

Project #2

Dylan Robertson
PHSX 815

May 11, 2023

Abstract

This project compared two hypotheses and attempted to distinguish between the two. The Null Hypothesis assumed that the rate parameter λ for a Poisson Distribution was constant, while the Alternative Hypothesis assumed that the rate parameter λ was distributed according to a Gamma Distribution. Data was generated for each hypothesis with $N_{meas} = 1,000$ measurements per experiment, and $N_{exp} = 10,000$ experiments total. A 95% Confidence Level was used for the Hypothesis comparison, which was done using the Log Likelihood Ratio. The power to distinguish between the two hypotheses was found to be 39.4%. Reasons for the low power of test will be discussed.

1 Null Hypothesis

The null hypothesis assumes that the rate parameter λ in a Poisson Distribution is constant. The Poisson Distribution is,

$$P(X|\lambda) = \frac{\lambda^X e^{-\lambda}}{X!}, \quad (1)$$

where λ is a positive, real number, representing the average rate that an event occurs at, and X is a whole number that represents the number of events that occurred within some fixed time interval.

The Poisson Distribution can be seen for various values of λ in Fig. 1. As λ gets sufficiently large, the distribution approaches that of a normal distribution, as predicted by the Central Limit Theorem.

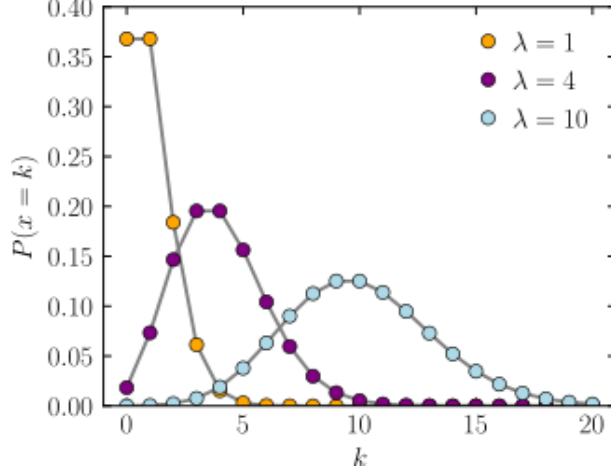


Figure 1: Poisson distributions with various values of λ . The expected peak in the distribution (the mean) occurs at the value of λ .

25 The likelihood of the a distribution is found by Baye's Theorem.

$$P(\lambda|X) \approx \prod_{i=1}^{N_{meas}} P(X_i|\lambda) = \prod_{i=1}^{N_{meas}} \frac{\lambda^{X_i} e^{-\lambda}}{X_i!} \quad (2)$$

26 where X_i is each individual measurement made in a single experiment and
 27 N_{meas} is the total number of measurements in the given experiment.

28 A more useful quantity is the logarithm of the likelihood, as it allows us to
 29 turn the product into a summation using properties of logarithmic functions.
 30 Going through the calculation, and using the Stirling approximation for the
 31 factorial results in the following expression.

$$\ln(P(\lambda|X)) = \sum_{i=1}^{N_{meas}} [x_i \ln(\lambda) - \frac{1}{2} \ln(2\pi x_i) - x_i \ln(x_i) + x_i - \lambda] \quad (3)$$

32 Which is the expression needed to calculate the likelihood of the Null Hy-
 33 pothesis given a sample data set.

34 For the simulation, random measurements from the Poisson Distribution
 35 were done using a standard function from the numpy library. The rate param-
 36 eter was chosen to be $\lambda = 5.0$, the number of measurements per experiment
 37 was $N_{meas} = 1,000$, and the number of experiments was $N_{exp} = 10,000$. The
 38 data for the Null Hypothesis can be seen in Fig. 2.

```

39
40 #Poisson Distribution
41 def Poisson(self, lamb):
42     return np.random.poisson(lam=lamb, size=1)
43 
```

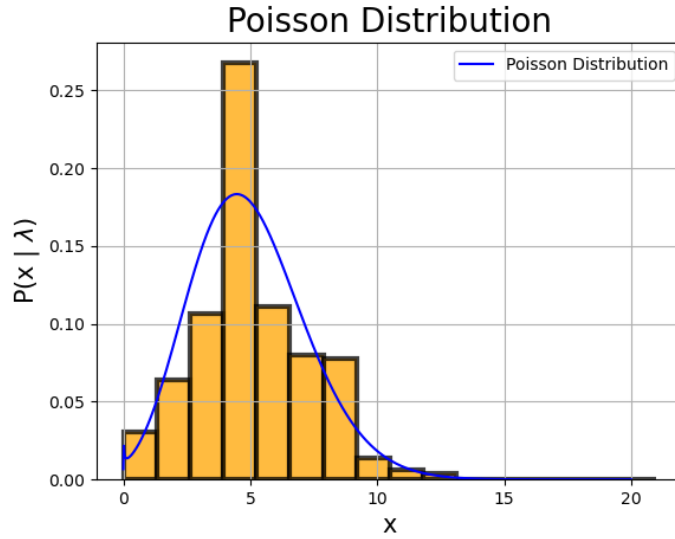


Figure 2: The data generated by the Null Hypothesis assuming a constant rate parameter $\lambda = 5.0$. The curve in blue is the actual distribution.

2 Alternative Hypothesis

The Alternative Hypothesis assumed that the rate parameter λ was distributed according to a Gamma Distribution. The Gamma Distribution is,

$$P(\lambda|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}, \quad (4)$$

where α is a positive, real number representing the shape of the distribution, β is a positive, real number representing the width of the distribution, and $\Gamma(\alpha)$ is the Gamma function.

Fig. 3 shows the Gamma Distribution for various values of α and β . As the value $\frac{\alpha}{\beta}$ gets sufficiently large, the distribution approaches a normal distribution as predicted by the Central Limit Theorem.

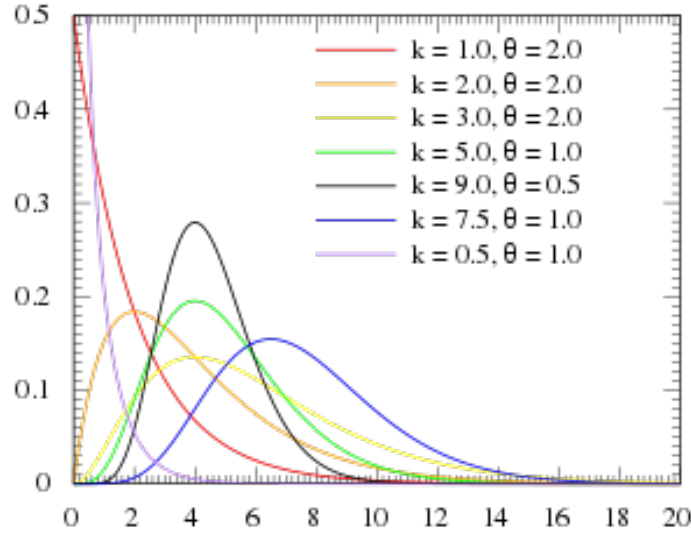


Figure 3: Gamma distributions with various values of $\alpha = k$ and $\beta = \frac{1}{\theta}$. The expected peak in the distribution (the mean) occurs at the value of $\frac{\alpha}{\beta}$.

The simulation for the Alternative Hypothesis worked as follows. The parameters for the Gamma Distribution were chosen such that the most likely λ value remains the same as the Null Hypothesis: $\alpha = 6.0$, and $\beta = 1.2$. For each new experiment, a new value of the rate parameter λ was generated from this Gamma Distribution. Then a data set X was generated from a Poisson Distribution for that value of λ . The number of measurements per experiment and the number of experiments remains the same as that for the Null Hypothesis.

Fig. 4 shows the distribution of λ values according to the Gamma Distribution, and Fig. 5 shows the data generated under this multi-step hypothesis.

Random measurements of λ from the Gamma Distribution were done using a standard function from the numpy library.

```
#Gamma Distribution
def Gamma(self, alpha, beta):
```

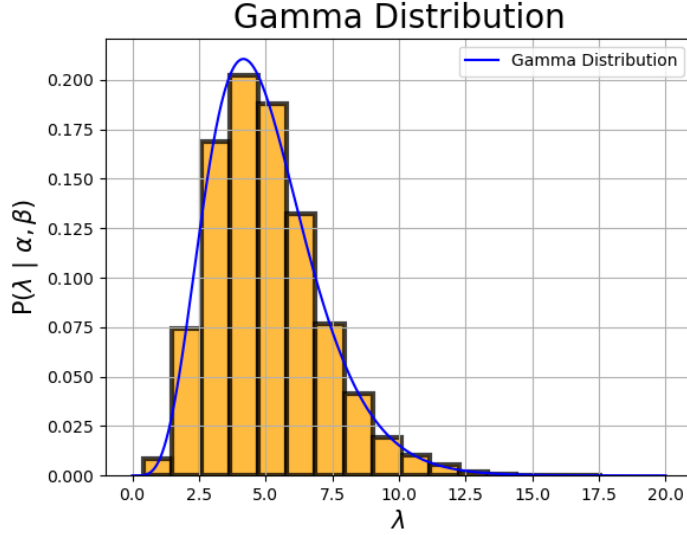


Figure 4: The distribution of λ values according to a Gamma distribution with parameters $\alpha = 6.0$ and $\beta = 1.2$. The expected peak in the distribution (the mean) occurs at the value of $\lambda = \frac{\alpha}{\beta} = 5.0$. The curve in blue is the actual distribution.

```

68     #both alpha and beta must be positive
69     k = alpha
70     theta = 1/beta
71
72     return np.random.gamma(shape = k, scale = theta, size = 1)
73

```

3 Hypothesis Comparison

Hypothesis Comparison was done using the Log Likelihood Ratio (LLR).

$$LLR = \ln\left(\frac{P(X|H_0)}{P(X|H_1)}\right) = \ln(P(X|H_0)) - \ln(P(X|H_1)) \quad (5)$$

where H_0 is the Null Hypothesis described by the parameter λ , and H_1 is the Alternative Hypothesis described by the parameters α and β .

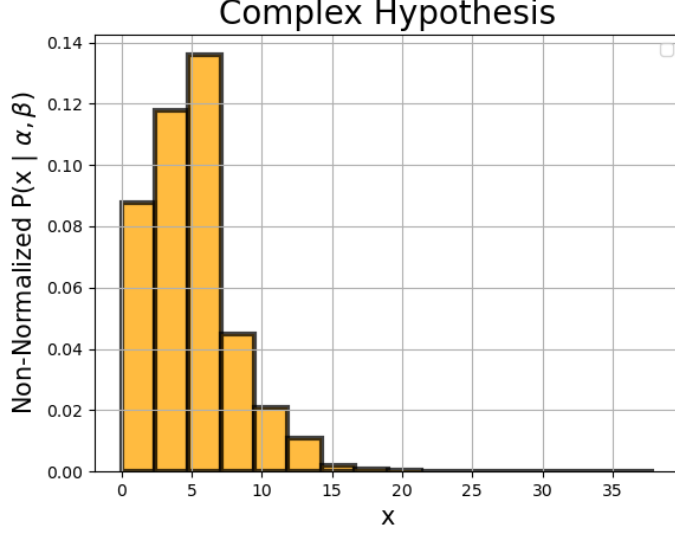


Figure 5: The data generated under the Alternative Hypothesis where the rate parameter of a Poisson Distribution λ was distributed according to a Gamma Distribution.

78 The log likelihood of a data set for the Null Hypothesis is described
79 succinctly by Eqn. 3. The Alternative Hypothesis is more complex due
80 to it's multi-step nature. The probability distribution shown in Fig. 5 is
81 described by the following equation.

$$P(X|\alpha, \beta) = \int P(X|\lambda)P(\lambda|\alpha, \beta)d\lambda \quad (6)$$

82 where $P(X|\lambda)$ is the Poisson Distribution, and $P(\lambda|\alpha, \beta)$ is the Gamma Dis-
83 tribution, and all λ values are summed over.

84 The likelihood is once again given by Baye's Theorem shown in Eqn.
85 2. While this likelihood can be calculated analytically, we instead took a
86 numerical approach using the histogram shown in Fig. 5. While useful to
87 visualize the data under the Alternative Hypothesis, the histogram was really
88 plotted in order to save two lists. The first list, n_{save} , contains the height of
89 each bin. The second list, $bins_{save}$, has the location of each bin edge. Once
90 the histogram has been normalized (by ensuring that the sum of all of the bin
91 heights is equal to one), the height of each bin is the probability of measuring
92 that data point (number of events).

93 Overall, the log likelihood function for the complex hypothesis is as fol-
94 lows.

```
95 #Log Likelihood for Complex Hypothesis
96 logprob_min = np.log(1/len(hist_complex))
97
98
99 def ComplexLikelihood(data): #takes in one experiment
100     g = 0 #initialize value
101
102     for d in data: #measurements in an experiment
103         bin_current = 0
104
105         #protect against going over the farthest right bin edge
106         if d > bins_save[len(bins_save) - 1]:
107             logprob = logprob_min
108
109         #find what bin the measurement falls in, starting from
110         the left
111         else:
112             while d > bins_save[bin_current + 1]:
113                 bin_current += 1
114
115             #protect against bins w/ no counts
116             if n_save[bin_current] <= 0:
117                 logprob = logprob_min
118
119             else:
120                 logprob = np.log(n_save[bin_current]) #once
121                 normalized, height of the histogram =
122                 probability of measurement
123
124         g += logprob
125
126     return g
127
```

128 Generating the Log Likelihood Ratio, as shown in Eqn. 5, for each Hy-
129 pothesis requires that you compute the likelihood of each hypothesis for each
130 data set. The code is shown below.

```
131 #Log Likelihood Ratios
132
```

```

133 LL_simple_simple = [] #Log Likelihood of the Simple Hypothesis
134     using the data made from the Simple Hypothesis
135 LL_complex_simple = [] #Log Likelihood of the Complex
136     Hypothesis using the data made from the Simple Hypothesis
137
138 LL_simple_complex = [] #Log Likelihood of the Simple Hypothesis
139     using the data made from the Complex Hypothesis
140 LL_complex_complex = [] #Log Likelihood of the Complex
141     Hypothesis using the data made from the Complex Hypothesis
142
143 LLR_simple = []
144 LLR_complex = []
145
146 #Log Likelihood of the Simple Hypothesis using the data made
147     from the Simple Hypothesis
148 for exp in data_simple: #each experiment
149     LL = LogLikelihoodPoisson(lamb, exp)
150     LL_simple_simple.append(LL)
151
152 #Log Likelihood of the Complex Hypothesis using the data made
153     from the Simple Hypothesis
154 for exp in data_simple: #each experiment
155     LL = ComplexLikelihood(exp)
156     LL_complex_simple.append(LL)
157
158 #Log Likelihood of the Simple Hypothesis using the data made
159     from the Complex Hypothesis
160 for exp in data_complex: #each experiment
161     LL = LogLikelihoodPoisson(lamb, exp)
162     LL_simple_complex.append(LL)
163
164 #Log Likelihood of the Complex Hypothesis using the data made
165     from the Complex Hypothesis
166 for exp in data_complex: #each experiment
167     LL = ComplexLikelihood(exp)
168     LL_complex_complex.append(LL)
169
170 #LLR for simple hypothesis
171 for i in range(len(LL_simple_simple)):
172     LLR_simple.append(LL_simple_simple[i] -

```



```

173         LL_complex_simple[i])
174
175     #LLR for complex hypothesis
176     for i in range(len(LL_simple_complex)):
177         LLR_complex.append(LL_simple_complex[i] -
178                             LL_complex_complex[i])
179

```

180 The last portion of the code sorted the LLR lists in ascending order and
181 then found the power of the test by finding the critical LLR value (determined
182 by the confidence level) in both lists.

```

183
184     #Sort the LLRs
185     LLR_simple = Sorter.DefaultSort(LLR_simple)
186     LLR_complex = Sorter.DefaultSort(LLR_complex)
187
188     #Define Confidence Level and find critical LLR value
189     alpha_CL = 0.05 #Confidence Level = 95%
190
191     LLR_alpha_CL = LLR_simple[math.ceil(Nexp * alpha_CL)] #critical
192                     LLR value
193
194     #Find Beta and Power of Test
195     for i in range(len(LLR_complex)):
196         if LLR_complex[i] >= LLR_alpha_CL:
197             LLR_Beta = LLR_complex[i]
198             LLR_Beta_Position = i
199             break
200
201     Beta_CL = (len(LLR_complex) - LLR_Beta_Position)/Nexp #Beta_CL
202               = percent of entries in LLR_complex above LLR_alpha_CL
203
204     Power = 1 - Beta_CL
205

```

206 4 Results

207 Results are summarized in Fig. 6. The power of the test to discern
208 between the two hypothesis was found to be 39%, which is quite low for
209 a number of reasons. The first is that the number of measurements per

210 experiment, N_{meas} was only 1,000, which is quite low. Using a higher value
 211 for N_{meas} would improve the power of the test. The second reason is that the
 212 two hypothesis had the same Most Likely Estimate (MLE) for λ . Changing
 213 the parameters of the Alternative Hypothesis, α and β to create a different
 214 MLE for the value of λ would make the two hypotheses more distinguishable.

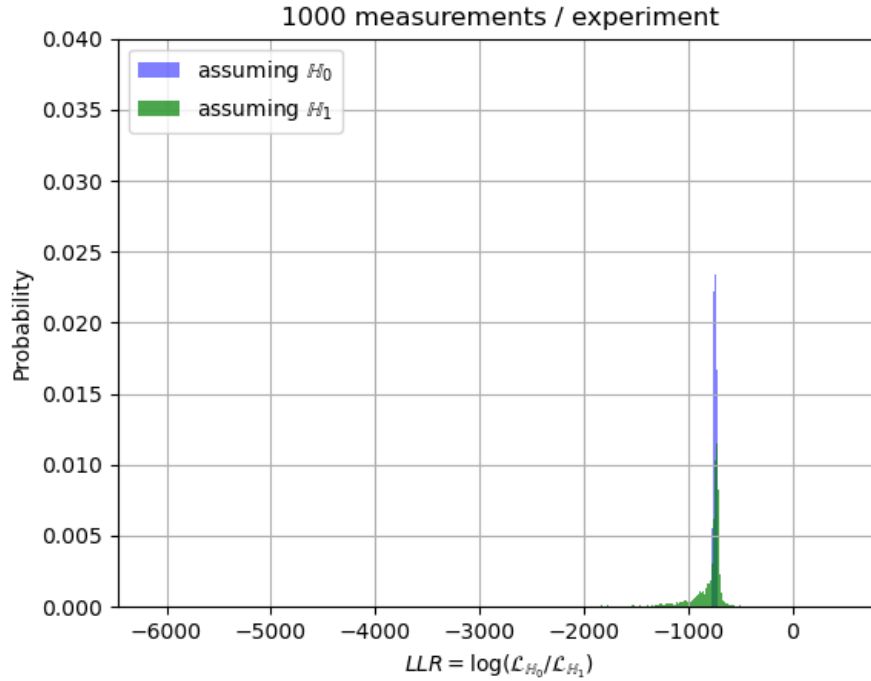


Figure 6: A graph of the log likelihood ratio for the two Hypotheses. The Null Hypothesis (blue) assumes a constant rate parameter λ , while the Alternative Hypothesis assumes that the rate parameter λ is distributed according to a Gamma Distribution. Parameter values were chosen such that both hypotheses had the same MLE for the value λ , which makes it hard to differentiate between the two hypotheses.

215 5 Summary

216 This paper described the process for how to use the Log Likelihood Ratio
217 to distinguish between two hypotheses. The Null Hypothesis was quite sim-
218 ple, assuming a constant rate parameter of $\lambda = 5$ for the Poisson Distribution.
219 The Alternative Hypothesis utilized a new rate parameter for each experi-
220 ment, being randomly generated from a Gamma Distribution with $\alpha = 6.0$,
221 and $\beta = 1.2$. Due to the similarity of the two hypotheses, the power of the
222 test was quite low, found to be 39.4%.

223 Improving the power of test should be done by simulating more measure-
224 ments per experiment (which my computer could not handle for the given
225 number of experiments), and by choosing values of α and β such that the
226 Maximum Likelihood Estimate of the most-likely rate parameter λ is distinct
227 between the two hypotheses.

228 6 Hypothesis.py

229 The code file **Hypothesis.py** was used to do the generation of the data
230 and the analysis of the data for this simulation. Since most of the analysis
231 was shown in Section 3, only the data generation will be shown here. The
232 entire code file can be viewed in the GitHub Repository.

```
233  
234 #import packages  
235 import math  
236 import numpy as np  
237 import matplotlib.pyplot as plt  
238 import sys  
239  
240 #import Random class  
241 sys.path.append(".")  
242 import Random as rng  
243  
244 #import Sorting class  
245 sys.path.append(".")  
246 import MySort as mys  
247  
248 # main function  
249 if __name__ == "__main__":  
250     # if the user includes the flag -h or --help print the options
```

```

251 if '-h' in sys.argv or '--help' in sys.argv:
252     print ("Usage: %s [options]" % sys.argv[0])
253     print (" options:")
254     print (" --help(-h)          print options")
255     print (" -seed [integer number] seed")
256     print (" -Nmeas [integer number] number of measurements
257           per experiments")
258     print (" -Nexp [integer number] number of experiments")
259     print (" -lambda [float number] Parameter for the Null
260           Hypothesis")
261     print (" -alpha [float number] Parameter for the
262           Alternative Hypothesis")
263     print (" -beta [float number] Parameter for the
264           Alternative Hypothesis")
265     print
266     sys.exit(1)
267
268 #Initialize
269 seed = 5555
270
271 Nmeas = 1 #number of measurements per experiment
272 Nexp = 1 #number of experiments
273
274 lamb = 1.0 #must be positive
275 alpha = 1.0 #must be positive
276 beta = 1.0 #must be positive
277
278 #System Inputs
279 if '-seed' in sys.argv:
280     p = sys.argv.index('-seed')
281     seed = sys.argv[p+1]
282
283 if '-Nmeas' in sys.argv:
284     p = sys.argv.index('-Nmeas')
285     ptemp = int(sys.argv[p+1])
286     Nmeas = ptemp
287
288 if '-Nexp' in sys.argv:
289     p = sys.argv.index('-Nexp')
290     ptemp = int(sys.argv[p+1])

```

```

291     Nexp = ptemp
292
293     if '-lambda' in sys.argv:
294         p = sys.argv.index('-lambda')
295         ptemp = float(sys.argv[p+1])
296         lamb = ptemp
297
298     if '-alpha' in sys.argv:
299         p = sys.argv.index('-alpha')
300         ptemp = float(sys.argv[p+1])
301         alpha = ptemp
302
303     if '-beta' in sys.argv:
304         p = sys.argv.index('-beta')
305         ptemp = float(sys.argv[p+1])
306         beta = ptemp
307
308     #class instance of Random and Sorting class
309     random = rng.Random(seed)
310     Sorter = mys.MySort()
311
312     #initialize data
313     data_simple = [] #[[exp1], [exp2], ...] each experiment in the
314                     simple hypothesis
315     data_graph = [] #[meas1, meas1, ...] used to plot data from
316                     simple hypothesis
317
318     data_complex = [] #[[exp1], [exp2], ...] each experiment in the
319                      complex hypothesis
320     hist_complex = [] #[meas1, meas1, ...] every measurement in the
321                      complex hypothesis from all lambdas
322     lamb_graph = [] #[lamb1, lamb2, ...] used to plot the
323                     distribution of lambdas
324
325     #Generate Data for Simple Hypothesis (fixed lambda)
326     for e in range(0, Nexp): #each experiment
327         data_exp_simple = [] #all measurements in a given experiment
328
329         for m in range(0, Nmeas): #each measurement
330             measurement_simple = float(random.Poisson(lamb))

```

```

331         data_exp_simple.append(measurement_simple)
332         data_graph.append(measurement_simple)
333
334     data_simple.append(data_exp_simple)
335
336     #Generate Data for Complex Hypothesis (lambda comes from a
337     #Gamma Distribution)
338     for e in range(0, Nexp): #each experiment
339         data_exp_complex = [] #all measurements in a given
340         experiment
341
342         lamb_complex = float(random.Gamma(alpha, beta)) #new lambda
343         every experiment
344
345         lamb_graph.append(lamb_complex)
346
347         for m in range(0, Nmeas): #each measurement
348             measurement_complex = float(random.Poisson(lamb_complex))
349             hist_complex.append(measurement_complex)
350             data_exp_complex.append(measurement_complex)
351
352         data_complex.append(data_exp_complex)
353

```
