## Parametric Families

We have seen several families of probability distributions whose behavior is determined by a finite set of adjustable parameters $\theta$ :

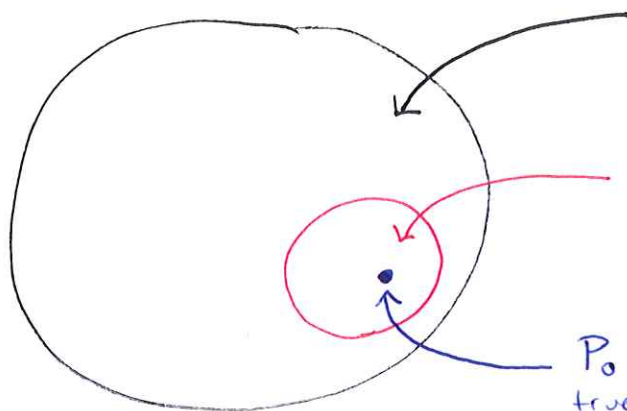| Family | Parameters |
|---|---|
| $X \sim \text{Bern}(p)$ | $\theta = \{p\}$ |
| $X \sim \text{Geometric}(p)$ | $\theta = \{p\}$ |
| $X \sim \text{Poisson}(\lambda)$ | $\theta = \{\lambda\}$ |
| $X \sim N(\mu, \sigma^2)$ | $\theta = \{\mu, \sigma^2\}$ |
| $X \sim \text{Gamma}(\alpha, \beta)$ | $\theta = \{\alpha, \beta\}$ |
| $X \sim t_\nu$ | $\theta = \{\nu\}$ |
| $\vdots$ | $\vdots$ |

## Parametric statistical inference

In statistics we observe data from some unknown probability distribution and use data to try to learn the distribution: $X_1, X_2, \ldots, X_n \overset{iid}{\sim} P_0$.

In __parametric__ statistics, we make the massive simplifying assumption that the unknown distribution $P_0$ belongs to some familiar parametric family. In that case, in order to learn $P_0$, all you have to do ~~is~~ is learn the parameters. Then you're done.

$$\underline{\text{Data}}: X_1, X_2, \ldots, X_n \overset{iid}{\sim} f_\theta \quad \longleftarrow \quad \text{pdf or pmf}$$

$$\underline{\text{Estimator}}: \hat{\theta}_n = \hat{\theta}(X_1, X_2, \ldots, X_n)$$

$\underline{\text{Assumption}}:$



The set of all possible probability distributions that could have generated the data.

The special lil' parametric family you have chosen to restrict yourself to.

$P_0$. We assume that the true data generating distribution belongs to the family. This is probably wrong in reality, but the approximation might be __good enough__.

Example: $X_1, X_2, \ldots, X_n \overset{iid}{\sim}$ Exponential $(\lambda)$.

$\lambda > 0$ is unknown. How do we use the data to estimate it? What is the joint pdf of the data?

$$f(x_1, x_2, \ldots, x_n \mid \lambda) = \overbrace{f_1(x_1) f_2(x_2) f_3(x_3) \cdots f_n(x_n)}^{\text{independence}}$$

independence

identically distributed

$$= f(x_1) f(x_2) f(x_3) \cdots f(x_n)$$

$$= \left(\lambda e^{-\lambda x_1}\right)\left(\lambda e^{-\lambda x_2}\right)\left(\lambda e^{-\lambda x_3}\right) \cdots \left(\lambda e^{-\lambda x_n}\right)$$

$$= \lambda^n e^{-\lambda \sum_{i=1}^{n} x_i}.$$

This is a function of $n$ arguments: $X_1, X_2, \ldots, X_n$. $\lambda$ is treated as a fixed constant. But stare at this formula and flip a switch in your brain:

$$\lambda^n e^{-\lambda \sum_{i=1}^{n} x_i}.$$

Instead of treating it as a function of the $x_i$ with $\lambda$ fixed, think of it as a function of $\lambda$ with the $x_i$ fixed. This is called the **likelihood function**:

$$L(\lambda \mid X_1, \ldots, X_n) = \lambda^n e^{-\lambda \sum_{i=1}^{n} X_i}.$$

It describes how the joint probability mass/density of a fixed dataset varies for different choices of the unknown parameter. For computational convenience, it will help to also define the **log likelihood function**:

$$\ell(\lambda \mid X_1, \ldots, X_n) = \ln L(\lambda \mid X_1, \ldots, X_n)$$

$$= \ln\left(\lambda^n e^{-\lambda \sum_{i=1}^{n} X_i}\right)$$

$$= \ln\left(\lambda^n\right) + \ln\left(e^{-\lambda \sum_{i=1}^{n} X_i}\right)$$
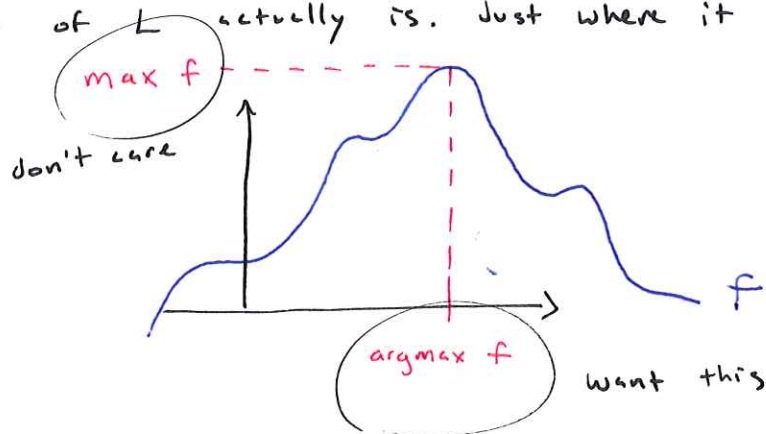
$$= n \ln \lambda + (-\lambda) \sum_{i=1}^{n} X_i.$$

# Maximum likelihood estimation:

To estimate the unknown parameter $\lambda > 0$, we pick the value that makes the likelihood of our observed data as large as possible. This is the <span style="color:red">Maximum likelihood estimator (MLE)</span>:

$$\hat{\lambda}_n^{(MLE)} = \underset{\lambda > 0}{argmax} \; L(\lambda | X_{1:n})$$

$$= \underset{\lambda > 0}{argmax} \; \ln L(\lambda | X_{1:n})$$

$$= \underset{\lambda > 0}{argmax} \left[ n \ln \lambda - \lambda \sum_{i=1}^{n} X_i \right].$$

natural log is an order-preserving function.

"argmax" means "the argument that does the maximizing." We want the location of the maximum value. We don't care what the maximum value of $L$ actually is. Just where it happens.



To find the MLE, we're back in calc I:

$$\frac{d \ln L}{d \lambda} = \frac{d}{d \lambda}\left( n \ln \lambda - \lambda \sum_{i=1}^{n} X_i \right) = \frac{d}{d\lambda}(n \ln \lambda) - \frac{d}{d\lambda}\left( \lambda \sum_{i=1}^{n} X_i \right)$$

$$= \frac{n}{\lambda} - \sum_{i=1}^{n} X_i .$$

$$\frac{d \ln L}{d \lambda} = 0 \implies \frac{n}{\lambda} - \sum_{i=1}^{n} X_i = 0 \implies \frac{n}{\lambda} = \sum_{i=1}^{n} X_i$$

$$\frac{n}{\sum_{i=1}^{n} X_i} = \lambda \quad \left( \begin{array}{l} \text{yada yada} \\ \text{check the} \\ \text{second} \\ \text{derivative} \end{array} \right)$$

$$\hat{\lambda}_n^{(MLE)} = \frac{n}{\sum_{i=1}^{n} \hat{X}_i} = \frac{1}{\bar{X}_n}.$$

Given data, this is our best guess for the unknown parameter $\lambda > 0$. Like all estimators, $\hat{\lambda}_n^{(MLE)}$ is a function of the data, and the data are random. As such, the MLE is a random variable with a <span style="color:red">sampling distribution</span> that describes how it varies across alternative random samples. We itemize the statistical properties of the estimator by itemizing the properties of its sampling distribution: when is it centered, how spread out is it, what happens to it as $n \to \infty$, and so on. For example...

---

Whet is the bias of $\hat{\lambda}_n^{(MLE)}$

---

$$\text{bias}\left(\hat{\lambda}_n^{(MLE)}\right) = E\left(\hat{\lambda}_n^{(MLE)}\right) - \lambda$$

<span style="color:red">true but unknown value</span>

So, what is the mean of the sampling distribution? To compute, note that $X_1, X_2, \ldots, X_n \overset{iid}{\sim} \text{Gamma}(1, \lambda)$, so $\bar{X}_n \sim \text{Gamma}(n, n\lambda)$. To compute the mean of $\hat{\lambda}_n^{(MLE)} = 1/\bar{X}_n$, we apply LOTUS:

$$E\left(\hat{\lambda}_n^{(MLE)}\right) = E\left(1/\bar{X}_n\right) = \int_0^\infty \frac{1}{x} f_{\bar{X}_n}(x)\, dx$$

$$= \int_0^\infty \frac{1}{x} \frac{(n\lambda)^n}{\Gamma(n)} x^{n-1} e^{-n\lambda x}\, dx$$

$$= \frac{n^n \lambda^n}{\Gamma(n)} \int_0^\infty x^{n-1-1} e^{-n\lambda x}\, dx$$

$$= \frac{n^n \lambda^n}{(n-1)!} \frac{\Gamma(n-1)}{(n\lambda)^{n-1}}$$

$$= \frac{n^n \lambda^n}{(n-1)!} \frac{(n-2)!}{n^{n-1} \lambda^{n-1}} = \frac{n\lambda}{n-1}.$$

So $\quad$ bias$\left( \hat{\lambda}_n^{(MLE)} \right) = \dfrac{n\lambda}{n-1} - \lambda$

$$= \frac{n\lambda}{n-1} - \frac{(n-1)\lambda}{n-1}$$

$$= \frac{n\lambda - n\lambda + \lambda}{n-1}$$

$$= \frac{\lambda}{n-1} .$$

Note two things.

- $\quad$ bias$\left( \hat{\lambda}_n^{(MLE)} \right) = \dfrac{\lambda}{n-1} > 0$ , $\quad$ so

  we are always <u>overestimating</u>.

- $\quad$ bias$\left( \hat{\lambda}_n^{(MLE)} \right) = \dfrac{\lambda}{n-1} \longrightarrow 0$ $\quad$ as $\quad$ $n \to \infty$.

So for finite $n$, the MLE is <u>biased</u>, but
in the limit, it is <u>consistent</u>:

$$\hat{\lambda}_n^{(MLE)} \xrightarrow{\text{prob}} \lambda .$$

Let's investigate properties of the MLE via simulation:

$$P_0 = \text{Exponential}\left(\lambda = 2\right)$$

ground truth parameter value

simulate many different datasets of size $n$ from the true distribution

$X_1 \quad X_1 \qquad\qquad X_1 \quad X_1$

$X_2 \quad X_2 \qquad\qquad X_2 \quad X_2$

repeat
$m = 1000$ times

$X_n \quad X_n \qquad\qquad X_n \quad X_n$

each dataset gives a different point estimate

$\hat{\lambda}_n \quad \hat{\lambda}_n \qquad\qquad \hat{\lambda}_n \quad \hat{\lambda}_n$

how do these estimates vary across datasets

visualize the sampling distribution

REPEAT the entire process for different sample sizes $n$.

```r
set.seed(123)   # for reproducibility

# Parameters
lambda_true <- 2          # true rate parameter
n_sim <- 1000             # number of simulations per sample size
sample_sizes <- 2 * 2^(0:6)

# Container for results
mle_estimates <- vector("list", length(sample_sizes))
names(mle_estimates) <- paste0("n = ", sample_sizes)

# Simulate MLEs
for (i in seq_along(sample_sizes)) {
  n <- sample_sizes[i]
  estimates <- numeric(n_sim)
  for (j in 1:n_sim) {
    sample <- rexp(n, rate = lambda_true)
    estimates[j] <- 1 / mean(sample)  # MLE of lambda
  }
  mle_estimates[[i]] <- estimates
}

# Combine into data frame for plotting
mle_data <- stack(mle_estimates)

# Box plot
boxplot(values ~ ind, data = mle_data,
        main = "Sampling Distribution of MLE for Exponential(2)",
        ylab = "MLE",
        xlab = "",
        col = "lightblue",
        las = 2,
        ylim = c(0, 8),
        outline = FALSE)
abline(h = lambda_true, col = "red", lty = 2)  # true value of lambda
```
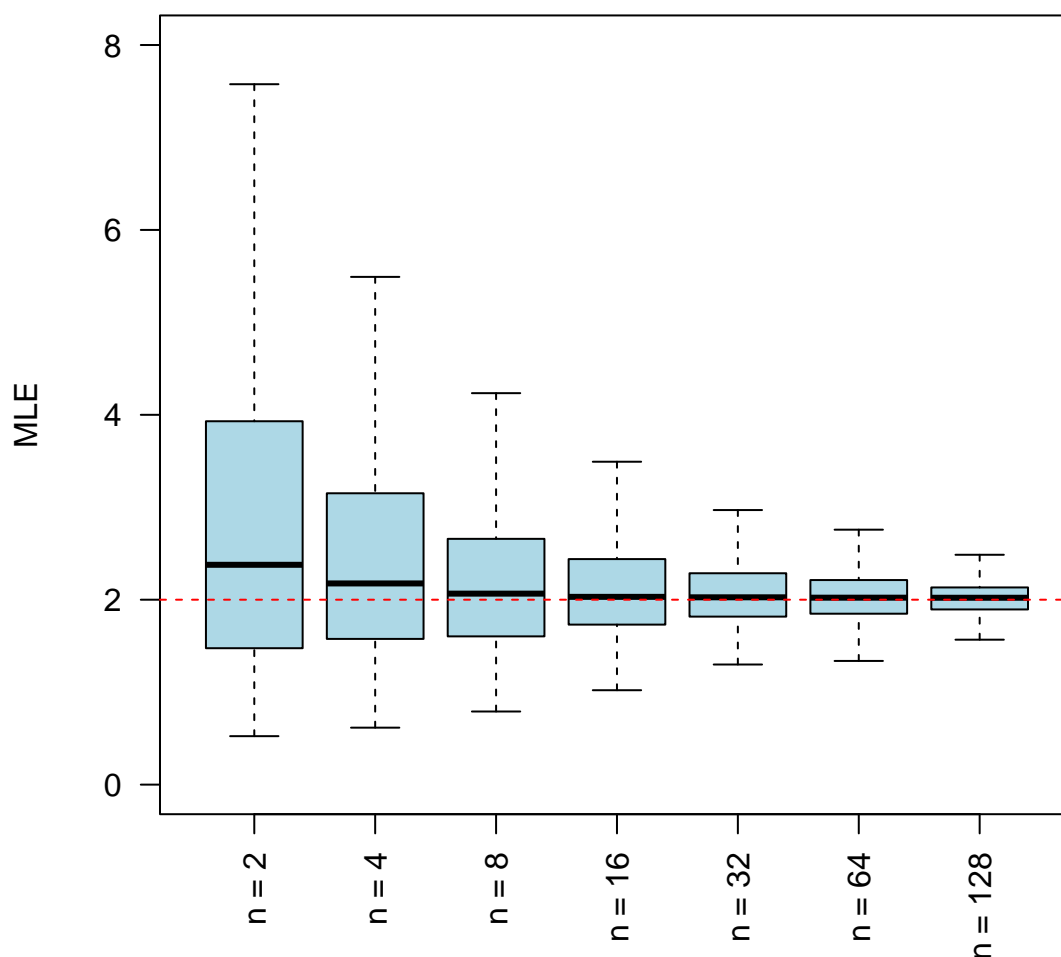
**Sampling Distribution of MLE for Exponential(2)**

Each box plot is displaying the sampling distribution of the MLE for a different sample size (based on $M = 1000$ simulations). For smaller sample sizes, we see that there is some positive bias, but as the sample size increases, the sampling distribution becomes more and more concentrated around the true value of $\lambda = 2$. So via simulation, we demonstrate that this estimator is biased but consistent.

2

## Maximum likelihood in general

**Data:** $X_1, X_2, \ldots X_n \overset{iid}{\sim} f(x \mid \theta)$

some pmf or pdf

true value of unknown parameter(s)

**Likelihood function:**

$$L(\theta \mid X_{1:n}) = f(X_1, X_2, \ldots, X_n \mid \theta)$$

$$= \prod_{i=1}^{n} f(X_i \mid \theta)$$

it's a product because of independence. It's the same $f$ in each factor because of identical distribution.

**Log-likelihood:**

$$\ell(\theta \mid X_{1:n}) = \ln L(\theta \mid X_{1:n})$$

$$= \sum_{i=1}^{n} \ln f(X_i \mid \theta)$$

**MLE:**

$$\hat{\theta}_n^{(MLE)} = \underset{\theta}{\arg\max} \; L(\theta \mid X_{1:n})$$

$$= \underset{\theta}{\arg\max} \; \ln L(\theta \mid X_{1:n})$$

**Some properties:**

consistent:

$$\hat{\theta}_n^{(MLE)} \overset{prob}{\longrightarrow} \theta$$

asymptotic normality:

$$\frac{\hat{\theta}_n^{(MLE)} - \theta}{se(\hat{\theta}_n)} \overset{dist}{\longrightarrow} N(0,1)$$

may or may not be biased:

$$bias\left(\hat{\theta}_n^{(MLE)}\right) = ?$$

may or may not be zero. no guarantees in general.