



ESTIMATION FOR COMPLETE DATA

Introduction

- Recap from chapter 4 – Data dependent distribution versus parametric distribution
 - **Definition 11.1 (13.1)** – A **data-dependent distribution** is at least as complex as the data or knowledge that produced it, and the number of “parameters” increases as the number of data points or amount of knowledge increases.
 - **Definition 11.2 (13.2)** – A **parametric distribution** is a set of distribution functions, each member of which is determined by specifying one or more values called *parameters*. The number of parameters is **fixed** and **finite**.



- Usually we deal with parametric distributions. However two data-dependent distributions are considered: the **empirical distribution** and the kernel smoothed distribution.
 - **Definition 11.3 (13.3)** – The **empirical distribution** is obtained by assigning probability $1/n$ to each data point *in the sample*.
 - **Definition 11.4 (13.4)** – A **kernel smoothed distribution** is obtained by replacing each data point *in the sample* with a continuous random variable and then assigning probability $1/n$ to each random variable. The random variable used must be identical except for a location or scale change that is related to its associated data point (see chapter 14).
 - **Observation:** The empirical distribution is a special case of kernel smoothed distribution in which the random variable assigns probability 1 to the data point (and 0 elsewhere).
- **Data sets:** When observations are collected the “ideal” situation is to have the *exact* value for each observation (“complete individual data” as in data set B and data set D1). However, complete individual data are not always available: one reason is grouping (data set C or data set A for drivers with 5 or more claims); another reason is **censoring** and/or **truncation**.



- 4 data sets are repeatedly used:
 1. Data set A – Number of accidents by one driver in one year (data presented in Dropkin, 1959).
 2. Data set B – Amounts paid on workers compensation medical benefits: Random sample (artificial data) of 20 payments (full amount of the loss).
 3. Data set C – Random sample of payments from 227 claims from a general liability insurance. Data classified by payment range.
 4. Data set D – Time at which a five-year term insurance policy terminates (artificial data). For some policyholders termination is by death, for some others it is by surrender (cancellation of the insurance contract) and for the remainder it is at the expiration of the five-years period. Two versions of this data set are presented. The first one (data set D1) with **full information** (time of death and time of surrender when applicable) and in the second one (data set D2) only the first event is recorded.

Data sets A and B will be presented in Example 11.1 (13.1).



Data Set C			Data set D1			Data set D2			
Payment range		Number	Policyholder	Time of death	Time of surrender	Policyholder	First observed	Last Observed	Event
Linf	Lsup	payment							
0	7500	99	1		0.1	1	0	0.1	s
7500	17500	42	2	4.8	0.5	2	0	0.5	s
17500	32500	29	3		0.8	3	0	0.8	s
32500	67500	28	4	0.8	3.9	4	0	0.8	d
67500	125000	17	5	3.1	1.8	5	0	1.8	s
125000	300000	9	6		1.8	6	0	1.8	s
300000	Infinity	3	7		2.1	...			
			8		2.5	15	0	4.1	s
			9		2.8	16	0	4.8	d
			10	2.9	4.6	17	0	4.8	s
			11	2.9	4.6	18	0	4.8	s
			12		3.9	19 -30	0	5.0	e
			13	4.0		31	0.3	5.0	e
			14		4.0	32	0.7	5.0	e
			15		4.1	33	1	4.1	d
			16	4.8		34	1.8	3.1	d
			17		4.8	35	2.1	3.9	s
			18		4.8	36	2.9	5.0	e
			19 -30			37	2.9	4.8	s
						38	3.2	4.0	d
						39	3.4	5.0	e
						40	3.9	5.0	e
Total number of observations		227							



- As we can notice the information given by data sets C to D is incomplete.
 - Data set C – grouped data
 - Data set D1 – censoring: For some observations, we only know that the time of death is greater than a given value (the time of surrender)
 - Data set D2 – censoring and truncation: Some observations are first observed at time $c > 0$
- **Censoring** and **truncation** are problems that will be discussed in the next chapter (Estimation for modified data)

The empirical distribution for **complete individual data**

- Let us consider a sample of size n , (x_1, x_2, \dots, x_n) and let us also define the indicator function of a set A

$$\text{by } I_A(x) = I(x \in A) = \begin{cases} 0 & x \notin A \\ 1 & x \in A \end{cases}$$

- Definition 11.5 (13.5)** – The empirical distribution function (also known as empirical cumulative distribution function or ecdf) is

$$F_n(x) = \frac{\text{number of obs } \leq x}{n} = \frac{\sum_{i=1}^n I(x_i \leq x)}{n}$$

- Whatever the type (discrete, continuous, mixed) of the random variable in the “theoretical” model, the empirical distribution function behaves as a distribution function of a discrete random variable.
- Klugman *et al* (*Loss Models*) introduce the concept of empirical probability function as

$$f_n(x) = \frac{\text{number of obs } = x}{n} = \frac{\sum_{i=1}^n I(x_i = x)}{n}.$$



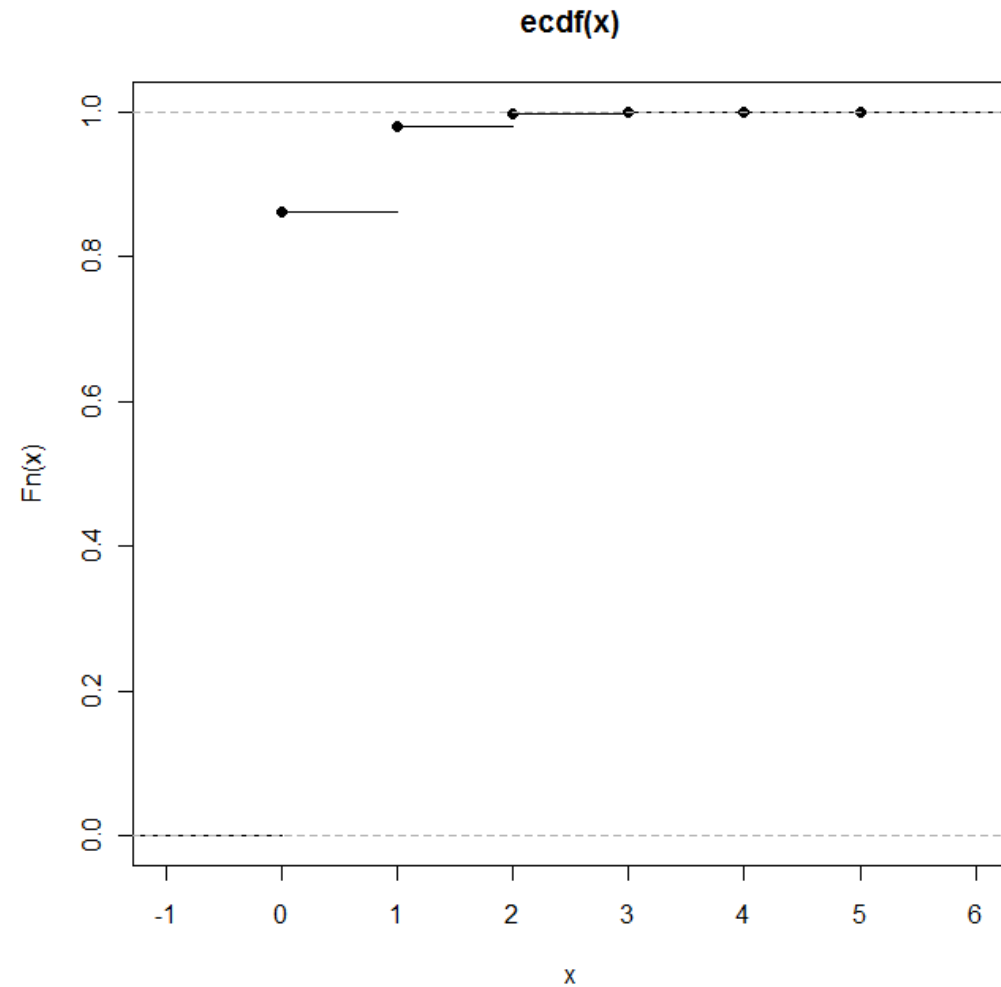
- If we are sampling from a continuous random variable, the probability that we observe a tie is 0 (some exceptions can arise due to the rounding of the observed values) and consequently in most situations $f_n(x) = 1/n$;
- The empirical distribution function is a much more important concept in statistical inference than the empirical probability function.
- **Example 11.1 (13.1)** – (changed – distribution instead of probability and vice-versa) Provide the empirical distribution functions for the data in data set A and B. For data set A also provide the empirical probability function. For data set A assume that all seven drivers who had five or more accidents had exactly five accidents.



Data Set A

Number of Accidents	Number of drivers
0	81714
1	11306
2	1618
3	250
4	40
5 or more	7

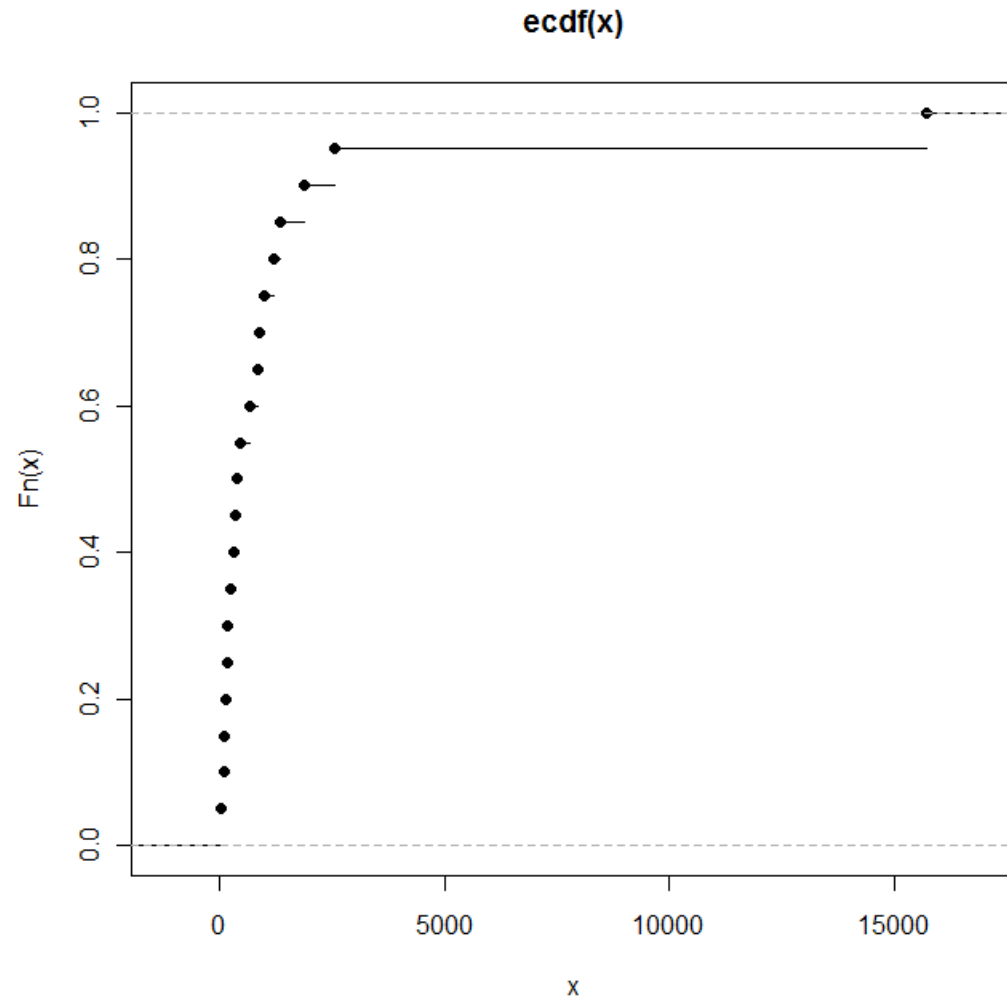
Total number of observations: 94935





Data Set B

27	82	115	126	155
161	243	294	340	384
457	680	855	877	974
1193	1340	1884	2558	15743





Data set B and R - Empirical distribution function

```
> # read data - Data set B
> x=c(27,82,115,126,155,161,243,294,340,384,457,680,855,877,974,
1193,1340,1884,2558,15743)
> F20=ecdf(x)
> summary(F20)
Empirical CDF:      20 unique values with summary
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  27.0  159.5   420.5 1424.0 1029.0 15740.0
> plot(F20)
```



Data Set A and R - Empirical distribution function

```
> # read data
>x=c(rep(0,81714),rep(1,11306),rep(2,1618),rep(3,250),rep(4,40),
rep(5,7))
> length(x)
[1] 94935
> F94935=ecdf(x)
> summary(F94935) # Be very careful with the results!!!!
F94935 is treated as an array with 6 observations equally distributed
Empirical CDF:      6 unique values with summary
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  0.00   1.25   2.50   2.50   3.75   5.00
# To get the empirical quartiles (all equal to 0 in this example) do
> quantile(x,c(0.25,0.5,0.75))
25% 50% 75%
  0   0   0
>plot(F94935)
```



```
> # Empirical probability function
> z=rep(1,length(x)); zz=tapply(z,x,sum)
> zz
      0      1      2      3      4      5
81714 11306  1618   250   40      7
> values=as.numeric(names(zz))
> values
[1] 0 1 2 3 4 5
> EmpProb=as.numeric(zz)/sum(as.numeric(zz))
> EmpProb
[1] 8.607363e-01 1.190920e-01 1.704324e-02 2.633381e-03 4.213409e-04
[6] 7.373466e-05
> F=cumsum(EmpProb)
> F
[1] 0.8607363 0.9798283 0.9968715 0.9995049 0.9999263 1.0000000
```



Risk set and cumulative hazard rate

- Consider a sample of size n , (x_1, x_2, \dots, x_n) , and let $y_1 < y_2 < \dots < y_k$ be the k unique values that appear in the sample ($k \leq n$). Let s_j be the number of times the observation y_j appears, $j = 1, 2, \dots, k$.

Obviously $\sum_{j=1}^k s_j = n$.

- Let us define the **risk set** as the observations that are greater than or equal to a given value (most of the time we will use “risk set” to refer the cardinal of the risk set) and let $r_j = \sum_{i=j}^k s_i$ be the risk set for the value y_j .
- Notice that, for $j = 2, 3, \dots, k$,

$$r_j = r_{j-1} - s_{j-1}; \quad n - r_j = \sum_{i=1}^{j-1} s_i; \quad r_j = n - \sum_{i=1}^{j-1} s_i, \text{ i.e. } r_1 = n; \quad r_2 = n - s_1; \quad \dots$$



- The empirical distribution can be written as

$$F_n(x) = \begin{cases} 0 & x < y_1 \\ 1 - \frac{r_j}{n} = \frac{\sum_{i=1}^{j-1} S_i}{n} & y_{j-1} \leq x < y_j \quad j = 2, 3, \dots, k \\ 1 & y_k \leq x \end{cases}$$

Note that $F_n(y_1) = 1 - \frac{r_2}{n}$, $F_n(y_2) = 1 - \frac{r_3}{n}$, ..., $F_n(y_k) = 1$.

- Illustrate using previous example.
- Definition 11.6 (13.6)** – The **cumulative hazard rate function** is defined as $H(x) = -\ln S(x)$
- Recall that $S(x) = 1 - F(x) = P(X > x)$
- Note that, if $H(x)$ is differentiable, $H'(x) = -S'(x) / S(x) = f(x) / S(x) = h(x)$ and then $H(x) = \int_{-\infty}^x h(y) dy$ where $h(x) = f(x) / S(x)$ is the hazard rate function.
- From previous definition, we get $F(x) = 1 - e^{-H(x)} \Leftrightarrow S(x) = e^{-H(x)}$

- **Definition 11.7 (13.7)** – The **Nelson-Aalen estimate** of the cumulative hazard rate function is

$$\hat{H}(x) = \begin{cases} 0 & x < y_1 \\ \sum_{i=1}^{j-1} \frac{S_i}{r_i} & y_{j-1} \leq x < y_j \quad j = 2, 3, \dots, k \\ \sum_{i=1}^k \frac{S_i}{r_i} & y_k \leq x \end{cases}$$

- **Comment:** Although the Nelson-Aalen estimator can be used with complete individual data, it has been established in a different framework, i.e. to be used with censored (and truncated) data. We shall return to this problem latter.
- **Examples 11.2 and 11.3 (13.2 and 13.3)** – Consider a data set containing the numbers
1.0, 1.3, 1.5, 1.5, 2.1, 2.1, 2.1, 2.8.
Determine the quantities described in the previous paragraph and then obtain the empirical distribution function. Determine the Nelson-Aalen estimate of the cumulative hazard function. Solve the example using EXCEL and R. Use the Nelson-Aalen estimate of the cumulative hazard function to estimate the distribution function
Nelson-Aalen estimate using R and following definition 11.7 ◻



```
> # Examples 11.2 and 11.3 following definition 11.7
> x=c(1.0, 1.3, 1.5, 1.5, 2.1, 2.1, 2.1, 2.8) # The sample
> z=rep(1,length(x)) # To provide a count using tapply
> zz=tapply(z,x,sum)
> zz
 1 1.3 1.5 2.1 2.8
 1  1  2  3  1
> y=as.numeric(names(zz)) # y_j
> s=as.numeric(zz) # s_j
> r=rep(length(x),length(s))
> r=r-c(0,cumsum(s)[1:length(s)-1]) # r_j
> y
[1] 1.0 1.3 1.5 2.1 2.8
> s
[1] 1 1 2 3 1
> r
[1] 8 7 6 4 1
> F=c(1-r/length(x),1)
> F # Example 11.2 finished
[1] 0.000 0.125 0.250 0.500 0.875 1.000
> H=c(0,cumsum(s/r)) # Nelson-Aalen estimate # Example 11.3
> H
[1] 0.0000000 0.1250000 0.2678571 0.6011905 1.3511905 2.3511905
```




```
> F_NA=1-exp(-H)
> F_NA                                     # another estimate of F_n
[1] 0.0000000 0.1175031 0.2349829 0.4518413 0.7410682 0.9047443
```

- **Empirical survival function**

- Using the empirical cumulative distribution function, $F_n(x)$, it is straightforward to get the empirical survival function $S_n(x) = 1 - F_n(x)$ which can act as an estimate of the survival function.

$$S_n(x) = \begin{cases} 1 & x < y_1 \\ \frac{r_j}{n} & y_{j-1} \leq x < y_j \quad j = 2, 3, \dots, k \\ 0 & y_k \leq x \end{cases}$$

- As we will see in the next chapter the case $S_n(x)$ for $x \geq y_k$ deserves some comments.
- We can also get an estimate of the survival function using the Nelson-Aalen estimate $\hat{S}(x) = e^{-\hat{H}(x)}$



Empirical distribution for grouped data

- For grouped data it is not possible to construct the empirical distribution function. The main idea is to approximate the empirical distribution by means of 2 points:
 - Wherever it is possible (at the groups boundaries) obtain the value of the empirical distribution;
 - Connect those points using a linear interpolation (other interpolation methods are possible)
- Let the group boundaries be $c_0 < c_1 < \dots < c_k$, i.e. group j is limited by c_{j-1} and c_j (often $c_0 = 0$ and $c_k = \infty$) and let us denote by n_j the number of observations in group j . Obviously $\sum_{j=1}^k n_j = n$.
- It is straightforward to see that $F_n(c_j) = (1/n) \sum_{i=1}^j n_i$, $j = 1, 2, \dots, k$ and that $F_n(c_0) = 0$.
- Treatment of the group boundaries: No rule is given. If the underlying variable is continuous, as it is generally the case, there is no real problem. For other situations, the best solution is to use group boundaries such that we can guarantee that the observed values are not equal to group boundaries. Technically, in order to guarantee that $F_n(x)$ is a distribution function, the value c_{j-1} should be excluded from group j and c_j included.



- **Definition 11.8 (13.8)** – For grouped data the distribution function obtained by connecting the values of the empirical distribution function at the group boundaries with straight lines is called the **ogive**. The formula is

$$F_n(x) = \frac{c_j - x}{c_j - c_{j-1}} F_n(c_{j-1}) + \frac{x - c_{j-1}}{c_j - c_{j-1}} F_n(c_j), \quad c_{j-1} \leq x < c_j$$

- Comments:
 - As this function is differentiable at all points except group boundaries, the (empirical) density function can be obtained. To specify the density function at the boundaries it is arbitrarily made right continuous.
 - We can re-write the empirical distribution function as

$$F_n(x) = \frac{c_j F_n(c_{j-1}) - c_{j-1} F_n(c_j)}{c_j - c_{j-1}} + \frac{F_n(c_j) - F_n(c_{j-1})}{c_j - c_{j-1}} x, \quad c_{j-1} \leq x < c_j$$

$$S_n(x) = 1 - F_n(x) = 1 - \frac{c_j F_n(c_{j-1}) - c_{j-1} F_n(c_j)}{c_j - c_{j-1}} - \frac{F_n(c_j) - F_n(c_{j-1})}{c_j - c_{j-1}} x$$

$$= \frac{c_j S_n(c_{j-1}) - c_{j-1} S_n(c_j)}{c_j - c_{j-1}} - \frac{S_n(c_{j-1}) - S_n(c_j)}{c_j - c_{j-1}} x, \quad c_{j-1} \leq x < c_j$$



- **Definition 11.9 (13.9)** – For grouped data the empirical density function can be obtained by differentiating the ogive. The resulting function is called a **histogram**. The formula is

$$f_n(x) = \frac{F_n(c_j) - F_n(c_{j-1})}{c_j - c_{j-1}} = \frac{n_j}{n(c_j - c_{j-1})}, \quad c_{j-1} \leq x < c_j$$

- Histograms and computer programs – be careful when classes do not have equal length
- **Example 11.5 (13.5)** – Construct the ogive and histogram for data set C.

Use EXCEL to define the empirical distribution function

You can also use R taking advantage of the **actuar** library or you can write your own solution



Using library actuar

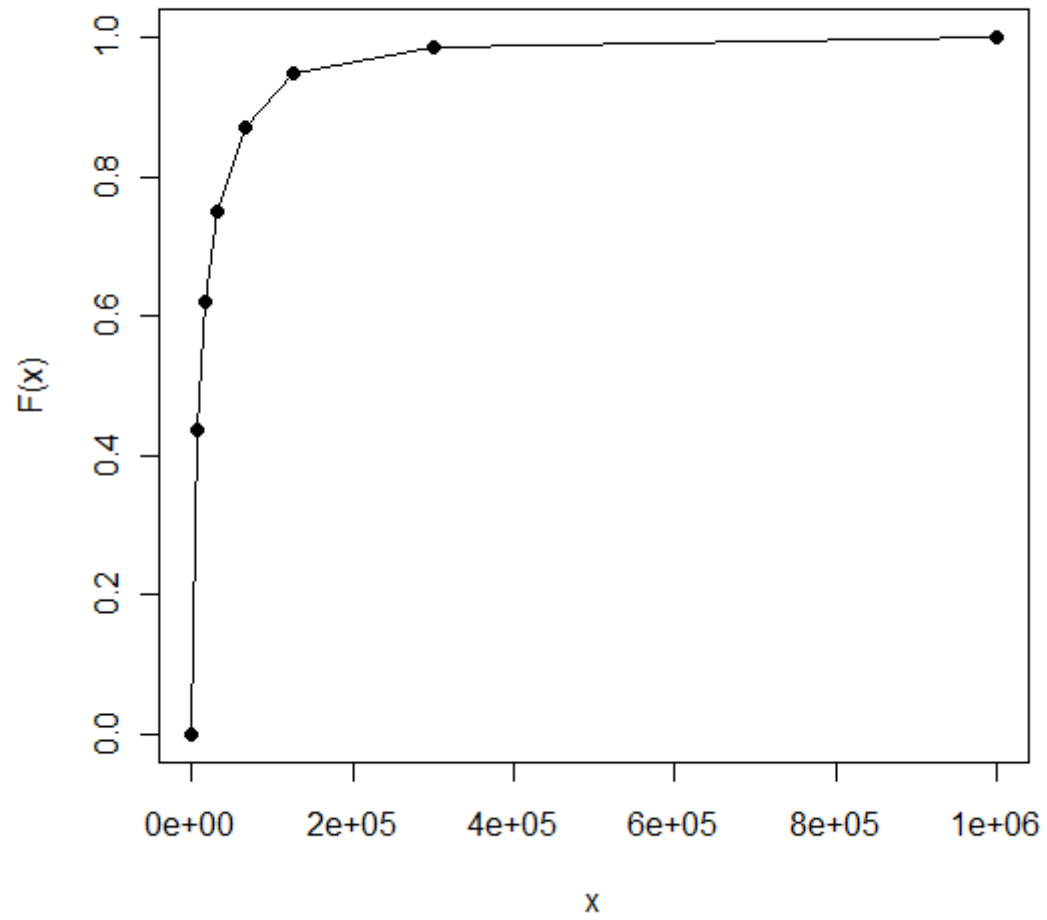
```
> library(actuar)
Attaching package: 'actuar'
```

The following object(s) are masked from package:grDevices : cm

```
> # 1000000 chosen arbitrarily
> x=c(0,7500,17500,32500,67500,125000,300000,1000000) # breaks
> y=c(99,42,29,28,17,9,3) # counts
> a=ogive(x,y)
> a
Ogive for grouped data
Call: ogive(x, y)
      x =      0,      7500,     17500, ..., 3e+05, 1e+06
  F(x) =      0, 0.43612, 0.62115, ..., 0.98678,      1
> plot(a)
> a(1000)
[1] 0.05814978
> a(7500)
[1] 0.4361233
```



ogive(x, y)

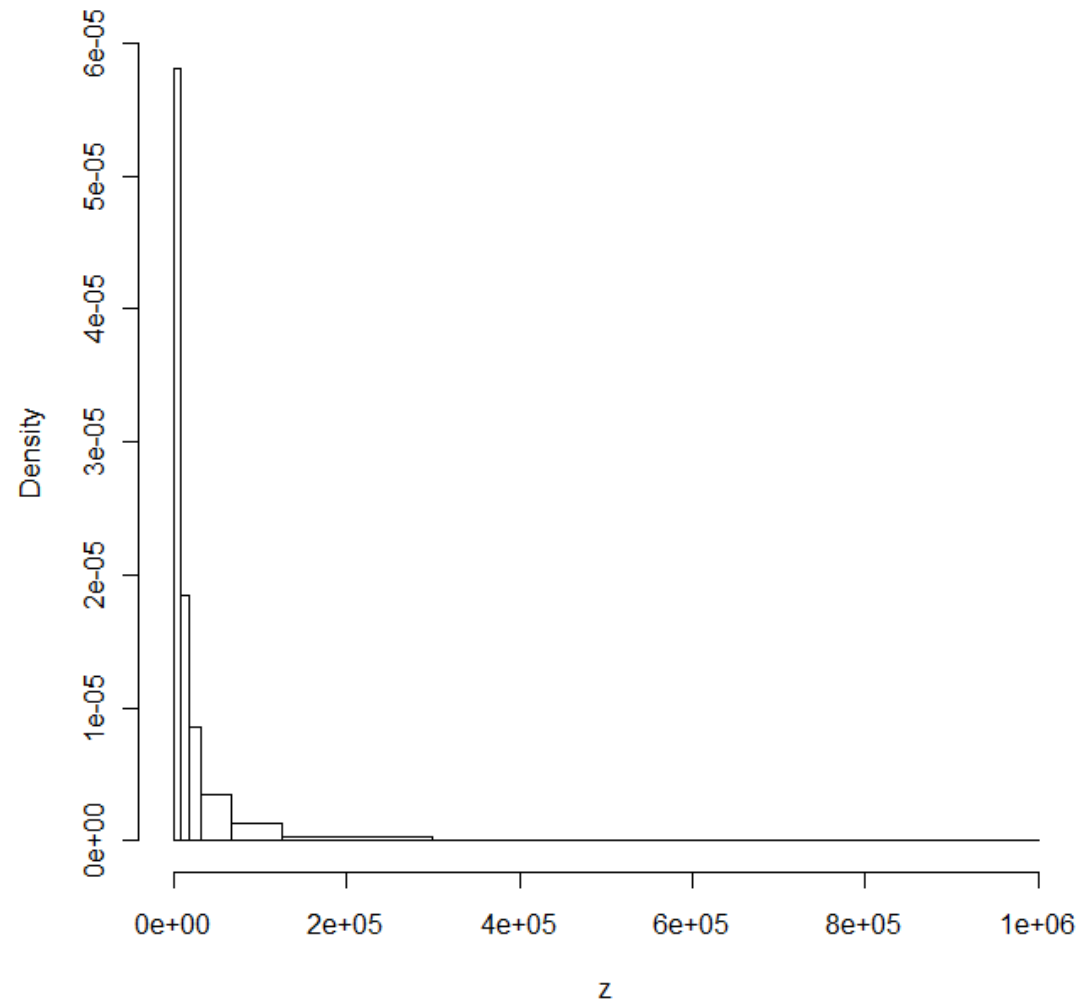




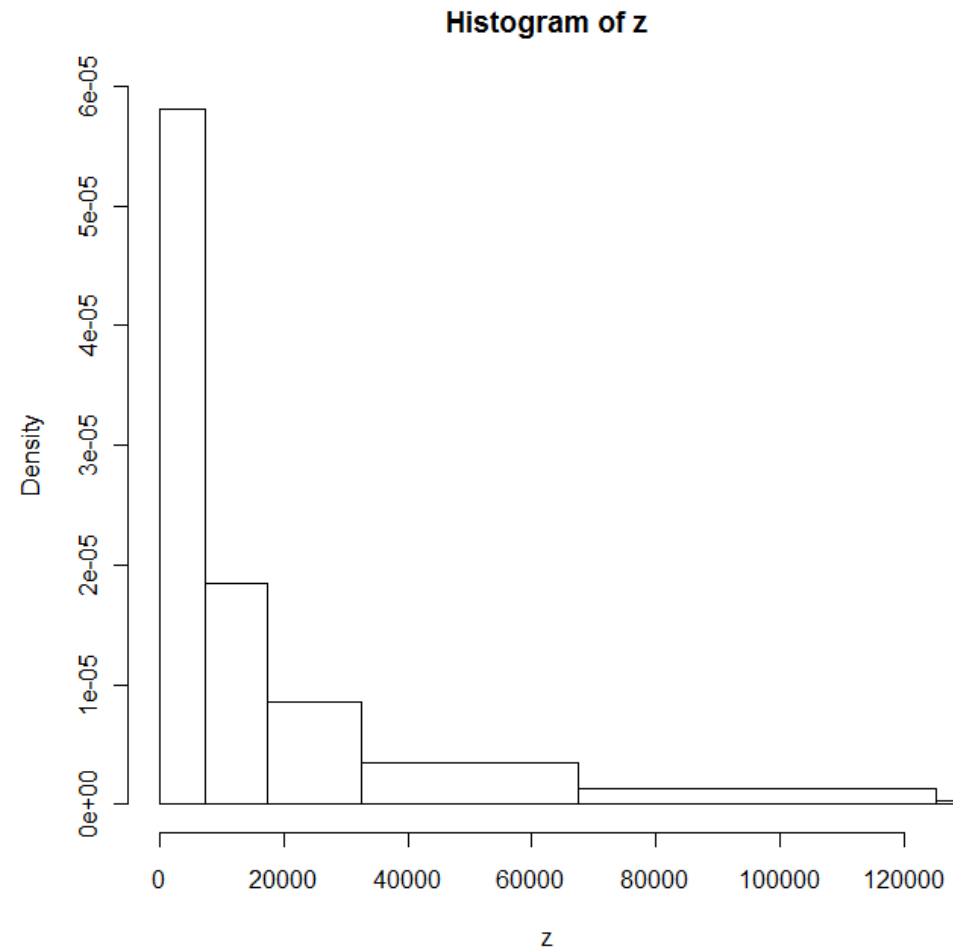
```
> lb=x[1:(length(x)-1)]; ub=c(lb[2:length(lb)],NA)
> a=cumsum(y)/sum(y);
> la=c(0,a[1:(length(a)-1)]); ua=a[1:length(a)]
> const=(ub *la-lb*ua)/(ub-lb)
> xcoef=(ua-la)/(ub-lb)
> ogive_table=data.frame(lower_bound=lb,upper_bound=ub,
constant=const,x_coef=xcoef)
> ogive_table
  lower_bound upper_bound  constant      x_coef
1           0         7500 0.0000000 5.814978e-05
2          7500        17500 0.2973568 1.850220e-05
3         17500        32500 0.4720999 8.516887e-06
4         32500        67500 0.6343612 3.524229e-06
5         67500       125000 0.7843325 1.302432e-06
6        125000       300000 0.9188169 2.265576e-07
7       300000            NA          NA          NA
> # empirical density in column 4 of ogive_table (x_coef)
> # To build array z choose an arbitrarily value in each class
> z=c(rep(5000,99),rep(10000,42),rep(20000,29),rep(50000,28),
rep(70000,17),rep(150000,9),rep(400000,3))
> b=hist(z,breaks=x)
```



Histogram of z




```
> hist(z, breaks=x, xlim=c(0, 125000))
```





The empirical survival function (from chapter 12 (14))

- Let us consider a random sample (X_1, X_2, \dots, X_n) and let us define the **estimator** of the empirical survival function

$$S_n^*(x) = \frac{1}{n} \#\{X_i > x\} = \frac{1}{n} \sum_{i=1}^n I(X_i > x) = \frac{N_x}{n}, \quad x > 0,$$

where $N_x = \#\{X_i > x\} = \sum_{i=1}^n I(X_i > x)$. It is straightforward to see that $N_x \sim b(n; S(x))$. If we consider an observed sample the corresponding estimate is

$$S_n(x) = \frac{1}{n} \#\{x_i > x\} = \frac{1}{n} \sum_{i=1}^n I(x_i > x) = \frac{n_x}{n}, \quad x > 0.$$

Following *Loss Models*, from now on **we will use the same notation for the estimator, $S_n^*(x)$, and the estimate, $S_n(x)$** . Both will be denoted by $S_n(x)$.



- **Problem 1** – How to estimate an unconditional probability like $\Pr(a < X \leq b)$?

Noting that $\Pr(a < X \leq b) = \Pr(X > a) - \Pr(X > b) = S(a) - S(b)$ a possible **estimator** is given by

$$\hat{\Pr}(a < X \leq b) = S_n(a) - S_n(b) = \frac{N_a - N_b}{n} = \frac{N_{(a,b]}}{n}.$$

Defining $N_{(a,b]}$ as the number of observations in the sample that fall in the interval $(a, b]$.

As $N_{(a,b]} \sim b(n; S(a) - S(b))$, it is straightforward to obtain the expected value and the variance of the estimator. Our estimate is

$$\hat{\Pr}(a < X \leq b) = S_n(a) - S_n(b) = \frac{n_a - n_b}{n} = \frac{n_{(a,b]}}{n}$$

- **Problem 2** – How to estimate a conditional probability like ${}_{y-x}q_x$

$${}_{y-x}q_x = \Pr(X \leq y - x + x | X > x) = \Pr(X \leq y | X > x) = \frac{\Pr(x < X \leq y)}{\Pr(X > x)} = \frac{S(x) - S(y)}{S(x)}$$

The “natural” estimate is ${}_{y-x}\hat{q}_x = \frac{S_n(x) - S_n(y)}{S_n(x)} = \frac{n_x - n_y}{n_x}$, assuming that $S_n(x) > 0$.

The corresponding estimator is ${}_{y-x}\hat{q}_x = \frac{N_x - N_y}{N_x}$.

Note that **this estimator do not have neither expected value nor variance** since $\Pr(N_x = 0) > 0$.



The usual solution

Assume that $S(x) = S_n(x)$ (or equivalently that $N_x = n_x$), given that $n_x > 0$.

Now the estimator is ${}_{y-x}\hat{q}_x = \frac{n_x - N_y}{n_x}$ but the distribution of N_y (and then the distribution of $S_n(y)$)

is conditioned by $S(x) = S_n(x)$.

The estimator is still unbiased and

$$\text{var}({}_{y-x}\hat{q}_x | S(x) = S_n(x)) = \frac{\text{var}(N_y | N_x = n_x)}{n_x^2} = \frac{1}{n_x^2} \times n_x \times \frac{n S(y)}{n_x} \times \left(1 - \frac{n S(y)}{n_x}\right) = \frac{1}{n_x^3} n S(y) (n_x - n S(y))$$

And the estimate of the variance is

$$\hat{\text{var}}({}_{y-x}\hat{q}_x | S(x) = S_n(x)) = \frac{1}{n_x^3} n_y (n_x - n_y)$$

How does it work?

Using the condition $S(x) = S_n(x)$ is equivalent to consider a sub-sample with all the observations greater than x and to estimate the probability of the variable being greater than y .

The sub-sample has n_x observations and we get the conditional estimator, ${}_{y-x}\hat{q}_x = \frac{n_x - N_y}{n_x} = 1 - \frac{N_y}{n_x}$.

Remember that, in this framework, $N_y \sim b(n_x, S(y)/S(x))$.

The variance of $\frac{N_y}{n_x}$, is estimated using the usual procedure applied to the sub-sample, i.e.

$$\hat{\text{vâr}}\left(\frac{N_y}{n_x}\right) = \frac{n_x \times \frac{n_y}{n_x} \times \left(1 - \frac{n_y}{n_x}\right)}{n_x^2} = \frac{n_y \times (n_x - n_y)}{n_x^3}.$$

As it is straightforward to see, $\hat{\text{vâr}}\left({}_{y-x}\hat{q}_x\right) = \hat{\text{vâr}}\left(1 - \frac{N_y}{n_x}\right) = \hat{\text{vâr}}\left(\frac{N_y}{n_x}\right)$.



- **Example 12.4 (14.5)** – Using the full information of data set D1, empirically estimate q_2 and estimate the variance of this estimator.

$$x = 2, y = 3, n = 30, n_2 = 29, n_3 = 27$$

$$\hat{q}_2 = \frac{29 - 27}{29} = \frac{2}{29} \approx 0.06897$$

$$\hat{\text{var}}(\hat{q}_2 | S(2) = 29/30) = \frac{27 \times (29 - 27)}{29^3} \approx 0.002214$$

- **Example 12.5 (14.6)** – Using data set B, empirically estimate the probability that a payment will be at least 1000 when there is a deductible of 250.

Let X be the value of a payment. Since there is a deductible of 250 we want to estimate $p = \Pr(X > 1250 | X > 250)$. Since there is a deductible we only have 13 observations

$$\hat{p} = \frac{S_n(1250)}{S_n(250)} = \frac{n_{1250}}{n_{250}} = \frac{4}{13} \approx 0.3077$$

$$\hat{\text{var}}(\hat{p}) = \frac{4 \times 9}{13^3} \approx 0.016386$$

Note that this variance is conditional to the existence of observations above the deductible.

Empirical estimation of probabilities

Let us consider a discrete random variable and let us assume that we want to estimate $p(x_j) = \Pr(X = x_j)$.

Let N_j be the number of times the value x_j was observed in a sample of size n . As it is straightforward to see $N_j \sim b(n; p(x_j))$.

The empirical estimator is $p_n(x_j) = N_j / n$. Consequently

$E(p_n(x_j)) = p(x_j)$, the estimator is unbiased

$\text{var}(p_n(x_j)) = \frac{p(x_j) \times (1 - p(x_j))}{n}$. The estimator is consistent.

The estimate of the variance is given by $\hat{\text{var}}(p_n(x_j)) = \frac{n_j \times (n - n_j)}{n^3}$

Note that the usual approximation from the binomial to the normal distribution can be used to get a confidence interval for $p(x_j)$.

Note also that similar results can be obtained for a continuous random variable when considering the probability of a particular event.



- **Example 12.7 (14.10)** – For Data Set A determine the empirical estimate of $p(2)$ and estimate the variance of the estimator.

$$n = 94935 \quad p_n(2) = 1618/94935 \approx 0.017043$$

$$\hat{\text{var}}(p_n(2)) = \frac{1618 \times (94935 - 1618)}{94935^3} \approx 1.76466 \times 10^{-7}$$

- **Example 12.8 (14.11)** – Use (10.3) and (10.4) – (12.3) and (12.4) – to construct approximate 95% confidence intervals for $p(2)$ using Data Set A

First approximation using (10.4): $\frac{p_n(2) - p(2)}{\sqrt{p_n(2) \times (1 - p_n(2)) / n}} \overset{\circ}{\sim} n(0;1)$

Confidence interval: $p_n(2) \pm 1.96 \times \sqrt{p_n(2) \times (1 - p_n(2)) / n}$, i.e. (0.01622; 0.01789)

Second approximation using (10.3): $\frac{p_n(2) - p(2)}{\sqrt{p(2) \times (1 - p(2)) / n}} \overset{\circ}{\sim} n(0;1)$

□ Confidence interval: $\frac{2n p_n(2) + 1.96^2 \pm 1.96 \sqrt{1.96^2 + 4n p_n(2) - 4n p_n(2)^2}}{2(n + 1.96^2)}$, i.e. (0.01624; 0.01789)



Empirical survival distribution for grouped data

Let Y be the number of observations in the sample (size n) whose values are less than or equal to c_{j-1} and let Z be the number of observations whose value are less than or equal c_j but greater than c_{j-1} .

- Then, for $c_{j-1} \leq x < c_j$, we have $S_n(x) = 1 - \frac{(c_j - c_{j-1})Y + (x - c_{j-1})Z}{n(c_j - c_{j-1})}$

Remember that, from definition 12.8, $F_n(x) = \frac{c_j - x}{c_j - c_{j-1}} F_n(c_{j-1}) + \frac{x - c_{j-1}}{c_j - c_{j-1}} F_n(c_j)$. Using the new

setup $F_n(c_{j-1}) = \frac{Y}{n}$ and $F_n(c_j) = \frac{Y+Z}{n}$.

- Now the marginal distributions of Y and Z are still binomial – $Y \sim b(n; 1 - S(c_{j-1}))$ and $Z \sim b(n; S(c_{j-1}) - S(c_j))$ – but the joint distribution is a multinomial (trinomial) distribution (Y and Z are not independent). Then

$$E(Y) = n(1 - S(c_{j-1})); \text{var}(Y) = n(1 - S(c_{j-1}))S(c_{j-1});$$

$$E(Z) = n(S(c_{j-1}) - S(c_j)); \text{var}(Z) = n(S(c_{j-1}) - S(c_j))(1 - S(c_{j-1}) + S(c_j));$$

$$\text{cov}(Y, Z) = -n(1 - S(c_{j-1}))(S(c_{j-1}) - S(c_j))$$



- The Expected value and variance of the estimator are given by

$$E(S_n(x)) = \frac{(c_j - x)}{(c_j - c_{j-1})} S(c_{j-1}) + \frac{(x - c_{j-1})}{(c_j - c_{j-1})} S(c_j)$$

$$\text{var}(S_n(x)) = \frac{(c_j - c_{j-1})^2 \text{var}(Y) + (x - c_{j-1})^2 \text{var}(Z) + 2(c_j - c_{j-1})(x - c_{j-1})\text{cov}(Y, Z)}{n^2(c_j - c_{j-1})^2}$$

- For the density estimate we get

$$f_n(x) = \frac{Z}{n(c_j - c_{j-1})}$$

Then

$$E(f_n(x)) = \frac{E(Z)}{n(c_j - c_{j-1})} = \frac{n(S(c_{j-1}) - S(c_j))}{n(c_j - c_{j-1})} = \frac{S(c_{j-1}) - S(c_j)}{c_j - c_{j-1}}$$

$f_n(x)$ is a biased estimator for $f(x)$. The variance is

$$\text{var}(f_n(x)) = \frac{\text{var}(Z)}{n^2(c_j - c_{j-1})^2} = \frac{(S(c_{j-1}) - S(c_j))(1 - S(c_{j-1}) + S(c_j))}{n(c_j - c_{j-1})^2}$$



Example 12.6 (14.8) – For data set C, estimate $S(10000)$, $f(10000)$ and the variance of your estimators.

Estimates

$$S_n(10000) = 1 - \frac{99 \times 10000 + 42 \times 2500}{227 \times 10000} \approx 0.51762$$

$$f_n(x) = \frac{42}{227 \times 10000} \approx 1.85022 \times 10^{-5}$$

Estimates of the variance of the estimators

$$\hat{\text{var}}(Y) = 227 \times \frac{128}{227} \times \frac{99}{227} = \frac{12672}{227} = 55.82379$$

$$\hat{\text{var}}(Z) = 227 \times \frac{42}{227} \times \frac{185}{227} = \frac{7770}{227} = 34.22907$$

$$\hat{\text{cov}}(Y, Z) = -227 \times \frac{42}{227} \times \frac{99}{227} = -\frac{4158}{227} = -18.31720$$

$$\hat{\text{var}}(S_n(x)) = \frac{10000^2 \times \frac{12672}{227} + 2500^2 \times \frac{7770}{227} - 2 \times 10000 \times 2500 \times \frac{4158}{227}}{227^2 \times 10000^2} \approx 0.000947127$$

$$\sqrt{\hat{\text{var}}(S_n(x))} \approx 0.030775$$

A 95% confidence interval for $S(10000)$ is given by (0.45730 ; 0.57794)