

CS 475 Machine Learning: Homework 1 Analytical

(50 points)

Assigned: Friday, September 6, 2024

Due: Friday, September 20, 2024, 11:59 pm US/Eastern

Trevor Black (tblack20)

Instructions

We have provided this L^AT_EX document for turning in this homework. We give you one or more boxes to answer each question. The question to answer for each box will be noted in the title of the box. You can change the size of the box if you need more space.

Other than your name, do not type anything outside the boxes. Leave the rest of the document unchanged.

Do not change any formatting in this document, or we may be unable to grade your work. This includes, but is not limited to, the height of textboxes, font sizes, and the spacing of text and tables. Additionally, do not add text outside of the answer boxes. Entering your answers are the only changes allowed.

We strongly recommend you review your answers in the generated PDF to ensure they appear correct. We will grade what appears in the answer boxes in the submitted PDF, NOT the original latex file.

1 Maximum Likelihood [10 pts]

1. (5 pts) Given a dataset of n independently and identically distributed data points $\{x_i\}_{i=1}^n$, where each x_i is drawn from a Gaussian distribution, $x_i \sim \mathcal{N}(\mu, \sigma^2)$. Show that the MLE estimator of the mean is:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad (1)$$

Likelihood

$$\mathcal{L}(\mu, \sigma^2 | \{x_i\}_{i=1}^n) = \prod_{i=1}^n p(x_i | \mu, \sigma^2)$$

$$p(x_i | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

$$\mathcal{L}(\mu, \sigma^2 | \{x_i\}_{i=1}^n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

Log likelihood

$$\log \mathcal{L}(\mu, \sigma^2 | \{x_i\}_{i=1}^n) = \ell(\mu, \sigma^2 | \{x_i\}_{i=1}^n) = \log\left(\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right)\right)$$

$$\ell(\mu, \sigma^2 | \{x_i\}_{i=1}^n) = n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

Maximize log likelihood for μ

$$\frac{d}{d\mu} \ell(\mu, \sigma^2 | \{x_i\}_{i=1}^n) = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0$$

$$0 = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = \sum_{i=1}^n (x_i - \mu) = -n\mu + \sum_{i=1}^n (x_i)$$

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

$$\text{So... } \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

2. (5 pts) Show that the MLE estimator of the variance is:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 \quad (2)$$

Likelihood (see previous question for details)

$$\mathcal{L}(\mu, \sigma^2 | \{x_i\}_{i=1}^n) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

Log likelihood (see previous question for details)

$$\ell(\mu, \sigma^2 | \{x_i\}_{i=1}^n) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

Maximize log likelihood for σ^2

$$\frac{d}{d\sigma^2} \ell(\mu, \sigma^2 | \{x_i\}_{i=1}^n) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

$$0 = -n\sigma^2 + \sum_{i=1}^n (x_i - \mu)^2$$

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

So...

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

2 Conditional Independence [20 pts]

A large group of people were surveyed on their recent health. Of these, 0.20 had a fever and 0.05 had pneumonia. Among the people who had pneumonia, 0.70 had cough as a symptom and 0.50 had fever as a symptom. Among the people who had a fever, 0.40 had cough as a symptom.

Let us create a probabilistic model where the presence/absence of each of these two symptoms, cough and fever, are conditionally independent given the presence/absence of pneumonia. Using this data for the empirical probabilities of our model, answer the following questions.

1. (5 pts) Find the probability that someone has both a cough and a fever.

Let $C = \text{cough}$, $F = \text{fever}$, and $P = \text{pneumonia}$

$$p(C \cap F) = p(C \cap F | P)p(P) + p(C \cap F | \neg P)p(\neg P)$$

$$p(C \cap F | P) = p(C | P)p(F | P) = 0.7 \cdot 0.5 = 0.35$$

$$p(C \cap F | \neg P) = p(C | \neg P)p(F | \neg P) = p(C | F) \cdot p(F | \neg P)$$

Find $p(F | \neg P)$...

$$p(F) = p(F | P) \cdot p(P) + p(F | \neg P) \cdot p(\neg P)$$

$$0.2 = 0.5 \cdot 0.05 + p(F | \neg P) \cdot 0.95 \implies p(F | \neg P) \approx 0.1842$$

$$p(C \cap F | \neg P) = .4 \cdot 0.1842 \approx 0.7368$$

$$p(C \cap F) = 0.35 \cdot 0.05 + 0.7368 \cdot 0.95 \approx 0.0875 = 8.75\%$$

2. (5 pts) Find the probability that someone has pneumonia given that they have a fever but no cough.

Let $C = \text{cough}$, $F = \text{fever}$, and $P = \text{pneumonia}$

$$p(P | F, \neg C) = \frac{p(F \cap \neg C | P)p(P)}{p(F \cap \neg C)}$$

$$p(F \cap \neg C | P) = p(F | P)p(\neg C | P) = 0.5 \cdot (1 - 0.7) = .15$$

$$p(F \cap \neg C) = p(F \cap \neg C | P)p(P) + p(F \cap \neg C | \neg P)$$

$$p(F \cap \neg C | \neg P) = p(F | \neg P)p(\neg C | F) \approx 0.1842 \cdot (1 - 0.4) = 0.1105$$

Note that $p(F | \neg P)$ was found in the previous question.

$$p(F \cap \neg C) = 0.15(0.05) + 0.11052(0.95) \approx 0.1135$$

$$p(P | F \cap \neg C) = \frac{p(F \cap \neg C | P)p(P)}{p(F \cap \neg C)} = \frac{0.15 \cdot 0.05}{0.1135} \approx 0.0661 = 6.61\%$$

3. (10 pts) Given assumptions described above, how many parameters do we need to specify the joint distribution $p(\text{fever}, \text{cough}, \text{pneumonia})$?

Let $C = \text{cough}$, $F = \text{fever}$, and $P = \text{pneumonia}$

The total number of combinations of these variables is $2^3 = 8$.

$$p(F, C, P) = p(P)p(F | P)p(C | F, P)$$

$$p(C | F, P) = p(C | P) \text{ Due to conditional independence.}$$

Counting the parameters...

$p(P) \rightarrow 1$ parameter (it is a binary variable)

$p(F | P) \rightarrow 2$ parameters (P may be $\neg P$ as well, giving two binary variables)

$p(C | P) \rightarrow 2$ parameters for identical reason as above.

The sum of these results in 5 total parameters.

3 Bayesian Reasoning [20 pts]

1. (5 pts) Define what conjugate priors are, and explain why they are useful.

A conjugate prior is a prior distribution that belongs to the same probability distribution of a posterior distribution for a likelihood function. $p(\theta | x) \propto p(x | \theta)p(\theta)$

Using a conjugate prior is useful as the posterior distribution is more closely related to the prior distribution. This makes it easier for calculations, updating the posterior by adjusting the prior, and providing an easier interpretation for data.

2. (10 pts) Show that the Gamma distribution is a conjugate prior of the exponential distribution. That is, show that if $x \sim \text{Exp}(\lambda)$ and $\lambda \sim \text{Gamma}(\alpha, \beta)$, then $p(\lambda|x) \sim \text{Gamma}(\alpha^*, \beta^*)$ for some α^*, β^* .

Likelihood of x :

$$p(x | \lambda) = \lambda e^{-\lambda x} \text{ for } x \geq 0$$

Gamma prior:

$$p(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \text{ for } \lambda > 0$$

Bayes' theorem:

$$\begin{aligned} p(\lambda | x) &\propto p(x | \lambda)p(\lambda) = (\lambda e^{-\lambda x}) \cdot \left(\frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \right) \\ &= \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^\alpha e^{-\lambda(x+\beta)} \end{aligned}$$

We can see that $\alpha^* = \alpha + 1$ and $\beta^* = \beta + x$. So...

$$p(\lambda | x) \sim \text{Gamma}(\alpha + 1, \beta + x)$$

3. (5 pts) Derive the maximum a posteriori (MAP) of λ under the $\text{Gamma}(\alpha, \beta)$ prior.

Maximum likelihood estimate of $\hat{\lambda}$

$$\hat{\lambda}_{MAP}(x) = \arg \max_{\lambda} f(\lambda | x)$$

From last question, we see $f(\lambda | x) \propto \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^\alpha e^{-\lambda(x+\beta)}$

$\log p(\lambda | x) = \alpha \log(\lambda) - \lambda(x + \beta) + C$ where C is the constant term.

$$\frac{d}{d\lambda} \log p(\lambda | x) = \frac{\alpha}{\lambda} - (x + \beta) \quad 0 = \frac{\alpha}{\lambda} - (x + \beta) \implies \lambda = \frac{\alpha}{x+\beta}$$

$$\hat{\lambda}_{MAP} = \frac{\alpha}{x+\beta}$$