

# Comparing Effective Estimators of Large Wildfires in California Using the Pareto Distribution

Grace Hanson, Simon Hempel-Costello, and Ammy Lin

March 15, 2024

## 1 Introduction

The frequency of wildfires has increased in the past few years. For example, fire alerts around the world increased by 13% in 2020 compared to 2019, which was already a record-setting year; this is seen as an effect of climate change [PA]. The cost of wildfires extends to issues in public health, wildlife management, and the economy, among others, and their impact is projected to continue growing [PA].

In response, there is a growing focus on improving fire management strategies. Benefits to this include saving lives, minimizing risks to the ecosystem, and reducing economic losses [PA]; the US Forest Service writes that such strategies allow for more efficient “operational management of and strategic planning for wildland firefighting assets, such as firefighters, aircraft, and engines” [Agr]. As such, we are interested in using data to predict the occurrence of large wildfires in California.

Our data includes information about the area of spread and frequency of fires from 2017 to 2022 in California; the minimum area of spread is 50,000 acres [FP]. We are interested in finding a method to estimate the number of fires that have a large area of spread (greater than or equal to 300,000 acres). However, we note that the number of fires that meet this criterion in our dataset is small.

Given this, we note that an empirical average may not provide an accurate estimate of the true expected number of fires covering more than or equal to 300,000 acres in California during a given year. We introduce a parametric model, which may be more accurate when working with the tails of a distribution. Using the Pareto model, we explore whether this estimated distribution fits the data well and then estimate the expected number of fires that have an area distribution greater than or equal to 300,000 acres.

## 2 Methods

Considering the impact of outlying data, we use a Pareto distribution to model our data for the area of spread and frequency of fires in California between 2017 and 2022. In the Pareto model, the shape parameter  $\alpha$  tells us how heavy the tail of the distribution is (i.e. the distribution of the more extreme values); when  $\alpha \leq 1$ , the heavy tail of the distribution causes it not to have a mean. We are interested in finding an estimator for  $\alpha$  and observing how it varies by changing values of  $\alpha$  and the sample size  $n$ . Moreover, we find the 95% confidence interval for  $\alpha$ , which gives us a better understanding of the interval in which the true value of  $\alpha$  may lie.

In our derivations (see Section 5), we found the maximum likelihood estimator (MLE) and method of moments estimator (MOM) of  $\alpha$ , our shape parameter. The maximum likelihood estimator for  $\alpha$  (see Appendix 5.2 for our derivation process) is

$$\hat{\alpha}_{MLE} = \frac{1}{\frac{1}{n} \sum_{i=1}^n \log(Y_i/y_m)}.$$

Likewise, the method of moments estimator of  $\alpha$  (see Appendix 5.6 for our derivation process) is

$$\hat{\alpha}_{MOM} = \frac{\bar{Y}}{\bar{Y} - y_m}.$$

In the context of our problem, these estimators of  $\alpha$  allow us to estimate the value of the shape parameter,  $\alpha$ , which tells us how heavy the tail of the distribution is. We discuss how these estimators compare in the next section (3).

Additionally, we found that the exact 95% confidence interval using the pivotal method for  $\alpha$  (full details in Appendix 5.4) is

$$0.95 = Pr(\hat{\alpha}_{MLE} \times q_{g0.025} < \alpha < \hat{\alpha}_{MLE} \times q_{g0.975}).$$

The approximate 95% confidence interval using the  $t_{n-1}$  distribution (full details in Appendix 5.5) is

$$\hat{\alpha}_{MLE} \pm \frac{\sqrt{n}}{sq_{t0.975}}.$$

A 95% confidence interval tells us the range of values for which we can be confident that the true value of  $\alpha$  lies. The exact method is used when more precision is necessary, while the approximate method is typically more efficient because it makes assumptions or simplifications at the cost of precision. Choosing between the two methods depends on factors such as the sample size and other characteristics of the data. We present and discuss which method is closer to the nominal 95% coverage in Section 3.

To explore the properties of our point estimates and confidence intervals for  $\alpha$ , we looked at how they were influenced by different values of  $\alpha$  and sample size  $n$ . Expanding on this, the values of  $\alpha$  that we chose were 0.5, 1.5, and 2.5. We chose 0.5 in order to see how the point estimates and confidence intervals are affected by values of  $\alpha \leq 1$  (since the heavy tail would affect the existence of a mean), and the values 1.5 and 2.5 to note how an increase in  $\alpha$  may impact our conclusions. Likewise, we simulated with sample sizes  $n = 15, 50, 100$ , and 1000 in order to understand how increasing the sample size may change the outcome of the simulation, especially with the approximate confidence interval method (which relies on the Central Limit Theorem).

The code for our simulation can be found in Section 6. Using the results from our derivations, we first sampled  $n$  i.i.d. Pareto random variables, then calculated the MLE and MOM estimators. Following this, we calculated the upper and lower bounds of the exact 95% confidence interval, and likewise for the approximate 95% confidence interval.

To decide which of these estimators and confidence intervals are the most accurate estimates of our data, we simulated sampling with replacement from a Pareto distribution with several different “true” values of  $\alpha$  and several different sample sizes. We did this by **bootstrapping** values from our generated Pareto distribution data and calculating the estimators and confidence intervals for each, then generated a collection of point estimators and confidence intervals that we used to examine several values that tell us more about which estimator is most appropriate.

For example, we calculated the **root mean squared errors (RMSEs)** for the MLE and MOM estimators. The RMSE allows us to understand and measure the average magnitude of error between the estimated and true value of  $\alpha$ , our parameter of interest. This tells us how far we would expect each estimator to fall from the true value of  $\alpha$ .

We also computed the **coverage rates** for our exact 95% confidence interval, approximate 95% confidence interval, and bootstrap 95% confidence intervals for the MLE and MOM estimators. Since we are looking at 95% confidence intervals, the ideal coverage rate should be as close to 0.95 as possible, meaning that the interval estimate is expected to contain the true value of  $\alpha$  approximately 95% of the time. In other words, it tells us the true likelihood that each confidence interval contains our true value of  $\alpha$ .

## 2.1 Additional Exploration

Another trait that we will explore about the MLE and MOM estimators of  $\alpha$  is the **asymptotic efficiency** of each estimator. MLEs are asymptotically efficient, which means that they become more and more accurate as sample sizes increase. On the other hand, MOM estimators are not guaranteed to be asymptotically efficient and depend on certain conditions in order to be classified as such. We want to understand, in our context, if the MOM estimator is asymptotically efficient, and if so, whether the MLE or MOM estimator is *more* asymptotically efficient (i.e. which becomes more accurate more quickly). This could help us infer which estimator would be better for a large sample size without having to test them, as the estimator that is more asymptotically effective would tend to be better for very large samples.

To test this, we will generate 1000 95% confidence intervals for each estimator at a single value of  $\alpha$  ( $\alpha = 1.5$ ) and sample sizes of  $n = 10, 15, 20, 25, 50, 100, 200, 1,000$ , and  $10,000$ . Through graphing this, we will be able to observe which interval grows smaller more quickly.

### 3 Results

#### 3.1 Coverage rates for the exact and approximate 95% confidence intervals

The results of our simulated study on the coverage rates of the exact and approximate 95% confidence intervals are displayed in Figure 1. The charts are divided by  $\alpha$  values 0.5, 1.5, and 2.5.

The coverage rates for the exact confidence interval method appear to stay around 0.95 regardless of the values of  $\alpha$  and  $n$ . This suggests that there is a high level of confidence in this method and interval estimate, as, with repeated sampling, it is expected to contain the true value of  $\alpha$  approximately 95% of the time.

For the approximate 95% confidence interval using the  $t_{n-1}$  distribution, we see that as  $\alpha$  increases, the coverage rate for the approximate confidence interval method decreases. For example, for  $\alpha = 0.5$  and sample size  $n = 15$ , the coverage rate is 0.852. Keeping  $n = 15$  and looking at  $\alpha = 1.5$  and  $2.5$ , we obtain the coverage rates 0.138 and 0.055, respectively. Likewise, keeping  $\alpha$  constant and increasing  $n$ , we observe that the coverage rate for the approximate confidence interval method increases: setting  $\alpha = 0.5$  and observing the coverage rate when  $n = 15$  and  $n = 1000$ , we notice that the value is approximately 0.1 greater (from 0.852 to 0.953) from  $n = 15$  to  $n = 1000$ . This observation is supported by the Central Limit Theorem, which works more effectively with a greater sample size.

In nearly every instance, the exact confidence interval method had a higher coverage rate than the approximate, suggesting that it is the best choice for our calculations. However, for values of  $\alpha$  equal to or below 1 and large sample sizes, the approximate method would be a fine choice as well.

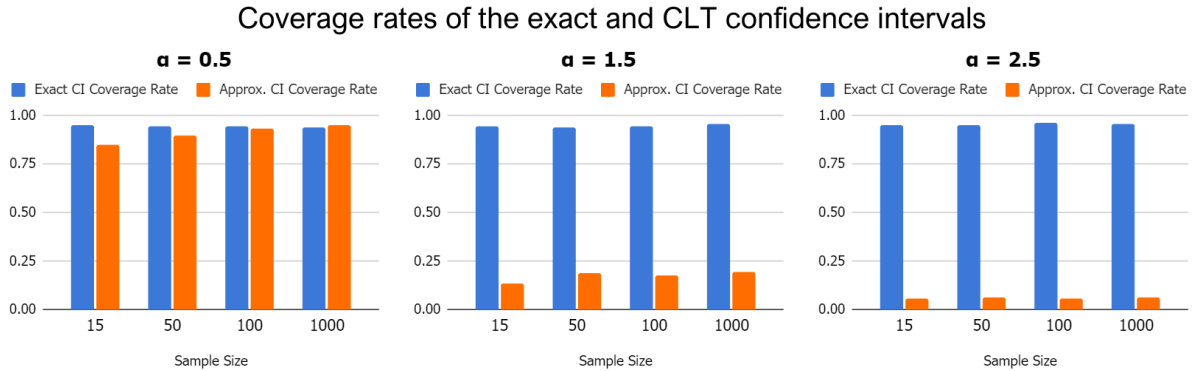


Figure 1: Coverage rates for the exact and approximate 95% confidence intervals, varying  $\alpha$ .

#### 3.2 Coverage rates for the bootstrap MLE and MOM 95% confidence intervals

Our simulation outcomes on the coverage rates for the exact and approximate 95% confidence intervals are displayed in Figure 2.

For the bootstrap 95% confidence interval applied to the maximum likelihood estimator, we note that increasing  $n$  increases the coverage rate; increasing  $\alpha$  similarly increases the coverage rate. The MLE is asymptotically unbiased, meaning that as the sample size increases, the average value of the estimator gets closer and closer to the actual value. As such, the MLE's asymptotic properties can be used to explain the greater coverage rates as  $n$  increases.

On the other hand, the bootstrap 95% confidence interval applied to the method of moments estimator shows that when  $\alpha = 0.5$ , we obtain a coverage rate of 0—this can be explained by the fact that when  $\alpha \leq 1$ ,

the distribution is extremely heavy-tailed (to the point where the mean does not exist), which decreases the certainty of the interval estimate. This is also reiterated in our derivations: we found the method of moments estimator using an expected value of  $Y$  that was only for values of  $\alpha$  greater than 1. As  $\alpha$  increases, the coverage increases from 0.547 for  $\alpha = 0.5, n = 15$  to 0.808 for  $\alpha = 2.5, n = 15$ .

Despite the MOM estimator working better for higher values of  $\alpha$ , the MLE confidence intervals have coverage rates much closer to 95% in every scenario, suggesting that the MLE generates a value close to the true value more often.

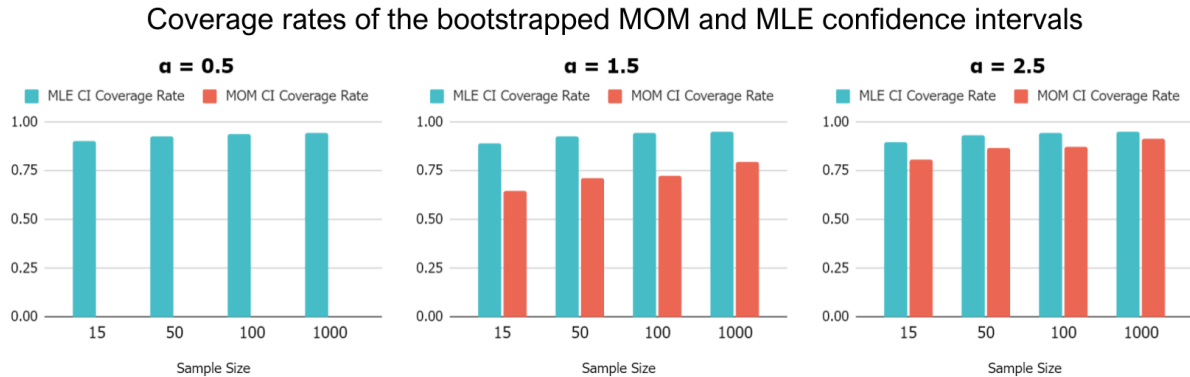


Figure 2: Coverage rates for the bootstrap MOM and MLE 95% confidence intervals, varying  $\alpha$ .

### 3.3 Root mean squared error for the MLE and MOM estimators

We also calculated the root mean squared errors for the MLE and MOM estimators; the results are shown in Figure 3.

We found that as  $\alpha$  increases and the sample size is kept constant, the RMSE for the MLE increases. Additionally, keeping  $\alpha$  constant and increasing  $n$ , we note that the RMSE decreases; again, this can be explained by the asymptotic characteristics of the MLE.

Likewise, we notice that the RMSE for the MOM estimator decreases as  $n$  increases, keeping  $\alpha$  constant. We did not make any significant observations while varying values of  $\alpha$ , keeping  $n$  constant, and observing the RMSE.

This measure shows that the MLE tends to generate values closer to the true value, just as we saw in the coverage rates of bootstrapped confidence intervals.

### 3.4 Asymptotic efficiency of the MLE and MOM estimators

We found that when we calculated and graphed the confidence intervals of each estimator, both estimators were asymptotically efficient and their intervals became narrower as the sample size grew. However, the MLE was more asymptotically efficient and grew smaller at a faster rate. Specifically, the MLE became narrower than the MOM estimator at around a sample size of  $n = 200$ , and was much narrower by  $n = 10,000$ .

Notably, the MLE begins much wider than the MOM estimator, which could suggest that the MOM estimator is more precise than the MLE at smaller sample sizes. However, because we know from Figure 2 that the coverage of this confidence interval is not truly 95%, we would need to account for this lack of coverage to make any strong claims about the MOM being more precise.

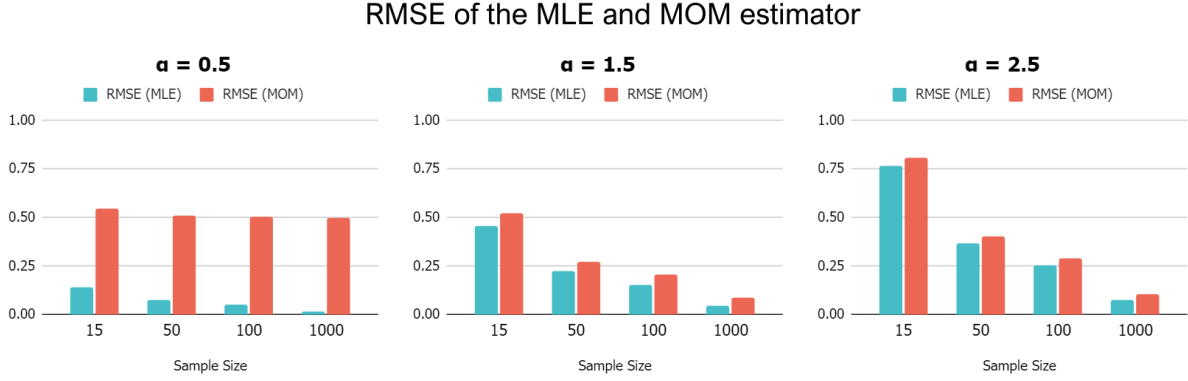


Figure 3: Root mean squared error of the MLE and MOM estimators, varying  $\alpha$ .

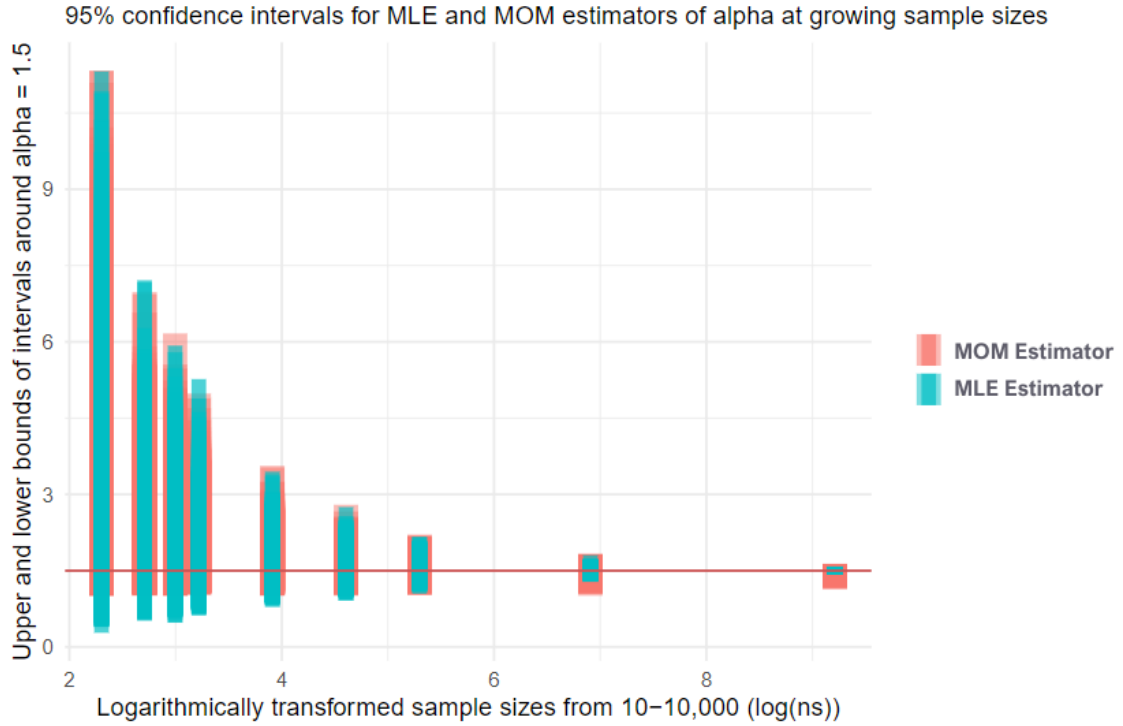


Figure 4: Effect of varying the sample size on the bootstrap MOM and MLE 95% confidence intervals.

### 3.5 Determining the best choice for a confidence interval method and estimator

To determine the best choice for a confidence interval method, we took into consideration the values of our data. We do not know the value of  $\alpha$ , but the sample size of our dataset is 45, so it might be helpful to examine coverage rates when  $\alpha = 1.5$  and  $n = 50$ . We see in Figure 5 that the coverage rate is the highest for the exact confidence interval method at 94.1%, followed by the bootstrap MLE confidence interval at 92.6%.

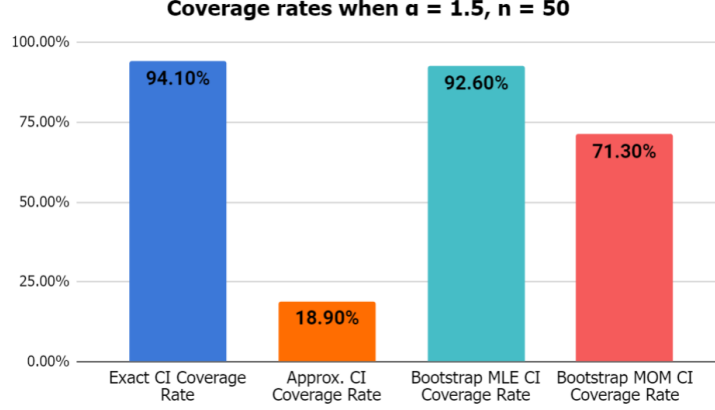


Figure 5: Coverage rates when  $\alpha = 1.5, n = 50$ .

To that end, we determined that the best confidence interval method is the **exact confidence interval**.

Finding the best choice for an estimator involved considering the RMSEs of each. From Figure 3, we know that when the sample size is 50, the MLE has a lower RMSE for all values of  $\alpha$ . Therefore, we can determine that the most accurate point estimator method is the **MLE**.

### 3.6 Applying our estimates to find the expected number of large wildfires

Through our previous work, we determined that the optimal point estimator was the MLE, while the exact confidence interval method was the most accurate. From here, our goal is to use these findings to produce an estimate of the rate of large wildfires that are observed in California during a given year. We define large wildfires to be those that cover more than or equal to 300,000 acres in a given year. Given that these events are relatively rare, we will attempt to estimate their rate by applying our previous work to the population of all fires covering at least 50,000 acres between 2017 and 2022.

Toward that end, we argue that the expected value can be given by

$$E[\text{area} \geq 300,000] = E[\text{area} \geq 50,000] * (1 - F(300,000)),$$

where  $F$  is the CDF of the Pareto distribution. Through simply averaging, we argue that  $E[\text{area} \geq 300,000] = 9$ , and we know that the *CDF* of the Pareto distribution is given by  $(\frac{50,000}{300,000})^\alpha$ . This results in our estimation of fires being

$$E[\text{area} \geq 300,000] = 9 \times (\frac{1}{6})^{\hat{\alpha}_{MLE}}.$$

We then use the exact confidence interval formula to find our confidence interval, which results in the final estimate of

$$\hat{\alpha}_{MLE} = 1.08 \text{ with a 95\% confidence interval of } (0.788, 1.41).$$

This estimate for  $\hat{\alpha}_{MLE}$ , applied to our expected number of fires formula, gives us the final result of

$$E[\text{area} \geq 300,000] = 1.30 \text{ with a 95\% confidence interval of } (0.709, 2.19).$$

Thus, the expected number of fires with an area greater than or equal to 300,000 acres in a given year somewhere in California is 1.3, and the 95% confidence interval is estimated to be between 0.71 and 2.19 fires. That is, we are 95% confident that the true number of fires meeting the stated criteria in a given year will fall within this range.

Confirming this result with our data, we generate a sample of 100,000 values corresponding to our Pareto distribution and get quantiles for this sample corresponding to each of the data points from our population, giving us Figure 6.

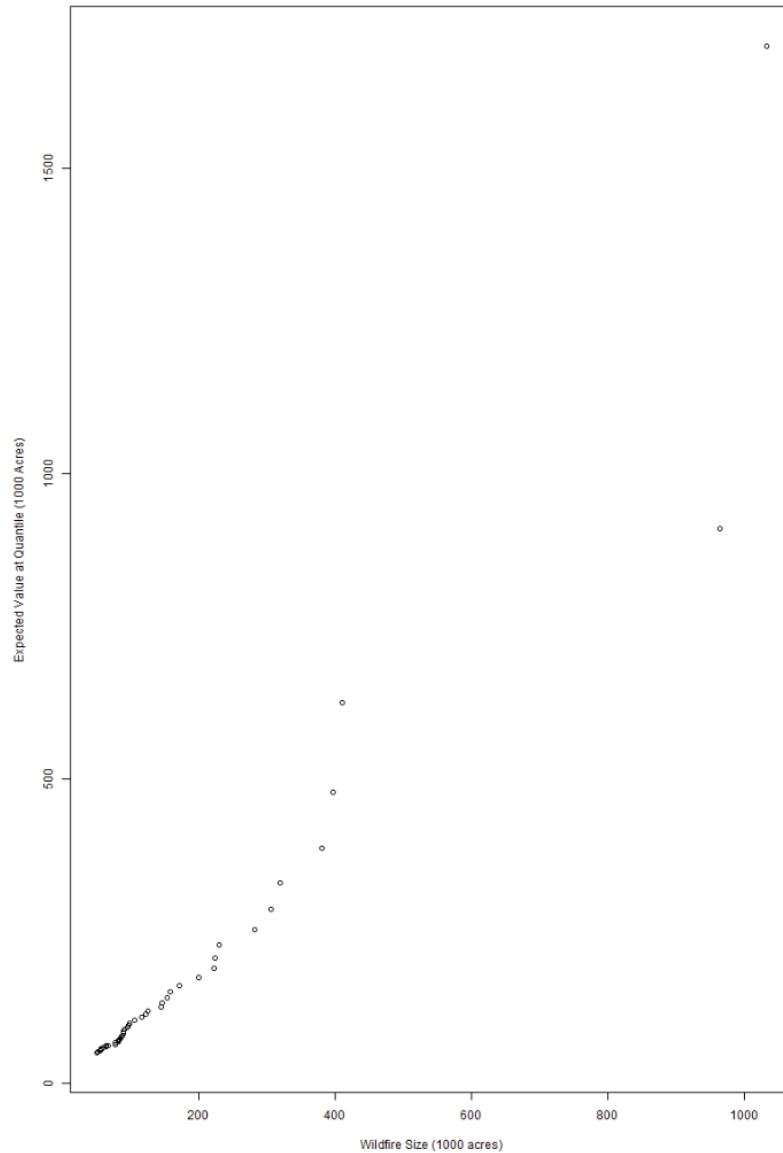


Figure 6: QQ plot between the representative quantiles of the expected Pareto distribution for our data, and the actual quantiles from the population of large fires.

As we can see, the distribution provides a strong description of the rate of relatively small fires, giving a nice linear correlation in the bottom left. As the scale of the fires increases, and the number of samples decreases, the linear relationship starts to fail, with an almost exponential curve occurring beyond 400,000 acres. This reflects a strong relationship between our distribution and the data up to the largest wildfires that California experiences.

## 4 Discussion

### 4.1 Summary and Interpretation of Findings

Our study was interested in exploring how well the estimated Pareto distribution fits the data of California fires between 2017-2022 with an area of at least 50,000 acres and estimating the expected number of large fires ( $\geq 300,000$  acres).

Through our simulation, we determined, by looking at coverage rates, that the exact confidence interval method was the best since its coverage rate was closest to 95% across varying sample sizes  $n$  and values of our shape parameter  $\alpha$ . Figure 5 displays a comparison of the four confidence interval methods that we considered when  $\alpha = 1.5, n = 50$ . Furthermore, we concluded that the best choice of an estimator of  $\alpha$  was the MLE because it had a lower RMSE for all values of  $\alpha$  in our simulation compared to the MOM estimator. Figure 3 showcases our findings.

By using the MLE and the exact confidence interval to estimate our data, we found that we are 95% certain that the number of large fires we would expect to happen in California in a given year falls between **0.709** and **2.19**.

### 4.2 Limitations and Future Directions

Given the scope and circumstances of our study, we acknowledge the limitations that may affect the ability for our results to be generalized and applied.

First, we limited our sample sizes for the simulation to  $n \leq 1000$  due to time constraints with simulating with a greater sample size. Although we did explore the asymptotic efficiency of the MLE and MOM estimators with a sample size of  $n = 10,000$ , we did not go into detail with calculating their coverage rate and RMSE at that size, for example. Similarly, we ran our simulation with three  $\alpha$  values: 0.5, 1.5, and 2.5. The scope of our analysis can be extended with a greater variety of  $\alpha$  values.

While implementing our bootstrap procedure, we were limited to 1000 bootstrap iterations per simulation, which may have impacted our results.

Likewise, we could have experimented with other confidence interval methods and estimators, such as the bootstrap percentile interval. Additionally, we could have implemented a more rigorous analysis and comparison of our estimators by looking at how they are affected by bias, for example.

To extend this project, it would be useful to examine and attempt to estimate other data we have for California wildfires. For example, we may be interested in what time of the year we would expect large wildfires to occur, or the part of the state where they are most frequent.

Additionally, our data would be more impactful if it were grounded in context about the increase or decrease in fires of this size. A natural extension of this project would be a comparison of our data to California fire data from five-year periods in the past and into the future, examining whether the number of expected large fires has increased.

## 5 Appendix

We are given that our data  $Y_1, \dots, Y_n$  are an i.i.d. sample from the  $\text{Pareto}(y_m, \alpha)$  distribution, with PDF

$$f(y) = \frac{\alpha y_m^\alpha}{y^{\alpha+1}}, y > y_m$$

and CDF

$$F(y) = 1 - \left(\frac{y_m}{y}\right)^\alpha, y > y_m.$$

The mean of this distribution is  $E[Y] = \frac{\alpha y_m}{\alpha - 1}$  if  $\alpha > 1$ . If  $0 < \alpha \leq 1$ , then  $E[Y] = \infty$ .

### 5.1 Showing that if $Y \sim \text{Pareto}(y_m, \alpha)$ then $\log(Y/y_m) \sim \text{Exp}(\alpha)$ .

Starting with the fact that  $Y \sim \text{Pareto}(y_m, \alpha)$ , we know that

$$\Pr(Y < k) = 1 - \left(\frac{y_m}{k}\right)^\alpha, k > y_m.$$



Then, if we let  $X = \log(\frac{Y}{y_m})$ , we can restructure the expression to be

$$Pr(X < k) = Pr(\log(\frac{Y}{y_m}) < k) = Pr(\frac{Y}{y_m} < e^k) = Pr(Y < e^k y_m).$$

Plugging this into  $Pr(Y < k) = 1 - (\frac{y_m}{k})^\alpha, k > y_m$  gives us:

$$Pr(Y < e^k y_m) = 1 - (\frac{y_m}{e^k y_m})^\alpha, \{e^k y_m < y_m\} = 1 - e^{-\alpha k}, e^k > 1.$$

Simplifying the support gives us our final distribution:

$$Pr(\log(\frac{Y}{y_m}) < k) = 1 - e^{-\alpha k}, k > 0.$$

This is the CDF of an exponential distribution with rate  $\alpha$ . Thus if  $Y \sim \text{Pareto}(y_m, \alpha)$ , then

$$\log(y/y_m) \sim \text{Exp}(\alpha).$$

## 5.2 Showing that the maximum likelihood estimator of $\alpha$ is $\hat{\alpha}_{\text{MLE}} = \frac{1}{\frac{1}{n} \sum_{i=1}^n \log(Y_i/y_m)}$ .

Since  $\log(Y/y_m) \sim \text{Exp}(\alpha)$ , we can estimate  $\alpha$  by first sending  $\log(Y_i/y_m) = X_i$ . We then set up a likelihood function to estimate  $\alpha$ :

$$L(\alpha) = \prod_{i=0}^n \alpha e^{-\alpha X_i}.$$

Taking the log of this gives us

$$\log(L(\alpha)) = \sum_{i=0}^n \log(\alpha) - \alpha X_i,$$

which we then take the derivative of:

$$\frac{\partial \log(L(\alpha))}{\partial \alpha} = \sum_{i=0}^n \frac{1}{\alpha} - X_i.$$

This can be simplified to

$$\frac{\partial \log(L(\alpha))}{\partial \alpha} = \frac{n}{\alpha} - \sum_{i=0}^n X_i.$$

Taking the extreme condition that  $\frac{\partial \log(L(\alpha))}{\partial \alpha} = 0$ , we get that

$$\frac{n}{\alpha} = \sum_{i=1}^n X_i,$$

which can be equivalently stated as

$$\alpha = \frac{n}{\sum_{i=1}^n \log(Y_i/y_m)}.$$

Therefore,

$$\hat{\alpha}_{\text{MLE}} = \frac{1}{\frac{1}{n} \sum_{i=1}^n \log(Y_i/y_m)}.$$

### 5.3 Showing that $\frac{1}{n} \sum_{i=1}^n \log \left( \frac{Y_i}{y_m} \right) \sim \text{Gamma}(n, n\alpha)$ .

We showed in 5.1 that  $\log \left( \frac{Y_i}{y_m} \right)$  is an exponential random variable, and from Appendix B.9 in [CH22] we can say that  $\text{Exp}(\alpha)$  is equivalent to  $\text{Gamma}(1, \alpha)$ . We are then left with  $\frac{1}{n} \sum_{i=1}^n X_i$ , where  $X_i \stackrel{iid}{\sim} \text{Gamma}(1, \alpha)$ . Adding two Gamma random variables of the form  $\text{Gamma}(1, \alpha)$  together yields a Gamma distribution of  $\text{Gamma}(1 + 1, \alpha)$ , or  $\text{Gamma}(2, \alpha)$ . Therefore by summing the random variables  $n$  times, we yield the distribution  $\text{Gamma}(n, \alpha)$ . We are left with  $\frac{1}{\hat{\alpha}_{MLE}} \sim \frac{1}{n} \text{Gamma}(n, \alpha)$ . We can use the fact that if  $c$  is some constant, then  $c\text{Gamma}(r, \lambda) = \text{Gamma}(r, \frac{1}{c}\lambda)$  for a Gamma distribution. Therefore,

$$\frac{1}{n} \text{Gamma}(n, \alpha) \sim \text{Gamma}(n, n\alpha).$$

### 5.4 Deriving an exact 95% confidence interval for $\alpha$ using the pivotal method.

First, we must calculate our pivotal statistic.

From 5.3, we know that  $\frac{1}{\hat{\alpha}_{MLE}} \sim \text{Gamma}(n, n\alpha)$ . However, in order to calculate a pivotal statistic, we want the unknown parameter to be in our estimate, not in our distribution. We again use the observation that if  $c$  is some constant, then  $c\text{Gamma}(r, \lambda) = \text{Gamma}(r, \frac{1}{c}\lambda)$  to manipulate our distribution. We can calculate that  $\text{Gamma}(n, n\alpha) = \frac{1}{\alpha} \text{Gamma}(n, n)$ . Therefore we can calculate the following:

$$\frac{1}{\alpha} \text{Gamma}(n, n) \sim \frac{1}{\hat{\alpha}_{MLE}},$$

so

$$\text{Gamma}(n, n) \sim \frac{\alpha}{\hat{\alpha}_{MLE}}.$$

This is our pivotal statistic!

To calculate our confidence interval, we will write  $q_{90.975}$  and  $q_{90.025}$  to represent the 97.5 and 2.5 percentiles of the  $\text{Gamma}(n, n)$  distribution. We then can write that

$$0.95 = \Pr \left( q_{90.025} < \frac{\alpha}{\hat{\alpha}_{MLE}} < q_{90.975} \right)$$

and then calculate our final confidence interval for  $\alpha$ , which is

$$0.95 = \Pr (\hat{\alpha}_{MLE} \times q_{90.025} < \alpha < \hat{\alpha}_{MLE} \times q_{90.975}).$$

### 5.5 Deriving an approximate 95% confidence interval for $\alpha$ using the $t_{n-1}$ distribution.

We have shown in 5.2 that  $\frac{1}{\hat{\alpha}_{MLE}} = \frac{1}{n} \sum_{i=1}^n X_i$ , where  $X_i = \log \left( \frac{Y_i}{y_m} \right)$ .

Likewise, we know that  $\bar{x} \pm qt_{0.975} \frac{s}{\sqrt{n}}$  is an approximate 95% CI for  $E(X_i) = \frac{1}{\alpha}$  by the Central Limit Theorem. Given that, we find that the 95% CI for  $\frac{1}{\alpha}$  is

$$\frac{1}{\hat{\alpha}_{MLE}} \pm qt_{0.975} \frac{s}{\sqrt{n}}.$$

Then, finding the approximate 95% confidence interval for  $\alpha$  gets us

$$\hat{\alpha}_{MLE} \pm \frac{\sqrt{n}}{sqt_{0.975}}.$$

## 5.6 Deriving the method of moments estimator of $\alpha$ , assuming that $\alpha > 1$ .

We are given that  $E[Y] = \frac{\alpha y_m}{\alpha - 1}$  if  $\alpha > 1$ .

We set our sample mean  $\bar{Y} = \frac{\hat{\alpha} y_m}{\hat{\alpha} - 1}$ . Then, solving for  $\hat{\alpha}$ :

$$\begin{aligned}(\hat{\alpha} - 1)(\bar{Y}) &= \hat{\alpha} y_m \\ \bar{Y} \hat{\alpha} - \bar{Y} &= \hat{\alpha} y_m \\ -\bar{Y} &= \hat{\alpha} y_m - \bar{Y} \hat{\alpha} \\ -\bar{Y} &= \hat{\alpha}(y_m - \bar{Y}),\end{aligned}$$

thus

$$\hat{\alpha}_{MOM} = \frac{\bar{Y}}{\bar{Y} - y_m}.$$

## 6 Appendix

Code and PDFs showing results for additional exploration, simulation, and determination of values can be found <https://github.com/simonhc24/2024FinalProjectStat250>.

Files named **final\_project** include the code for our simulations to determine the best estimator and confidence interval, **wildfire\_exploration** includes the qqplot and calculations of our expected number of large wildfires, and **AdditionalExploration** includes calculations and visualization of our exploration of asymptotic effectiveness.

## References

- [CH22] Laura Chihara and Tim Hesterberg. *Mathematical Statistics with Resampling and R*. Wiley, 2022.
- [Agr] US Department of Agriculture. *Fire Forecasting*. URL: <https://www.fs.usda.gov/science-technology/fire/forecasting>. (accessed: 03.09.2024).
- [FP] California Department of Forestry and Fire Protection. *California Fire Parameters (all)*. URL: <https://hub-calfire-forestry.hub.arcgis.com/maps/e3802d2abf8741a187e73a9db49d68fe/about>. (accessed: 03.09.2024).
- [PA] Kshitij Purwar and Saheel Ahmed. *Predicting Wildfires: Techniques and Tools for Monitoring and Mitigating Risk*. URL: <https://blueskyhq.io/blog/predicting-wildfires-techniques-and-tools-for-monitoring-and-mitigating-risk>. (accessed: 03.09.2024).