**MA 354: Data Analysis – Fall 2021 – Due 10/8 at 5p**
**Homework 2:**

*Complete the following opportunities to use what we've talked about in class. These questions will be graded for correctness, communication and succinctness. Ensure you show your work and explain your logic in a legible and refined submission.*

   The starting jobs will be applied in alphabetical order (last name) for question two.

1. **Solver:** provide a solution, if possible, and reasoning for the solution. **Due to group 10/5 or earlier.**

2. **Code Checker:** provides a first check of the solver's worked solutions and ensures they are correct with a solid interpretation.

3. **Checker** checks the solution for completeness, proposes and implements changes if agreed upon by the group. Provides a short paragraph summarizing the discussion of proposals and their reason for acceptance or non-acceptance.

4. **Double Checker** checks the solution for completeness, communication and polish. The Double Checker ensures that the solution is correct and highly polished for submission.

For subsequent questions student roles will move down one position. The rolls change as follows.

1. **Solver $\implies$ Code Checker**

2. **Code Checker $\implies$ Checker**

3. **Checker $\implies$ Double Checker**

4. **Double Checker $\implies$ Solver**

While students have assigned jobs for each question I encourage students to help each other complete the homework in collaboration.

1. Select a continuous distribution (Not the uniform or exponential). It does not have to be one that we cover in the notes! To explore the PDF of your distribution, specify two sets of parameter(s) for your distribution.

    (a) **History** Discuss what types of random variables are modeled with your distribution. Be sure to include a discussion about the support and ensure to provide the density function, and CDF. This requires some internet research – what's the history of the distribution, why was it created and named? What are some exciting applications of this distribution?

    Cite all of your sources in LaTeX by adding a BibTeX citation to the .bib file. To help, I've cited R (R Core Team, 2021) in parentheses here. R Core Team (2021) provides helpful tools for the rest of the questions below. BibTeX citations are available through Google Scholar by clicking the cite button below the article of interest and selecting the BibTeX option.

    Solution: **The Gamma distribution models the probability that $\alpha$ events occur in a Poisson process with mean arrival time $\beta = \frac{1}{\theta}$, $\theta$ a scale parameter. (Analytica Wiki, 2021). The probability density function, which is a function of the random variable $x$, and our parameters $\alpha$ and $\beta$ is given in the equation below: (provided by Gamma Function Wikipedia (2021))**

    $$f(x; \alpha, \beta) = \frac{\beta^{\alpha} x^{\alpha-1} e^{-\beta x}}{\Gamma(\alpha)} \quad \text{for } x > 0 \quad \alpha, \beta > 0, \tag{1}$$

    **Where $\Gamma(\alpha)$ is the Gamma function. The restrictions on the parameters and random variable comes from the need for the PDF to have positive values only. $x$, $\beta > 0$ ensures the numerator is positive, and $\alpha > 0$ ensures the denominator is positive, since $\Gamma(x) > 0 \ \forall \ x \ \in \mathbb{R}^{+}$. $x > 0$ is the support of the distribution - The random variable must be a positive real number. The cumulative distribution function is given below: (provided by Gamma Function Wikipedia (2021))**

    $$F(x; \alpha, \beta) = \int_{0}^{x} f(u; \alpha, \beta) \, du = \frac{\gamma(\alpha, \beta x)}{\Gamma(\alpha)}, \tag{2}$$

    **where $\gamma()$ is the lower incomplete Gamma function, given by: (provided by Incomplete gamma function Wikipedia (2021))**

    $$\gamma(s, x) = \int_{0}^{x} t^{s-1} e^{-t} \, dt. \tag{3}$$

    **The Gamma distribution is named after the Gamma function, which is principle in it's derivation. The Gamma function is a function that extends the idea of the factorial to all real numbers and even complex numbers, provided that any real input is not a non-negative integer. For positive real inputs, (which is the support for the Gamma distribution), the Gamma function draws a continuous line through all points in the $(x, y)$ plane whose $x$ value is an integer and $y$ value equals $x!$.**

    **The Gamma function is good at modeling any continuous random variable that takes on positive real values, is uni-modal, and is positively skewed. The Gamma distribution is frequently used in climatology to model random whether variables such as rainfall (Thom, Herbert CS, 1958). It has also been used to model the size of insurance claims (Boland, Philip J, 2007), multi-path fading of of signal power in wireless**

**communication (Gamma Function Wikipedia, 2021), and the age distribution of cancer incidence (Belikov, Aleksey V, 2017), to name just a few of its applications.**

(b) Show that you have a valid PDF. You will find the `integrate()` function in R helpful.

(c) Find the median for your two sets of parameter(s). Conduct some research to find the median based on our PDF to confirm that your numerical approach is correct.

(d) Graph the PDF for several values of the parameter(s) including the two sets you specified. What does changing the parameter(s) do to the shape of the PDF?

(e) Graph the CDF for the same values of the parameter(s) as you did in Question 1d. What does changing the parameter(s) do to the shape of the CDF? Comment on the aspects of the CDFs that show that the CDF is valid.

(f) Generate a random sample of size $n = 10, 25, 100$, and 1000 for your two sets of parameter(s). In a $4 \times 2$ grid, plot a histogram of each set of data and superimpose the true density function at the specified parameter values. Interpret the results.

2. Continue with the continuous distribution you selected for Question 1.

(a) Provide the mean, standard deviation, skewness, and kurtosis of the PDF. Ensure to interpret each.

(b) Generate a random sample of size $n = 10, 25, 100$, and 1000 for your two sets of parameter(s). Calculate the sample mean, standard deviation, skewness, and kurtosis. Interpret the results.

(c) Generate a random sample of size $n = 10$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a $1 \times 2$ grid, plot a histogram of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

(d) Generate a random sample of size $n = 25$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a $1 \times 2$ grid, plot a histogram of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

(e) Generate a random sample of size $n = 100$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a $1 \times 2$ grid, plot a histogram of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

(f) Generate a random sample of size $n = 100$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a $1 \times 2$ grid, plot a histogram of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

(g) Comment on the results of parts (c)-(f).

3. Select a discrete distribution (not the Poisson). It does not have to be one that we cover in the notes! To explore the PMF of your distribution, specify two sets of parameter(s) for your distribution.

   (a) **History** Discuss what types of random variables are modeled with your distribution. Be sure to include a discussion about the support and ensure to provide the mass function, and CDF. This requires some internet research – what's the history of the distribution, why was it created and named? What are some exciting applications of this distribution? Cite all of your sources.

   (b) Show that you have a valid PMF. You can show this approximately by calculating the series in a repeat loop until probability mass evaluations are infinitesimally small.

   (c) Find the median for your two sets of parameter(s). Conduct some research to find the median based on our PMF to confirm that your numerical approach is correct.

   (d) Graph the PMF for several values of the parameter(s) including the two sets you specified. What does changing the parameter(s) do to the shape of the PMF?

   (e) Graph the CDF for the same values of the parameter(s) as you did in Question 3d. What does changing the parameter(s) do to the shape of the CDF? Comment on the aspects of the CDFs that show that the CDF is valid.

   (f) Generate a random sample of size $n = 10, 25, 100$, and 1000 for your two sets of parameter(s). In a $4 \times 2$ grid, plot a histogram (with bin size 1) of each set of data and superimpose the true mass function at the specified parameter values. Interpret the results.

4. Continue with the discrete distribution you selected for Question 3.

   (a) Provide the mean, standard deviation, skewness, and kurtosis of the PMF. Ensure to interpret each.

   (b) Generate a random sample of size $n = 10, 25, 100$, and 1000 for your two sets of parameter(s). Calculate the sample mean, standard deviation, skewness, and kurtosis. Interpret the results.

   (c) Generate a random sample of size $n = 10$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a $1 \times 2$ grid, plot a histogram (with bin size 1) of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

   (d) Generate a random sample of size $n = 25$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a $1 \times 2$ grid, plot a histogram (with bin size 1) of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

   (e) Generate a random sample of size $n = 100$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a $1 \times 2$ grid, plot a histogram (with bin size 1) of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

   (f) Generate a random sample of size $n = 100$ for your two sets of parameter(s). Calculate the method of moments estimator(s) and maximum likelihood estimator(s). In a $1 \times 2$ grid, plot a histogram (with bin size 1) of each set of data with (1) the method of moments estimated distribution, (2) the maximum likelihood estimated distribution, and superimpose the true distribution in both.

   (g) Comment on the results of parts (c)-(f).

# References

Analytica Wiki (2021). Gamma distribution.

Belikov, Aleksey V (2017). The number of key carcinogenic events can be predicted from cancer incidence. *Scientific reports*, 7(1):1–8.

Boland, Philip J (2007). *Statistical and probabilistic methods in actuarial science.* CRC Press.

Gamma Function Wikipedia (2021). Gamma function.

Incomplete gamma function Wikipedia (2021). Incomplete gamma function.

R Core Team (2021). *R: A Language and Environment for Statistical Computing.* R Foundation for Statistical Computing, Vienna, Austria.

Thom, Herbert CS (1958). A note on the gamma distribution. *Monthly weather review*, 86(4):117–122.