

# Likelihoods

History of the course:

Taught since 2017 or 2018 by

Haoran } Thanks for notes/help preparing class.  
Dan }

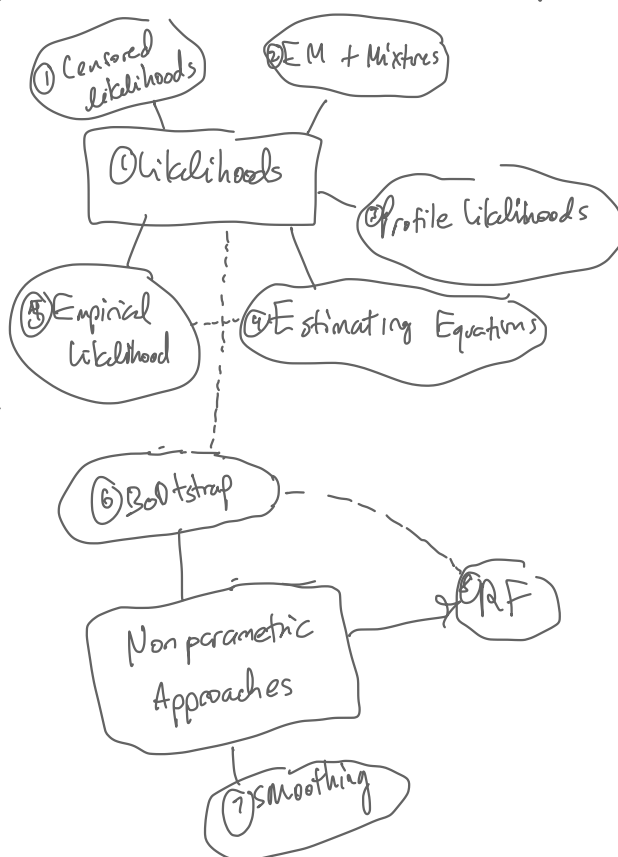
Now: Me

Idea: Many important topics not taught in other courses but should be "core" to a PhD in statistics. Not linear!

## 0.1 Outline

- ① likelihoods
- ② EM, k-means
- ③ Profile likelihood
- ④ Estimating Equations (M estimation).
- ⑤ Empirical likelihoods.
- ⑥ Bootstrap
- ⑦ Smoothing methods
- ⑧ Random Forests (?)

midterm



# 1 Likelihood Construction and Estimation

Likelihood-based Methods

MLE

LRT

likelihood-based uncertainty (CI's).

Why do Statisticians love likelihood-based estimation?

1. Invariance property of MLE: If a distribution is parametrized by  $\theta$  but interest is in a function of  $\theta$ ,  $\psi(\theta)$ , if  $\hat{\theta}$  is MLE of  $\theta$ ,  $\psi(\hat{\theta})$  is MLE of  $\psi(\theta)$ .
2. Asymptotically unbiased; consistent  
$$\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| > \epsilon) = 0$$
3. Asymptotically efficient; variance achieve Cramer-Rao Lower Bound.  
estimator has all the information.
4. Relationship w/ Fisher Information matrix allows for construction of CI's  
(based on asymptotic properties).

Downsides?

1. Very model based! You are assuming know entire distribution
2. It often require numerical optimization.

Still... we ♥ it!

## 1.1 Introduction

**Definition:** Suppose random variables  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  has joint density or probability mass function  $f_{\mathbf{Y}}(\mathbf{y}, \boldsymbol{\theta})$  where  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_b)$ . Then the *likelihood function* is

$$L(\boldsymbol{\theta}|\mathbf{Y}) = f_{\mathbf{Y}}(\mathbf{Y}, \boldsymbol{\theta}). \leftarrow \text{in general, likelihood} = \text{joint}$$

$\uparrow$  likelihood is random.

(because it depends on the data!).

Know MLE is random + we quantify its uncertainty.

Given a vector of observations  $\mathbf{y}$ , the likelihood is a function of  $\boldsymbol{\theta}$ :  
For any (valid) value of  $\boldsymbol{\theta}$ , it returns a number (the likelihood).

$\hat{\boldsymbol{\theta}}_{MLE}$  obtained by finding value of  $\boldsymbol{\theta}$  which yields max likelihood value.

**Key concept:** In all situations, the likelihood is the joint density of the observed data to be analyzed.

Comments:

(1) "density" can mean continuous density or pmf.

(2) "observed data" will be generalized. E.g, censored data.

### 1.1.1 Notation

Given  $\mathbf{y}$ , note that  $L(\boldsymbol{\theta}|\mathbf{y}) : \mathbb{R}^b \rightarrow \mathbb{R}$ .

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_b)^\top$$

Likelihood  $L(\boldsymbol{\theta}|\mathbf{y})$  is scalar valued!

Generally, we optimize  $\ell(\boldsymbol{\theta}) = \log L(\boldsymbol{\theta}|\mathbf{y})$ .

$\underbrace{\hspace{10em}}$   
monotone increasing

$$\arg \max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}|\mathbf{y}) = \arg \max_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}).$$

How? Take derivatives, set to zero, solve.

Generally convention is the derivative of a function (i.e.  $\ell(\boldsymbol{\theta})$ ) wrt a vector  $\boldsymbol{\theta}$ ,  
is a row vector  $\ell'(\boldsymbol{\theta}) = \frac{\partial \ell(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \left( \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_1}, \dots, \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_b} \right)$ .

Define score function

$$\begin{aligned} S(\boldsymbol{\theta}) &= \ell'(\boldsymbol{\theta})^\top \\ &= \begin{pmatrix} \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_1} \\ \vdots \\ \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_b} \end{pmatrix} \end{aligned}$$

$b \times 1$   
column vector.

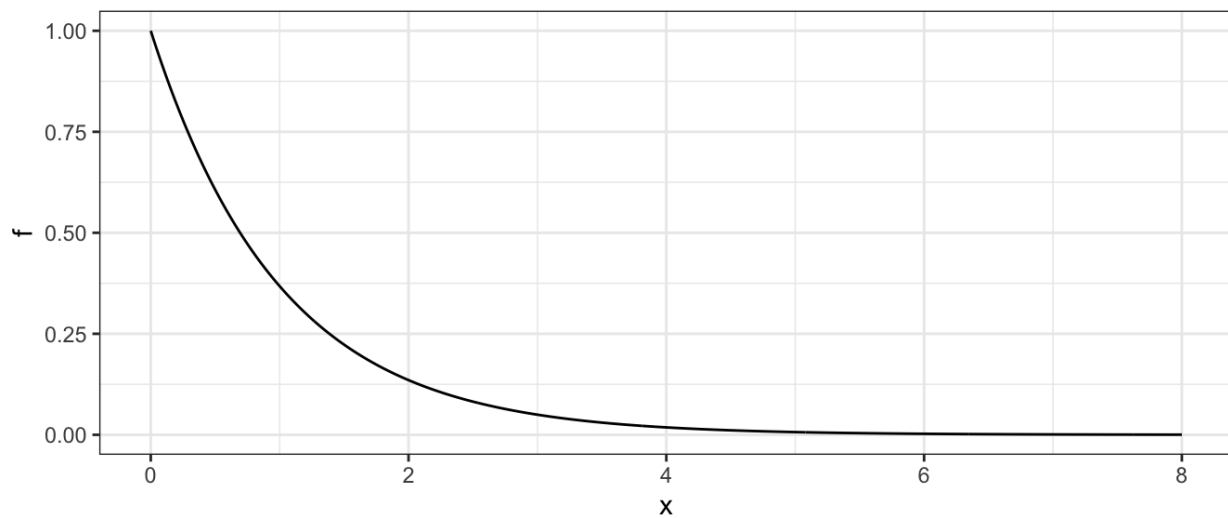


**Example:** Suppose we have  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Exp}(\lambda)$ . The likelihood function is defined as

$$\begin{aligned} L(\lambda | Y) &= f_Y(Y, \lambda) \\ &= \prod_{i=1}^n f_Y(y_i; \lambda) \\ &= \prod_{i=1}^n \lambda e^{-\lambda y_i} = \lambda^n e^{-\lambda \sum_{i=1}^n y_i} \Rightarrow \ell(\lambda) = n \log \lambda - \lambda \sum_{i=1}^n y_i \end{aligned}$$

```
# likelihood simulation
n <- 10
lambda <- 1

# plot of exponential(lambda) density
data.frame(x = seq(0, 8, .01)) |>
  mutate(f = dexp(x, rate = lambda)) |>
  ggplot() +
  geom_line(aes(x, f))
```



```

# define likelihood
loglik <- function(lambda, data)
{
  lik <- prod(dexp(data, rate = lambda))
  loglik <- sum(dexp(data, rate = lambda, log = T))
  out <- data.frame(lik = lik, loglik = loglik)
  return(out)
}

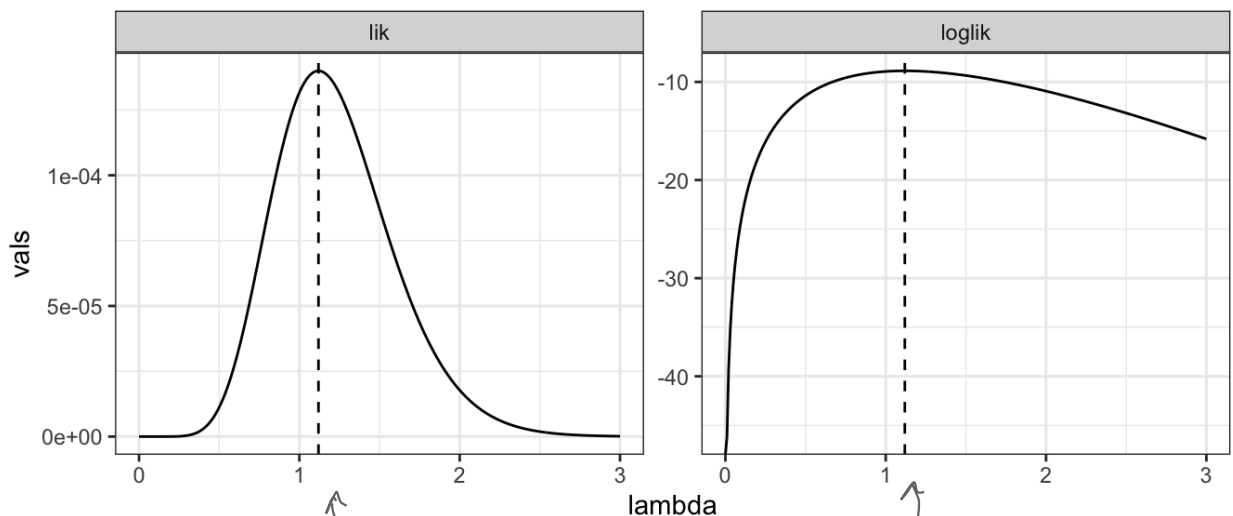
# simulate data
data <- rexp(n = n, rate = lambda)

# plot likelihood and loglikelihood
data.frame(lambda = seq(0, 3, by = .01)) |>
  rowwise() |>
  mutate(loglik = loglik(lambda, data)) |>
  unnest(cols = c(loglik)) |>
  pivot_longer(-lambda, names_to = "func", values_to = "vals") |>
  ggplot() +
  geom_vline(aes(xintercept = 1 / mean(data)), lty = 2) + # max ← MLE
  geom_line(aes(lambda, vals)) +
  facet_wrap(~func, scales = "free")

```

now we have realized data!

likelihood estimate is  $1/\text{mean}$



max happens at same place!

The likelihood function is random!

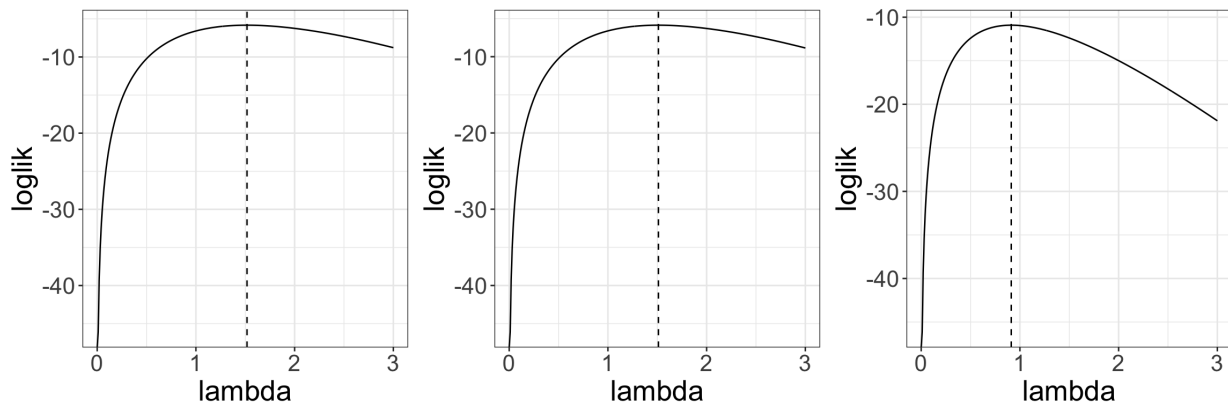
```

for(i in seq_len(3)) {
  # simulate data
  data <- rexp(n = n, rate = lambda)

  # plot likelihood and loglikelihood
  data.frame(lambda = seq(0, 3, by = .01)) |>
    rowwise() |>
    mutate(loglik = loglik(lambda, data)) |>
    unnest(cols = c(loglik)) |>
    ggplot() +
    geom_vline(aes(xintercept = 1 / mean(data)), lty = 2) + # max
      likelihood estimate is 1/mean
    geom_line(aes(lambda, loglik)) +
    theme(text = element_text(size = 20)) -> p ## make legible in
      notes

  print(p)
}

```



Question: What does likelihood integrate to?

Answer: not 1!  $\rightarrow$  function of  $\theta$ , not  $\psi$ .

**Your Turn:** What is the effect of sample size on the log-likelihood function? Make a plot showing the log-likelihood function that results from  $n = 10$  vs.  $n = 100$  with corresponding MLE.

## 1.2 Construction

The use of the likelihood function in parameter estimation is easiest to understand in the case of discrete iid random variables.

### 1.2.1 Discrete IID Random Variables

Suppose each of the  $n$  <sup>random</sup> variables in the sample  $Y_1, \dots, Y_n$  have probability mass function  $f(y; \theta) = P_\theta(Y_1 = y), y = y_1, y_2, \dots$ . The likelihood is then defined as:

$L(\theta | \mathbf{Y}) =$  joint density of observed random variables

$\stackrel{\text{iid}}{=} \text{product of univariate "densities"}$

$$= \prod_{i=1}^n \underbrace{f(y_i; \theta)}_{\text{pmf}}$$

$$= \prod_{i=1}^n P_\theta(Y_i^* = y_i | y_i)$$

where  $Y_1^*, \dots, Y_n^*$  are iid RV's w/ same distribution as  $Y_1, \dots, Y_n$   
(but mutually independent of  $Y_1, \dots, Y_n$ )

In other words,

the likelihood is the probability of getting sample actually obtained for a given value of  $\theta$ .

- ① In discrete case, can be thought of as a probability. (over what domain?)
- ② Will the likelihood sum to 1 over the parameter space? No.
- ③ Probability of finding a particular realization for a given  $\theta$ .

**Example (Fetal Lamb Movements):** Data on counts of movements in five-second intervals of one fetal lamb ( $n = 240$  intervals:)

No. of Movements	0	1	2	3	4	5	6	7
Count	182	41	12	2	2	0	0	1

$= 240.$

Assume a <sup>iid</sup> Poisson model:  $P(Y = y) = f_Y(y; \lambda) = \frac{\exp(-\lambda)\lambda^y}{y!}$ . Then the likelihood is

$$L(\lambda | \mathcal{Y}) = \prod_{i=1}^n f_Y(y_i; \lambda) = \prod_{i=1}^n \frac{\lambda^{y_i} e^{-\lambda}}{y_i!} = \lambda^{\sum_{i=1}^n y_i} e^{-n\lambda} \left( \prod_{i=1}^n y_i! \right)^{-1}$$

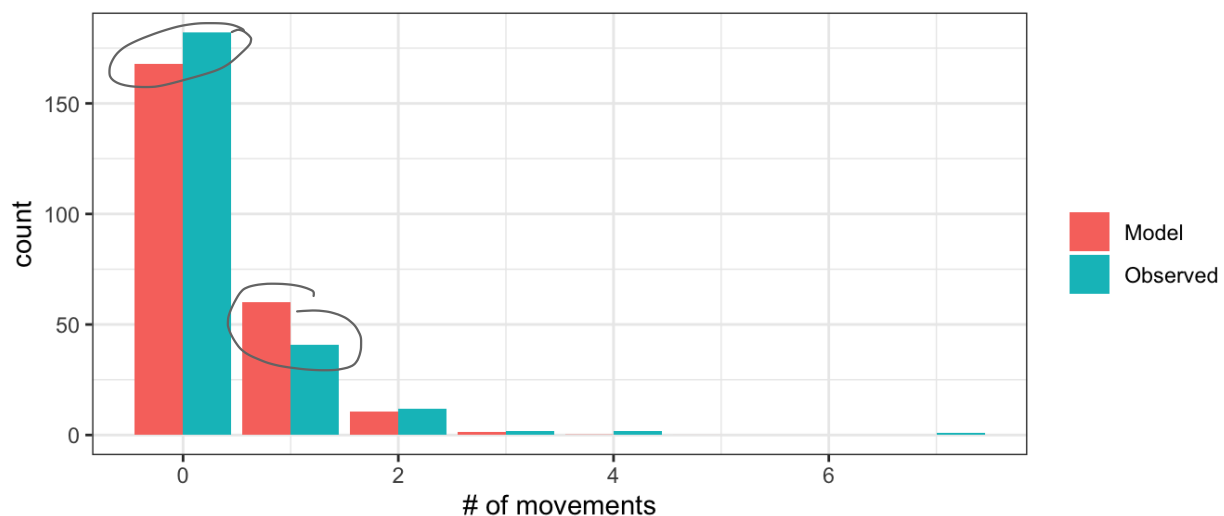
$$\ell(\lambda) = \sum_{i=1}^n y_i \log \lambda - n\lambda - \log \left( \prod_{i=1}^n y_i! \right)$$

$$\ell'(\lambda) = \frac{\sum y_i}{\lambda} - n \stackrel{\text{set}}{=} 0.$$

Equating the derivative of the loglikelihood with respect to  $\lambda$  to zero and solving results in the MLE

$$\hat{\lambda}_{\text{MLE}} = \frac{\sum y_i}{n} = \frac{86}{240} = 0.358.$$

This is the best we can do with this model. But is it good?



$\chi^2$  GOF test returns p-value of 0.00025, NOT good.

Illustration of a disadvantage of likelihood-based methods:  
very model based!

pg 31-32 of ESI extends to zero-inflated Poisson.

## 1.2.2 Multinomial Likelihoods

*more interesting discrete likelihoods.*

The multinomial distribution is a generalization of the binomial distribution where instead of 2 outcomes (success or failure), there are now  $k \geq 2$  outcomes.

Consider independently tossing  $n$  balls into  $k$  urns, where  $p_i$  is proba of the ball landing in  $i^{\text{th}}$  urn on each toss  $i=1, \dots, k$ .

$$\Rightarrow N_i \text{ balls in } i^{\text{th}} \text{ urn and } \sum_{i=1}^k N_i = n.$$

The probability mass function is ( $n$  total trials,  $k$  categories).

$$\begin{aligned} P(N_1 = n_1, \dots, N_k = n_k) &= p(n_1, \dots, n_k; p_1, \dots, p_k) \\ &= \frac{n!}{n_1! \dots n_k!} p_1^{n_1} p_2^{n_2} \dots p_k^{n_k} \quad \text{where } 0 \leq p_i \leq 1 \text{ and } \sum_{i=1}^k p_i = 1 \end{aligned}$$

For  $N_1, \dots, N_k$ ,  $N_i$  = the number of balls in  $i^{\text{th}}$  urn,  $\sum_{i=1}^k N_i = n$  (total balls thrown).  
*note  $p_k = 1 - \sum_{i=1}^{k-1} p_i$*

(to  $k$ ) observed data

$$L(p | N_1, \dots, N_k) = \frac{n!}{N_1! \dots N_k!} p_1^{N_1} \dots p_k^{N_k}$$

$\Rightarrow N_i$ 's not independent!

$$= \frac{n!}{N_1! \dots N_k!} p_1^{N_1} \dots p_{k-1}^{N_{k-1}} \left(1 - \sum_{i=1}^{k-1} p_i\right)^{N_k}$$

$$\Rightarrow \ell(p) = \text{const} + N_1 \log p_1 + N_2 \log p_2 + \dots + N_k \log \left(1 - \sum_{i=1}^{k-1} p_i\right)$$

*only have  $k-1$  params to estimate.*

The maximum likelihood estimator of  $p_i$ :

$$\frac{\partial \ell(p)}{\partial p_j} = \frac{N_j}{p_j} - \frac{N_k}{1 - \sum_{i=1}^{k-1} p_i} \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \begin{pmatrix} N_k & 0 & \dots & 0 & -N_1 \\ 0 & N_k & 0 & \dots & 0 & -N_2 \\ 0 & 0 & N_k & \dots & 0 & -N_3 \\ \vdots & 0 & \dots & N_k & N_{k-1} & \vdots \\ 0 & 0 & \dots & 0 & 0 & 0 \\ 1 & 1 & \dots & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ \vdots \\ p_k \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix}$$

$$\Rightarrow N_k p_j - N_j p_k = 0$$

$$\Rightarrow \hat{p}_{j \text{ MLE}} = \frac{N_j}{n} \quad (\text{what you would think})$$

More interesting multinomial likelihoods arise when the  $p_i$  are modeled as a function of a lesser number of parameters  $\theta_1, \dots, \theta_m$ ,  $m < k - 1$ .

**Example (Capture-Recapture):** To estimate fish survival during a specific length of time (e.g., one month), a common approach is to use a removal design.

$S$  = prob of fish surviving one month.

Time 0: catch and tag  $n$  fish.

Time 1: catch and remove some # of tagged fish,

$p$  = prob tagged fish is caught and removed.

$N_1$  = # tagged fish removed at time 1.

prob a tagged fish is caught at time 1 =

prob survival  
prob caught  
 $S \cdot p = P_1$

Time 2: repeat

$N_2$  = # tagged fish removed at time 2

prob a tagged fish is caught at time 2 =  $S^2(1-p)p = P_2$

⋮

Time  $k-1$ : repeat

$N_{k-1}$  = # tagged fish removed at time  $k-1$

prob a tagged fish is caught at time  $k-1$   $S^{k-1}(1-p)^{k-2}p = P_{k-1}$

$k^{\text{th}}$ : tagged fish is not removed

$$N_k = n - \sum_{i=1}^{k-1} N_i$$

$$P_k = 1 - sp - s^2(1-p)p - s^3(1-p)^2p - \dots - s^{k-1}(1-p)^{k-2}p = 1 - \sum_{i=1}^{k-1} P_i$$

Goal: estimate  $p$  and  $S$

Say you catch and remove  $N_1, \dots, N_k$  fish @  $k$  times,  $\sum_{i=1}^k N_i = n$  known.

The likelihood is the probability of catching  $N_1, \dots, N_k$  w/  $n$  total tagged.

$\Rightarrow$  multinomial w/  $p_i = s^i(1-p)^{i-1}p$   $i=1, \dots, k-1$  and  $p_k = 1 - \sum_{i=1}^{k-1} P_i$ .

$$\text{Recall: } L(f | N_1, \dots, N_k) = \frac{n!}{N_1! \dots N_k!} p_1^{N_1} \dots p_k^{N_k}$$

$$\text{(substitute } p_i \text{ 's)} = \frac{n!}{N_1! \dots N_k!} (sp)^{N_1} (s^2(1-p)p)^{N_2} \dots (s^{k-1}(1-p)^{k-2}p)^{N_{k-1}} P_k^{N_k}$$

What now? take log + partial derivatives wrt  $s$  &  $p$ , solve?

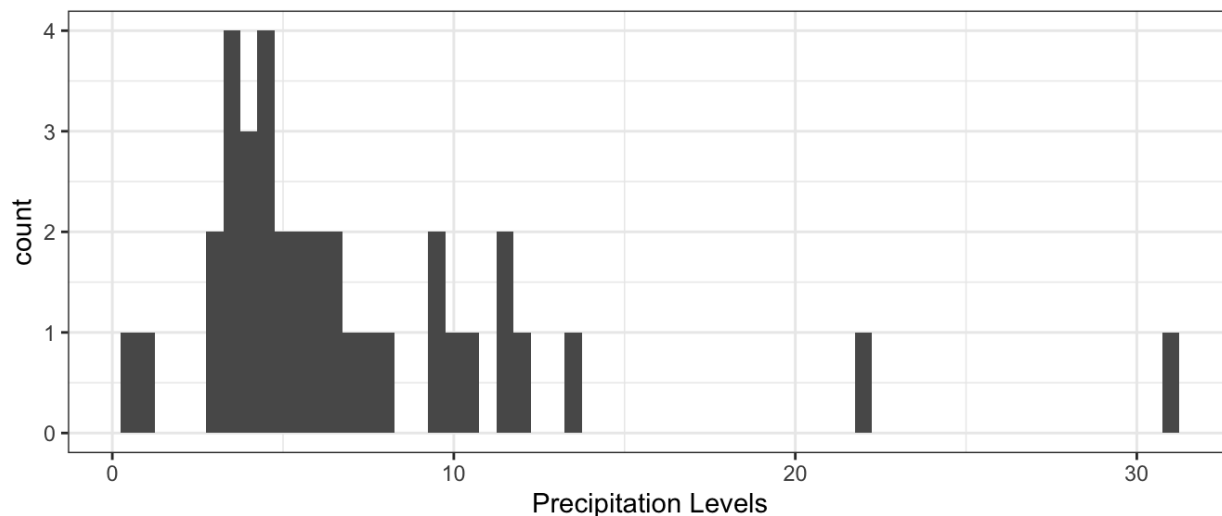
use a computer

where  $P_k = 1 - sp - s^2(1-p)p - \dots - s^{k-1}(1-p)^{k-2}p$   
↑ complicated, not log friendly.

### 1.2.3 Continuous IID Random Variables

Recall: the likelihood is the joint density of data to be analyzed.

**Example (Hurricane Data):** For 36 hurricanes that had moved far inland on the East Coast of the US in 1900-1969, maximum 24-hour precipitation levels during the time they were over mountains.



We model the precipitation levels with a gamma distribution, which has density

$$f(y; \alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} y^{\alpha-1} \exp(-y/\beta), \quad y > 0, \alpha, \beta > 0.$$

This leads to the likelihood

$$L(\theta | \mathcal{Y}) = \prod_{i=1}^n \frac{1}{\Gamma(\alpha)\beta^\alpha} y_i^{\alpha-1} e^{-y_i/\beta} = \{\Gamma(\alpha)\}^{-n} \beta^{-n\alpha} \left\{ \prod_{i=1}^n y_i \right\}^{\alpha-1} e^{-\sum_{i=1}^n y_i/\beta}$$

Of course, this cannot be interpreted as a probability because

$$P(Y = y) = 0 \quad \forall y!$$



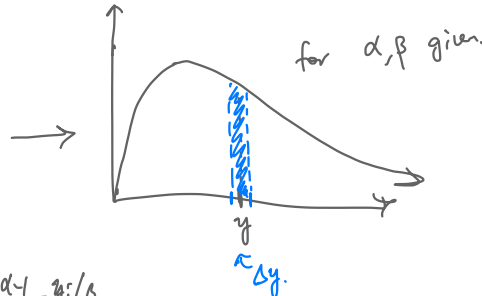
To get a probability, need to go from a density to a measure. i.e. integrate!

But likelihood not necessarily integrate to 1 (won't return value in  $[0,1]$  necessarily).

But it may be useful to think of the value of the likelihood as being proportional to a probability.

Given data  $y_1, \dots, y_n$

$$L(\alpha, \beta | y) = \prod_{i=1}^n \frac{1}{\Gamma(\alpha)\beta^\alpha} y_i^{\alpha-1} e^{-y_i/\beta}$$



prob is approximately  $\prod_{i=1}^n \frac{1}{\Gamma(\alpha)\beta^\alpha} y_i^{\alpha-1} e^{-y_i/\beta} \cdot \Delta y_i$

More formally, begin with the definition of a derivative

$$g'(x) = \lim_{h \rightarrow 0^+} \frac{g(x+h) - g(x-h)}{2h}$$

Let  $F$  be the cumulative distribution function of a continuous random variable  $Y$ , then (if the derivative exists)

$$f(y) = \lim_{h \rightarrow 0^+} \frac{F(y+h) - F(y-h)}{2h} = \lim_{h \rightarrow 0^+} \frac{P(Y \in (y-h, y+h])}{2h}$$

If we substitute this definition of a density into the definition of the likelihood

$$\begin{aligned} L(\theta | Y) &= \prod_{i=1}^n f(y_i; \theta) \\ &= \prod_{i=1}^n \lim_{h \rightarrow 0^+} \frac{F_\theta(y_i+h) - F_\theta(y_i-h)}{2h} \\ &= \lim_{h \rightarrow 0^+} \prod_{i=1}^n \frac{1}{2h} (F_\theta(y_i+h) - F_\theta(y_i-h)) \\ &= \lim_{h \rightarrow 0^+} \left( \frac{1}{2h} \right)^n \prod_{i=1}^n P_\theta(Y_i^* \in (y_i-h, y_i+h]) \end{aligned}$$

for small  $h$ , likelihood is proportional to a probability!

as  $h \rightarrow 0^+$  prob  $\downarrow 0$   
and  $(\frac{1}{2h})^n \uparrow$  to balance.

proportionality  
"constant" bc cause it doesn't depend on  $\theta$   
 $Y_i^*$  indep of  $Y_i$  w/ same dist as  $Y_i$

$\Rightarrow$  likelihood is proportional to the probability of obtaining a new sample that is close to the sample we obtained.

Compare this to the iid discrete case:

$$L(\underline{\theta} | \underline{y}) = \prod_{i=1}^n f(y_i; \theta)$$

$$= \lim_{h \rightarrow 0^+} \prod_{i=1}^n \{F(y_i + h; \theta) - F(y_i - h; \theta)\}$$

*now we don't need the proportionality constant.*

So likelihoods for discrete RVs get weighted differently than likelihoods for continuous RVs.

This has to do w/ underlying dominating measure of these RVs

(discrete: counting measure  
continuous: Lebesgue measure)

Thought for later: what about mixtures?

**Example (Hurricane Data, Cont'd):** Recall with a gamma model, the likelihood for this example is

$$L(\theta|Y) = \{\Gamma(\alpha)\}^{-n} \beta^{-n\alpha} \left\{ \prod Y_i \right\}^{\alpha-1} \exp\left(-\sum Y_i/\beta\right),$$

and log-likelihood

$$l(\theta) = -n \log \Gamma(\alpha) - n\alpha \log \beta + (\alpha-1) \sum \log Y_i - \frac{\sum Y_i}{\beta}$$

take derivs, set = 0 solve?

$$\frac{\partial \log \Gamma(\alpha)}{\partial \alpha} = -\gamma + \sum_{k=1}^{\infty} \left( \frac{1}{k} - \frac{1}{k+\alpha} \right)$$

Euler-Mascheroni constant.

## loglikelihood function

```
neg_gamma_loglik <- function(theta, data) {
  -sum(log(dgamma(data, theta[1], scale = theta[2])))
}
```

so we can use nlm (minimizes).

## maximize

partial derivs of Gamma likelihood do not result in linear system of eqns => use numerical optimization (common!).

```
mle <- nlm(neg_gamma_loglik, c(1.59, 4.458), data = hurr_rain)
mle$estimate
```

optim

$\hat{\alpha}_{MLE}$     $\hat{\beta}_{MLE}$

```
## [1] 2.187214 3.331862
```

quantiles based on Gamma fitted model.

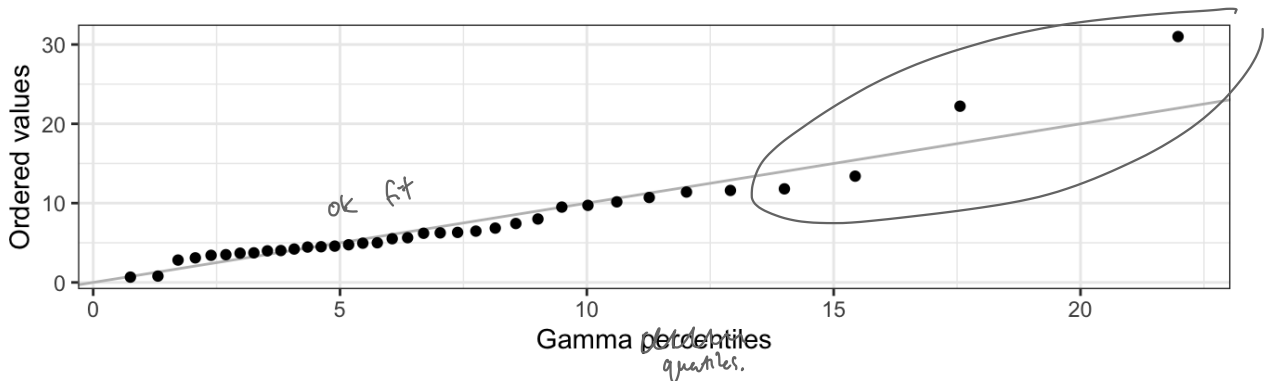
## Gamma QQ plot

```
data.frame(theoretical = qgamma(ppoints(hurr_rain), mle$estimate[1],
  scale = mle$estimate[2]),
  actual = sort(hurr_rain)) |>
ggplot() +
  geom_abline(aes(intercept = 0, slope = 1), colour = "grey") +
  geom_point(aes(theoretical, actual)) +
  xlab("Gamma percentiles") + ylab("Ordered values")
```

actual data (ordered)

percentiles

not great fit here.



### 1.2.4 Mixtures of Discrete and Continuous RVs

Some data  $Y$  often have a number of zeros and the amounts greater than zero are best modeled by a continuous distribution.

Ex: Rainfall, snowfall, Amt of time I climb in a day.

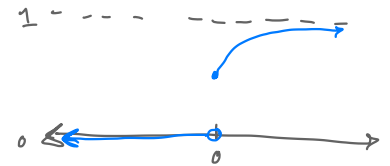
In other words, they have positive probability of taking a value of exactly zero, but continuous distribution otherwise.

- different from ZIP, which is discrete.

- can be generalized to other cases, point mass doesn't have to be at zero (e.g. record linkage)

A sensible model would assume  $Y_i$  are iid with cdf

$$F_Y(y; p, \theta) = \begin{cases} 0 & y < 0 \\ p & y = 0 \\ p + (1-p)F_T(y; \theta) & y > 0 \end{cases}$$



where  $0 < p \leq 1$  is  $P(Y = 0)$  and  $F_T(y; \theta)$  is a distribution function for a continuous positive random variable.

Another way to write this:

$$F_Y(y; p, \theta) = P(Y \leq y) = p\mathbb{I}(0 \leq y) + (1-p)F_T(y; \theta).$$

No problems w/ cdf.

How to go from here to get a likelihood?

What is the "density" here?

*practical.*

One approach: let  $n_0$  be the number of zeroes in the data and  $m = n - n_0$  be the number of non-zero  $Y_i$ . This leads to an intuitive way to construct the likelihood for iid  $Y_1, \dots, Y_n$  distributed according to the above distribution: *recalling the construction of likelihoods for cts & discrete i.i.d. r.v.s.*

$$\begin{aligned}
 L(\theta | \mathbf{Y}) &= \lim_{h \rightarrow 0^+} \left( \frac{1}{2h} \right)^m \prod_{i=1}^n \{ F_Y(Y_i + h; p, \theta) - F_Y(Y_i - h; p, \theta) \} \\
 &= \lim_{h \rightarrow 0^+} \left\{ F_Y(h; p, \theta) - F_Y(-h; p, \theta) \right\}^{n_0} \times \lim_{h \rightarrow 0^+} \prod_{Y_i > 0} \left\{ \frac{F_Y(Y_i + h; p, \theta) - F_Y(Y_i - h; p, \theta)}{2h} \right\} \\
 &= \lim_{h \rightarrow 0^+} \left\{ p + (1-p) F_T(h; \theta) \right\}^{n_0} \times \lim_{h \rightarrow 0^+} \prod_{Y_i > 0} \left\{ \frac{p + (1-p) F_T(Y_i + h; \theta) - p - (1-p) F_T(Y_i - h; \theta)}{2h} \right\} \\
 &= \underbrace{p^{n_0}}_{\text{Bernoulli component for } n_0 = n - m \text{ zeros.}} \times (1-p)^m \prod_{Y_i > 0} \underbrace{f_T(Y_i; \theta)}_{\text{continuous component for } m \text{ nonzero values.}}
 \end{aligned}$$

$$\ell(p, \theta) = (n - m) \log p + m \log(1 - p) + \sum_{Y_i > 0} \log f_T(Y_i; \theta)$$

Notice:  $\hat{p}_{MLE} = \frac{n - m}{n}$  &  $\hat{\theta}_{MLE}$  obtained in usual way from  $m$  obs. w/ density  $f_T(y; \theta)$  will only use  $Y_i > 0$

Kind of similar to law of total probability:

$$\begin{aligned}
 L(p, \theta | \mathbf{Y}) &\propto \prod_{i=1}^n P(Y_i = y_i) \\
 &= \prod_{i=1}^n \{ P(Y_i = y_i | Y_i = 0) P(Y_i = 0) + P(Y_i = y_i | Y_i \neq 0) P(Y_i \neq 0) \} \\
 &= \prod_{i=1}^n \{ \delta_{Y_i=0} \cdot p + \delta_{Y_i \neq 0} f_T(y_i; \theta) \cdot (1-p) \} = p^{n-m} (1-p)^m \prod_{Y_i > 0} f_T(y_i; \theta)
 \end{aligned}$$

Feels a little arbitrary in how we are defining different weights on our likelihood for discrete and continuous parts.

*(mostly)*

Turns out, it doesn't matter! (Need some STAT 630/720 to see why.)  
not much!

Small aside:

**Definition (Absolute Continuity)** On  $(\mathbb{X}, \mathcal{M})$ , a finitely additive set function  $\phi$  is *absolutely continuous* with respect to a measure  $\mu$  if  $\phi(A) = 0$  for each  $A \in \mathcal{M}$  with  $\mu(A) = 0$ . We also say  $\phi$  is *dominated* by  $\mu$  and write  $\phi \ll \mu$ . If  $\nu$  and  $\mu$  are measures such that  $\nu \ll \mu$  and  $\mu \ll \nu$  then  $\mu$  and  $\nu$  are *equivalent*.

Ex: continuous dsn is dominated by the Lebesgue measure.

Ex: a discrete dsn is dominated by the counting measure.

**Theorem (Lebesgue-Radon-Nikodym)** Assume that  $\phi$  is a  $\sigma$ -finite countably additive set function and  $\mu$  is a  $\sigma$ -finite measure. There exist unique  $\sigma$ -finite countably additive set functions  $\phi_s$  and  $\phi_{ac}$  such that  $\phi = \phi_{ac} + \phi_s \ll \mu$ ,  $\phi_s$  and  $\mu$  are mutually singular and there exists a measurable extended real valued function  $f$  such that

$$\phi_{ac}(A) = \int_A f d\mu, \quad \text{for all } A \in \mathcal{M}.$$

$\exists$  a set  $B$  s.t.  
 $\phi_s(B) = 0$  and  $\mu(B^c) = 0$   
 ↙ R-N derivative "density"

If  $g$  is another such function, then  $f = g$  a.e. wrt  $\mu$ . [If  $\phi \ll \mu$  then  $\phi(A) = \int_A f d\mu$  for all  $A \in \mathcal{M}$ .]

Think about a  $\phi$  w/ pos value at 0 and continuous  $> 0$ . Let  $\mu = \text{Lebesgue}$   
 $\nu = \text{counting measure on } \mathbb{Z}^3$ .  
 $\mu(\{0\}) = 0$ , but  $\nu(\{0\}^c) = 0 \Rightarrow \phi_s = \nu$  and  $\phi_{ac}$  is the rest.

**Definition (Radon-Nikodym Derivative)**  $\phi = \phi_{ac} + \phi_s$  is called the *Lebesgue decomposition*. If  $\phi \ll \mu$ , then the density function  $f$  is called the *Radon-Nikodym derivative* of  $\phi$  wrt  $\mu$ .

So what?

let  $\mu = \text{Lebesgue measure over } \mathbb{R}^t$   
 $\nu = \text{counting measure over } \mathbb{Z}^3$ .

$\hat{=}$  sum over elements in  $A \cap \mathbb{Z}^3$ .

$$\Rightarrow P(Y \in A | \theta) = \int_A f(y; \theta) d\mu(y) + \sum_A p(y; \theta) d\nu(y).$$

let  $\lambda = \mu + \nu$  and  $f_{\lambda}(y; \theta) = \mathbb{I}(y \neq \mathbb{Z}^3) f(y; \theta) + \mathbb{I}(y = \mathbb{Z}^3) p(y; \theta)$ .

$$\Rightarrow P(Y \in A | \theta) = \int_A f_{\lambda}(y; \theta) d\lambda(y).$$

↖ R-N derivative of prob measure of  $Y \Rightarrow$  valid density.

Let  $L_{\lambda}(\theta | \mathcal{Y}) = \text{joint density of observed data!}$

Claim: We can scale the continuous and discrete parts of likelihood and still have a valid likelihood.

See: Let use dominating measure  $\lambda_{**} = \alpha \cdot \mu + \beta \cdot \nu$  for  $\alpha, \beta > 0$ .

Then corresponding R-N derivative  $f_{**}(y; \theta) = \frac{\mathbb{I}(y \notin \mathcal{E}_0)}{\alpha} f(y; \theta) + \frac{\mathbb{I}(y \in \mathcal{E}_0)}{\beta} p(y; \theta)$

$$\Rightarrow P(Y \in A | \theta) = \int_A f_{**}(y; \theta) d\lambda_{**}(y).$$

A valid likelihood would be  $L_{**}(\theta | Y) = f_{**}(Y | \theta)$ .

$\Rightarrow$  we can scale the continuous and discrete parts of the likelihood however we like and its still valid.

Implications: We can scale the discrete and continuous components however we like.  
What to do? Mostly doesn't matter.

Let's say we have a sample w/  $n_0$   $Y_i = 0$  and  $m = n - n_0$   $Y_i > 0$  iid.

$$\begin{aligned} L_{**}(\theta | Y) &= \prod_{i=1}^n f_{**}(Y_i; \theta) \\ &= \prod_{Y_i=0} \frac{1}{\beta} p(Y_i; \theta) \prod_{Y_i>0} \frac{1}{\alpha} f(Y_i; \theta) \\ &= \frac{1}{\beta^{n_0} \alpha^m} \prod_{Y_i=0} p(Y_i; \theta) \prod_{Y_i>0} f(Y_i; \theta). \\ &= \frac{1}{\beta^{n_0} \alpha^m} \prod_{i=1}^n f_{**}(Y_i; \theta) \propto L_{*}(\theta | Y). \end{aligned}$$

$\Rightarrow$  Scaling can be ignored in MLE applications.

### 1.2.5 Proportional Likelihoods

Likelihoods are equivalent for point estimation as long as they are proportional and the constant of proportionality does not depend on unknown parameters.

Why?

Consider if  $Y_i, i = 1, \dots, n$  are iid continuous with density  $f_Y(y; \theta)$  and  $X_i = g(Y_i)$  where  $g$  is increasing and continuously differentiable. Because  $g$  is one-to-one, we can construct  $Y_i$  from  $X_i$  and vice versa.

knows  $\rightarrow$

$$y_i = g^{-1}(x_i)$$

$\Rightarrow \{y_1, \dots, y_n\}$  and  $\{x_1, \dots, x_n\}$  are "equivalent" because they contain exactly the same information

(intuition)  $\Rightarrow$  likelihood inference based on  $\{y_1, \dots, y_n\}$  should be identical to inference based on  $\{x_1, \dots, x_n\}$ .

More formally, the density of  $X_i$  is  $f_X(x; \theta) = f_Y(h(x); \theta) h'(x)$ , where  $h = g^{-1}$ , and

← from 530(?)

$$h(x_i) = g^{-1}(x_i) = y_i$$

$$L(\theta | \mathbf{X}) = \prod_{i=1}^n f_Y(h(x_i); \theta) h'(x_i)$$

$$= \prod_{i=1}^n f_Y(y_i; \theta) h'(g(y_i))$$

$$= \prod_{i=1}^n f_Y(y_i; \theta) \frac{1}{g'(y_i)}$$

$$= L(\theta | \mathbf{Y}) \left\{ \prod_{i=1}^n \frac{1}{g'(y_i)} \right\}$$

doesn't depend on  $\theta$ .

$$\frac{dh(x)}{dx} = \frac{dg^{-1}(x)}{dx} = \frac{1}{g'(g^{-1}(x))}$$

$x = f(f^{-1}(x))$  take derivative of both sides, solve.

$\Rightarrow$  MLE are identical when derived from  $L(\theta | \mathbf{Y})$  or  $L(\theta | \mathbf{X})$ .



**Example (Likelihood Principle):** Consider data from two different sampling plans:

1. A Binomial experiment w/  $n=12$ . Let  $Y_i = 1$  if  $i^{th}$  trial is successful and 0 otherwise.

\$\$

$L_1(p | \mathbf{Y}) = \binom{12}{S} p^S (1-p)^{12-S}$ , where  $S = \sum_{i=1}^n Y_i$

\$\$

$$L_1(p | \mathbf{Y}) = \binom{12}{S} p^S (1-p)^{12-S} \text{ where } S = \sum_{i=1}^n Y_i$$

2. A negative binomial experiment, i.e. run the experiment until three zeroes are obtained.

$$L_2(p | \mathbf{Y}) = \binom{S+2}{S} p^S (1-p)^3.$$

The ratio of these likelihoods is

$$\frac{L_1(p | \mathbf{Y})}{L_2(p | \mathbf{Y})} = \frac{\binom{12}{S} p^S (1-p)^{12-S}}{\binom{S+2}{S} p^S (1-p)^3} = \frac{\binom{12}{S}}{\binom{S+2}{S}} (1-p)^{9-S}$$

Suppose  $S = 9$ . Is all inference equivalent for these likelihoods? Debatable.

Then  $\frac{L_1(p | \mathbf{Y})}{L_2(p | \mathbf{Y})} = \frac{\binom{12}{9}}{\binom{11}{9}}$  doesn't depend on  $p$ !

Arguing:

YES: both cases,  $\hat{p}_{MLE} = \frac{9}{12}$

$$l_1(p) = \log \binom{12}{9} + 9 \log p + 3 \log (1-p)$$

$$\frac{d l_1(p)}{d p} = \frac{9}{p} - \frac{3}{1-p} \stackrel{set}{=} 0 \Rightarrow 9 - 3p = 3p \Rightarrow \hat{p}_{MLE} = \frac{9}{12}$$

$$l_2(p) = \log \binom{11}{9} + 9 \log p + 3 \log (1-p)$$

same exp constant  $\Rightarrow \hat{p}_{MLE} = \frac{9}{12}$

NO: Consider hypothesis test  $H_0: p = \frac{1}{2}$  vs.  $H_a: p > \frac{1}{2}$ .

If in experiment (1):  $\text{sum}(\text{dbinom}(c(9,10,11,12), \text{size}=12, p=\frac{1}{2})) = .0730$   
 (2):  $1 - \text{sum}(\text{dnbinom}(\text{seq}(0,8), \text{size}=3, \text{prob}=\frac{1}{2})) = 0.0327$  } p-values.

Bayesians: Discrepancy in p-values implies frequentist methods not logical because inference not based solely on likelihood.

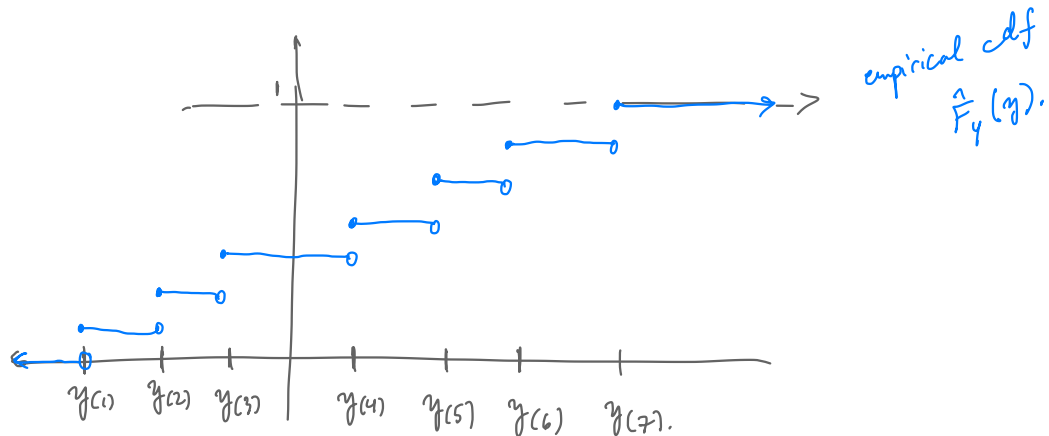
Maintain proportional likelihoods contain same information (formalized by Berger and Wolpert (1984)).  
**The likelihood principle** states all the information about  $\theta$  from an experiment is contained in the actual observation  $\mathbf{y}$ . Two likelihood functions for  $\theta$  (from the same or different experiments) contain the same information about  $\theta$  if they are proportional.

### 1.2.6 Empirical Distribution Function as MLE

Recall the empirical cdf:

Suppose  $y_{(1)} \leq y_{(2)} \leq \dots \leq y_{(n)}$  are the order statistics of an iid sample from an unknown distribution function  $F_Y$ . Our goal is to estimate  $F_Y$ .

$$\hat{F}_Y(y) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y \geq y_{(i)})$$



Is this a “good” estimator of  $F_Y$ ?

Maybe not... If you believe  $F_Y$  has support on  $\mathbb{R}$ ,

having  $\hat{F}_Y(y) = 0$  for  $y < y_{(1)}$  +  $\hat{F}_Y(y) = 1$  for  $y \geq y_{(n)}$ .

OTOH,  $F_Y(y)$  is likely to be pretty close to 0 for  $y < y_{(1)}$  and close to 1 for  $y \geq y_{(n)}$ .

↳ hopefully

otherwise talk to somebody who does “extremes”

Another answer:

Yes, because it's MLE.

→ require  $F(y)$  to have properties of dsn function

Suppose  $Y_1, \dots, Y_n$  are iid with distribution function  $F(y)$ . Here  $F(y)$  is the unknown parameter.

1.  $F(y)$  is nonnegative, non decreasing,
2.  $F(y)$  is right continuous
3.  $\lim_{y \rightarrow -\infty} F(y) = 0$  and  $\lim_{y \rightarrow \infty} F(y) = 1$ .

⇒ "parameter space" is set of all distribution functions.

An approximate likelihood for  $F$  is

(ignoring  $(2h)^{-n}$  factor).

$$L_h(F|Y) = \prod_{i=1}^n \{F(Y_i + h) - F(Y_i - h)\}$$

Assume no ties so that  $h$  is small enough such that  $[Y_i - h, Y_i + h]$  doesn't contain  $Y_j$  for any  $j \neq i$ .

let  $p_{i,h} = F(Y_i + h) - F(Y_i - h) \Rightarrow L_h(F|Y) = \prod_{i=1}^n p_{i,h}$

Note:  $\left. \begin{array}{l} \textcircled{1} \text{ increasing } p_{i,h} \text{ increases } L(F, Y). \\ \textcircled{2} L_h(F|Y) \text{ is maximized only if } p_{i,h} > 0 \text{ } i=1, \dots, n. \end{array} \right\} \Rightarrow$  want  $p_{i,h}$  to be as large as possible w/  $p_{i,h} > 0 \text{ } \& \sum_{i=1}^n p_{i,h} \leq 1$ . (3. above).  
max happens when = 1.

goal: maximize  $\prod_{i=1}^n p_{i,h}$  subject to  $p_{i,h} > 0$  and  $\sum_{i=1}^n p_{i,h} = 1$ .

Optimization problem to be solved using Lagrange multipliers (find stationary points of  $g$ )

$$g(p_{1,h}, \dots, p_{n,h}, \lambda) = \underbrace{\sum_{i=1}^n \log(p_{i,h})}_{\text{log likelihood}} + \lambda \left( \underbrace{\sum_{i=1}^n p_{i,h} - 1}_{\text{constraint}} \right).$$

get by solving:

$$\left. \begin{array}{l} \frac{\partial g(\dots)}{\partial p_{i,h}} = \frac{1}{p_{i,h}} + \lambda \stackrel{\text{set}}{=} 0 \quad i=1, \dots, n. \\ \frac{\partial g(\dots)}{\partial \lambda} = \sum_{i=1}^n p_{i,h} - 1 \stackrel{\text{set}}{=} 0 \end{array} \right\} \text{ implies } p_{i,h} = -\frac{1}{\lambda} \Rightarrow \lambda = -n.$$

⇒ any function  $\hat{F}_h(y)$  which satisfies  $\hat{F}_h(Y_i + h) - \hat{F}_h(Y_i - h) = \frac{1}{n} \quad i=1, \dots, n$  maximize  $L_h(F, Y)$ .

As  $h \rightarrow 0$ ,  $\hat{F}_h(y) \rightarrow \frac{1}{n} \sum_{i=1}^n \mathbb{1}(Y_i \leq y) \leftarrow$  this is the empirical dsn  $f_n^h \Rightarrow$  MLE.

### 1.2.7 Censored Data

Censored data occur when the value is only partially known. This is different from *truncation*, in which the data does not include any values below (or above) a certain limit.

*truncation:*

For example, we might sample only households that have an income above a limit,  $L_0$ . If all incomes have distribution  $F(x; \theta)$ , then for  $y > L_0$ ,

$$P(Y_1 \leq y | Y_1 > L_0) = \frac{P(Y_1 \leq y, Y_1 > L_0)}{P(Y_1 > L_0)} = \frac{F(y; \theta) - F(L_0; \theta)}{1 - F(L_0; \theta)}$$

The likelihood is then

$$L(\theta | y) = \prod_{i=1}^n \left\{ \frac{f(Y_i; \theta)}{1 - F(L_0; \theta)} \right\}$$

↑  
this is just an iid likelihood w/ densities adjusted to take into account that  $Y_i > L_0$ .

#### 1.2.7.1 Type I Censoring

*censoring:*

Suppose a random variable  $X$  is normally distributed with mean  $\mu$  and variance  $\sigma^2$ , but whenever  $X \leq 0$ , all we observe is that it is less than or equal to 0. If the sample is set to 0 in the censored cases, then define

$$Y = \begin{cases} 0 & \text{if } X \leq 0 \\ X & \text{if } X > 0. \end{cases} \quad \begin{array}{l} Y: \text{ observed,} \\ X: \text{ latent/hidden.} \end{array}$$

The distribution function of  $Y$  is  $X \sim N(\mu, \sigma^2)$ .

$$\begin{aligned} \text{at } y=0: F_Y(0) &= P(Y=0) = P(X \leq 0) \\ &= \Phi\left(-\frac{\mu}{\sigma}\right), \text{ where } \Phi \text{ is standard Normal cdf.} \end{aligned}$$

$$\text{for } y > 0: F_Y(y) = P(Y \leq y) = P(X \leq y) = \Phi\left(\frac{y-\mu}{\sigma}\right)$$

$$\text{and } F_Y(y) = 0 \text{ for } y < 0.$$

Suppose we have a <sup>random</sup> sample  $Y_1, \dots, Y_n$  and let  $n_0$  be the number of sample values that are 0. Then  $m = n - n_0$  and

$$L_h(\theta | Y) = \left(\frac{1}{2h}\right)^m \prod_{i=1}^m \{F_Y(Y_i + h; \theta) - F_Y(Y_i - h; \theta)\}$$

$$= \left\{ \Phi\left(\frac{h - \mu}{\sigma}\right) - 0 \right\}^{n_0} \times \prod_{Y_i > 0} \left\{ \frac{\Phi((Y_i + h - \mu)/\sigma) - \Phi((Y_i - h - \mu)/\sigma)}{2h} \right\}$$

$$L(\theta | Y) = \lim_{h \rightarrow 0^+} L_h(\theta | Y)$$

$$= \left\{ \Phi\left(-\frac{\mu}{\sigma}\right) \right\}^{n_0} \prod_{Y_i > 0} \left\{ \frac{1}{\sigma} \phi\left(\frac{Y_i - \mu}{\sigma}\right) \right\}$$

same as mixture!

$p_0(Y \text{ censored})^{n_0} \prod_{Y_i \text{ not censored}} f_Y(Y_i | \theta)$

More generally for left censoring @  $L_0$ :

$$L(\theta | Y) = \{F(L_0)\}^{n_0} \prod_{Y_i > L_0} f_Y(Y_i | \theta)$$

$L_0$  fixed = "type I censoring"

Type II censoring = "let's collect until first  $r$  censored values"

We might have censoring on the left at  $L_0$  and censoring on the right at  $R_0$ , but observe all values of  $X$  between  $L_0$  and  $R_0$ . Suppose  $X$  has density  $f(x; \theta)$  and distribution function  $F(x; \theta)$  and

$$Y_i = \begin{cases} L_0 & \text{if } X_i \leq L_0 \\ X_i & \text{if } L_0 < X_i < R_0 \\ R_0 & \text{if } X_i \geq R_0 \end{cases}$$

If we let  $n_L$  and  $n_R$  be the number of  $X_i$  values  $\leq L_0$  and  $\geq R_0$  then the likelihood of the observed data  $Y_1, \dots, Y_n$  is

$$L(\theta | Y) = \underbrace{\{F(L_0; \theta)\}^{n_L}}_{\text{left censored part}} \underbrace{\left\{ \prod_{L_0 < Y_i < R_0} f_Y(Y_i; \theta) \right\}}_{\text{observed}} \underbrace{\{1 - F(R_0; \theta)\}^{n_R}}_{\text{right censored}}$$

We could also let each  $X_i$  be subject to its own censoring values  $L_i$  and  $R_i$ . For the special case of right censoring, define  $Y_i = \min(X_i, R_i)$ . In addition, define  $\delta_i = \mathbb{I}(X_i \leq R_i)$ . Then the likelihood can be written as

$$= \begin{cases} 1 & \text{if observed } X_i \\ 0 & \text{if } X_i \text{ censored.} \end{cases}$$

$$L(\theta|Y) = \prod_{i=1}^n f(y_i; \theta)^{\delta_i} \underbrace{[1 - F(R_i; \theta)]^{1-\delta_i}}_{\text{right censored part.}}$$

**Example (Equipment failure times):** Pieces of equipment are regularly checked for failure (but started at different times). By a fixed date (when the study ended), three of the items had not failed and therefore were censored.

y	2	72	51	50	33	27	14	24	4	21
delta	1	0	1	0	1	1	1	1	1	0

Suppose failure times follow an exponential distribution  $F(x; \sigma) = 1 - \exp(-x/\sigma)$ ,  $x \geq 0$ . Then

$$L(\sigma|\mathbf{Y}) =$$

### 1.2.7.2 Random Censoring

So far we have considered censoring times to be fixed. This is not required.

This leads to random censoring times, e.g.  $R_i$ , where we assume that the censoring times are independent of  $X_1, \dots, X_n$  and iid with distribution function  $G(t)$  and density  $g(t)$ .

Let's consider the contributions to the likelihood:

which results in

$$L(\boldsymbol{\theta}|\mathbf{Y}, \boldsymbol{\delta}) =$$