

CSE 312: Foundations of Computing II

Section 9: Concentration Inequalities and Maximum Likelihood

0. Concentration Inequalities

Suppose $X \sim \text{Binomial}(6, 0.4)$. We will bound $\Pr(X \geq 4)$ using the tail bounds we've learned, and compare this to the true result.

- (a) Give an upper bound for this probability using Markov's inequality. Why can we use Markov's inequality?
- (b) Give an upper bound for this probability using Chebyshev's inequality. You may have to rearrange algebraically and it may result in a weaker bound.
- (c) Give an upper bound for this probability using the Chernoff bound.
- (d) Give the exact probability.

1. Laplace MLE

Suppose x_1, \dots, x_{2n} are iid realizations from the Laplace density (double exponential density): for $x \in \mathbb{R}$,

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

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:

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

2. 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 $\Pr(\theta | \mathcal{D})$. This is because it doesn't make sense to compute $\Pr(\theta)$, 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 Bernoulli(θ) 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/Binomial 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 $W \sim \text{Beta}(\alpha, \beta)$ 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 $\text{Beta}(k + \alpha, n - k + \beta)$ density, and find the MAP estimator for Θ . (Hint: use the mode given). Notice that $\text{Beta}(1, 1) \equiv \text{Uniform}(0, 1)$. 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/Binomial distribution. Interpret what the parameters α, β mean as to the prior.
- (g) Which do you think is "better", MLE or MAP?