

CSE 312: Foundations of Computing II

Quiz Section #9: Maximum Likelihood Estimation, Chernoff Bound

Review: Main Theorems and Concepts

Weak Law of Large Numbers (WLLN): Let X_1, \dots, X_n be iid random variables with common mean μ and variance σ^2 . Let $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ be the sample mean for a sample of size n . Then, for any $\epsilon > 0$, $\lim_{n \rightarrow \infty} \mathbb{P}(|\bar{X}_n - \mu| > \epsilon) = 0$. We say that \bar{X}_n converges in probability to μ .

Chernoff Bound: Let X_1, \dots, X_n be independent random variables with $X_i \sim \text{Ber}(p_i)$, and let $X = \sum_{i=1}^n X_i$ have $\mathbb{E}[X] = \mu$. Then, for any $\delta > 0$, we have the following

- $\mathbb{P}(X \geq (1 + \delta)\mu) \leq \exp\left(\frac{-\delta^2\mu}{2+\delta}\right)$
- $\mathbb{P}(X \leq (1 - \delta)\mu) \leq \exp\left(-\frac{\delta^2\mu}{2}\right)$
- $\mathbb{P}(|X - \mu| > \delta\mu) \leq 2 \exp\left(\frac{-\delta^2\mu}{2+\delta}\right)$

Sample/Realization: A sample (or realization) x of a random variable X is the value that is actually observed.

Likelihood: Let x_1, \dots, x_n be iid samples from probability mass function $p_X(x | \theta)$ (if X discrete) or density $f_X(x | \theta)$ (if X continuous), where θ is a parameter (or a vector of parameters). We define the likelihood function to be the probability of seeing the data.

If X is discrete:

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

If X is continuous:

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

Maximum Likelihood Estimator (MLE): We denote the MLE of θ as $\hat{\theta}_{\text{MLE}}$ or simply $\hat{\theta}$, the parameter (or vector of parameters) that maximizes the likelihood function (probability of seeing the data).

$$\hat{\theta}_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} L(x_1, \dots, x_n | \theta) = \underset{\theta}{\operatorname{argmax}} \ln L(x_1, \dots, x_n | \theta)$$

Log-Likelihood: We define the log-likelihood as the natural logarithm of the likelihood function. Since the logarithm is a strictly increasing function, the value of θ that maximizes the likelihood will be exactly the same as the value that maximizes the log-likelihood.

If X is discrete:

$$\ln L(x_1, \dots, x_n | \theta) = \sum_{i=1}^n \ln p_X(x_i | \theta)$$

If X is continuous:

$$\ln L(x_1, \dots, x_n | \theta) = \sum_{i=1}^n \ln f_X(x_i | \theta)$$

Steps to find the maximum likelihood estimator, $\hat{\theta}$:

1. Find the likelihood and log-likelihood of the data.
2. Take the derivative of the log-likelihood and set it to 0 to find a candidate for the MLE, $\hat{\theta}$.
3. Take the second derivative and show that $\hat{\theta}$ indeed is a maximizer, that $\frac{d^2L}{d\theta^2} < 0$ at $\hat{\theta}$. Also ensure that it is the global maximizer: check points of non-differentiability and boundary values.

Exercises

1. Suppose x_1, \dots, x_n are iid samples from a distribution with density

$$f_X(x | \theta) = \begin{cases} \frac{\theta x^{\theta-1}}{3^\theta}, & 0 \leq x \leq 3 \\ 0, & \text{otherwise} \end{cases}$$

Find the MLE for θ .

2. Suppose x_1, \dots, x_{2n} are iid samples from the Laplace density (double exponential density)

$$f_X(x | \theta) = \frac{1}{2} e^{-|x-\theta|}$$

Find the MLE for θ . For this problem, you need not verify that the MLE is indeed a maximizer. You may find the **sign** function useful:

$$\text{sgn}(x) = \begin{cases} +1, & x \geq 0 \\ -1, & x < 0 \end{cases}$$

3. Suppose $X \sim \text{Bin}(6, 0.4)$. We will bound $\mathbb{P}(X \geq 4)$ using the concentration inequalities we've learned, and compare this to the true result.
 - (a) Give an upper bound for this probability using Markov's inequality.
 - (b) Give an upper bound for this probability using Chebyshev's inequality.
 - (c) Give an upper bound for this probability using the appropriate Chernoff bound.
 - (d) Give the exact probability.

4. (**MAP Estimation**) Let x_1, \dots, x_n be iid realizations from a distribution with common pmf $p_X(x; \theta)$ where θ is an unknown but **fixed** parameter. Let's call the event $\{X_1 = x_1, \dots, X_n = x_n\} = \mathcal{D}$ for data. You may wonder why in MLE, we seek to maximize the likelihood $L(\mathcal{D} | \theta)$, rather than $\mathbb{P}(\theta | \mathcal{D})$. This is because it doesn't make sense to compute $\mathbb{P}(\theta)$, since θ is fixed. However, in **Maximum a Posteriori (MAP) estimation**, we assume the parameter is a random variable (denoted Θ), and attempt to maximize $\pi_\Theta(\theta | \mathcal{D})$, where π_Θ is the pmf or pdf of Θ , depending on whether Θ is continuous or discrete. Using Bayes Theorem, we get $\pi_\Theta(\theta | \mathcal{D}) = \frac{L(\mathcal{D} | \theta)\pi_\Theta(\theta)}{L(\mathcal{D})}$. To maximize the LHS with respect to θ , we may ignore the denominator on the RHS since it is constant with respect to θ . Hence MAP seeks to maximize $\pi_\Theta(\theta | \mathcal{D}) \propto L(\mathcal{D} | \theta)\pi_\Theta(\theta)$. We call $\pi_\Theta(\theta)$ the **prior** distribution on the parameter Θ , and $\pi_\Theta(\theta | \mathcal{D})$ the **posterior** distribution on Θ . MLE maximizes the likelihood, and MAP maximizes the product of the likelihood and the prior. If the prior is uniform, we will see that MAP is the same as MLE (since $\pi_\Theta(\theta)$ won't depend on θ).

- (a) Suppose we have the samples $x_1 = 0, x_2 = 0, x_3 = 1, x_4 = 1, x_5 = 0$ from the $Ber(\theta)$ distribution, where θ is unknown. Assume θ is unrestricted; that is, $\theta \in (0, 1)$. What is $\hat{\theta}_{MLE}$?
- (b) Suppose we impose that $\theta \in \{0.2, 0.5, 0.7\}$. What is $\hat{\theta}_{MLE}$?
- (c) Assume Θ is restricted as in part (b) (now a random variable for MAP). Assume a (discrete) prior of $\pi_{\Theta}(0.2) = 0.1, \pi_{\Theta}(0.5) = 0.01, \pi_{\Theta}(0.7) = 0.89$. What is $\hat{\theta}_{MAP}$?
- (d) Show that we can make the MAP estimator whatever we want it to be. That is, for each of the three candidate parameters above, find a prior distribution on Θ such that the MAP estimate is that parameter.
- (e) Typically, for the Bernoulli distribution, if we use MAP, we want to be able to get any value $\theta \in (0, 1)$ (not just ones in a finite set such as $\{0.2, 0.5, 0.7\}$). So we assign θ the **Beta distribution** with parameters $\alpha, \beta > 0$ and density $\pi_{\Theta}(\theta) = c\theta^{\alpha-1}(1-\theta)^{\beta-1}$ for $\theta \in (0, 1)$ and 0 otherwise as a prior, where c is a normalizing constant which has a complicated form. The **mode** of a $W \sim Beta(\alpha, \beta)$ random variable is given as $\frac{\alpha-1}{\alpha+\beta-2}$ (the mode is the value with the highest density = $\arg \max_{w \in (0,1)} f_W(w)$). Suppose x_1, \dots, x_n are iid samples from the Bernoulli distribution with unknown parameter, where $\sum_{i=1}^n x_i = k$. Recall that the MLE is k/n . Show that the posterior $\pi_{\Theta}(\theta | \mathcal{D})$ has a $Beta(k + \alpha, n - k + \beta)$ density, and find the MAP estimator for Θ . (Hint: use the mode given). Notice that $Beta(1, 1) \equiv Unif(0, 1)$. If we had this prior, how would the MLE and MAP estimates compare?
- (f) Since the posterior is also a Beta distribution, we call the Beta distribution the **conjugate prior** to the Bernoulli distribution. Interpret what the parameters α, β mean as to the prior.
- (g) Which do you think is “better”, MLE or MAP?