

DS-GA 3001.001 Applied Statistics: Homework #1

Due on Thursday, September 26, 2024

Please hand in your homework via Gradescope (entry code: DKYKGY) before 11:59 PM.

1. The Gamma distribution has a shape parameter $\alpha > 0$ and a scale parameter $\beta > 0$, with density given by

$$\Gamma_{\alpha,\beta}(y) = \frac{\beta^\alpha}{\Gamma(\alpha)} y^{\alpha-1} e^{-\beta y}, \quad y > 0.$$

Here $\Gamma(\alpha)$ is the Gamma function - you only need to know that this is a function of α and will not need any further properties.

- (a) Show that the family of Gamma distributions $\{\Gamma_{\alpha,\beta}(y)\}_{\alpha,\beta>0}$ belongs to the exponential family. Write down the expressions of $(\theta, T(y), A(\theta), h(y))$.
 - (b) Verify that the Gamma distribution is a conjugate prior for the Poisson family, i.e. if $\lambda \sim \Gamma_{\alpha,\beta}$ and $y \sim \text{Poi}(\lambda)$, then $\lambda \mid y \sim \Gamma_{\alpha(y),\beta(y)}$.
2. Recall from the lecture that for an exponential family $p_\theta(y) = \exp(\langle \theta, T(y) \rangle - A(\theta))h(y)$, the family of conjugate priors has two parameters $\xi \in \mathbb{R}^d$ and $\tau > 0$, with density

$$\pi_{\xi,\tau}(\theta) = \exp(\langle \xi, \theta \rangle - \tau A(\theta))b(\xi, \tau).$$

- (a) Using $\mathbb{E}_{\xi,\tau}[\nabla_\theta \log \pi_{\xi,\tau}(\theta)] = 0$ (you don't need to prove this), show that

$$\mathbb{E}_{\xi,\tau}[\nabla A(\theta)] = \frac{\xi}{\tau}.$$

- (b) Given i.i.d. observations $y_1, \dots, y_n \sim p_\theta(y)$, show that the posterior distribution takes the form

$$\pi_{\xi,\tau}(\theta \mid y_1, \dots, y_n) = \pi_{\xi + \sum_{i=1}^n T(y_i), \tau+n}(\theta).$$

- (c) Show that the posterior mean of $\mu_\theta = \nabla A(\theta)$ is

$$\mathbb{E}_{\xi,\tau}[\nabla A(\theta) \mid y_1, \dots, y_n] = \frac{\tau}{\tau+n} \cdot \mathbb{E}_{\xi,\tau}[\nabla A(\theta)] + \frac{n}{\tau+n} \cdot \frac{1}{n} \sum_{i=1}^n T(y_i).$$

How would you interpret this result?

3. Recall Fisher's 2×2 table in class, but this time we use a multinomial model $(X_1, \dots, X_4) \sim \text{Multi}(N; (\pi_1, \dots, \pi_4))$ to fit the data. It is easy to verify that the MLE $\hat{\pi} = (\hat{\pi}_1, \dots, \hat{\pi}_4)$ is given by $\hat{\pi}_i = X_i/N$, for all $i = 1, 2, 3, 4$.

- (a) Based on the definition of $\hat{\pi}_i$, directly verify that

$$\text{Cov}(\hat{\pi}_i, \hat{\pi}_j) = \begin{cases} \frac{\pi_i(1-\pi_i)}{N} & \text{if } i = j, \\ -\frac{\pi_i\pi_j}{N} & \text{if } i \neq j. \end{cases}$$

(Hint: Recall how to compute the variance of the Binomial distribution.)

- (b) For the log odds estimator $\hat{\theta} = \log \frac{\hat{\pi}_1 \hat{\pi}_4}{\hat{\pi}_2 \hat{\pi}_3}$, using the delta method and the plug-in approach, show that

$$\text{Var}(\hat{\theta}) \approx \sum_{i=1}^4 \frac{1}{X_i}.$$

- (c) If $(X_1, X_2, X_3, X_4) = (9, 12, 7, 17)$, compute $\hat{\theta}$ and the approximation of $\text{Var}(\hat{\theta})$ in (b). Compare your results with Fisher's hypergeometric modeling in class.
4. Coding I: we will verify numerically that in Poisson models, the deviance residual looks more normally distributed than the Pearson residual. Based on the inline instructions, fill in the missing codes in <https://tinyurl.com/nhez6cu>. Be sure to submit a pdf with your codes, outputs, and colab link.
5. Coding II: in this problem, we investigate if a newly discovered poem is indeed written by Shakespeare, replicating the paper “Did Shakespeare write a newly discovered poem?” by Thisted and Efron in 1987. To this end, we collect 884,647 total words of known Shakespeare, and count the following:
- there are 9 words in the poem which never appears in known Shakespeare;
 - there are 7 words in the poem which exactly appears once in known Shakespeare;
 - there are 5 words in the poem which exactly appears twice in known Shakespeare;
 - ...

We use a data vector $y = (9, 7, 5, 8, 11, 10, 21, 16, 18, 8, 5)$ to record these numbers.

On the other hand, using a theory based on empirical Bayes (we will cover it in Lecture 7), a statistician may predict the number of new words in the poem if Shakespeare indeed wrote it, and similarly the number of words which appear once in known Shakespeare, etc. These predictions are presented by a vector

$$\theta = (6.97, 4.21, 3.33, 5.36, 10.24, 13.96, 10.77, 8.87, 13.77, 9.99, 7.48).$$

The key assumption here is that, if Shakespeare indeed wrote this poem, then

$$y_i \stackrel{\text{ind.}}{\sim} \text{Poi}(\theta_i), \quad \forall i = 0, 1, \dots, 10.$$

This problem aims to test this null hypothesis using the tests we learned in class. Based on the inline instructions, fill in the missing codes in <https://tinyurl.com/yn7tjtxp>. Be sure to submit a pdf with your codes, outputs, and colab link.