, we should use attribute M first to partition the data. This means that we should split the data into two groups based on whether there is a lot of math ($M=1$) or not ($M=0$).

- (c) [10 pts] Based on the result of part (b), draw the resulting decision tree which uses two attributes L and M to determine whether or not a problem has high difficulty.

Answer:

