

---

# SurPyval

Apr 12, 2023



<b>1</b>	<b>Contents:</b>	<b>3</b>
1.1	Quickstart . . . . .	3
1.2	Types of Data . . . . .	4
1.3	Conventions . . . . .	7
1.4	Handy References - Aide-mémoire . . . . .	8
1.5	Data Wrangling Examples . . . . .	10
1.6	Datasets . . . . .	12
1.7	Non-Parametric Estimation . . . . .	12
1.8	Parametric Estimation . . . . .	17
1.9	Regression Analysis . . . . .	23
1.10	Non-Parametric SurPyval Modelling . . . . .	26
1.11	Parametric SurPyval Modelling . . . . .	33
1.12	Regression Modelling with SurPyval . . . . .	62
1.13	Example Applications . . . . .	66
1.14	API . . . . .	79
1.15	Changelog . . . . .	80
1.16	Support . . . . .	82
1.17	Contributing . . . . .	83
1.18	Acknowledgements . . . . .	83
1.19	Installation . . . . .	83
<b>2</b>	<b>Indices and tables</b>	<b>85</b>
	<b>Bibliography</b>	<b>87</b>



*surpyval* is an implementation of survival analysis in Python. The intent of this was to see if I could actually make it, and therefore learn a lot about survival analysis along the way, but also so that each time a model is created, it can be reused by other planned projects for monte carlo simulations (used in reliability engineering) and optimisations.

One feature of *surpyval* that separates it from other survival analysis packages is the intuitive way with which you can pass data to the fit methods. There are many different formats that can be used for survival analysis; *surpyval* handles many of the conceivable ways you can have your data stored. This is discussed in the data format tab.

*Surpyval* is also unique in the way in which it lets you estimate the parameters. With *surpyval*, you can use any of the following methods to estimate the parameters of you distribution of interest:

Table 1: SurPyval Modelling Methods

Method	Para/Non-Para	Observed	Censored	Truncated
Maximum Likelihood (MLE)	Parametric	Yes	Yes	Yes
Probability Plotting (MPP)	Parametric	Yes	Yes	Limited
Mean Square Error (MSE)	Parametric	Yes	Yes	Limited
Method of Moments (MOM)	Parametric	Yes	No	No
Maximum Product Spacing (MPS)	Parametric	Yes	Yes	No (planned)
Kaplan-Meier	Non-Parametric	Yes	Right only	Left only
Nelson-Aalen	Non-Parametric	Yes	Right only	Left only
Fleming-Harrington	Non-Parametric	Yes	Right only	Left only
Turnbull	Non-Parametric	Yes	Yes	Yes

Most other survival analysis packages focus on just using the MLE, or maybe the Probability Plotting. This package grew out of replicating the historically used probability plotting method from engineering, and as it progressed, it was discovered that there are many many ways parameters of distributions can be estimated. The product spacing estimator is particularly useful for offset distributions or finitely bounded distributions.

*SurPyval* attempts to use the combination of these methods to make parameter estimation possible for any distribution with arbitrary combinations of observations, censoring, and truncation.

Becoming a competent survival analyst depends strongly on having a very strong understanding of censoring, truncation, and observations in conjunction with a solid understanding of different types of distributions. Knowing and being able to identify situations as being censored or truncated in real applications will ensure you do not make an errors in your analysis. This can be very difficult to do. This documentation can be used as a reference to understand the types of censoring and truncation so that you can identify these situations in your work. Further, having a deep understanding of the types of distributions used in survival analysis will allow you to identify the process that is generating your data. This will then allow you to select an appropriate distribution, if any, to solve your problem. Survival analysis is an extremely powerful, and thoroughly interesting tool, so don't give up, or if you do give up, do the survival statistics on it.



## 1.1 Quickstart

So, you know what survival analysis is and you just want to see what this can do.

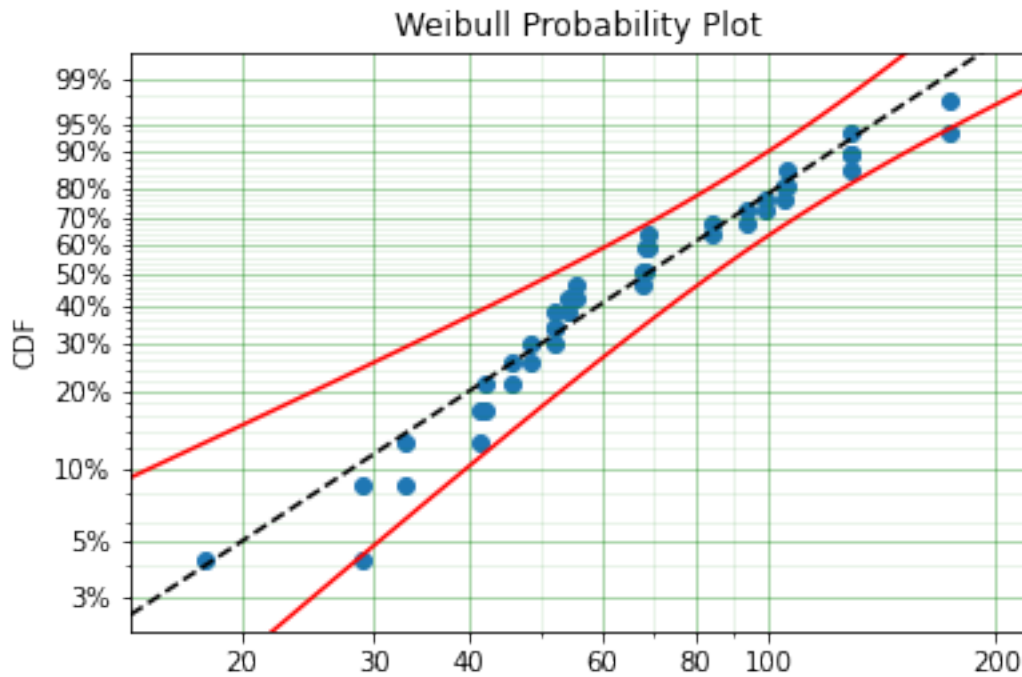
Alas

```
import surpyval as surv

x = [17.88, 28.92, 33, 41.52, 42.12, 45.6, 48.4, 51.84,
     51.96, 54.12, 55.56, 67.8, 68.64, 68.64, 68.88, 84.12,
     93.12, 98.64, 105.12, 105.84, 127.92, 128.04, 173.4]

model = surv.Weibull.fit(x)
model.plot()
```

This gives us the Weibull plot



## 1.2 Types of Data

Survival analysis is the statistics about durations. To understand durations, or time to events, we must have data that captures how long something lasts. This is the start of survival analysis where we have data in the form of a list of durations of some time to event. This time to event can be engineering failure data, health data on time to death from a given disease, economic data on the duration of a recession or time between recessions, or it could be race times for a group of athletes in a triathlon.

Survival analysis is unique in statistics because of the types of data that we encounter, specifically censoring and truncation. The purpose of this section is to explain these types of data and the scenarios under which they are generated so that you can understand when you might need to use the different flags in `surpyval` in your analysis.

### 1.2.1 Exactly Observed

The first type of data is exactly observed data. This is the type of data where we know exactly when the death or failure occurred. For example, if I run a test on how long some light bulbs will last, I get 5 and turn them on and watch them continuously. Then as each fail I record their failure times 983, 1321, 1889, 1923, and 2932 hours. Each of these times is exact because I saw the exact moment at which they failed. So the task is then understanding the distribution of these exact failures.

### 1.2.2 Censored Data

Say that I got bored of sitting and looking at light bulbs. And because of this boredom I stopped looking at the light bulbs at 1900 hours. I would therefore not have seen two of the light bulbs fail, the failures that would have occurred at 1923 and 2932 hours. All we would know about these two light bulbs is that they failed sometime *after* 1,900 hours. That is, we know that these two light bulbs would have failed if they continued the test but that this failure



time is greater than 1,900. This is to say that the failure time has been *censored*. Specifically the failure has been right censored. ‘Right’ is used because if we consider a (horizontal) timeline with time progressing along the line from left to right, we know that the failure would have occurred to the right of the time at which we stopped our observation. Hence, the observation is right censored.

If on the other hand, I also knew it would take some time for the bulbs to start failing so instead of waiting from the very start of the test I did not sit there for the first 1,000 hours. That is, the test continues to run but is not being observed for the first 1,000 hours, then after the first 1,000 hours I return to the test and start my observations. When I return I find that a bulb has failed. From the original data, I see that there was a bulb that failed at 983 hours. But if I was not observing for the first 1,000 hours all I would know about this failure is that it occurred sometime *before* 1,000 hours. Using the timeline concept again, I know that the failure would have occurred to the left of the 1,000 hour mark. Therefore, we say that the failure is *left* censored.

Finally, had I not been patient enough to sit down for any extended period of time and instead inspected the light bulbs at different times to see if any had failed. So say I inspect the bulbs every 100 hours from 1000 hours till 2,000 hours. The first and last failures would be left and right censored. But the middle failures would be known to fail between inspections. So the second failure would have occurred between the 1300 and 1400 hours inspections, the third between 1800 and 1900, and the second last failure would have happened between the 1900 and 2000 hours inspections. These failures are said to be *intervally* censored. That is because they are known to have happened in a given interval on a timeline.

Survival analysis has several methods for handling censored data in the parametric and non-parametric analysis. Surpyval is able to handle an input that has an arbitrary combination of observed and left, right, and intervally censored failure data. Although, not all methods can handle all types of data. This is covered in the sections on each of the estimation and fitting methods.

Surpyval uses a convention regarding censoring. Specifically, surpyval takes as input, with a list of failure times ‘x’, an optional censoring flag array ‘c’. If no flagging array is provided, it is assumed that all the data are exact observations, i.e. that they are not censored. But if the ‘c’ array is provided, it must have a value for each value in the x input. That is, they must be the same length. The possible values of c are -1, 0, 1, and 2. The convention tries to illustrate the concept of left, right, and interval censoring on the timeline. That is, -1 is the flag for left censoring because it is to the left of an observed failure. With an observed failure at 0. 1 is used to flag a value as right censored. Finally, 2 is used to flag a value as being intervally censored because it has 2 data points, a left and right point. In practice this will therefore look like:

```
import surpyval

x = [3, 3, 3, 4, 4, [4, 6], [6, 8], 8]
c = [-1, -1, -1, 0, 0, 2, 2, 1]

model = surpyval.Weibull.fit(x=x, c=c)
```

This example shows the flexibility surpyval offers. It allows users to analyse data that has any arbitrary combination of the different types of censoring. The surpyval format is even more powerful, because the above example can be condensed even further through using the ‘n’ value.

```
import surpyval

x = [3, 4, [4, 6], [6, 8], 8]
c = [-1, 0, 2, 2, 1]
n = [3, 2, 1, 1, 1]

model = surpyval.Weibull.fit(x=x, c=c, n=n)
```

The first step of the fit method actually wrangles the input data into the densest form possible. So internally, the example without the n value, will be condensed to be the second example without you seeing it. But it shows the capability of how data can be input to surpyval if you have different formats. But we are getting away from data

types...

### 1.2.3 Truncated Data

For my light bulb test, let's say I test a different manufacturers bulbs. This time, I know that the bulbs from this manufacturer have been tested for 500 hours prior to shipping them. This situation needs to be treated differently because we know that in this circumstance we only have the bulbs because they survived more than 500 hours. If there were any failures prior to 500 hours the bulb would not have been shipped and therefore would not be being tested by me. This is to say, that my observation of the distribution of the light bulb failures has been *truncated*. In this regime there is no way I can have any observation below 500 hours because of the testing then discarding done by the manufacturer. The astute reader might have observed that this data is in fact *left* truncated. This is because the truncation occurs to the left of the observation on a timeline. In this example, all the bulbs are left truncated at the 500 hour mark.

In biostatistics left truncation is known as 'late-entry', this is because in clinical trials a participant can enter a trial later than other participant. Therefore this participant was at risk of not being present in the trial. This is because they could have died prior to entering the trial. Morbid, yes, but the estimate of the distribution needs to account for this risk otherwise the estimate will overestimate the true risk of the event.

Right truncated data is when you only observe a value because it happened below some time. For example, in the light bulb experiment, I received some of the bulbs that passed the burn in test. That is, I received some of the bulbs that survived the original 500 hours of testing. But if the failed bulbs were then given to an engineering team to investigate possible design changes that will improve reliability; they will have a series of failure times that must be below 500 hours. That is, from their perspective, they have data that is right truncated. There is one condition to this situation, they must not know how many other bulbs were tested. If they knew how many other bulbs were tested, they would know how many would fail after 500 hours. That is, they would know that all the other bulbs are right censored. So for our engineers investigating the failed bulbs, they must be ignorant of how many other bulbs were actually tested for the right truncation to work for them. In many applications we do know how many were under test and therefore right truncation become right censoring, but from our engineers circumstance, we can see that they are right censored.

Parametric and non-parametric analysis can both handle left truncated data. This is explained further in the estimation methods for both these methods. Right truncation can only be handles in surpyval with parametric analysis, specifically, with Maximum Likelihood Estimation and in limited cases with Minimum Product Spacing. This is also explained in their respective sections of these notes.

In surpyval, passing truncated data to the fitting method looks like:

```
import surpyval

x = [674, 792, 1153, 1450, 1555, 1923, 2019]
tl = 500

model = surpyval.Weibull.fit(x=x, tl=tl)
```

### 1.2.4 Concluding Points

Having read through the above explanation you might be thinking how often these scenarios appear in real data, if ever. The vast majority of data used in survival analysis is observed or right censored. This is what happens when you observe a whole population but finish the observation before the event happens on all the items being observed.

Right truncation is extremely rare because it only happens if you do not know the size of the whole population under test. It can happen with scientific instruments where say, a camera is limited in the frequencies of light it can capture. So if we were to try capture a distribution of light of an object, say a star, this distribution could be truncated above and below certain frequencies. Meeker and Escobar provide an example in their book on reliability statistics for warranty analysis, similar to the contrived example provided above. If you have some returns of products from the field, these

are right-truncated because you do not know what has been bought and used in the field. A more realistic example could be the estimation of race finish times at a triathlon or marathon. If I arrive at the finish line of a race and record the times of participants as they cross the line during that window I will have truncated data. I do not know how many people started the race (presumably) and I only stay and watch for a given period of time, therefore all the observations I make are truncated within the window of my observation time. In conclusion though, right truncation in survival analysis is rare.

Left truncation is common in insurance studies. If an insurance company wants to estimate the distribution of losses due to property crime based on policy payouts they need to consider the impact of 'excess'. Excess is the cost of making a claim on an insurance policy. So if I have an insurance policy with an excess of \$500, if I lose \$20,000 worth of property in a robbery I will have to pay \$500 to be paid \$20,000. Because of this, it is clear that if I lost \$400 in a robbery I would not pay the \$500 excess to make a claim. Therefore the distribution of property crime will be truncated by the value of the excesses on the policies. Actuaries need to consider this in their calculations of policy fees.

Insurance is also a good example of right censoring. An insurance policy will also have a maximum payout. So if calculating the distribution of the value of property crime an analyst will need to consider that those payouts that are at the maximum of that policy value are in fact censored. That is, the value of the loss or damage was greater than the actual payout and therefore the payout is a censored value. In the classic Boston housing pricing data there is censored data! A histogram of the values of houses shows that there is a large number of houses at the highest price. This can be understood because a limit was set on the highest possible value, therefore these house prices are actually censored, not exact observations.

## 1.3 Conventions

### 1.3.1 Variable Names

Before discussing the formats, the conventions for variable names need to be clarified. These conventions are:

- $x$  = The random variable (time, stress etc.) array
- $x_l$  = The random variable (time, stress etc.) array for the left interval of interval censored data.
- $x_r$  = The random variable (time, stress etc.) array for the right interval of interval censored data.
- $c$  = Refers to the censoring data array associated with  $x$
- $n$  = The count array associated with  $x$
- $t$  = the truncation values for the left and right truncation at  $x$  (must be two dim, or use  $tl$  and  $tr$  instead)
- $r$  = risk set at  $x$
- $d$  = failures/deaths at  $x$
- $t$  = two dimensional array of the truncation interval at  $x$
- $tl$  = one dimensional array of the value at which  $x$  is left truncated
- $tr$  = one dimensional array of the value at which  $x$  is right truncated

### 1.3.2 Data Formats

The conventional formats use in surpyval are:

- $xcnt$  =  $x$  variables, with  $c$  as the censoring scheme,  $n$  as the counts, and  $t$  as the truncation
- $xrd$  =  $x$  variables, with the risk set,  $r$ , at  $x$  and the deaths,  $d$ , also at  $x$

All functions in `surpyval` have default handling conditions for `c` and `n`. That is, if these variables aren't passed, it is assumed that there was one observation and it was a failure for every `x`.

`Surpyval` `fit()` functions use the `xcnt` format. But the package has handlers for other formats to rearrange it to the needed format. Other formats are:

wrangers for formats:

- `fs` = failure time array, `f`, and right censored time array, `s`
- `fsl` = `fs` format plus an array for left censored times.

### 1.3.3 Censoring conventions

For the censoring values, `surpyval` uses the convention used in Meeker and Escobar, that is:

- `-1` = left
- `0` = failure / event
- `1` = right
- `2` = interval censoring. Must have left and right value in `x`

This convention gives an intuitive feel for the placement of the data on a timeline.

### 1.3.4 Function Conventions

The conventions for `SurPyval` are that each object returned from a `fit()` call has the ability to compute the following:

- `df()` - The density function
- `ff()` - The CDF
- `sf()` - The survival function, or reliability function
- `hf()` - The (instantaneous) hazard function
- `Hf()` - The cumulative hazard function

These functions can be used to plot or even in optimisers so that you can optimize decisions that you are guiding with your survival analysis.

## 1.4 Handy References - Aide-mémoire

### 1.4.1 Relationship between functions of a probability distribution

There exists a relationship between each of the functions of a distribution and the others. This can be very useful to keep in mind when understanding how `surpyval` works. For example, the Nelson-Aalen estimator is used to estimate the cumulative hazard function (`Hf`), the below relationships is how distribution for this can be used to estimate the survival function, or the cdf.

	$f(t)$	$F(t)$	$R(t)$	$h(t)$	$H(t)$
$f(t)$	-	$F'(t)$	$-R'(t)$	$h(t)e^{-\int h(t)}$	$H'(t)e^{-H(t)}$
$F(t)$	$\int_{-\infty}^t f(\tau) d\tau$	-	$1 - R(t)$	$1 - e^{-\int h(t)}$	$1 - e^{-H(t)}$
$R(t)$	$1 - \int_{-\infty}^t f(\tau) d\tau$	$1 - F(t)$	-	$e^{-\int h(t)}$	$e^{-H(t)}$
$h(t)$	$\frac{f(t)}{1 - \int_{-\infty}^t f(t) dt}$	$\frac{F'(t)}{1 - F(t)}$	$\frac{-R'(t)}{R(t)}$	-	$H'(t)$
$H(t)$	$-\ln\left(1 - \int_{-\infty}^t f(\tau) d\tau\right)$	$-\ln(1 - F(t))$	$-\ln R(t)$	$\int_{-\infty}^t h(\tau) d\tau$	-

The above table shows how the function on the left, can be described by the function along the top row (I leave out the function describing itself as it is simply itself. . .). So, an interesting one is that the reliability or survival function,  $R(t)$ , is simply the exponentiated negative of the cumulative hazard function! This relationship holds for **every** distribution.

### 1.4.2 AFT, AL, or PH?

What is the difference, if any, between an Accelerated Failure Time model, an Accelerated Life model, and a Proportional Hazard model? SurPyval uses the distinctions defined in [Bagdonavicius]. The explanation of these are:

- ALT is an accelerated life model. That is, a model where the ‘characteristic life’ of the distribution is a function of the stress or stresses applied to the system. Another way to describe it is that, for two different stresses and two different times,  $t_1$  and  $t_2$ , if the probability of failure at the times is the same.
- AFT is an Accelerated Failure Time model. This is simply a distribution where the time is multiplied by a function of covariates. This has the effect of ‘accelerating’ the time. Concretely, for a function  $f(t)$  it can be accelerated with a function to give  $f(\phi(x)t)$ .
- PH is a proportional hazard model. In a proportional hazard model, the hazard function is multiplied by some function of covariates. Hence if a function has a hazard rate of  $h(x)$  then the proportional hazard model will give simply  $\phi(x)h(t)$ .

SurPyval has implementations, and even a general constructor, for AFT, AL, and PH models. Each of which can handle arbitrary censoring (truncation coming).

### 1.4.3 How an AFT and PH Model Relate to a regular distribution

An AFT, or accelerated failure time, model does exactly that. It ‘accelerates’ the actual time by multiplying the time in the hazard function by a function of factors,  $\phi(x)$ . This factor can be any function. A Proportional Hazard model also does exactly what it says, if changes the hazard rate by a particular proportion.

Regular Distribution	$h(t) = h(t)$
Proportional Hazard	$h(t x) = \varphi(x).h(t)$
AFT	$h(t x) = h(\varphi(x).t)$

Given the relationship between variables and a distribution with either the PH or AFT models, you can see, using the above relationships that the survival, failure, and density functions can all be determined. This relationship is good to know to understand how AFT and PH models work.

#### 1.4.4 References

### 1.5 Data Wrangling Examples

Lets just say we have a list of right censored data and a list of failures. How can we wrangle these into data for the `fit()` method to accept?

```
import surpyval as surv

# Failure data
f = [2, 3, 4, 5, 6, 7, 8, 8, 9]
# 'suspended' or right censored data
s = [1, 2, 10]

# convert to xcn format!
x, c, n = surv.fs_to_xcn(f, s)
print(x, c, n)

model = surv.Weibull.fit(x, c, n)
print(model)
```

```
[ 1  2  2  3  4  5  6  7  8  9 10] [1 0 1 0 0 0 0 0 0 0 1] [1 1 1 1 1 1 1 2 1 1]
Parametric Surpyval model with Weibull distribution fitted by MLE yielding parameters_
→ (7.200723109183674, 2.474773882227539)
```

You can even bring in your left censored data as well:

```
# Failure data
f = [2, 3, 4, 5, 6, 7, 8, 8, 9]
# 'suspended' or right censored data
s = [1, 2, 10]
# left censored data
```

(continues on next page)

(continued from previous page)

```
l = [7, 8, 9]

# convert to xcn format!
x, c, n = surv.fsl_to_xcn(f, s, l)
print(x, c, n)

model = surv.Weibull.fit(x, c, n)
print(model)
```

```
[ 1  2  2  3  4  5  6  7  7  8  8  9  9 10] [ 1  0  1  0  0  0  0 -1  0 -1  0 -1  0  1]
→1] [1 1 1 1 1 1 1 1 1 1 2 1 1 1]
Parametric Surpyval model with Weibull distribution fitted by MLE yielding parameters
→(6.814750943874994, 2.4708983791967163)
```

Another common type of data that is provided is in a simple text list with “+” indicating that the observation was censored at that point. Using some simple python list comprehensions can help.

```
# Example provided data
data = "1, 2, 3+, 5, 6, 8, 10, 3+, 5+"

f = [float(x) for x in data.split(',') if "+" not in x]
s = [float(x[0:-1]) for x in data.split(',') if "+" in x]

data = surv.fs_to_xcn(f, s)

model = surv.Weibull.fit(*data)
```

```
Parametric Surpyval model with Weibull distribution fitted by MLE yielding parameters
→(6.737537377506333, 1.9245506420162473)
```

Again, this can be extended to left censored data as well:

```
data = "1, 2, 3+, 5, 6, 8, 10, 3+, 5+, 15-, 16-, 17-"
split_data = data.split(',')

f = [float(x) for x in split_data if ("+" not in x) & ("-" not in x)]
s = [float(x[0:-1]) for x in split_data if "+" in x]
l = [float(x[0:-1]) for x in split_data if "-" in x]

# Create the x, c, n data
data = surv.fsl_to_xcn(f, s, l)

model = surv.Weibull.fit(*data)
```

Surpyval also offers the ability to use a pandas DataFrame as an input. All you need to do is tell it which columns to look at for x, c, n, and t. Columns for c, n, and t are optional. Further, if you have interval censored data you can use the ‘xl’ and ‘xr’ column names instead. If you have mixed interval and observed or censored data, just make sure the value in the ‘xl’ column is the value of the observation or left or right censoring.

```
xr = [2, 4, 6, 8, 10]
xl = [1, 2, 3, 4, 5]
df = pd.DataFrame({'xl' : xl, 'xr' : xr})

model = surv.Weibull.fit_from_df(df)
print(model)
```

```
Parametric Surpyval model with Weibull distribution fitted by MLE yielding parameters_
↳ (4.694329418712716, 2.4106930022962714)
```

## 1.6 Datasets

### 1.6.1 Ball Bearing Failures

```
>>> from surpyval.datasets import Bearing

>>> Bearing.df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 23 entries, 0 to 22
Data columns (total 1 columns):
#   Column                                Non-Null Count  Dtype
---  ---                                -
0   Cycles to Failure (millions)         23 non-null     float64
dtypes: float64(1)
memory usage: 312.0 bytes
```

	Cycles to Failure (millions)
0	17.88
1	28.92
2	33
3	41.52
4	42.12

## 1.7 Non-Parametric Estimation

Non-parametric survival analysis is the attempt to capture the distribution of survival data without making any assumptions about the shape of the distribution. That is, non-parametric analysis, unlike parametric analysis, does not assume that the survival data was Weibull distributed or that it was Normally distributed etc. Concretely, non-parametric estimation does not attempt to estimate the parameters of a distribution, therefore “non-parametric.” Parametric analysis is covered in more detail in the section covering parametric estimation but it is important to contrast non-parametric estimation against what it is not. So what exactly is non-parametric analysis?

Survival analysis is using statistics to answer the question ‘what is the probability that the thing survived to a particular time?’ Non-parametric analysis answers this by estimating the probability from the proportion failed upto a given time. This can be done by either estimating the probability of surviving a particular segment or estimating the hazard rate.

Non-parametric estimation is well understood by appreciating the data format used to estimate the CDF. Specifically, the ‘xrd’ format and particularly understanding the r and d sets of that format.

The number of components at risk, r, at a given time, x, is the number of things at risk just prior to time x. The number of deaths, d, is the number of the at risk items that died (or failed) at time x. So for completely observed data the number at risk counts down for every death. So r would count down, e.g. 6, 5, 4, 3... for each death, 1, 1, 1, 1, ... So in this example there were 6 items at risk at one death at the first time. Then, because there was 1 death at the first time the number of items at risk has decreased to 5, therefore for the next death there are only 5 at risk. This continues further until there are no more items at risk because they have all died, i.e. there is 1 at risk and 1 death.

This can be extended to more than one death. For example, the risk set could be 8, 6, 5, 3, 2, 1. with an accompanying death set of 2, 1, 2, 1, 1, 1. In this example there were times where there were 2 deaths and therefore the number at



risk decreased by 2 after that number of deaths.

So a complete example of this format is:

```
x = [1, 2, 3, 4, 5, 6]
r = [7, 5, 4, 3, 2, 1]
d = [2, 1, 1, 1, 1, 1]
```

This format for data is not how survival data is usually provided in text books or papers. Survival data is usually displayed with the simple list of failure times such as “1, 3, 6, 7, 10, 16”. The first step surpyval does for non-parametric analysis is to transform data into the xrd format. All the `fit()` methods for surpyval take as input the xcnt format, see more at the data types docs. So if you provide surpyval with the data “1, 2, 3, 4, 5, 6” it will assume that each of them are one death, and then create the risk set from the death counts resulting in the xrd format from above.

Given we now understand the format of the data we can estimate the probability of survival to some time with non-parametric methods. The first method we will visit is the Kaplan-Meier.

### 1.7.1 Kaplan-Meier Estimation

Kaplan-Meier [KM] is a very popular method for estimating survival curves for populations. The insight for this method is that for each time there is a death, we can estimate the probability of having survived since the previous deaths. Using the data from above as an example, at time 1, there are 7 items at risk and there are 2 deaths. We can therefore say that the probability of surviving this period was  $(7 - 2)/7$ , i.e.  $5/7$ . Then the next time there is a death, the probability of having survived that extra time is  $(5 - 1)/5$ , i.e.  $4/5$ .

To be clear, this is the chance of survival between each death. Therefore the chance of surviving up to a given time is the chance of surviving each segment. Therefore the probability of surviving up to any given time is the probability of surviving through all the previous segments. The probability of surviving multiple outcomes is the multiplication of each of the survival probabilities. Surviving through three sections is equal to the probability that I survive the first, then multiply this by the probability of surviving the second, then multiplying this result with the probability of surviving the third. So continuing our example from above, the probability of surviving the first two segments is  $(5/7) \times (4/5) = 4/7$ .

Therefore using the at risk count,  $r$ , and the death count,  $d$ , can be used to estimate the segment survival probabilities and the survival probability to any point can be found by multiplying these probabilities. Formally, this has the following formula:

$$R(x) = \prod_{i: x_i \leq x} \left(1 - \frac{d_i}{r_i}\right)$$

### 1.7.2 Nelson-Aalen Estimation

The Nelson-Aalen estimator [NA] (also known as the Breslow estimator), instead of finding the probability, estimates the cumulative hazard function, and given that we know the relationship between the cumulative hazard function and the reliability function, the Nelson-Aalen cumulative hazard estimate can be converted to a survival curve.

The first step in computing the NA estimate is to convert your data to the  $x, r, d$  format. Once in this format the instantaneous hazard rate is found by:

$$h(x) = \frac{d_x}{r_x}$$

This estimate of the instantaneous hazard rate is the proportion of deaths/failures at a value,  $x$ . Then to find the cumulative hazard rate for any  $x$  we simply take the sum of the instantaneous hazard rates for all the values below  $x$ .

Mathematically:

$$H(x) = \sum_{i: x_i \leq x} \frac{d_i}{r_i}$$

Then, since we know that the reliability, or survival function, is related to the cumulative hazard function, we can easily compute it.

$$R(x) = e^{-H(x)}$$

So we now have the survival/reliability function. One benefit of the Nelson-Aalen estimator is that it does not estimate a probability of 0 for the highest value (in a completely observed data set). This means that for a completely observed data set the whole estimation can be plotted on a transformed y-axis. For this reason SurPyval uses the Nelson-Aalen as the default plotting position.

### 1.7.3 Fleming-Harrington Estimation

The Fleming-Harrington estimator [FH], uses the same principal as the Nelson-Aalen estimator. That is, it finds the cumulative hazard function and then converts that to the reliability/survival estimate. However, the NA estimate assumes, for any given step that the number of items at risk is equal for each death, the FH estimate changes this. Mathematically, the hazard rate is calculated with:

$$h(x) = \frac{1}{r_x} + \frac{1}{r_x - 1} + \frac{1}{r_x - 2} + \dots + \frac{1}{r_x - (d_x - 1)}$$

Which can be summarised as:

$$h(x) = \sum_{i=0}^{d_x-1} \frac{1}{r_x - i}$$

The cumulative hazard rate therefore becomes:

$$H(x) = \sum_{i: x_i \leq x} \sum_{i=0}^{d_x-1} \frac{1}{r_x - i}$$

You can see that the cumulative hazard rate will be slightly higher than the NA estimate since:

$$\frac{1}{r_x} + \dots + \frac{1}{r_x} \leq \frac{1}{r_x} + \dots + \frac{1}{r_x - (d_x - 1)}$$

The above is less than or equal for the case where there is one death/failure. The Fleming-Harrington and Nelson-Aalen estimates are particularly useful for small samples, see [FH].

### 1.7.4 Turnbull Estimation

The Turnbull estimator is a remarkable non-parametric estimation method for data that can handle arbitrary censoring and truncation [TB]. The Turnbull estimator can be found with a procedure of finding the most likely survival curve from the data, for that reason it is also known as the Non-Parametric Maximum Likelihood Estimator. The Kaplan-Meier is also known as the Maximum Likelihood estimator, so is there a contradicton? No, the Turnbull estimator is the same as the Kaplan-Meier for fully observed data.

The Turnbull estimate is really an estimate of the observed failures given censoring, and then the ‘ghost’ failures (as Turnbull describes it) due to truncation. Turnbull’s estimate converts all failures to interval failures regardless of the censoring. This is because a left censored point is equivalent to an intervally censored observation in the interval -Inf to x, and a right censored point is equivalent to an intervally censored observation in the interval x to Inf. Then for all

the intervals between negative infinity we find how many failures happened in that interval. This value need not be a whole number since a single observation could have failed across several intervals. To estimate the failures, we use:

$$\mu_{ij}(s) = \frac{\alpha_{ij}s_j}{\sum_{k=1}^m \alpha_{ik}s_k}$$

Where  $\mu_{ij}$  is the probability of the i-th observation failing in the j-th interval,  $\alpha_{ij}$  is a flag to indicate if the i-th failure was at risk in interval j, (1 if at risk and 0 if not), and  $s_j$  is the probability of failure in an interval. That is,  $s_j$  is the survival function we are trying to estimate.

If an observation is truncated, it was only a possible observation among others that would have been seen had the observation not been limited. To estimate the additional at risk items outside of the domain for which an observation is truncated we use:

$$\nu_{ij}(s) = \frac{(1 - \beta_{ij})s_j}{\sum_{k=1}^m \alpha_{ik}s_k}$$

Where  $\nu_{ij}$  is the probability of the i-th observation failing in the j-th interval and  $\beta_{ij}$  is a flag to indicate if the i-th failure was observable in interval j, (1 if at risk and 0 if not).

This formula then finds the number of failures outside the truncated interval for a given observation.

We can then estimate the probability of failure in a given interval using the total failures in each interval divided by the total number of failures:

$$s_j = \frac{\sum_{i=1}^N \mu_{ij} + \nu_{ij}}{M(s)}$$

where

$$M(s) = \frac{\sum_{i=1}^N \sum_{j=1}^m \mu_{ij} + \nu_{ij}}{M(s)}$$

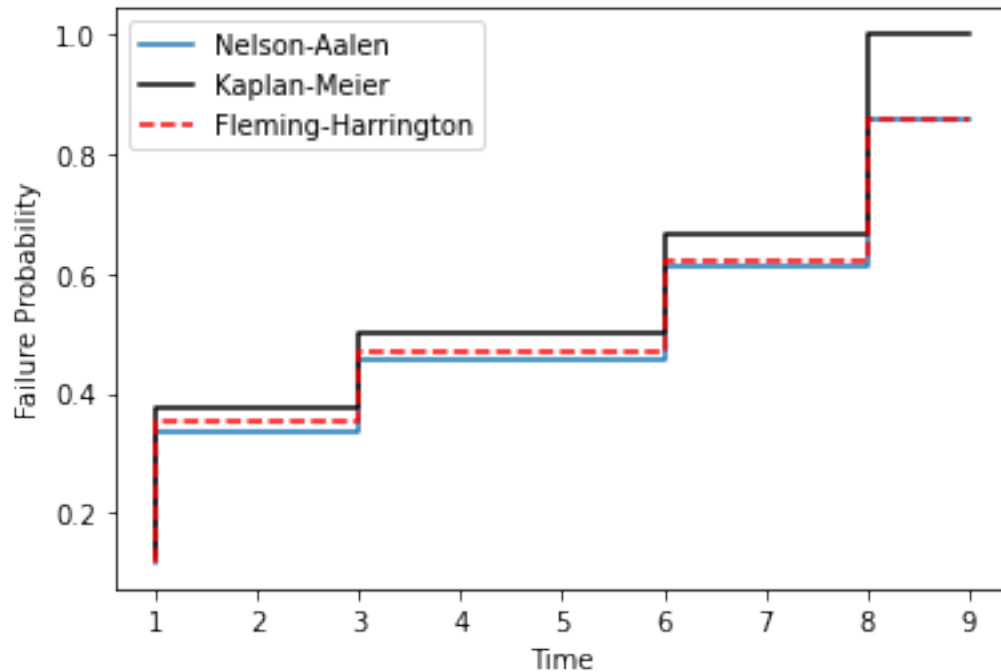
Using this estimate of the survival function, it can be input to the start of this procedure and it done again. This can then be repeated over and over until the values do not change. At this point we have reached the NPMLE estimate of the survival function!

The Turnbull estimation is the only non-parametric method that can be used with truncated and left censored data. Therefore it must be used when using the plotting methods in the parametric package when you have truncated or left censored data.

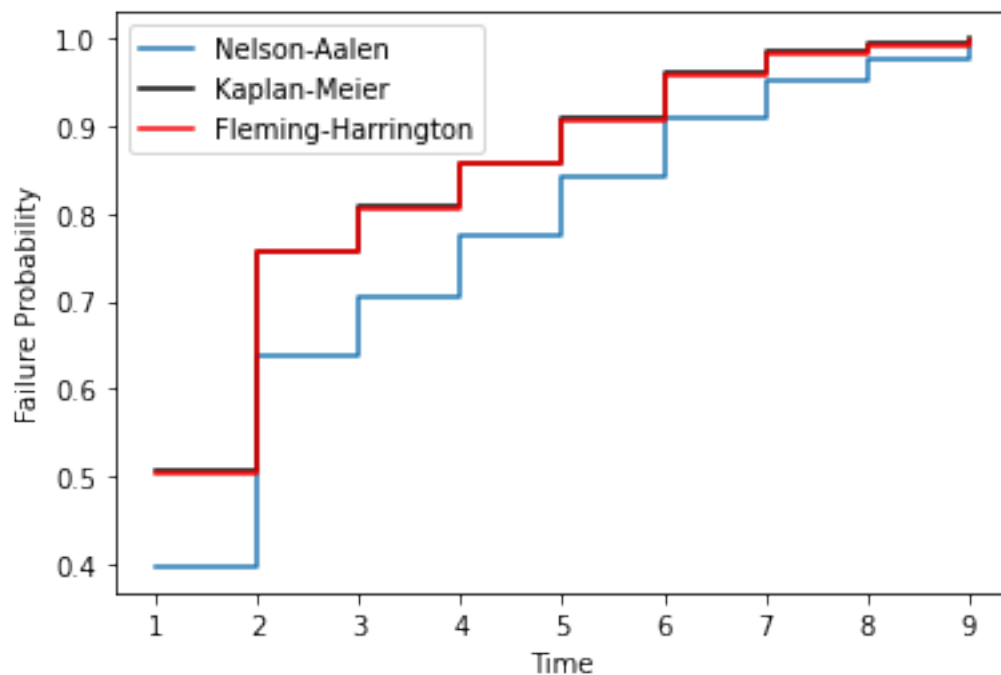
### 1.7.5 On Surpyval's Default

Surpyval uses the Fleming-Harrington estimator as the default. The rationale for this is because it has optimal behaviour. That is, it performs well where the Kaplan-Meier and the Nelson-Aalen behave poorly.

The Kaplan-Meier, since it tends to 1, results in cases where it overstates the probability of failure. It is because of this that the Kaplan-Meier should not be used in circumstances of competing risks. As an example, a comparison between a Nelson-Aalen and Kaplan-Meier estimate over time (I have plotted the Fleming-Harrington estimate for later discussion):



On the contrary, the Nelson-Aalen estimate performs poorly with lots of ties. This results in the Nelson-Aalen estimator overstating the risk for lower failure times. This is in contrast to the Kaplan-Meier estimator which does well with lots of tied values. For example:



The Fleming-Harrington, plotted in red in the above two charts, optimises between these two estimators. The Fleming-Harrington estimate approaches the Nelson-Aalen under the conditions of where the Nelson-Aalen estimate performs well and the Kaplan-Meier does poorly. Fleming-Harrington also does well where the Nelson-Aalen estimate does poorly but the Kaplan-Meier does well. Although the two examples provided are in the extreme, it is worth using the Fleming-Harrington by default since it is more flexible; it is therefore, for this reason, that surpyval does exactly that.

This is not to say not to use KM or NA, but only when you are sure you are making the correct assumptions about what you are doing!

### 1.7.6 References

## 1.8 Parametric Estimation

Parametric modelling is the process of estimating the parameters of a particular distribution from a set of data. This is distinct from non-parametric modelling where we make no assumptions about the shape of the distribution. In parametric modelling we make some assumptions, explicit or implied, about the shape of the data that we have.

For this segment I will use the Weibull distribution as the example distribution. The Weibull distribution is a very useful distribution for one interesting reason. It is the distribution for the ‘weakest link.’ As the normal distribution is the limiting distribution of averages, the Weibull distribution is the limiting distribution for minimums. What does that mean? If we have a large number of sets of samples from something that is normally distributed the average of these sets will also be normally distributed but the minimums of these sets of samples will be Weibull distributed. This is analogous to a chain. It is common wisdom that a chain is only as strong as its weakest link. The Weibull distribution enables us to model the strength of a chain based on the strength of the links.

The Weibull distribution can then be used in scenarios where we assume that the shape of the distribution will be due to a weakest link effect. This assumption holds in many scenarios, the strength of materials, the fielded life of equipment, the lifetime of animals, the time until another recession, or the time until germination of seeds. This example makes clear the assumption that we can make when using the Weibull distribution. Other distributions have differing processes that can result in their generation. If we know and understand these processes we can check them against the scenario we are analysing and choose a distribution from them. For example, a lognormal distribution can arise due to the combined effect of the product of random variables so in petroleum engineering the total recoverable oil is a product of the height, width, depth, features of the rock and an infinitude of other variables of the field. Therefore fields can be lognormally distributed. Similar considerations can be applied for many other types of distributions. Finally, If we don’t know, or mind, what distribution we have, we can simply find the best fit amongst a set of distributions.

SurPyval offers users several methods for estimating parameters, these are:

- Method of Moments (MOM)
- Method of Probability Plotting (MPP)
- Mean Square Error (MSE)
- Maximum Likelihood Estimation (MLE)
- Minimum Product Spacing (MPS)

There are other methods that can be used, e.g. L-moments or generalised method of moments. These are interesting, and may be added in future, but for now surpyval offers the above estimation methods. Surpyval is unique in the capability to provide the estimation technique. Most other survival analysis methods do not allow for specifying different methods. The advantage of this flexibility will become apparent.

### 1.8.1 Method of Moments (MOM)

This method is the simplest (and least accurate) method to find parameters of a distribution. The intent of the Method of Moments (MOM) is to find the closest match of a distribution’s moments, to those of the moments of a sample of data.

For a given data set, or sample, the  $k$ th moment is defined as:

$$M_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

If the distribution has only one parameter, like the exponential distribution, then the method of moments is simply equates the sample moment to the distribution moment. For a continuous distribution the  $k$ th moment is defined as:

$$M_k = \int_{-\infty}^{\infty} x^k f(x) dx$$

Where  $f(x)$  is the density function of that distribution. Therefore, for the exponential distribution, the moments can be computed (with some working) to be:

$$E[X^k] = \frac{k!}{\lambda^k}$$

Because there is only one parameter of the exponential distribution, we need to only match the first moment of the distribution ( $k=1$ ) to the first moment of the sample. Therefore we get:

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

This is to say that the method of moments solution for the parameter of the exponential is simply the inverse of the average. This is an easy result. When we extend to other distributions with more than one parameter, such simple analytical solutions are not available, so numeric optimisation is needed. SurPyval uses numeric optimisation to compute the parameters for these distributions.

The method of moments, although interesting, can produce incorrect results, and it can only be used with observed data, so it cannot account for truncation or censoring. But it is good to understand as it is one of the oldest methods used to estimate the parameters of a distribution.

## 1.8.2 Method of Probability Plotting (MPP)

Probability plotting is an extremely simple way to find the parameters of a distribution. This method has a long history because it is a simple activity to do while providing an easy to understand graphic. Further, probability plotting produces a good estimate for the parameters even with few data points. All this combined with the fact that probability plotting can be used for all types of data, observed, censored, and truncated, it is easy to understand why it is widely used.

SurPyval uses the MPP method as an initial guess, when not provided, because it is the only method that does not require an initial guess of the parameters. This is because numeric optimisers require an initial guess, however, when using a probability plotting method, an initial guess is not needed. It therefore provides an excellent method to get an initial guess for subsequent optimisation. But the method itself can be sufficient enough for the majority of applications.

So how does it work?

Probability plotting works of the idea that a distributions CDF can be made into a straight line if the data is transformed. This can be shown by rearranging the CDF of a distribution. For the Weibull:

$$F(x) = 1 - e^{-\left(\frac{x}{\alpha}\right)^\beta}$$

If we negate, add one, and then take the log of each side we get:

$$\ln(1 - F(x)) = -\left(\frac{x}{\alpha}\right)^\beta$$

Then take the log again:

$$\ln(-\ln(1 - F(x))) = \beta \ln(x) - \beta \ln(\alpha)$$

From here, we can see that there is a relationship between the CDF and  $x$ . That is, the log of the log of  $(1 - \text{CDF})$  has a linear relationship with the log of  $x$ . Therefore, if we take the log of  $x$ , and take the log of the negative log of 1 minus the CDF and plot these, we will get a straight line. To make this work, we therefore need a method to estimate the CDF empirically. Traditionally, there have been heuristics used to create the CDF. However, we can also use the non-parametric estimate as discussed in the non-parametric session. Concretely, we can use the Kaplan-Meier, the Nelson-Aalen, Fleming-Harrington, or Turnbull estimates to approximate the CDF,  $F(x)$ , transform it, plot, and then do the linear regression. SurPyval uses as a default, the Nelson-Aalen estimator for the plotting point.

Other methods are available. The simplest estimate, for complete data, is the empirical CDF:

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n 1_{X_i \leq x}$$

This equation says, that (for a fully observed data set) for any given value,  $x$ , the estimate of the CDF at that value is simply the sum of all the observations that occurred below that value divided by the total number of observations. This is a simple percentage estimate that has failed at any given point. This equation will therefore make a step function that increases from 0 to 1.

One issues with this is that the highest value is always 1. But if this is transformed as above, this will be an undefined number. As such, you can adjust the value with a simple change:

$$\hat{F}(x) = \frac{1}{n+1} \sum_{i=1}^n 1_{X_i \leq x}$$

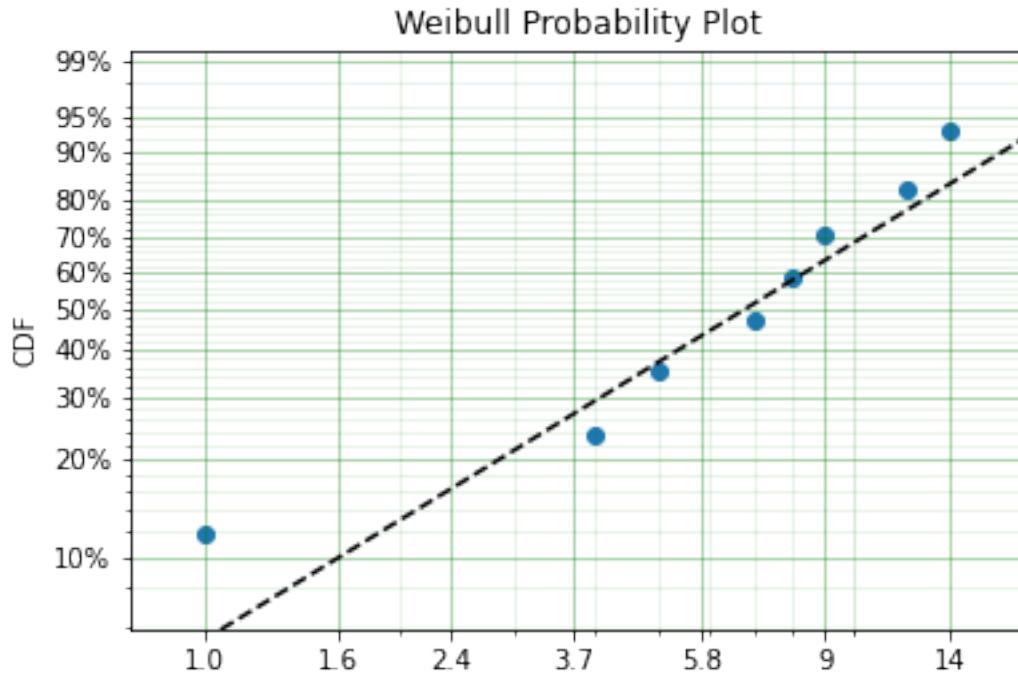
By using this simple change, the highest value will not be 1, and will therefore be plottable, and not undefined. There are many different methods used to adjust the simple ECDF to be used with a plotting method to estimate the parameters of a distribution. For example, consider Blom's method:

$$\hat{F}_k = (k - 0.375)/(n + 0.25)$$

Where  $k$  is the rank of an observation  $k$  is in  $(1, 2, 3, 4, \dots, n)$  for  $n$  observations. Using these methods we can therefore plot the linearised version above.

Combining this all together is simple witht surpyval.

```
x = [1, 4, 5, 7, 8, 9, 12, 14]
model = surv.Weibull.fit(x, how='MPP', heuristic='Blom')
model.plot()
```



In this example we have used the probability plotting method with the Blom heuristic to estimate the parameters of the distribution. SurPyval has the option to use many different plotting methods, including the regular KM, NA, and FH non-parametric estimates. All you need to do is change the ‘heuristic’ parameter; SurPyval includes:

Table 1: SurPyval Modelling Methods

Method	A	B
Blom	0.375	0.25
Median	0.3	0.4
ECDF	0	0
ECDF_Adj	0	1
Mean	0	1
Weibull	0	1
Modal	1	-1
DPW	1	0
Midpoint	0.5	0
Benard	0.3	0.2
Beard	0.31	0.38
Hazen	0.5	0
Gringorten	0.44	0.12
Larsen	0.567	-0.134
Larsen	1/3	1/3
None	0	0

Which is used with the general formula to estimate the plotting position heuristic:

$$\hat{F}_k = (k - A)/(n + B)$$

One final option available is that of the Filliben estimate:



### 1.8.3 Mean Square Error (MSE)

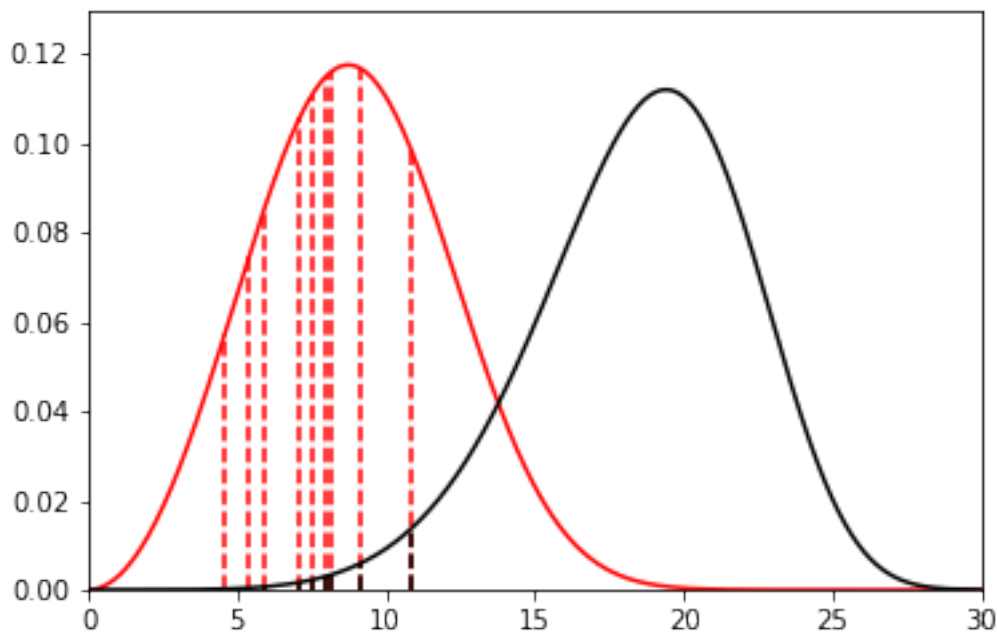
MSE is essentially the same as probability plotting. Instead of finding the minimum against the transformed data in the x and y axes. The parameters are found by minimising the distance to the non-parametric estimate without transforming the data to be linear. Mathematically, MSE find the parameters by minimising:

$$\Sigma \left( \hat{F} - F(x; \theta) \right)^2$$

This is the difference between the, untransformed, empirical estimate of the CDF and the parametric distribution.

### 1.8.4 Maximum Likelihood Estimation (MLE)

Maximum Likelihood Estimation (MLE) is the most widely used, and most flexible of all the estimation methods. It's relative simplicity (because of modern computing power) makes it the reasonable first choice for parametric estimation. What does it do? Essentially MLE asks what parameters of a distribution are 'most likely' given the data that we have seen. Consider the following data and distributions:



The solid lines are the densities of two different Weibull distributions. The dashed lines represent the data we have observed, their height is the density of the two distributions at the x value for each observation. Given the data and the two distributions, which one seems to explain the distribution of the data better? That is, which distribution is more likely to produce, if sampled, the dashed lines? It should be fairly intuitive that the red distribution is more likely to do so. For example, the observation just above 10, you can see the height to the black line and the height to the red line. The red line is taller than the black line, therefore this observation is more 'likely' to have come from the red distribution than the black one. Conversely, the value near 15 is more likely to have come from the black distribution than the red one because the height to the black line is greater than the height to the red line. To find the distribution of best fit then we need to find the parameters that best averages the height of all these lines.

MLE formalises this concept by saying that the most likely distribution is the one that has the highest (geometric) mean of the height of the density function for each sample of data. The height of the density at a particular observation is known as the likelihood. Mathematically, (for uncensored data) MLE then maximises the following:

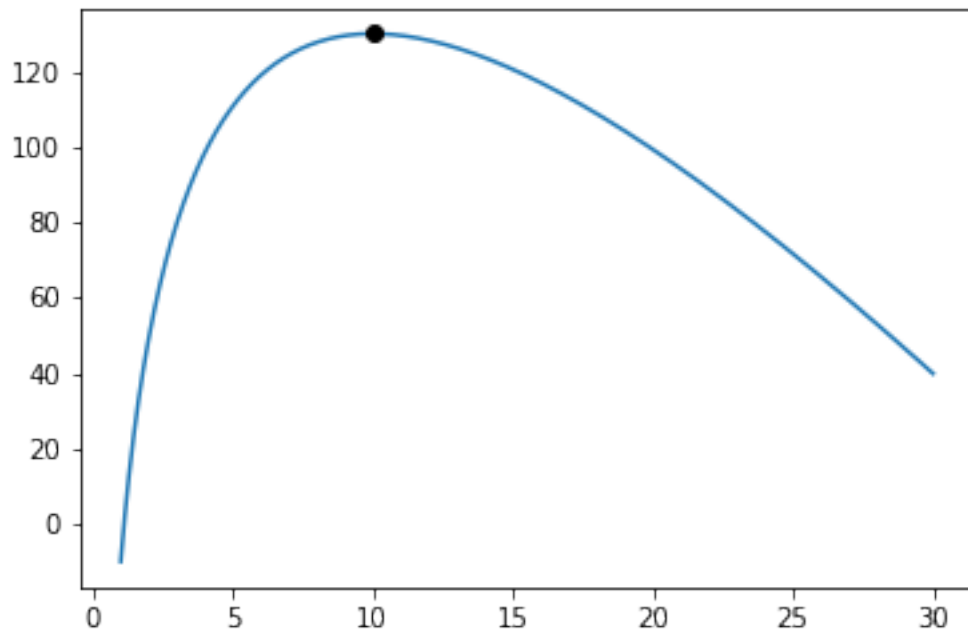
$$L = \left( \prod_{i=1}^n f(x_i|\theta) \right)^{1/n}$$

$f$  is the pdf of the distribution being estimated,  $x$  is the observed value,  $\theta$  is the parameter vector, and  $L$  is the geometric mean of all the values. This is complicated, but a simplification is available by taking the log of this product yielding:

$$l = \frac{1}{n} \sum_{i=1}^n \ln f(x_i|\theta)$$

Therefore MLE simply finds the parameters of the distribution that maximise the average of the log of the likelihood for each point. . . One final transform that is used in optimisers is that we take the negative of the above equation so that we find the minimum of the negative log-likelihood.

Armed with the log likelihood we can then search for the parameter where the log likelihood is maximised. Using an Exponential distribution as an example, we can see the change in the value of the log likelihood as the exponential parameter changes. The following is a random sample of 100 observations with a parameter of 10. Then changing the value of the parameter ‘lambda’ from low to high we can see what the log-likelihood is and find the value at which it is maximized.



On the chart above you can see that the maximum is near 10. As we would expect given that we know that the answer is 10. It is this simple and intuitive approach that allows the parameters of distributions are estimated with the MLE.

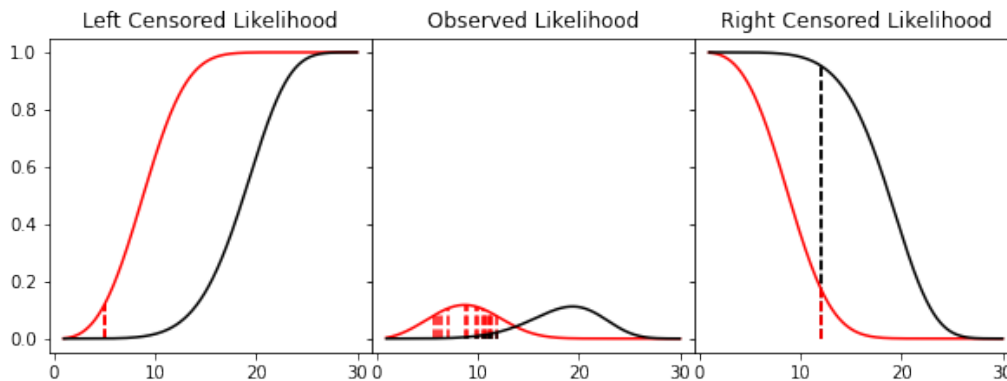
What about censored data?

All the equations above are for observed data. Handling the likelihood of censored data also has an intuitive understanding. What we know about the point when the data point is censored is that we know it is above or below the

value at which we observed. So for a right censored data point, we want to maximize the probability that we saw this observation, concretely we want a censored points contribution to the likelihood function is the probability that the point was left or right censored. This is simply the probability of failure (CDF) for left censored and the probability of surviving to that time (survival function). Formally:

$$l = \frac{1}{n_o} \sum_{i=1}^{n_o} \ln f(x_{o_i} | \theta) + \frac{1}{n_r} \sum_{i=1}^{n_r} \ln R(x_{r_i} | \theta) + \frac{1}{n_l} \sum_{i=1}^{n_l} \ln F(x_{l_i} | \theta)$$

An easy and intuitive way to understand this is to compare these two possibilities. With some randomly generated data with a few values made to be left censored, and a few to be right censored. We get:



In this example, again, we need to consider whether the red or black distribution is a more likely description of the observations, including some censored ones. Although the right censored point for the black distribution is very likely, this does not mean it is a good fit because the ‘average’ across all observations is poor. Therefore, it should be obvious that the red distribution is the better fit.

But what about truncated data

## 1.8.5 Maximum Product of Spacings (MPS)

Coming soon

## 1.9 Regression Analysis

The time until some event happens will, almost certainly, be impacted by factors. For example, when considering how long a machine will last before failure an engineer will want to account for the operational conditions. It may operate in a humid environment, or it may operate at a higher rate. The question is then, how do we account for these variations, or ‘covariates’, on the time until failure?

Regression analysis is the process of capturing the effect that covariates have on the item. That is, we use data on other factors to ‘regress’ onto the survival distribution. The purpose of this type of regression is so that you can ask, and answer, questions like “what effect will increasing X have on the survival time?”

**Surpyval covers three types of regression models. These are:**

- Proportional Hazards,
- Accelerated Time, and
- Accelerated Life.

There are special cases when these are the same, however, it is important to understand the difference between them in general. I detail the differences in the following sections.

### 1.9.1 Proportional Hazards Model

A proportional hazards model is one in which we change the hazard rate of the distribution by some proportional amount. You may recall that every distribution can be defined by a hazard rate or a cumulative hazard rate, see the “Handy References” section which shows that the density, CDF, and survival function can all be defined in terms of the hazard rate,  $h(t)$ .

So what we can do then is assume that the covariates will affect the survival time of the thing by having some effect on the hazard rate. The general definition for a proportional hazard model is:

$$h(t|X) = \phi(X)h_0(t)$$

This is to say that the hazard rate at time  $t$  is the function (of a vector) of covariates on a ‘baseline’ hazard rate. Let’s use a simple example, a proportional hazard model with covariates that affect a constant hazard rate. Let’s say that some factory produces one widget an hour. But this is only with one machine in operation, if we add a second machine, we can produce widgets at two per hours, if we had a third, it will be three per hour. In this case the base rate is 1 and the function linking  $X$  to the base rate is to simply multiply  $X$  by the baserate.

This is to say that for this example:

$$\begin{aligned}\phi(X) &= X \\ h_0(t) &= 1\end{aligned}$$

Therefore:

$$h(t|X) = X$$

This is an overly simple model, but it shows how we can construct a PH model.

In this case we have a simple proportional hazard model, also, it is limited to only an increasing hazard rate, but sometimes we need to capture a negative impact. Further, we may need a way to capture more covariates. For these reasons a very common selection for the function of covariates is an exponential function.

$$\phi(X) = e^{X \cdot \beta}$$

Where

$$X \cdot \beta = X_0\beta_0 + X_1\beta_1 + \dots + X_{n-1}\beta_{n-1} + X_n\beta_n$$

In this case the proportional term is the  $e$  raised to the power of the cross product of  $X$  and  $\beta$ . Using this as the covariate function is a very common choice. This is because it will not ever become negative. It can capture situations where a covariate will increase the hazard rate if it’s coefficient,  $\beta$ , is positive, and it will decrease the hazard rate if it’s coefficient is negative. Also, the dot product can capture a varying number of covariates with ease. For these reasons the Cox model is a widely used. Although you can choose any function for your covariates there is already likely literature about your problem which might indicate which function to use.

Survpyval uses MLE to estimate the parameters for proportional hazards models. This is a simple conversion from regular MLE since we know the relationship between a baseline distribution and the proportional hazards version. These relationships are:

$$f(t|X) = \phi(X)h_0(t)e^{-\phi(X)H(t)}$$

$$F(t|X) = 1 - e^{-\phi(X)H(t)}$$

$$S(t|X) = e^{-\phi(X)H(t)}$$

It is therefore relatively simple to adjust the MLE methods to accommodate proportional hazard models.

The details on fitting proportional hazards model is detailed more in the surpyval analysis section.

## Semi-Parametric

The previous sections covered ‘parametric’ and ‘non-parametric’ survival models, so what is ‘semi-parametric’? A semi-parametric model is a survival model with a non-parametric baseline and parametric function that affects that baseline. Recall that a proportional hazard model can be defined as:

$$h(t|X) = \phi(X)h_0(t)$$

It is interesting to note that the phi term must be parametric, however, the baseline hazard rate need not be parametric, it can be non-parametric! Therefore, what we have is a parametric relationship of the covariates to the baseline hazard rate, but a non-parametric baseline hazard rate, therefore, a ‘semi-parametric’ model.

By far the most common of any regression model of any kind (parametric, non-parametric, and semi-parametric of all the accelerated life, proportional hazard, and accelerated time) is the Cox Proportional Hazard model, it is a semi-parametric model.

The Cox model is used in a wide variety of fields. It has been used in criminology to study the recidivism of parolees, in engineering to understand the factors affecting tire reliability, and in medical science to understand factors affecting cancer and other diseases, among many many other applications. The wide use of the model shows the utility the model has and the broad applicability to solve problems.

## 1.9.2 Accelerated Time

An accelerated time model is very similar to a proportional hazards model. The difference is where the function is applied; instead of multiplying the hazard function, and accelerated time model multiplies the time by the function of covariates. The general definition is:

$$f(t|X) = f(\phi(X)t)$$

It is called an accelerated time since the time term is transformed by the covariates, i.e. time is ‘accelerated’ by the covariates.

$$t_a = \phi(X)t$$

Just like proportional hazards, there are simple transformations that apply

$$f(t|X) = f(\phi(X)t)$$

$$F(t|X) = F(\phi(X)t)$$

$$S(t|X) = S(\phi(X)t)$$

Given the simple transformation of the time term the MLE is feasible with an additional transformation step. This is how surpyval estimates the parameters.

## 1.9.3 Accelerated Life

An accelerated life model is, in many cases, simply the inverse of an accelerated time model. However, there are some cases where they are different. Consider an accelerated life model with a normal distribution:

$$F(t|X) = \Phi\left(\frac{\phi(X)t - \mu}{\sigma}\right)$$

Where  $\Phi$  is the CDF of the standard normal distribution. In this case  $\mu$  is the expected life of the model, however, we may instead be interested in determining what effect covariates have on the expected life of an item. In this case we can simply substitute the expected life:

$$F(t|X) = \Phi\left(\frac{t - \phi(X)}{\sigma}\right)$$

An accelerated life model is, therefore, simply a model where the life parameter of a distribution is substituted with a function of the covariates, that is, it ‘accelerates’ the expected life, as opposed to accelerating time as per an accelerated time model. For each of the distributions in Surpyval their life parameter that varies is as per the following table:

Distribution	Life Param
Weibull	alpha
Exponential	1./lambda
Normal	mu
LogNormal	mu
Gamma	alpha
Gumbel	mu
Logistic	mu
LogLogistic	alpha
ExpoWeibull	Not Avail
Uniform	Not Avail
Beta	Not Avail

Given the simple substitution into the life parameter, surpyval uses MLE to calculate the parameters.

For examples on how to do regression analysis, see the entry in the ‘SurPyval Analysis’ section of the docs.

## 1.10 Non-Parametric SurPyval Modelling

To get started, let’s import some useful packages, as such, for the rest of this page we will assume the following imports have occurred:

```
import surpyval as surv
import numpy as np
from matplotlib import pyplot as plt
```

Survival modelling with *surpyval* is very easy. This page will take you through a series of scenarios that can show you how to use the features of *surpyval* to get you the answers you need. The first example is if you simply have a list of event times and need to find the distribution of best fit.

In each of the examples below, each of the `KaplanMeier`, `NelsonAalen`, or `FlemingHarrington` can be substituted with any of the others. It is the choice of the analyst which should be used. The `Turnbull` estimator has additional capabilities that can be used when you have right truncated, left censored, or interval censored data.

### 1.10.1 Complete Data

Using data of the stress of Bofors steel from Weibull’s original paper we can estimate the reliability, that is, the probability that a sample of steel will survive up to a given applied stress. So what does that mean?

We can find when the steel will break. This is particularly useful when we know the application.

For this example, let’s say that the maximum tensile stress our design will see during use is 34 units. Let’s try and estimate the proportion that will fail during operation.

For this we can use the Nelson-Aalen estimator of the hazard rate, then convert it to the reliability. This is all done with one easy call.

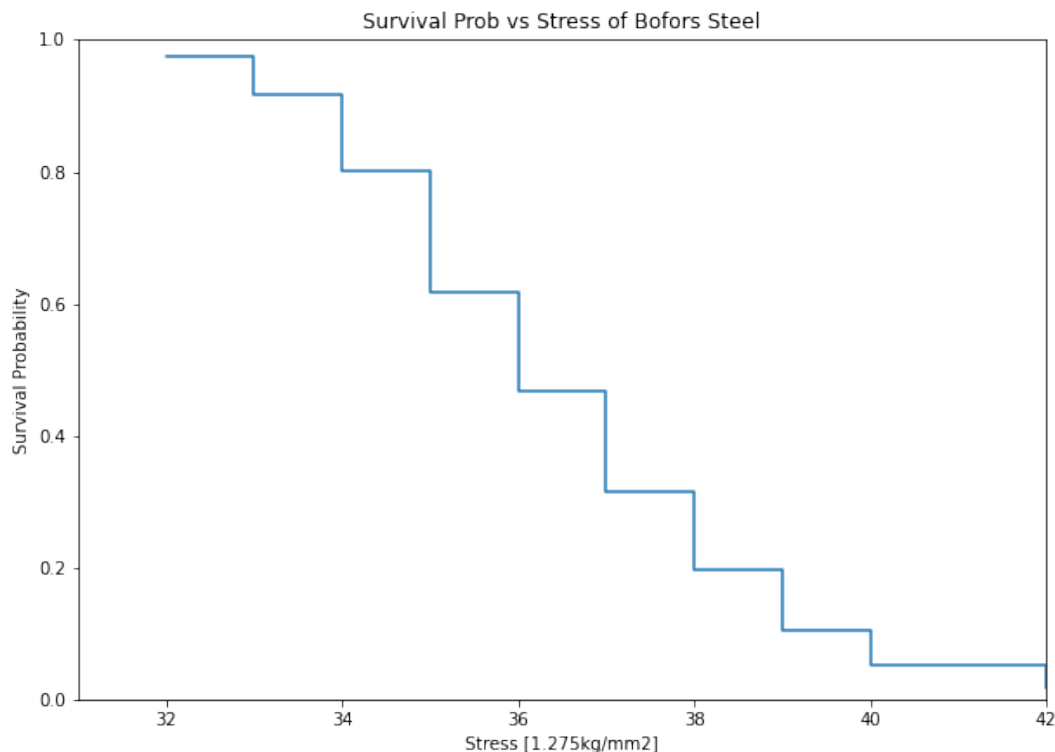
```
import surpyval as surv
import numpy as np
from matplotlib import pyplot as plt

x = np.array([32, 33, 34, 35, 36, 37, 38, 39, 40, 42])
n = np.array([10, 33, 81, 161, 224, 289, 336, 369, 383, 389])

# Weibull's measurements are cumulative so we need to transform them
n = np.concatenate([[n[0]], np.diff(n)])

bofors_steel_na = surv.NelsonAalen.fit(x, n=n)

plt.figure(figsize=(10, 7));
plt.ylabel('Survival Probability')
plt.xlabel('Stress [1.275kg/mm2]')
plt.ylim([0, 1])
plt.xlim([31, 42])
plt.step(bofors_steel_na.x, bofors_steel_na.R)
plt.title('Survival Prob vs Stress of Bofors Steel');
```



So what purpose is this?

With our non-parametric model of the Bofors steel. We can use this model to estimate the reliability in our application. Let's say that our application uses Bofors steel up to 34. What is our estimate of the number of failures?

```
print(str(bofors_steel_na.sf(34).round(4).item() * 100) + "%")
```

Which gives:

```
80.15%
```

The above shows that approximately 80% will survive up to a stress of 34. Therefore we will have an approximately 20% chance of our component failing in the design.

It is up to the designer to determine whether this is acceptable.

What if we want to take into account our uncertainty about the reliability. The non-parametric class automatically computes the Greenwood variance and uses that to compute the upper and lower confidence intervals. Let's plot the intervals to see.

```
plt.figure(figsize=(10, 7))
bofors_steel_na.plot(how='interp')
plt.xlabel('Stress [1.275kg/mm2]')
plt.ylabel('Survival Probability')
plt.ylim([0, 1])
plt.xlim([32, 42])
plt.title('Surv Prob vs Stress of Bofors Steel')
```



The confidence bounds can also be used to estimate the probability of survival up to some point with some degree of confidence. For example:

```
print(str(bofors_steel_na.R_cb(34, bound='lower', how='interp', confidence=0.95).
        .round(4).item() * 100) + "%")
```

(continues on next page)



(continued from previous page)

76.46%

Therefore we can be 95% confident that the reliability at 34 is above 76%. You can also see that the confidence interval stretches the entire span of the possible  $[0, 1]$  interval at the highest value. This is because the variance at the final value is infinite using the Greenwood confidence interval.

## 1.10.2 Right Censored Data

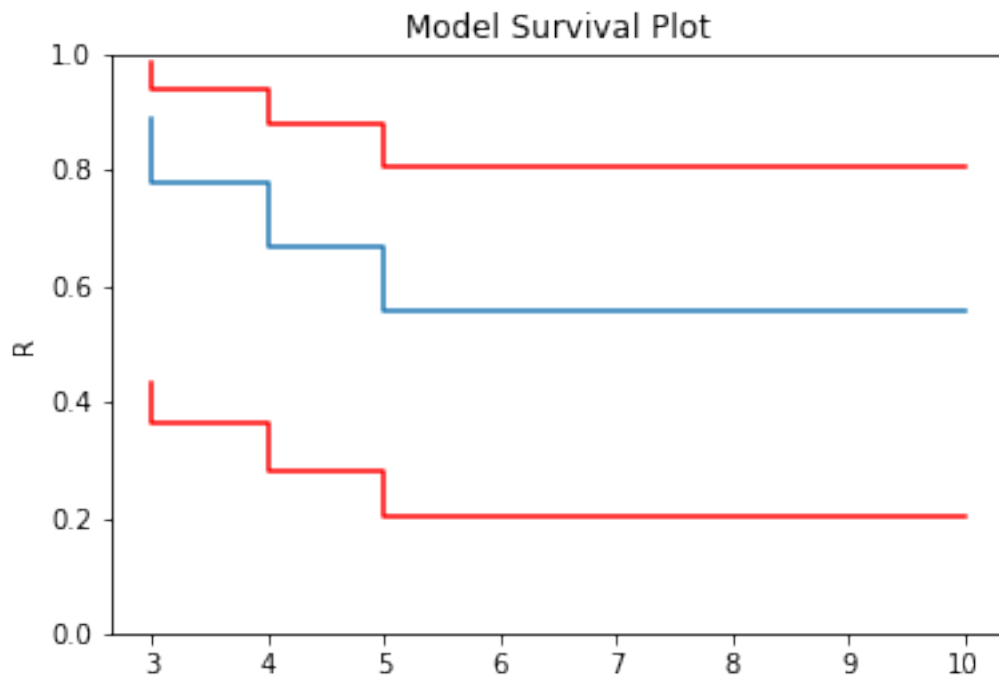
Non-Parametric estimation can handle right censored, this is possible because at the point of censoring the item is removed from the at risk group without counting a death/failure.

```
import numpy as np
from surpyval import KaplanMeier as KM

x = np.array([3, 4, 5, 6, 10])
c = np.array([0, 0, 0, 0, 1])
n = np.array([1, 1, 1, 1, 5])

model = KM.fit(x=x, c=c, n=n)
model.plot()
model.R
```

```
array([0.88888889, 0.77777778, 0.66666667, 0.55555556, 0.55555556])
```



In this example, we have included right censored data. This example can be done for the Nelson-Aalen, Fleming-Harrington, and Turnbull estimators as well.

### 1.10.3 Left Truncated Data

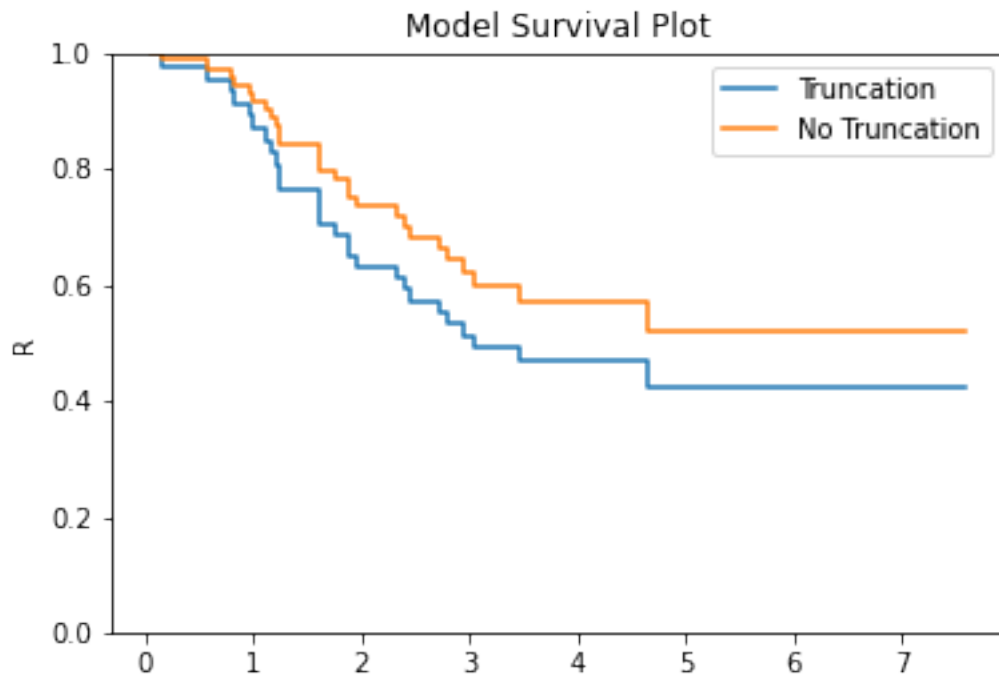
In some instances you will need to account for left truncated data. These data can be passed stright to the same KM, NA, and FH fitters, Using one (of the many) excellent data sets from the [lifelines](#) package:

```
from surpyval import KaplanMeier as KM
from lifelines.datasets import load_multicenter_aids_cohort_study
df = load_multicenter_aids_cohort_study()

x = df["T"].values
c = 1. - df["D"].values
tl = df["W"].values

model = KM.fit(x=x, c=c, tl=tl)
model_no_trunc = KM.fit(x=x, c=c)

model.plot(plot_bounds=False)
model_no_trunc.plot(plot_bounds=False)
plt.legend(['Truncation', 'No Truncation'])
```



The image above shows that if you fail to take into account the left truncation (using the `tl` keyword) you will overstate the survival probability. This can be used with any of the other non-parametric fitters.

### 1.10.4 Arbitrarily Truncated and Censored Data

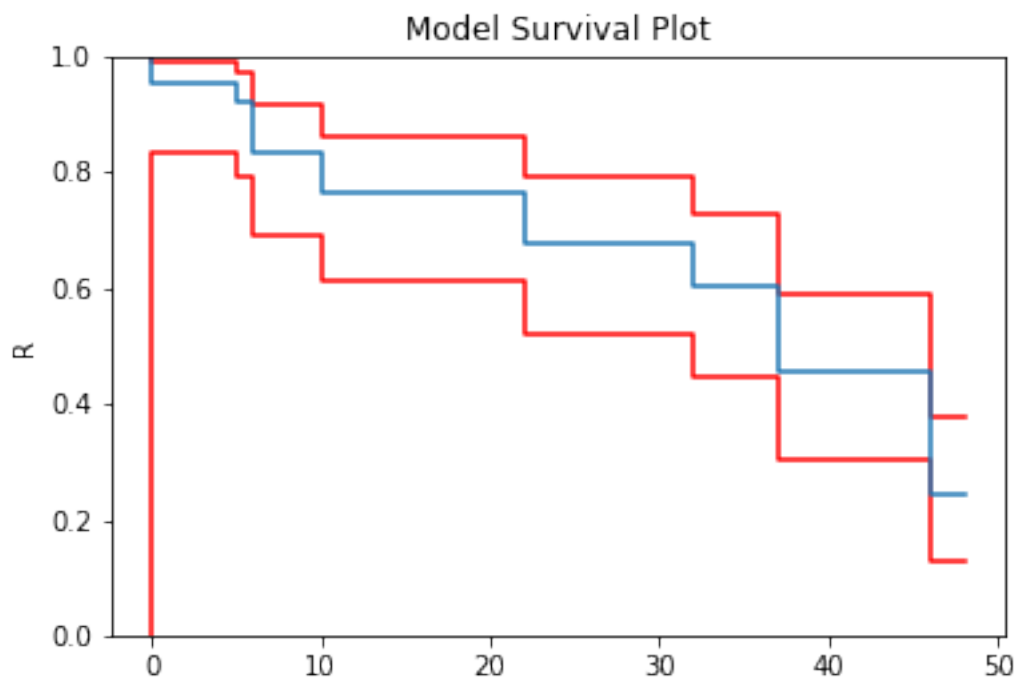
In the event you have data that has interval, left, or right censoring with no, left, or right truncation, the previous estimators will not work. Enter the Turnbull estimator. First an interval estimation example:

```

low = np.array([0, 0, 0, 4, 5, 5, 6, 7, 7, 11, 11, 15, 17, 17,
               17, 18, 19, 18, 22, 24, 24, 25, 26, 27, 32, 33,
               34, 36, 36, 36, 36, 37, 37, 37, 37, 38, 40, 45,
               46, 46, 46, 46, 46, 46, 46, 46])
upp = np.array([7, 8, 5, 11, 12, 11, 10, 16, 14, 15, 18, np.inf,
               np.inf, 25, 25, np.inf, 35, 26, np.inf, np.inf,
               np.inf, 37, 40, 34, np.inf, np.inf, np.inf, 44,
               48, np.inf, np.inf, 44, np.inf, np.inf, np.inf,
               np.inf, np.inf, np.inf, np.inf, np.inf, np.inf,
               np.inf, np.inf, np.inf, np.inf, np.inf])

x = np.array([low, upp]).T
model = TB.fit(x)
model.plot()

```



And finally, an example with arbitrary censoring and truncation:

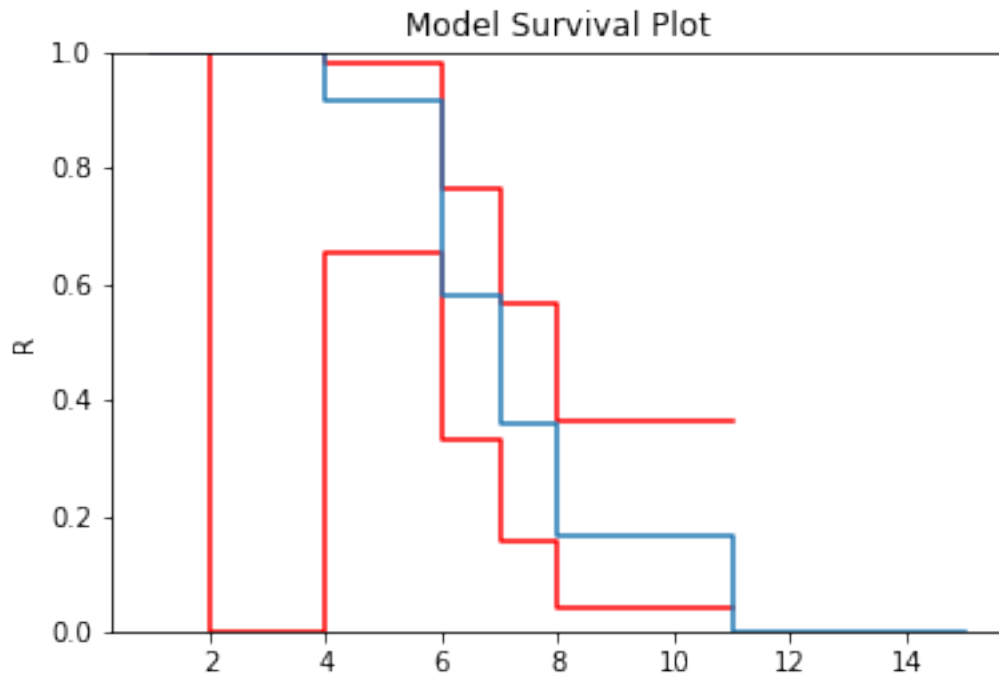
```

from surpyval import Turnbull as TB

x = [1, 2, [3, 6], 7, 8, 9, [5, 9], [4, 10], [7, 10], 11, 12]
c = [1, 1, 2, 0, 0, 0, 2, 2, 2, -1, 0]
n = [1, 2, 1, 3, 2, 2, 1, 1, 2, 1, 1]
tl = [0, 0, 0, 0, 0, 2, 3, 3, 1, 1, 5]
tr = [np.inf, np.inf, 10, 10, 10, 10, np.inf, np.inf, np.inf, 15, 15]

model = TB.fit(x=x, c=c, n=n, tl=tl, tr=tr)
model.plot()

```



With a completely arbitrary set of data we have created a non-parametric estimate of the survival curve that can be used to estimate probabilities.

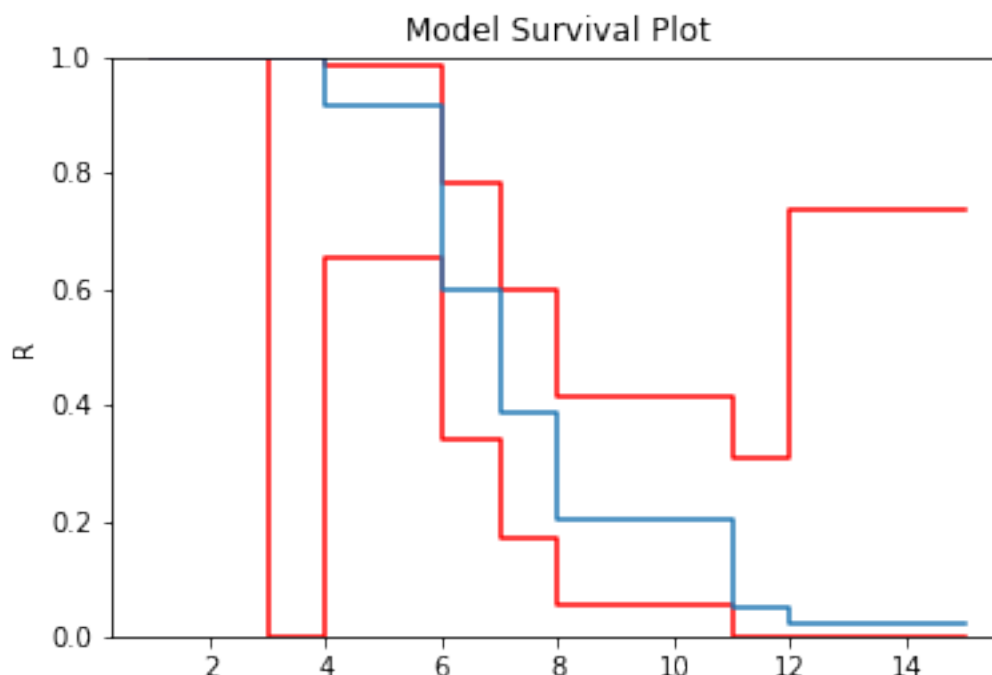
What is interesting about the Turnbull estimate is that it first finds the data in the 'xrd' format. This is done even though we might not have a complete failure occur in an interval. This can be seen by looking at the number of deaths/failures occur at each value.

```
model.d
```

```
array([0.00000000e+00, 0.00000000e+00, 0.00000000e+00, 2.76875496e-02,
       1.58808369e+00, 0.00000000e+00, 5.81471061e+00, 4.10951885e+00,
       3.54383160e+00, 7.67984832e-02, 3.93153047e-15, 3.09598691e+00,
       1.66794197e+00])
```

You can see that some values are 0 (or essentially 0) or that there is an interval where there were 4.1095188 failures. But because the Turnbull estimate finds the x, r, d format we can actually elect to use the Nelson-Aalen or Kaplan-Meier estimate with the Turnbull estimates of x, r, and d.

```
model = TB.fit(x=x, c=c, n=n, tl=tl, tr=tr, estimator='Nelson-Aalen')
model.plot()
```



The Greenwood confidence intervals do give us a strange set of bounds. But you can see that using the Nelson-Aalen estimator instead of the Kaplan-Meier gives us a better approximation for the tail end of the distribution.

### Some Issues with the Turnbull Estimate

Caution must be given when using the Turnbull estimate when all values are truncated by some left and/or right value. This will be shown below in the methods for estimating parameters with truncated values. But essentially the Turnbull method cannot make any assumptions about the probability by which the smallest value if left truncated should be adjusted. This is because there is no information available with the non-parametric method below this smallest value. The same is true for the largest value if it is also right truncated, there is no information available about the probability of its observation. Therefore the Turnbull method makes an implicit assumption that the first value, if left truncated has 100% chance of observation, and the highest value, if right truncated also has 100% chance of being observed.

The implications of this are detailed in the Parametric section, because the only way to gain an understanding of these situations is by assuming a shape of the distribution. That is, by doing parametric analysis. This is possible since if the distribution within the truncated ends has a shape that matches to a particular distribution you can then extrapolate beyond the observed values. Parametric analysis is therefore incredibly powerful for prediction / extrapolation.

### 1.10.5 Confidence Intervals

Coming soon

## 1.11 Parametric SurPyval Modelling

The parametric API is essentially the exact same as the non-parametric API. All models are fit by a call to the `fit()` method. However, the parametric models have more options that are only applicable to parametric modelling. The

inputs of  $x$  for the random variable,  $c$  for the censoring flag,  $n$  for count of each  $x$ ,  $x_l$  and  $x_r$  for intervally censored data (can't be used with  $x$ )  $t$  for the truncation matrix,  $t_l$  for the left truncation scalar or array, and  $t_r$  for the right truncation scalar or array all remain.

### 1.11.1 Complete Data

The easiest and simplest case is that when you have a dataset of exactly observed data. that is, you have one array of data with the values at which they failed. Fitting a parametric distribution to the data can be done with a simple call to the `fit()` method:

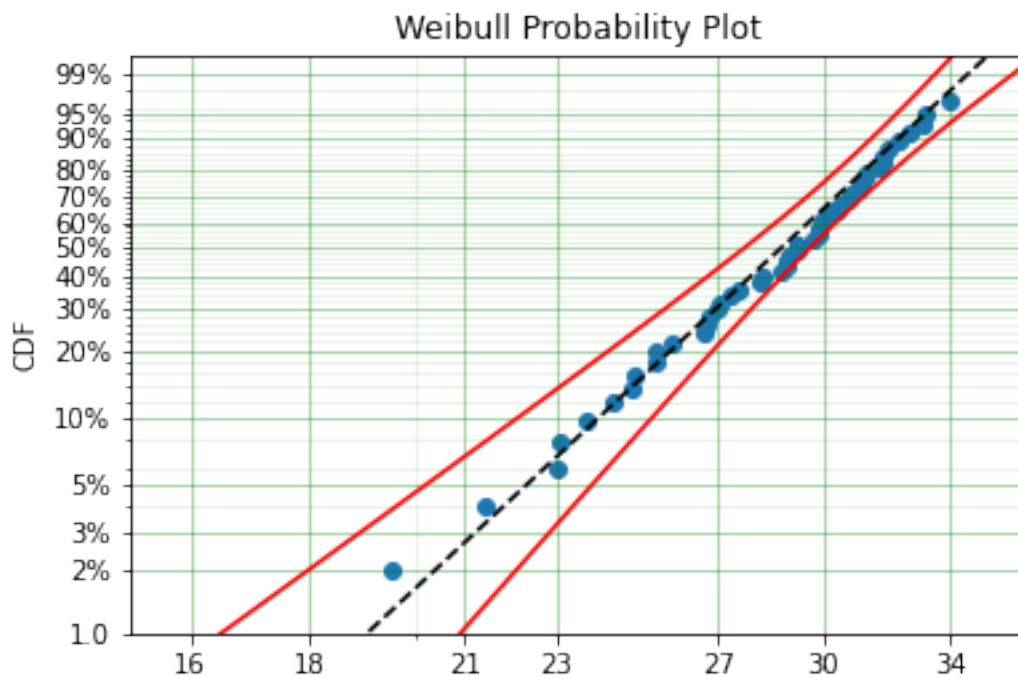
```
import surpyval as surv
import numpy as np

np.random.seed(10)
x = surv.Weibull.random(50, 30., 9.)
model = surv.Weibull.fit(x)
print(model)
model.plot();
```

```
Parametric SurPyval Model
=====
Distribution      : Weibull
Fitted by        : MLE
Parameters       :
    alpha: 29.805137406871953
    beta: 10.296037991991037
```

To visualise the outcome of this fit we can inspect the results on a probability plot:

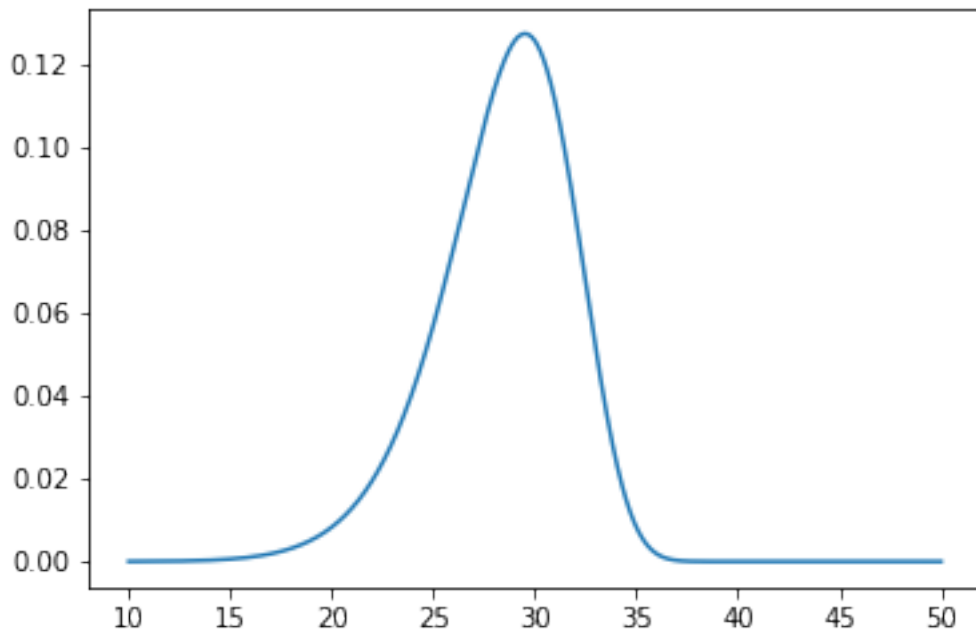
```
model.plot()
```



The `model` object from the above example can be used to calculate the density of the distribution with the parameters found with the best fit from above. This is very easy to do:

```
x = np.linspace(10, 50, 1000)
f = model.df(x)

plt.plot(x, f)
```



The CDF `ff()`, Survival (or Reliability) `sf()`, hazard rate `hf()`, or cumulative hazard rate `Hf()` can be computed as well. This functionality makes it very easy to work with `surpyval` models to determine risks or to pass the function to other libraries to find optimal trade-offs.

### 1.11.2 Using censored data

#### Right Censored

A common complication in survival analysis is that all the data is not observed up to the point of failure (or death). In this case the data is right censored, see the types of data section for a more detailed discussion, `surpyval` offers a very clean and easy way to model this. First, let's create a simulated data set:

```
import surpyval as surv
import numpy as np

np.random.seed(10)
x = surv.Weibull.random(50, 30, 2.)

observation_limit = 40
# Censoring flag
```

(continues on next page)

(continued from previous page)

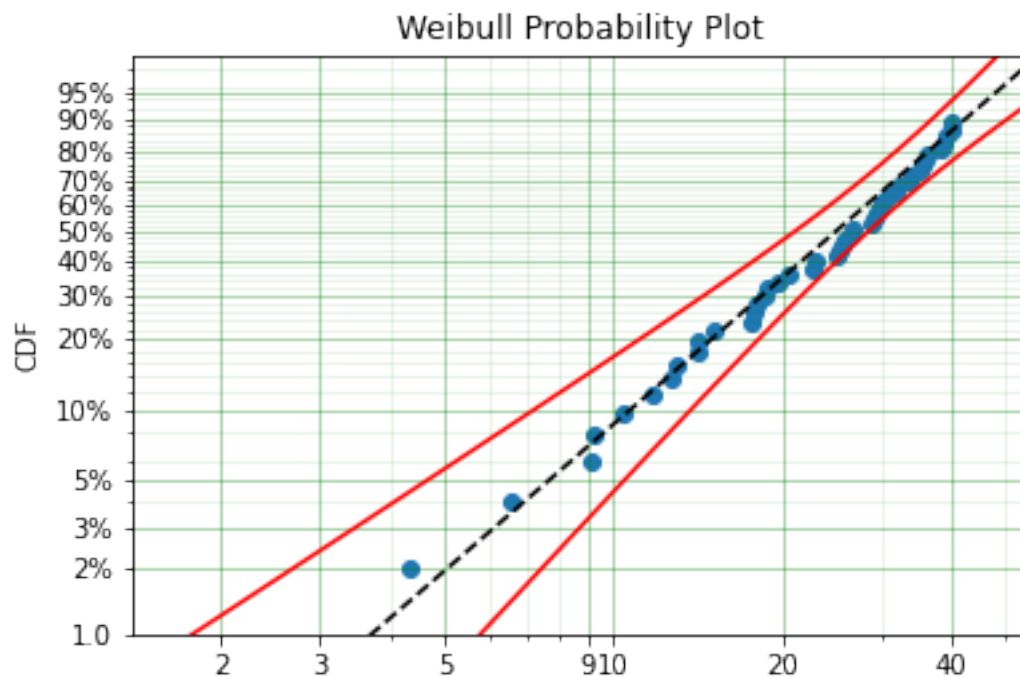
```
c = (x >= observation_limit).astype(int)
x[x >= observation_limit] = observation_limit
```

In this example, we created 50 random Weibull distributed values with  $\alpha = 30$  and  $\beta = 2$ . For this example the observation window has been set to 40. This value is where we stopped observing the events. For all the randomly generated values that are above this limit we create the censoring flag array  $c$ . This array has zeros where the event time was observed, and a 1 where the value is above the recorded value. For all the values in the data that are above 40 we set them to 40. This is a common occurrence in survival analysis and surpyval is designed to accept this input with a simple call:

```
model = surv.Weibull.fit(x, c)
print(model)
model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : Weibull
Fitted by        : MLE
Parameters       :
    alpha: 29.249243175049152
    beta: 2.2291485877426354
```

The plot for this can be seen to be:



The results from this model are very close to the data we input, and with only 50 samples.



## Left Censored

The above example can be extended to another kind of censoring; left censored data. This is the case where the values are known to fall below a particular value. We can change our example data set to have a start observation time for which we will left censor all the data below that:

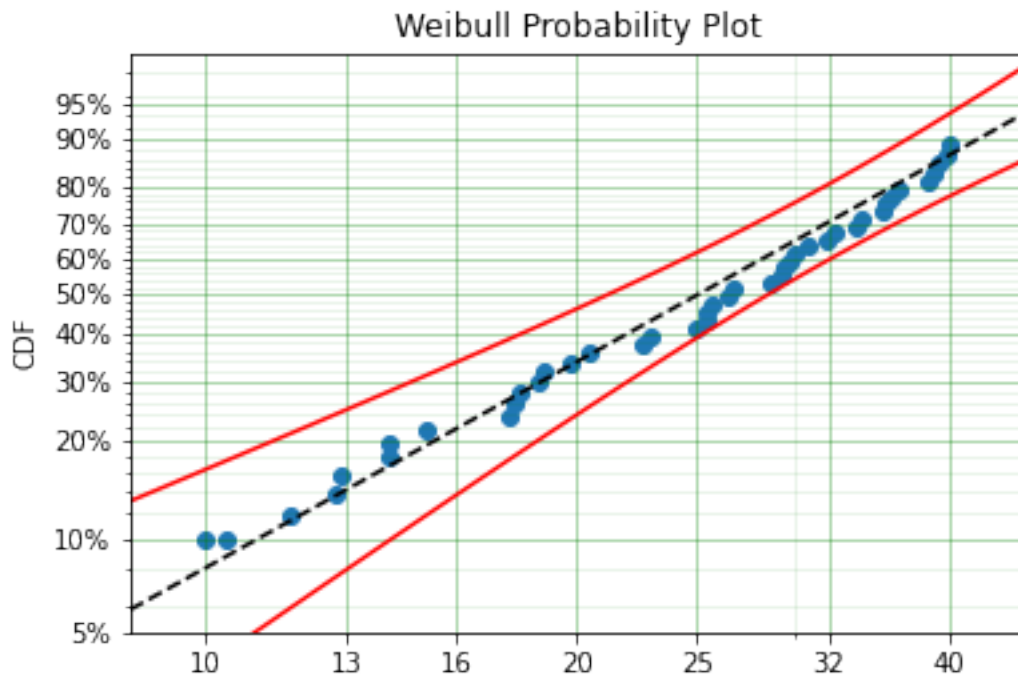
```
observation_start = 10
# Censoring flag
c[x <= observation_start] = -1
x[x <= observation_start] = observation_start
```

That is, we set the start of the observations at 10 and flag that all the values at or below this are left censored. We can then use the updated values of x and c:

```
model = surv.Weibull1.fit(x, c)
print(model)
model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : Weibull1
Fitted by        : MLE
Parameters       :
    alpha: 29.34709766238127
    beta:  2.3049027909575903
```

The values did not substantially change, although the plot does look different as there are no values below 10.



## Intervally Censored

The next type of censoring that is naturally handled by surpyval is interval censoring. Creating another example data set:

```
import surpyval as surv
import numpy as np

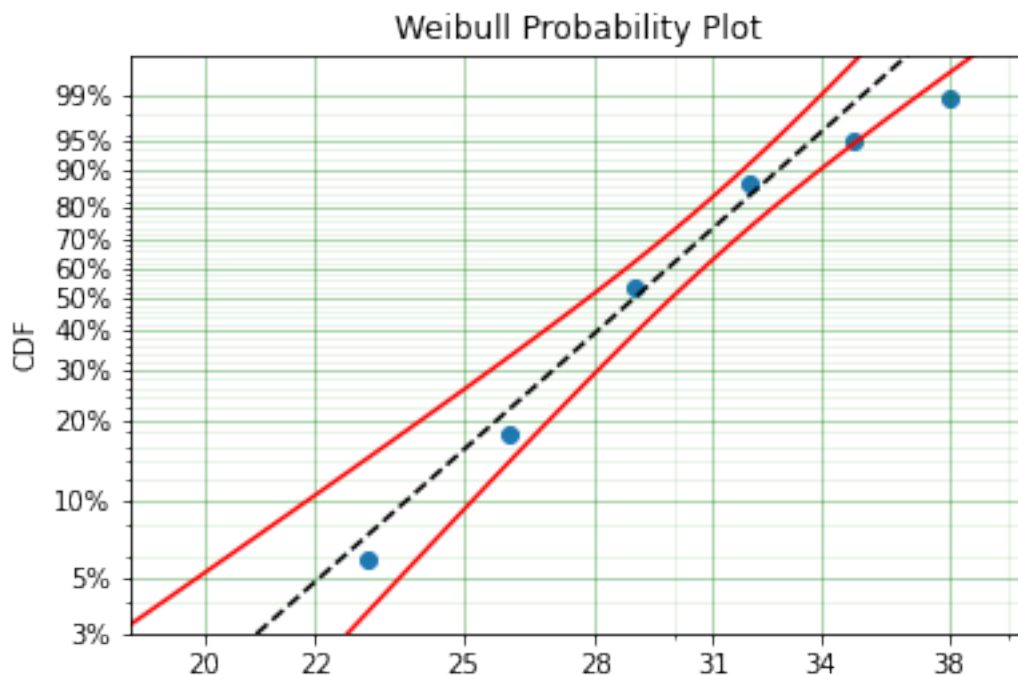
np.random.seed(30)
x = surv.Weibull.random(50, 30, 10.)
n, xx = np.histogram(x, bins=[20, 23, 26, 29, 32, 35, 38])
x = np.vstack([xx[0:-1], xx[1:]]).T
```

In this example we have created the variable `x` with a matrix of the intervals within which each of the observations have failed. That is each exact observation has been binned into a window and the `x` array has an entry `[left, right]` within which the event failed. We also have the `n` array that has the count of the failures within the window. With these two values we can make the simple surpyval call:

```
model = surv.Weibull.fit(x, n=n)
print(model)
```

```
Parametric SurPyval Model
=====
Distribution      : Weibull
Fitted by        : MLE
Parameters       :
    alpha: 30.074154903683105
    beta:  9.637405285678362
```

Again, we have a result that is very close to the original parameters. SurPyval can take as input an arbitrary combination of censored data. If we plot the data we will see:



This is a good fit! The data at the tails are a little bit off, but this is only 50 samples and the core of the model matches the data quite well.

## Mixed Censoring

Mixed censoring, or arbitrary censoring is easily handled by SurPyval. So no matter the combination of the data that you have, SurPyval will be able to fit a distribution to it.

```
import surpyval as surv

x = [0, 1, 2, [3, 4], [6, 10], [4, 8], 5, 19, 10, 13, 15]
c = [0, 0, 1, 2, 2, 2, 0, -1, 0, 1, 0]
surv.Gumbel.fit(x, c=c)
```

```
Parametric SurPyval Model
=====
Distribution      : Gumbel
Fitted by        : MLE
Parameters       :
    mu: 9.912232006272871
    sigma: 4.95952392045353
```

### 1.11.3 Using truncated data

#### Left truncated

Surpyval has the capacity to handle arbitrary truncated data. A common occurrence of this is in the insurance industry data. When customers make a claim on their policies they have to pay an ‘excess’ which is a charge to submit a claim for processing. If say, the excess on a set of policies in an area is \$250, then it would not be logical for a customer to submit a claim for a loss of less than that number. Therefore there will be no claims under \$250. This can also happen in engineering where a part may be tested up to some limit prior to be sold, therefore, as a customer you need to make sure you take into account the fact that some parts would have been rejected at the end of the line which you may not have seen. So a washing machine may run through 25 cycles prior to shipping. This is similar to, but distinct from censoring. When something is left censored, we know there was a failure or event below the threshold. Whereas with truncation, we do not see any variables below the threshold. A simulated example may explain this better:

```
import numpy as np
import surpyval as surv

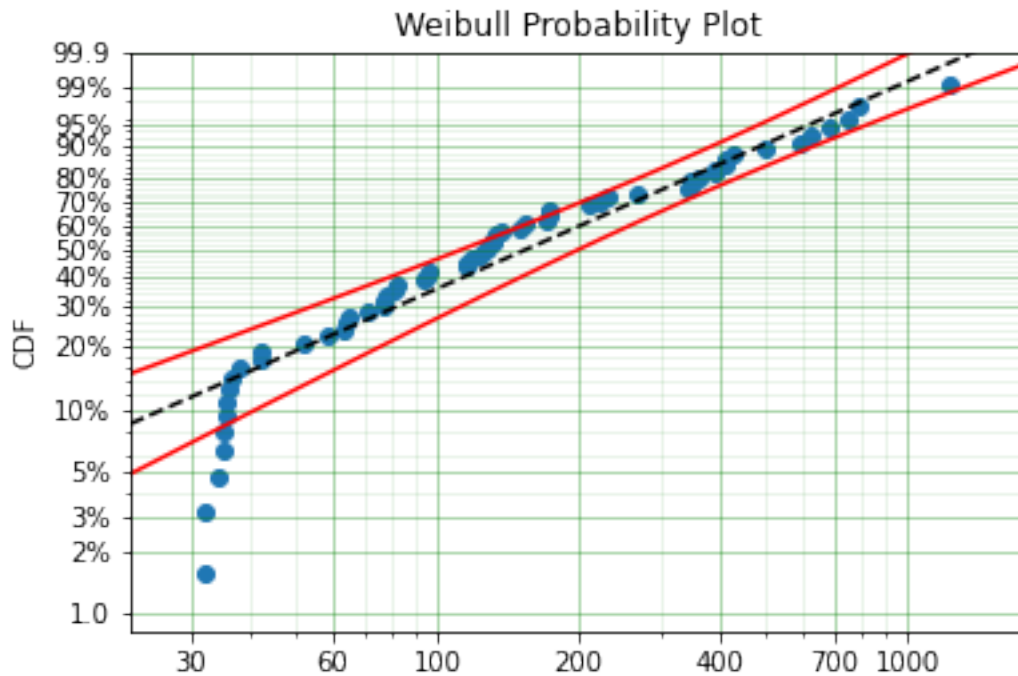
np.random.seed(10)
x = surv.Weibull.random(100, alpha=100, beta=0.6)
# Keep only those values greater than 250
threshold = 25
x = x[x > threshold]
```

We have therefore simulated a scenario where we have taken 100 random samples from a fat tailed Weibull distribution. We then filter to keep only those records that are above the threshold. In this case we assume we haven’t seen the data for the washing machines with less than 25 cycles. To understand what could go wrong if we ignore this, what do we get if we assume all the data are failures and there is no truncation?

```
model = surv.Weibull.fit(x=x)
print(model.params)
```

```
[218.39245675  1.0507186 ]
```

With a plot that looks like:



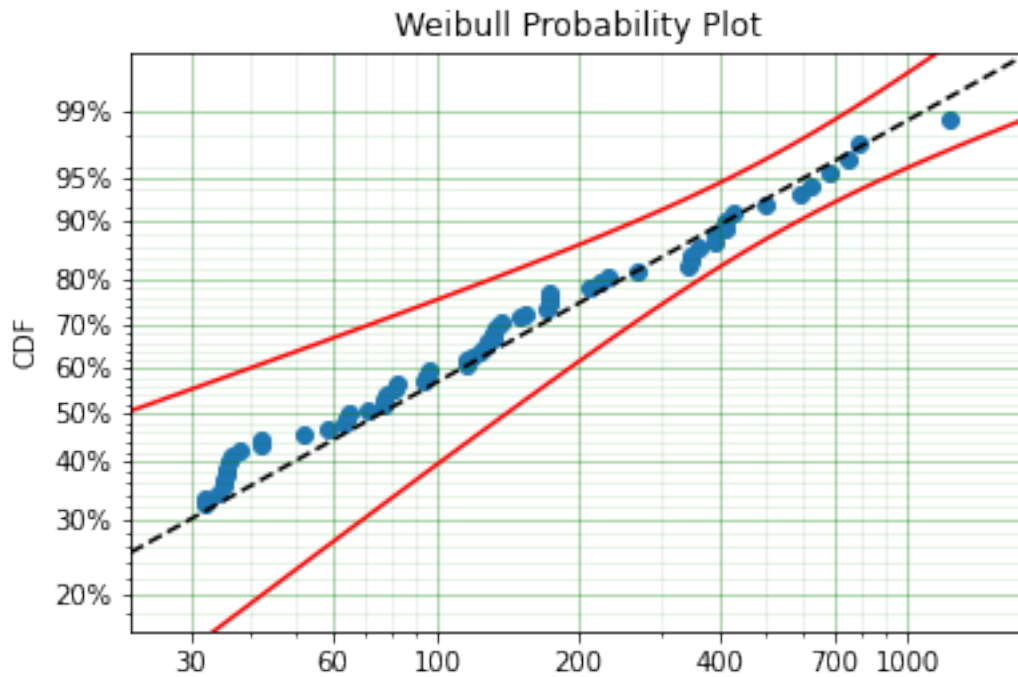
Looking at the parameters of the distribution, you can see that the beta value is greater than 1. Although only slightly, this implies that this distribution has an increasing hazard rate. If you were the operator of the washing machines (e.g. a hotel or a laundromat) and any downtime had a cost, you would conclude from this that replacing the machines after a fixed time would be a good policy.

But if you take the truncation into account:

```
model = surv.Weibull.fit(x=x, tl=threshold)
print(model.params)
```

```
[127.32704868  0.71053572]
```

With the plot:



You can see now that the model fits the data much better, but also that the beta parameter is actually below 1. This shows that ignoring the left-truncated data in parametric estimation can lead to errors in prediction.

### Right truncated

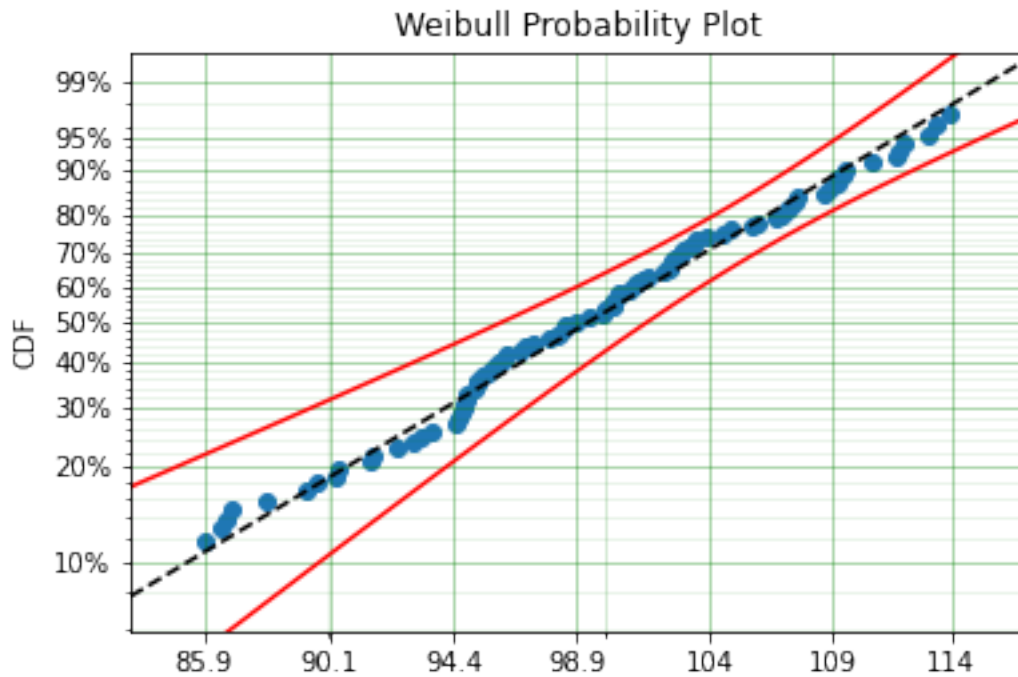
The example from above can be continued for right-truncated data as well.

```
import numpy as np
import surpyval as surv

np.random.seed(10)
x = surv.Normal.random(100, mu=100, sigma=10)
# Keep only those values greater than 250
tl = 85
tr = 115
# Truncate the data
x = x[(x > tl) & (x < tr)]

model = surv.Weibull.fit(x=x, tl=tl, tr=tr)
print(model.params)
```

```
[102.27078401 12.47906136]
```



From the output above, the number of data points we have has been reduced from the simulated 100, down to 87. Then with the 87 samples we now have we estimated the parameters to be quite close to the parameters used in the simulation. Further, the plot looks as though the parametric distribution fits the non-parametric distribution quite well.

In the cases above we used a scalar value for the truncation values. But some data has individual values for left truncation. This is seen in trials where someone may join the trial as a late entry. Therefore each data point as an entry time. For example:

```
import surpyval as surv

x = [3, 4, 6, 7, 9, 10]
tl = [0, 0, 0, 0, 5, 2]

model = surv.Weibull.fit(x, tl=tl)
print(model.params)
```

```
[7.05854717 2.70096672]
```

### Intervally and Arbitrarily truncated

Surpyval can even work with arbitrary left and right truncation:

```
import surpyval as surv

x = [3, 4, 6, 7, 9, 10]
tl = [0, 0, 0, 0, 5, 2]
tr = [10, 9, 8, 10, 15]

model = surv.Weibull.fit(x, tl=tl, tr=tr)
print(model.params)
```

```
[8.12377602 2.56917036]
```

In the above example we used both the `tl` and `tr`. However, `surpyval` has a flexible API where it can take the truncation data as a two dimensional array:

```
import surpyval as surv

x = [3, 4, 6, 7, 9, 10]
t = [[ 0, 10],
      [ 0, 9],
      [ 0, 8],
      [ 0, 10],
      [ 5, 15],
      [ 2, 15]]

model = surv.Weibull.fit(x, t=t)
print(model.params)
```

```
[8.12377602 2.56917036]
```

Which, obviously, gives the same result. This shows the flexibility of the `surpyval` API, you can use scalar, array, or matrix values for the truncations using the `t`, `tl`, and `tr` keywords with the `fit` method and `surpyval` does the rest.

### 1.11.4 Offsets

Another common feature in survival analysis is a requirement to fit a distribution with an offset. These distributions are sometimes referred to as the two-parameter (e.g. two parameter exponential) three parameter, (e.g., the three three parameter Weibull), or four parameter (e.g four parameter Exponentiated Weibull distribution). `SurPyval` however just uses an `offset` to increase the numbers of parameters and allow the distribution to be shifted.

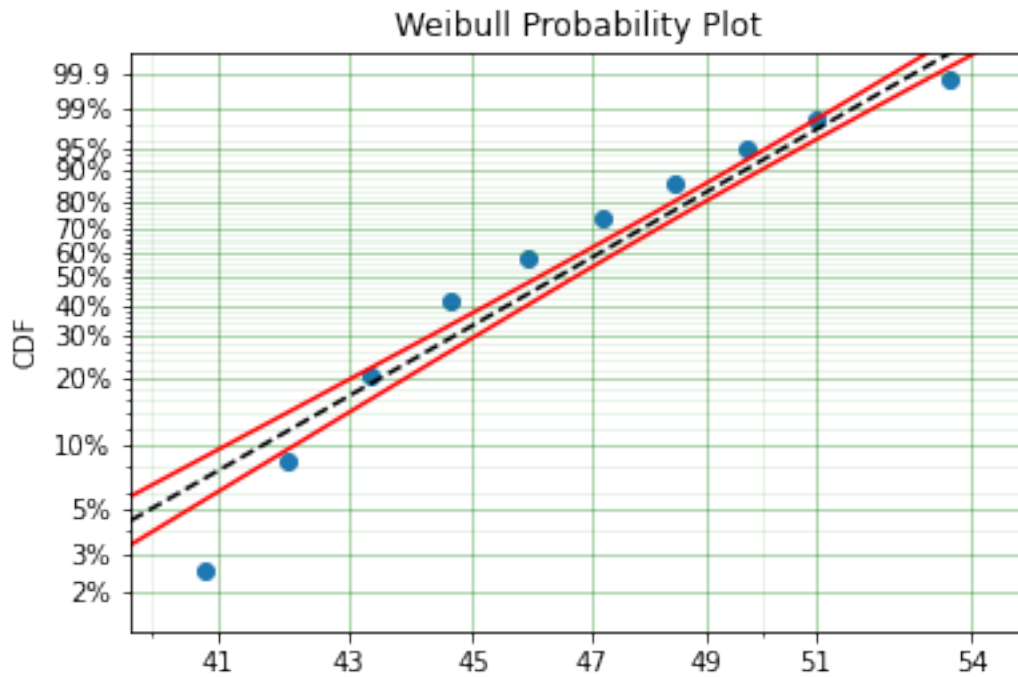
Using data from Weibull's original paper for the strenght of Bofor's steel shows when this might be necessary.

```
import surpyval as surv
from surpyval.datasets import BoforsSteel

df = BoforsSteel.df
x = df['x']
n = df['n']

model = surv.Weibull.fit(x=x, n=n)
print(model.params)
model.plot()
```

```
[47.36735846 17.5713195 ]
```



The above plot does not look to be a good fit. However, if we use an offset we can use the three parameter Weibull distribution to attempt to get a better fit. Using offset values with surpyval is very easy:

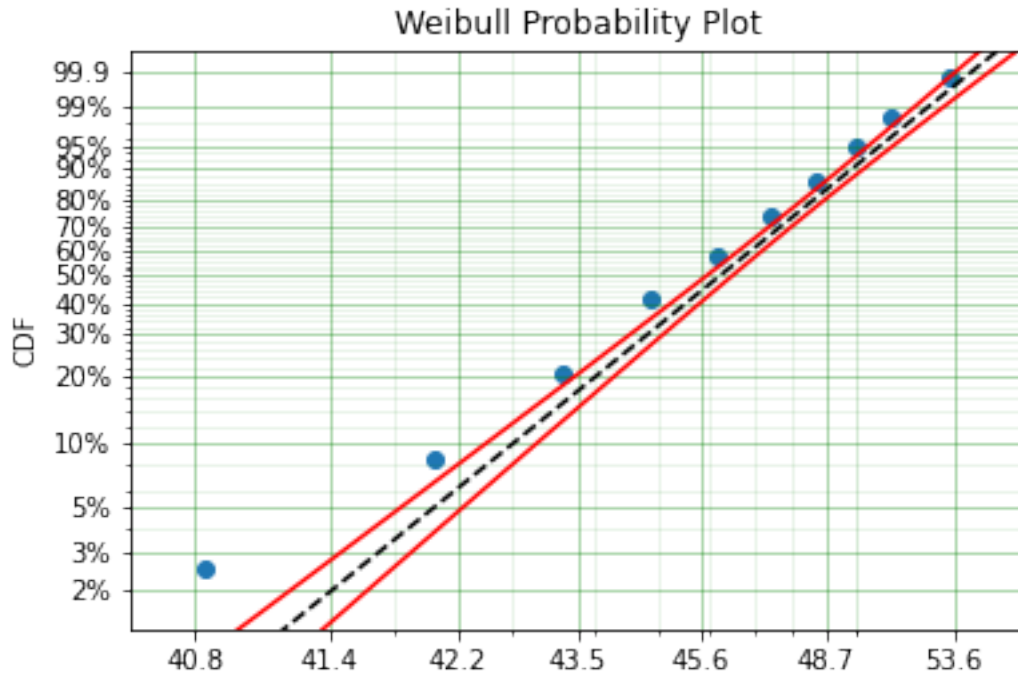
```
import surpyval as surv
from surpyval.datasets import BoforsSteel

df = BoforsSteel.df
x = df['x']
n = df['n']

model = surv.Weibull.fit(x=x, n=n, offset=True)
print(model)
model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : Weibull
Fitted by        : MLE
Offset (gamma)    : 39.76562962867477
Parameters       :
    alpha: 7.141925216146524
    beta: 2.6204524040137844
```



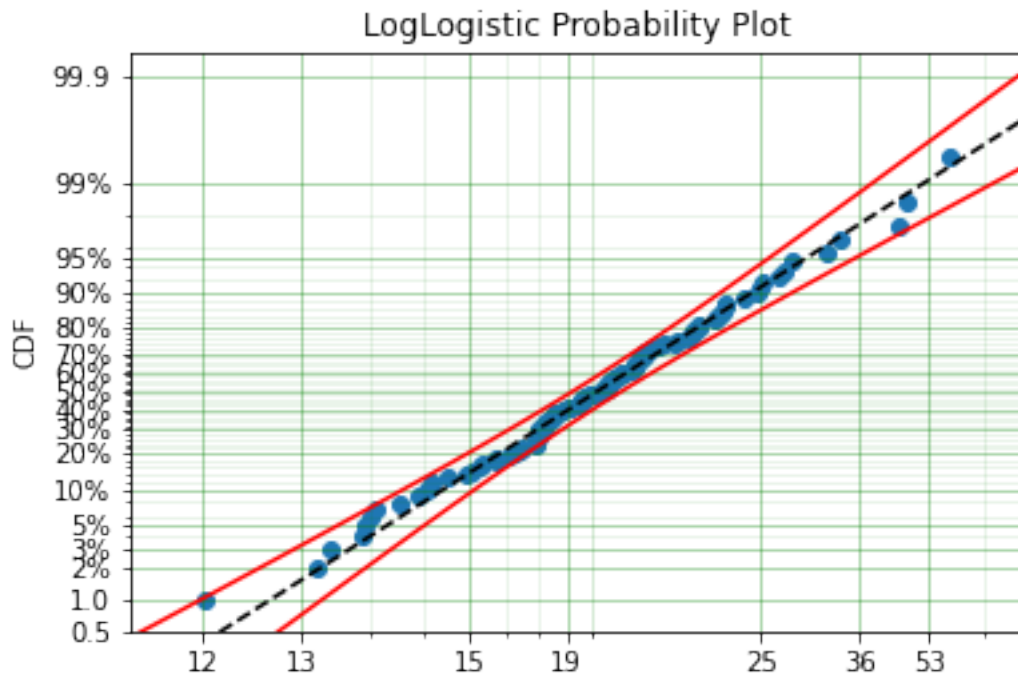


This is evidently a much better fit! The offset value for an offset distribution is saved as `gamma` in the model object. Offsets can be used for any distribution supported on the half real line. Currently, this is the Weibull, Gamma, LogNormal, LogLogistic, and Exponential. For example:

```
import surpyval as surv
import numpy as np

np.random.seed(10)
x = surv.LogLogistic.random(100, 10, 3) + 10
model = surv.LogLogistic.fit(x, offset=True, how='MLE')
print(model)
model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : LogLogistic
Fitted by        : MLE
Offset (gamma)    : 9.56270794050046
Parameters       :
    alpha: 10.18946967467503
    beta: 3.407325975660712
```

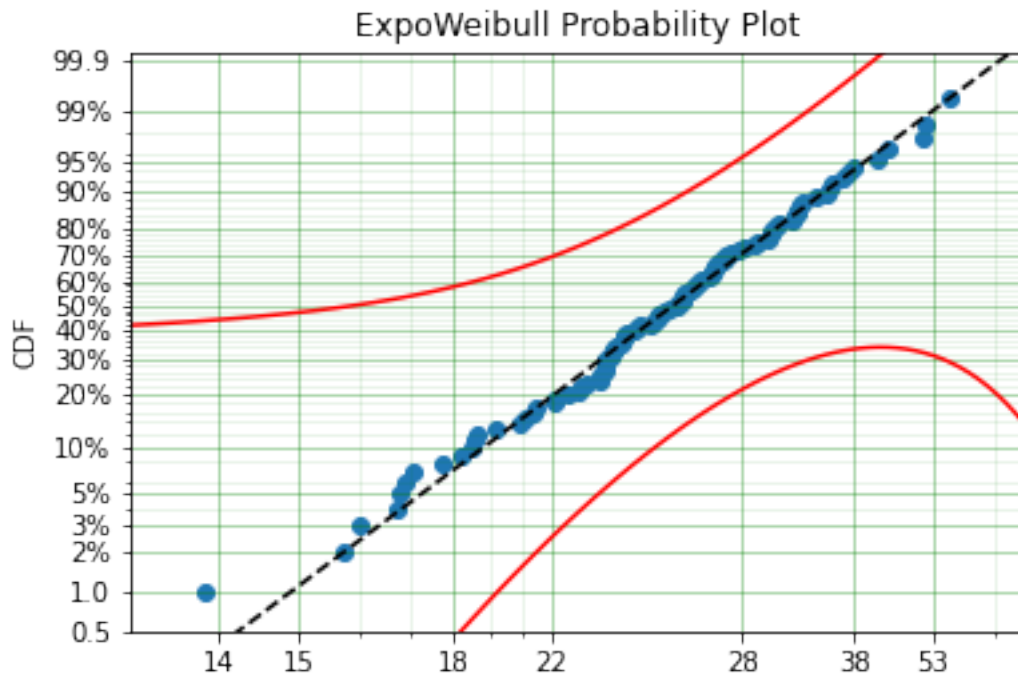


A four parameter exponentiated Weibull can also be found:

```
import surpyval as surv
import numpy as np

np.random.seed(10)
x = surv.ExpoWeibull.random(100, 10, 1.2, 4) + 10
model = surv.ExpoWeibull.fit(x, offset=True)
print(model)
model.plot(plot_bounds=False)
```

```
Parametric SurPyval Model
=====
Distribution      : ExpoWeibull
Fitted by        : MLE
Offset (gamma)   : 10.701280166551431
Parameters      :
    alpha: 11.47511146192537
    beta:  1.3969785125819283
    mu:    2.845307244239084
```



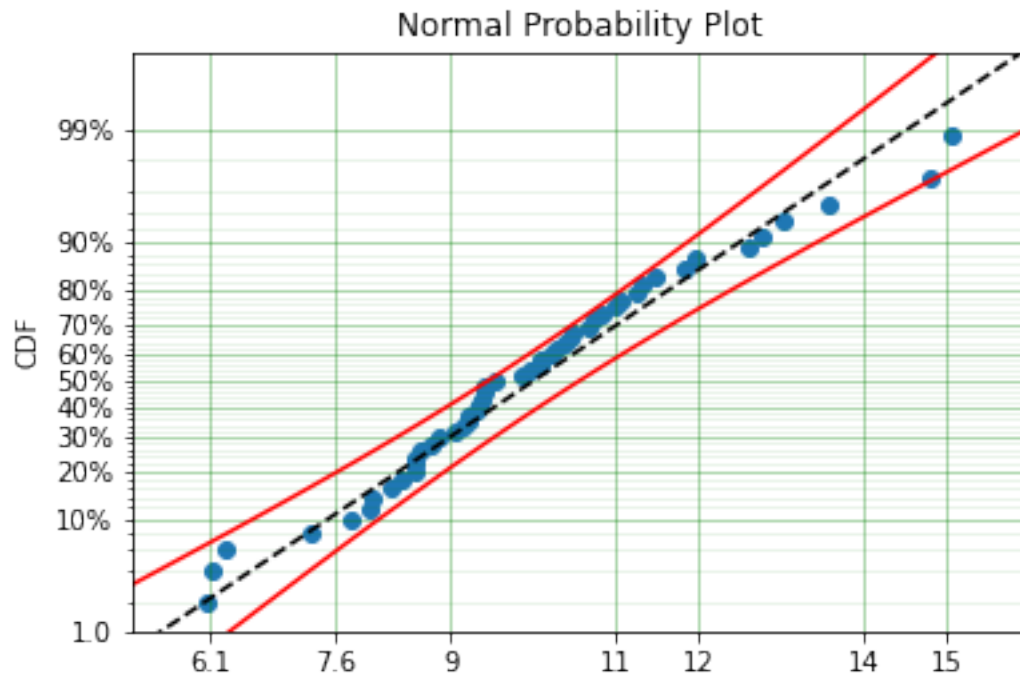
### 1.11.5 Fixing parameters

Another useful feature of surpyval is the ability to easily fix parameters. For example:

```
import surpyval as surv
import numpy as np

np.random.seed(30)
x = surv.Normal.random(50, 10., 2)
model = surv.Normal.fit(x, fixed={'mu' : 10})
print(model)
model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : Normal
Fitted by        : MLE
Parameters       :
    mu: 10.0
    sigma: 1.9353643871136006
```

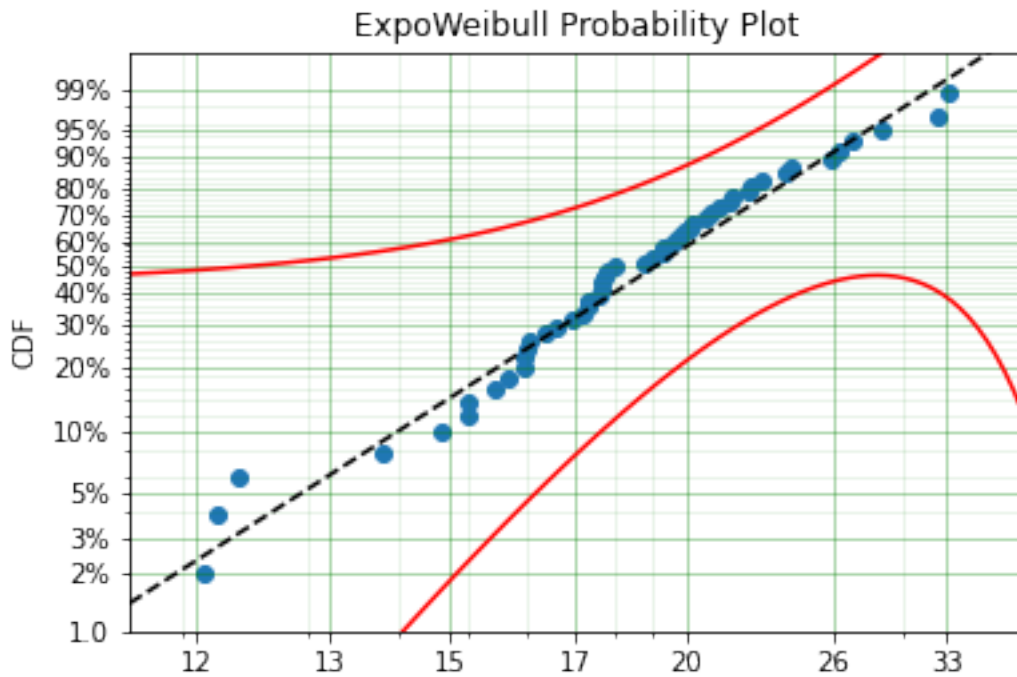


You can see that the  $\mu$  parameter has been fixed at 10. This can work for distributions with many more parameters, including the offset.

```
import surpyval as surv
import numpy as np

np.random.seed(30)
x = surv.ExpoWeibull.random(50, 10., 2, 4) + 10
model = surv.ExpoWeibull.fit(x, offset=True, fixed={'mu' : 4, 'gamma' : 10, 'alpha' : 10})
print(model)
model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : ExpoWeibull
Fitted by        : MLE
Offset (gamma)   : 10.0
Parameters      :
    alpha: 10.0
    beta: 1.9986073390210994
    mu: 1.2
```



We have fit three of the four parameters for an offset exponentiated-Weibull distribution!

### 1.11.6 Modelling with arbitrary input

The surpyval API is extremely flexible. All the unique examples provided above can all be used at once. That is, data can be censored, truncated, and directly observed with offsets and fixing parameters. The API is completely flexible. This makes surpyval an extremely useful tool for analysts where the data is gathered in a manner where it's cleanliness is not guaranteed.

```
import surpyval as surv

x = [0, 1, 2, [3, 4], [6, 10], [4, 8], 5, 19, 10, 13, 15]
c = [0, 0, 1, 2, 2, 2, 0, -1, 0, 1, 0]
tl = [-1, 0, 0, 0, 0, 0, 0, 2, 2, -np.inf, 0, 0]
tr = 25
model = surv.Normal.fit(x, c=c, tl=tl, tr=tr, fixed={'mu' : 1.})
print(model)
```

```
Parametric SurPyval Model
=====
Distribution      : Normal
Fitted by        : MLE
Parameters       :
    mu: 1.0
    sigma: 9.131202240846182
```

### 1.11.7 Using alternate estimation methods

Surpyval's API is very flexible because you can change which method is used to estimate parameters. This is useful when a more appropriate method is needed or the method you are using fails.

The default parametric method for surpyval is the maximum likelihood estimation (MLE), this is because it can take any arbitrary input. However, the MLE is not always the best estimator. Consider an example with the uniform distribution:

```
import surpyval as surv
import numpy as np

np.random.seed(5)
x = surv.Uniform.random(20, 5, 10)
print(x.min(), x.max())

mle_model = surv.Uniform.fit(x)
print(*mle_model.params)
```

```
5.9386061433062585 9.593054539689607
5.9386061433062585 9.593054539689607
```

You can see that the results are the same. This is because the maximum likelihood estimate of the parameters of a uniform distribution are just the smallest and largest values in the sample. If however we use the 'Maximum Product Spacing' method we get:

```
mps_model = surv.Uniform.fit(x, how='MPS')
print(*mps_model.params)
```

```
5.532556321486052 9.999104361509815
```

You can see that using the MPS method we have parameters that are closer to the real values. This is because the MPS method can 'look outside' the existing values to estimate where the real value lies. See the details of this method in the 'Parametric Estimation' section. But the MPS method is useful when you need to estimate the point at which a distribution's support starts or for any distribution that has unknown support. Concretely, this includes any offset distribution or a distribution with a finite upper and lower support (Uniform, Generalised Beta, Triangle)

The other important use case is when, for some reason, an alternate estimation method just does not work. For example:

```
import surpyval as surv
import numpy as np

np.random.seed(30)
x = surv.LogLogistic.random(10, 4., 2) + 10
model = surv.LogLogistic.fit(x, how='MLE', offset=True)
```

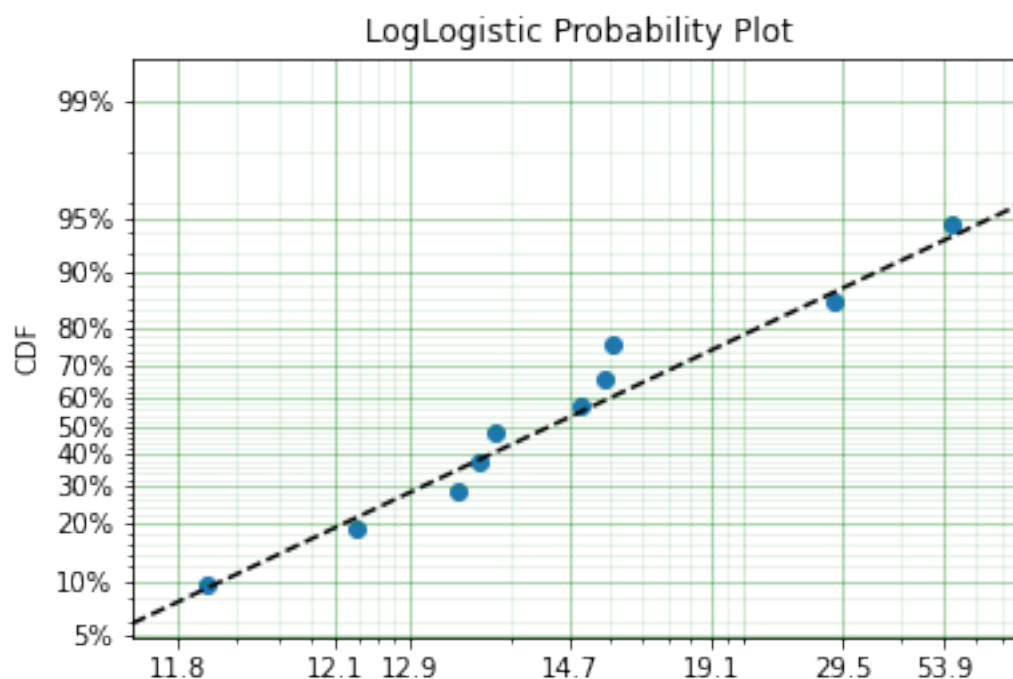
```
Precision was lost, try:
- Using alternate fitting method
- visually checking model fit
- change data to be closer to 1.
```

This shows, that the Maximum Likelihood Estimation may have failed for this data. However, because we have access to other methods, we can use an alternate estimation method:

```
import surpyval as surv
import numpy as np

np.random.seed(30)
x = surv.LogLogistic.random(10, 4., 2) + 10
model = surv.LogLogistic.fit(x, how='MPS', offset=True)
print(model)
model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : LogLogistic
Fitted by        : MPS
Offset (gamma)    : 11.524905733806891
Parameters       :
    alpha: 2.631868521887908
    beta: 0.9657662293516666
```



Our estimation has worked! Even though we used the MPS estimate for the parameters, we can still call all the same functions with the created variable to find the density `df()`, hazard `hf()`, CDF `ff()`, SF `sf()` etc. So regardless of the estimation method, we can still use the model.

This shows the power of the flexible API that `surpyval` offers, because if your modelling fails using one estimation method, you can use another. In this case, the MPS method is quite good at handling offset distributions. It is therefore a good approach to use when using offset distributions.

As stated in the Non-Parametric section, there is a risk that using the Turnbull estimator when all values are truncated by the same values. We will now show what happens. First, some example data:

```
import surpyval as surv
```

(continues on next page)

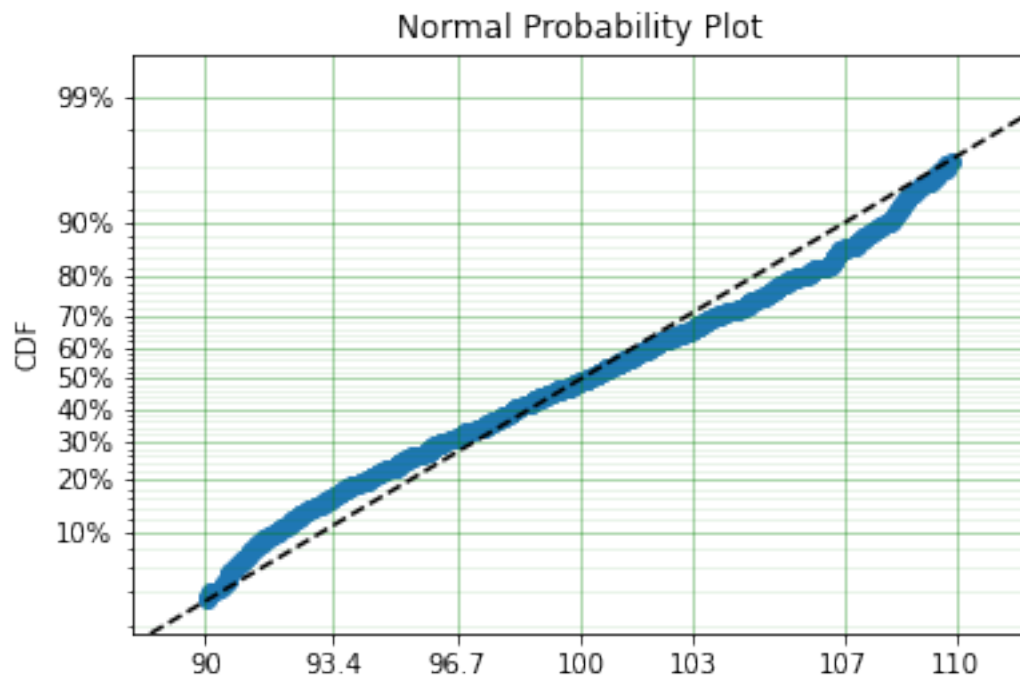
(continued from previous page)

```
import numpy as np

np.random.seed(1)
x = surv.Normal.random(1000, 100, 10)
tl = 90
tr = 110
x = x[x > tl]
x = x[x < tr]

mpp_model = surv.Normal.fit(x, tl=tl, tr=tr, how='MPP')
mpp_model.plot()
mpp_model
```

```
Parametric SurPyval Model
=====
Distribution      : Normal
Fitted by        : MPP
Parameters       :
    mu: 100.03108440743388
    sigma: 5.432878735738111
```



You can see that there is a strange match between the Turnbull estimate of the CDF and the parametric model. Also, you can see that the CDF at 90 is near 0% and the CDF at 110 is near 100%. This shows that it has not taken into account the truncation. Instead, if we use MLE we get:

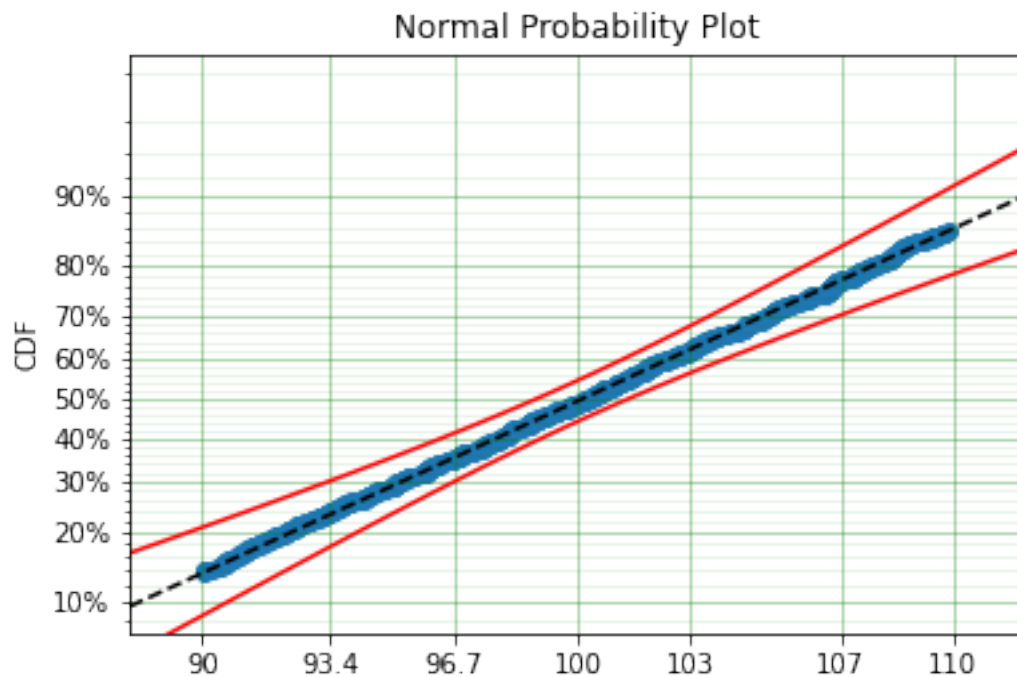
```
model = surv.Normal.fit(x, tl=tl, tr=tr, how='MLE')
model.plot()
model
```



```

Parametric SurPyval Model
=====
Distribution      : Normal
Fitted by        : MLE
Parameters       :
    mu: 100.13045397963812
    sigma: 9.17784957390746

```



We can see that the MLE method is a much better fit to this data, further, the MLE estimate of the  $\sigma$  parameter is much closer. The plotting points for the MLE plot have been adjusted in accordance with the truncation that the MLE model has estimated at the first entry. This is because it is known to be truncated and needs to be adjusted. This is not possible with the MPP method because the Turnbull estimator cannot adjust the truncation at the first and last value as it can make no assumptions about the truncation at those points.

This is just a word of warning for when using Truncation and the MPP method, make sure not all values are truncated by the same value, otherwise it will give a poor fit.

### 1.11.8 Mixture Models

On occasion, it can appear as though there are one, or two different distributions in the data you are using. On these occasions it can be useful to use a different type of distribution; or really, distributions. A mixture model is a distribution made from the partial combination of several distributions. Intuitively, it can be understood as a distribution where there is a proportion that fail for each kind of distribution. So 60% may come from a Weibull(3, 4) distribution but then another 40% come from a Weibull(19, 2) distribution.

SurPyval uses Expectation-Maximisation to

```

import surpyval as surv
import numpy as np

```

(continues on next page)

(continued from previous page)

```

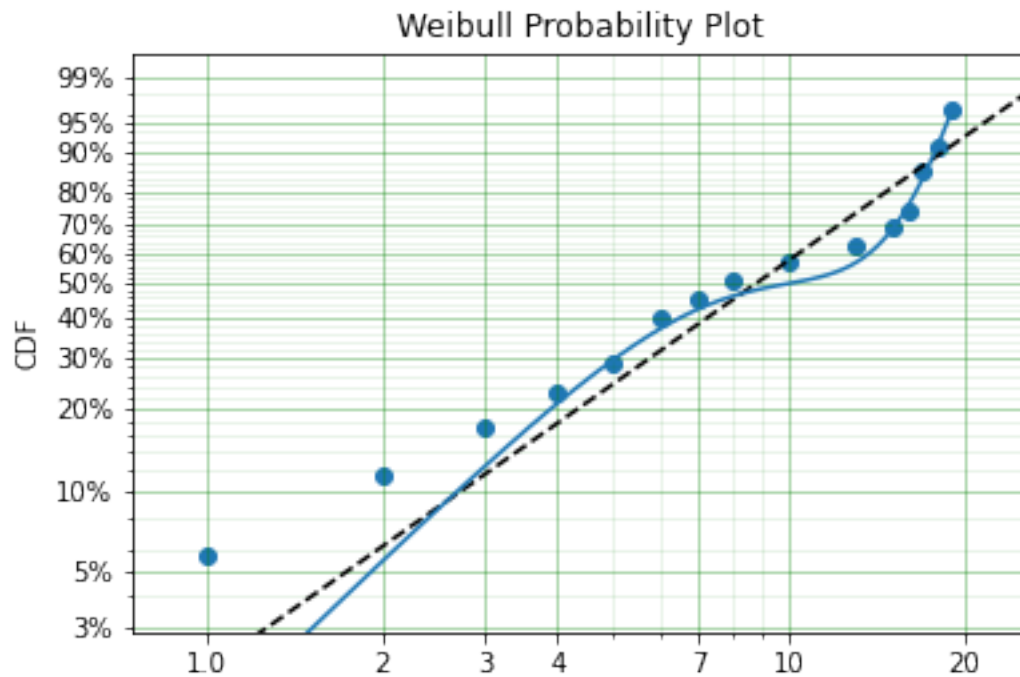
from matplotlib import pyplot as plt

x = [1, 2, 3, 4, 5, 6, 6, 7, 8, 10, 13, 15, 16, 17, 17, 18, 19]
x_ = np.linspace(np.min(x), np.max(x))

model = surv.Weibull.fit(x)
wmm = surv.MixtureModel(x=x, dist=surv.Weibull, m=2)

model.plot(plot_bounds=False)
plt.plot(x_, wmm.ff(x_))

```



You can see that the mixture model, in blue, tracks the data more closely than does the single model. SurPyval has incredible flexibility. The number of distributions can be changed by simply changing the value of `m`, and, the distribution passed to `dist` in the mixture can also be changed. Consider:

```

import surpyval as surv
import numpy as np
from matplotlib import pyplot as plt

np.random.seed(1)
x1 = surv.Normal.random(20, -10, 5)
x2 = surv.Normal.random(30, 10, 10)
x3 = surv.Normal.random(40, 50, 15)
x = np.concatenate([x1, x2, x3])
np.random.shuffle(x)
x_ = np.linspace(np.min(x), np.max(x))

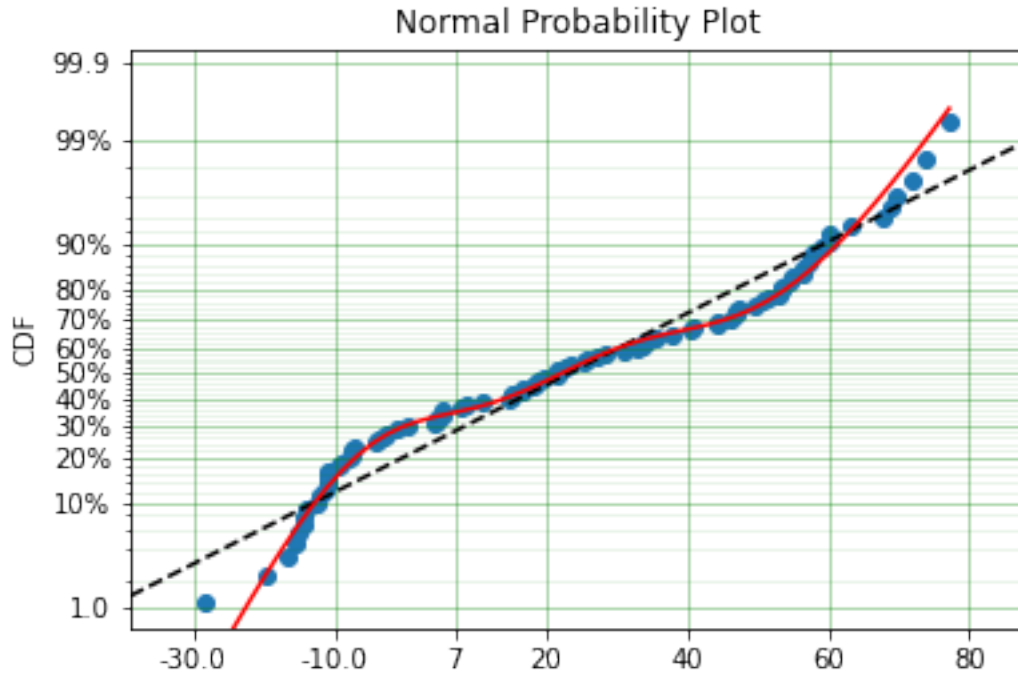
normal = surv.Normal.fit(x)
gmm = surv.MixtureModel(x=x, dist=surv.Normal, m=3)

```

(continues on next page)

(continued from previous page)

```
normal.plot(plot_bounds=False)
plt.plot(x_, gmm.ff(x_), color='red')
```



It was that simple to create a gaussian mixture model using `m=3` and the `dist=surv.Normal` parameters. SurPyval does default to 2 Weibull distributions if neither parameters are provided, but it can take any distribution in SurPyval as an input distribution.

Finally, mixture models can take counts and censoring flags as input (but not, yet, truncation). This makes SurPyval a truly powerful package for your survival analysis.

### 1.11.9 Limited Failure Population

Another kind of model that is useful in survival analysis is when a population has a limited number of items in the population that are susceptible to the failure. This is also known as a 'Defective Subpopulation' model. As such, no matter how long a test continues, it will not be possible for all items to fail (with the particular death/failure).

As an example, we can create a Defective Subpopulation Weibull, also known as a Limited Failure Population Model using a Weibull distribution:

```
import surpyval as surv
import numpy as np
from matplotlib import pyplot as plt

lfp_weibull = surv.Weibull.from_params([10, 2], p=0.6)
np.random.seed(10)
# LFP Model outputs x, c, and n from `random()`
x, c, n = lfp_weibull.random(100)
```

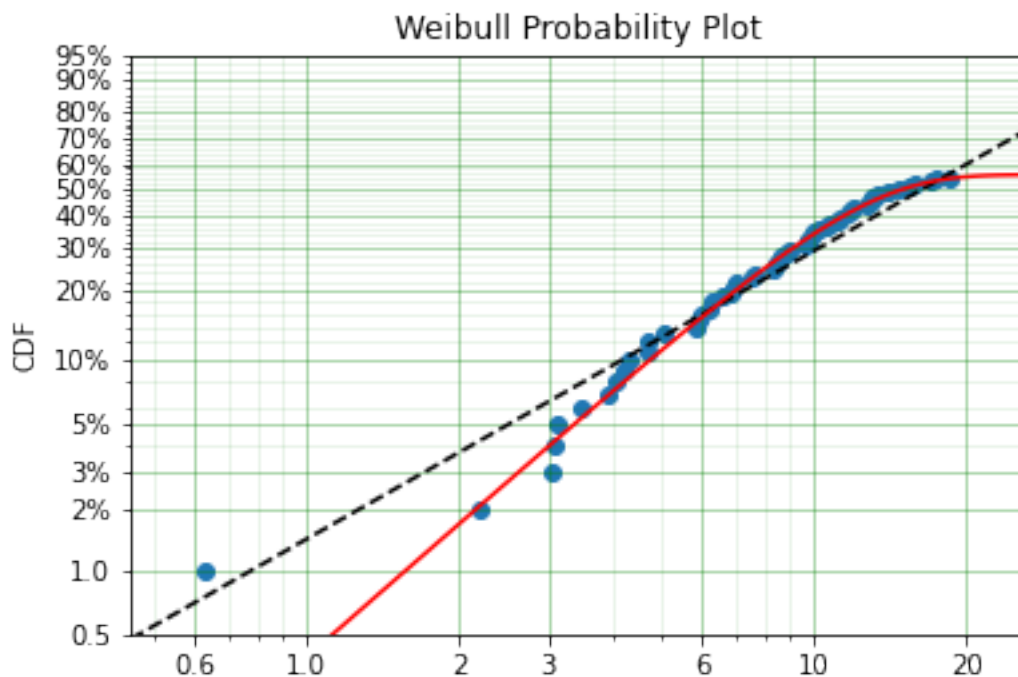
(continues on next page)

(continued from previous page)

```
# Fit regular Weibull
model = surv.Weibull.fit(x=x, c=c, n=n)
model.plot(plot_bounds=False)

# Set LFP to be `True`
lfp_model = surv.Weibull.fit(x=x, c=c, n=n, lfp=True)
print(lfp_model)
xx = np.linspace(np.min(x), np.max(x)*2)
plt.plot(xx, lfp_model.ff(xx), color='red')
```

```
Parametric SurPyval Model
=====
Distribution      : Weibull
Fitted by        : MLE
Max Proportion (p) : 0.5553951704157292
Parameters       :
    alpha: 10.180334244350309
    beta:  2.1358575854287265
```



This API works with any distribution so simply changing Weibull to Exponential would create a Defective Subpopulation Exponential / Limited Failure Population Exponential model. Further, if it was changed to Gamma it would create a Defective Subpopulation Gamma model / Limited Failure Population Gamma.

LFP models can only (as yet) work with MLE. It cannot (yet) work with the other estimation methods. The MSE is a good candidate for implementation.

### 1.11.10 Zero-Inflated Modelling

In survival analysis you might have the scenario where many failure times are 0, known as being dead on arrival. In this case we need a model that can account for the fact that many will be failed at 0, this is a situation that cannot be handled by regular distributions, since most have a 0% chance of failing at 0. Therefore what we need is something that is symmetrical to the LFP/DS case, where a proportion of the failures occur at 0 instead of there being a proportion that will never fail.

```
import surpyval as surv
from autograd import numpy as np

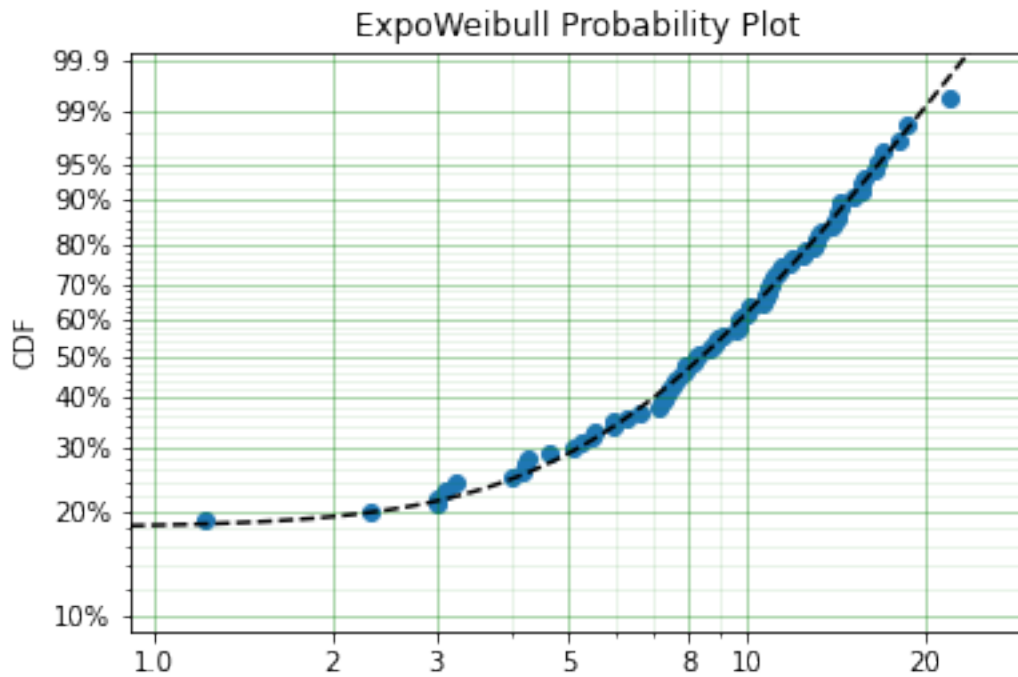
dist = surv.ExpoWeibull
model = dist.from_params([10.2, 2., 1.3], f0=0.15)
np.random.seed(10)
x = model.random(100)
model
```

```
Parametric SurPyval Model
=====
Distribution      : ExpoWeibull
Fitted by        : given parameters
Zero-Inflation (f0) : 0.15
Parameters       :
    alpha: 10.2
    beta: 2.0
    mu: 1.3
```

Using this random data, we can make a fitted model (with the added convenience not offered in the real world of knowing exactly what parameters we are aiming toward).

```
fitted_model = dist.fit(x, zi=True)
print(fitted_model)
fitted_model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : ExpoWeibull
Fitted by        : MLE
Zero-Inflation (f0) : 0.1799999522942094
Parameters       :
    alpha: 11.723925167019866
    beta: 2.769781748379123
    mu: 0.8437868556785479
```



We can see that we have made a good fit!

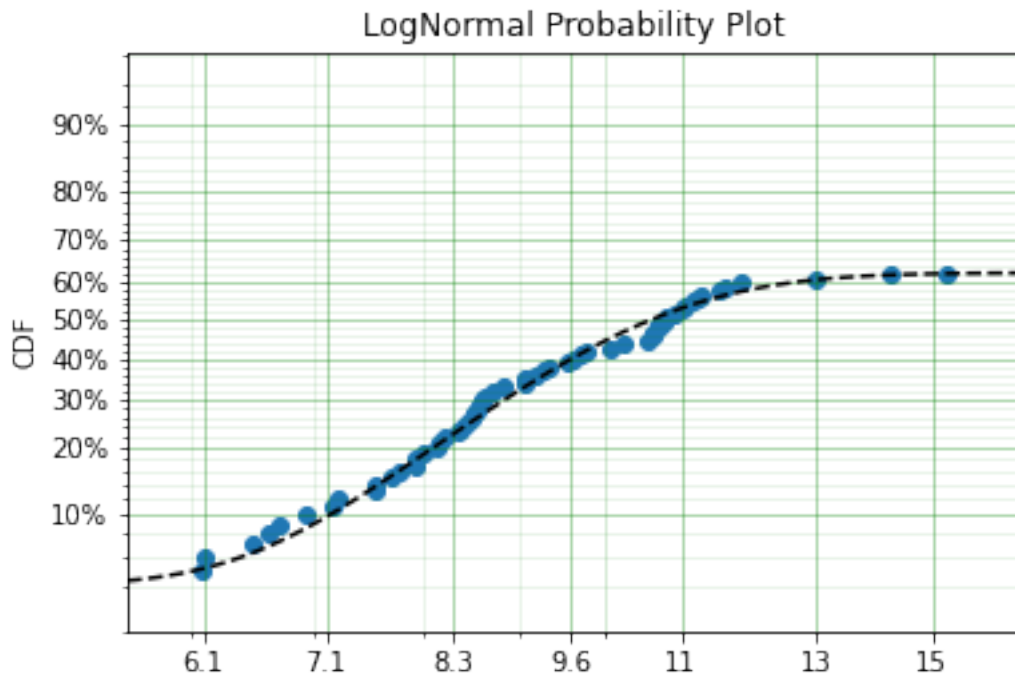
To showcase the SurPyval API again, and to demonstrate the flexibility, it is trivial to have Defective Subpopulation Zero Inflated (DSZI) model / Limited Failure Population and Zero Inflated model.

```
import surpyval as surv
import numpy as np

dist = surv.LogNormal
model = dist.from_params([2.2, .2], f0=0.05, p=0.6)
np.random.seed(10)
# Random values from LFP models come in xcn format!!!!
x, c, n = model.random(100)

fitted_model = dist.fit(x, c, n, zi=True, lfp=True)
print(fitted_model)
fitted_model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : LogNormal
Fitted by        : MLE
Max Proportion (p) : 0.6061204729747632
Zero-Inflation (f0) : 0.040000034963115105
Parameters       :
    mu: 2.2060270833195372
    sigma: 0.19060910628572927
```



Using a `LogNormal` distribution we were able to easily capture the DS/LFP and ZI behaviour of the data.

### 1.11.11 Confidence Intervals

*SurPyval* can be used to compute the confidence interval for any of the functions of a distribution. That is, *SurPyval* can compute the confidence interval for `ff()`, `sf()`, `hf()`, `Hf()`, and `df()`.

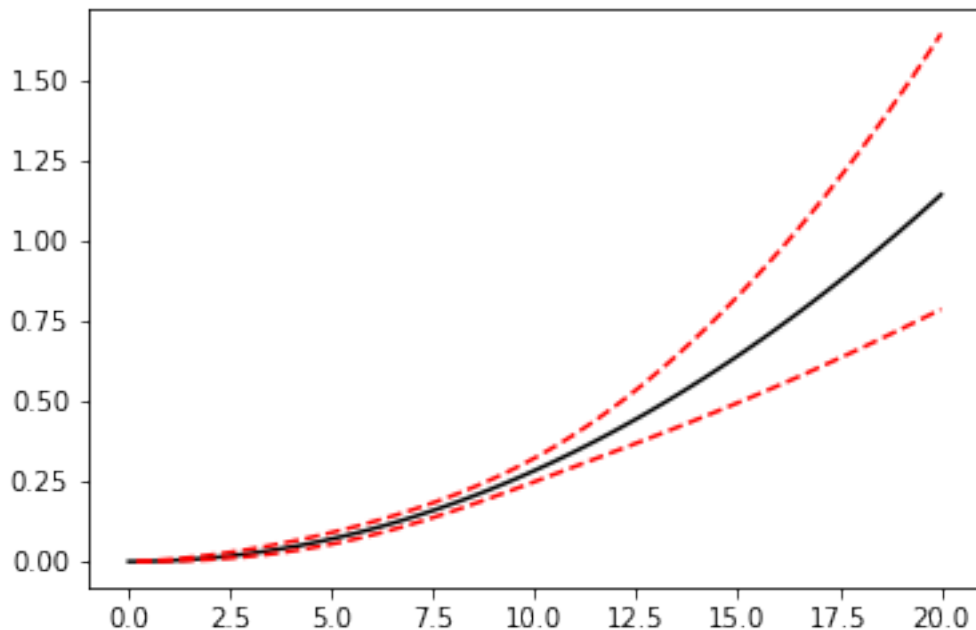
Once you have a model, this can easily be computed with the `cb()` method.

```
from surpyval import Weibull
import numpy as np
from matplotlib import pyplot as plt

x = Weibull.random(100, 10, 3)

model = Weibull.fit(x)

x_plot = np.linspace(0, 20, 100)
plt.plot(x_plot, model.sf(x_plot), color='black')
plt.plot(x_plot, model.cb(x_plot, on='sf', alpha_ci=0.1), color='red', linestyle='--')
```



This shows that we can change the confidence level with `alpha_ci` and that we can change the function for which we want the confidence interval. That is, the `on` keyword can be any of `sf`, `ff`, `df`, `hf`, or `Hf`. This will work with models that you create as well, so even a user defined Distribution will be able to have the confidence intervals computed. Creating these models is discussed in the section below.

### 1.11.12 Creating a custom Distribution

Given the implementation in SurPyval, it is possible to create a new distribution and use all the previously listed techniques. For example, the Gompertz distribution is not implemented in the `surpyval` API, this however can be quickly overcome. First, we set up a random number generator. Because SurPyval works based on the autograd numpy implementation, it is essential that you use the autograd numpy import to make this work.

```
import surpyval as surv
# IMPORTANT - Will not work with regular numpy
from autograd import numpy as np

def qf(p, mu, b):
    return (np.log((-np.log(p)/mu)) + 1)/b

# Generate random values from Gompertz distribution
np.random.seed(1)
x = qf(np.random.uniform(0, 1, 100), .3, 1.1)
```

Now that we have our random data set, we can fit a Gompertz distribution to it. To do so, we need to create a Gompertz distribution class, and to do this we need the cumulative hazard function, the names of the parameters, the bounds of the parameters, and the distribution support.

```
name = 'Gompertz'
```

(continues on next page)



(continued from previous page)

```
def Hf(x, *params):
    return params[0] * np.exp(params[1] * x - 1)

param_names = ['nu', 'b']
bounds = ((0, None), (0, None))
support = (-np.inf, np.inf)
Gompertz = surv.parametric.Distribution(name, Hf, param_names, bounds, support)
```

With this now created, all the calls to the regular surpyval API can be used.

```
Gompertz.fit(x)
```

```
Parametric SurPyval Model
=====
Distribution      : Gompertz
Fitted by        : MLE
Parameters       :
    nu: 1.15060014910275
    b: 1.8973107004872167
```

If we transform the data slightly, we can show that this can be used with censored and truncated data as well.

```
c = np.zeros_like(x)
# Right censor all values above 2
c[x > 2] = 1
x[x > 2] = 2
# Left truncate all values below 0
tl = 0
c = c[x > tl]
x = x[x > tl]

model = Gompertz.fit(x=x, c=c, tl=tl)
model
```

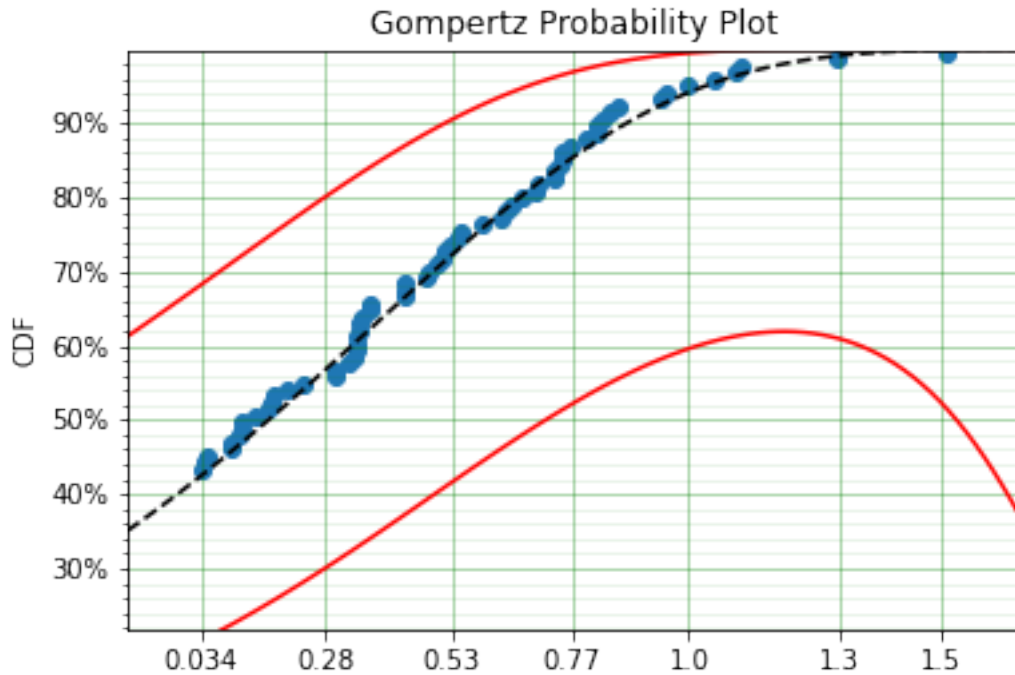
```
Parametric SurPyval Model
=====
Distribution      : Gompertz
Fitted by        : MLE
Parameters       :
    nu: 1.4228615499353794
    b: 1.688152800158132
```

This is extraordinary! We have created a new distribution using only the cumulative hazard function, but are able to handle arbitrary censoring and truncation. It shows the power of the SurPyval API and functionality.

Credit for this idea must be given to the creators of the *lifelines* package. *lifelines* is capable of receiving a cumulative hazard function that can then be used as a distribution to fit parameters. However, at the time of writing it could not handle arbitrarily censored or truncated data.

Even with a user defined `Hf()` we can still use the confidence bounds as well. The results of this can be seen by simply calling the plot function:

```
model.plot(alpha_ci=0.5)
```



You can see that the distribution is not linearised. This is because the  $H_f$  is not readily convertible into the transformation function needed to do the linearisation of the CDF. The defaults are a simple linear scale for both the x and y axis and it shows that the confidence bounds have worked nicely.

You can also see that the confidence bound expands quite widely above approximately 1.2. This is due to the heavy truncation and censoring, if using complete data the confidence bounds do not diverge. This shows the importance of inference when working with truncated and censored data, the uncertainty can be quite wide!

**Warning:** Due to the implementation of confidence bounds in *surpyval* it can result in numeric overflows which results in incredulous bounds. Please take caution when using the `cb` with non *surpyval* implemented distributions.

## 1.12 Regression Modelling with SurPyval

This section is about how we can understand the effect that covariates can have on survival times. As per the other entries in these docs, let's import some useful packages, as such, for the rest of this page we will assume the following imports have occurred:

```
import surpyval as surv
import numpy as np
from matplotlib import pyplot as plt
```

Regression survival modelling with *surpyval* is very easy. This page will take you through a series of scenarios that can show you how to use the features of *surpyval* to get you the answers you need.

### 1.12.1 Semi-Parametric - Cox Proportional Hazards Model

The first example is the Cox Proportional Hazards model. In this example we will use the data from Krivtsov et al. This data set is the results of testing tires time to failure with measurements about those tires. The authors of this paper intended to determine what factors affected tire reliability.

```
from surpyval.datasets import Tires
from surpyval import CoxPH

x = Tires.data['Survival']
c = Tires.data['Censoring']
Z = Tires.data[['Tire age', 'Wedge gauge', 'Interbelt gauge', 'EB2B', 'Peel force',
               'Carbon black (%)', 'Wedge gauge×peel force']]
model = CoxPH.fit(x=x, Z=Z, c=c)
model
```

```
Regression SurPyval Model
=====
Type                : Proportional Hazards
Kind                : Cox
Parameterization    : Semi-Parametric
Parameters          :
    beta_0 : 2.109496593744941
    beta_1 : -9.686078593632743
    beta_2 : -10.67809536747068
    beta_3 : -13.67594841333851
    beta_4 : -34.29448581473381
    beta_5 : -48.35286747450483
    beta_6 : 20.84037862912251
```

We can see that we have a mixture of coefficients. We can check the p-values:

```
print(model.p_values)
```

```
[0.13000628, 0.03677433, 0.02074066, 0.0918419 , 0.01200102, 0.14831534, 0.01867559]
```

We can see that it is only 1, 2, 4, and 6 that are significant at the 0.05 level.

We can redo the model using only those covariates:

```
from surpyval.datasets import Tires
from surpyval import CoxPH

x = Tires.data['Survival']
c = Tires.data['Censoring']
Z = Tires.data[['Wedge gauge', 'Interbelt gauge', 'Peel force', 'Wedge gauge×peel_
↪force']]
model = CoxPH.fit(x=x, Z=Z, c=c)
print(model.p_values)
model
```

```
[0.02207978 0.01368372 0.00956108 0.01030372]
```

```
Regression SurPyval Model
=====
Type                : Proportional Hazards
Kind                : Cox
```

(continues on next page)

(continued from previous page)

```

Parameterization      : Semi-Parametric
Parameters            :
  beta_0 : -9.313960179920473
  beta_1 : -7.069295556681021
  beta_2 : -27.413473066027667
  beta_3 : 18.105822313415462

```

All the coefficients can now be seen to be significant. It also shows that as the wedge gauge, interbelt gauge, and peel force increase, the hazard rate will decrease and the life will therefore increase. The opposite is the case for the wedge gauge x peel force coefficient.

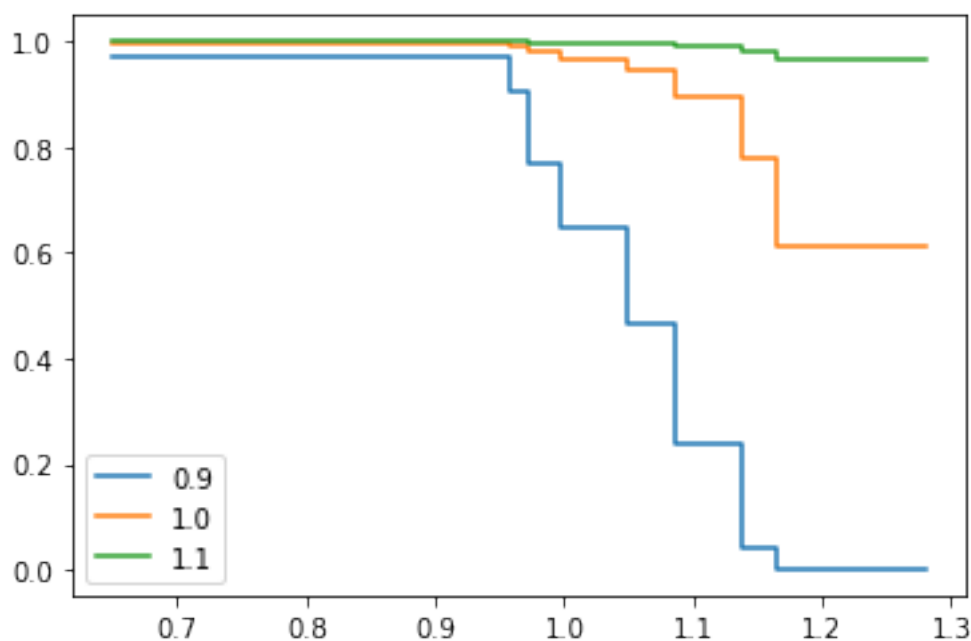
We can plot the survival curves of the average tire and the 10% above and 10% below average tire:

```

Z_mean = Tires.data[['Wedge gauge', 'Interbelt gauge', 'Peel force', 'Wedge_
↪gauge×peel force']].mean().values

plot_x = np.linspace(x.min(), x.max())
for f in [0.9, 1., 1.1]:
    plt.step(plot_x, model.sf(plot_x, Z=Z_mean * f), label=f)
plt.legend()

```



We can see that as the covariates increase there is a decrease in the probability of survival up to 1.2. The Semi-Parametric nature of the model can also be seen clearly in this plot. You can see that the baseline is non-parametric, but the baseline has been affected by the covariates.

### 1.12.2 Parametric Proportional Hazards Modelling

In the above example we used a semi-parametric model where the ‘baseline’ hazard rate was a non-parametric model but the hazard was multiplied by a parametric function of the covariates. We can use fully parametric models instead. These come with the advantages of parametric models, namely extrapolation, but are also disadvantaged by the assumption needed about the shape of the distribution. SurPyval has two Proportional Hazard models that are ready to

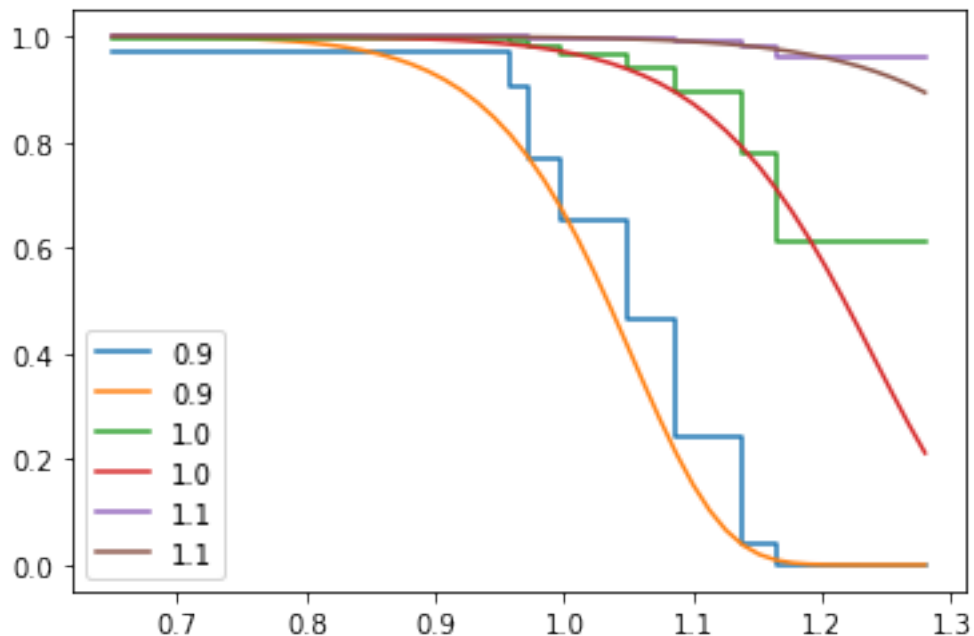
use with any number of covariate inputs (just like the CoxPH model); these are the *ExponentialPH* and the *WeibullPH* models. We will analyse the tires data using the Weibull Proportional hazards model.

```
from surpyval.datasets import Tires
from surpyval.regression import WeibullPH

x = Tires.data['Survival']
c = Tires.data['Censoring']
Z = Tires.data[['Wedge gauge', 'Interbelt gauge', 'Peel force', 'Wedge gauge×peel_
↪force']]
weibull_ph_model = WeibullPH.fit(x=x, Z=Z, c=c)
weibull_ph_model
```

```
Parametric Regression SurPyval Model
=====
Kind                : Proportional Hazard
Distribution         : Weibull
Regression Model    : Log Linear (Exponential)
Fitted by          : MLE
Distribution        :
    alpha: 0.24255057163126237
    beta: 16.057788534593193
Regression Model   :
    beta_0: -9.165054735694651
    beta_1: -7.998600691841399
    beta_2: -27.50328118957879
    beta_3: 18.38550168401127
```

You can see from the above that the coefficients for the covariates are very similar.

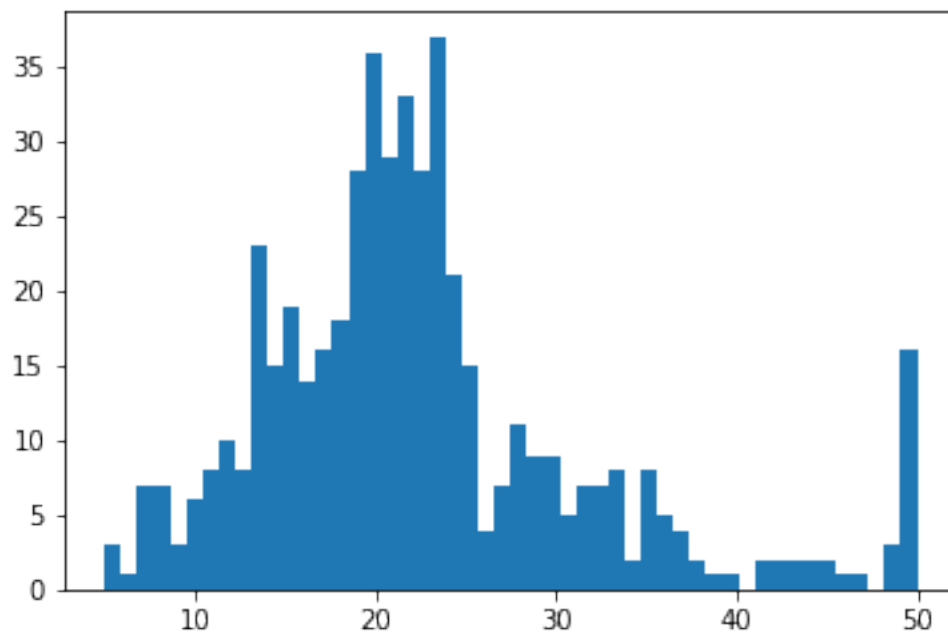


## 1.13 Example Applications

This section documents some of the applications that SurPyval as a survival analysis toolkit can be useful to you, no matter what discipline you need it for.

### 1.13.1 Boston House Prices

No statistical analysis package can avoid doing the ‘hello world’ task of analysing the ‘Boston House Pricing’ dataset. What might surprise some readers is that this would even be considered. . . The myriad blogs and Kaggle posts looking into this problem can not surely be improved upon. I agree, however, it is a good example of why one needs to be aware of censoring and how flexible the SurPyval API is when dealing with it. Looking at the boston house pricing dataset you can see that there is a suspicious number of houses at the top end that all have the same price:



On consideration, one can see that there are no houses above \$50,000 and that the density at that point is much higher than we would expect because we would expect some form of a ‘fat-tail.’ That is, we should expect a decreasing number of houses at the highest costs. It is therefore safe to conclude that all values above \$50,000 have been set to \$50,000; which is to say that the sale price is right censored! And because it is a censored observation we will need to use an analysis tool that can handle censored observations, lest we may be wrong in our estimates of the distribution of housing prices. So let's load the data and see what results we can get, starting with the raw data.

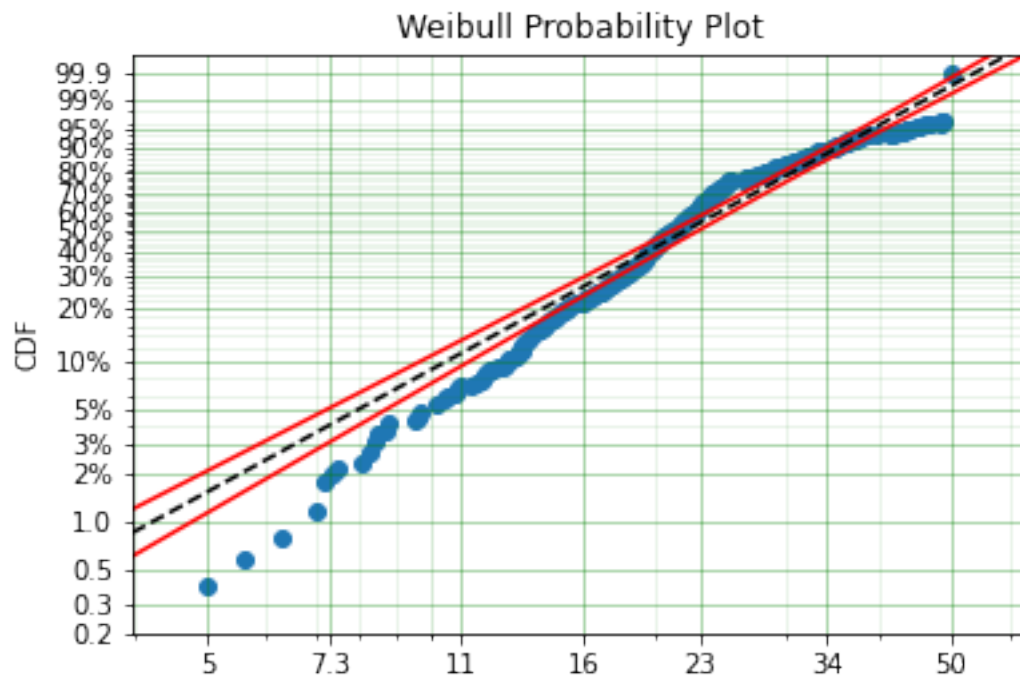
```
import surpyval as surv
x = Boston.df['medv'].values
x, c, n = surv.xcn_handler(x)

model = surv.Weibull.fit(x, c, n)
print(model)
model.plot()
```

```

Parametric SurPyval Model
=====
Distribution      : Weibull
Fitted by        : MLE
Parameters       :
    alpha: 25.386952832397032
    beta:  2.5651903209947684

```



From the above plot you can see that near 50, the parametric model diverges substantially from the actual data. So we can see that having not censored the highest values means that our model could be improved by doing so. Let's see:

```

import surpyval as surv
x = Boston.df['medv'].values
x, c, n = surv.xcn_handler(x)
# Right censor the highest value
c[-1] = 1

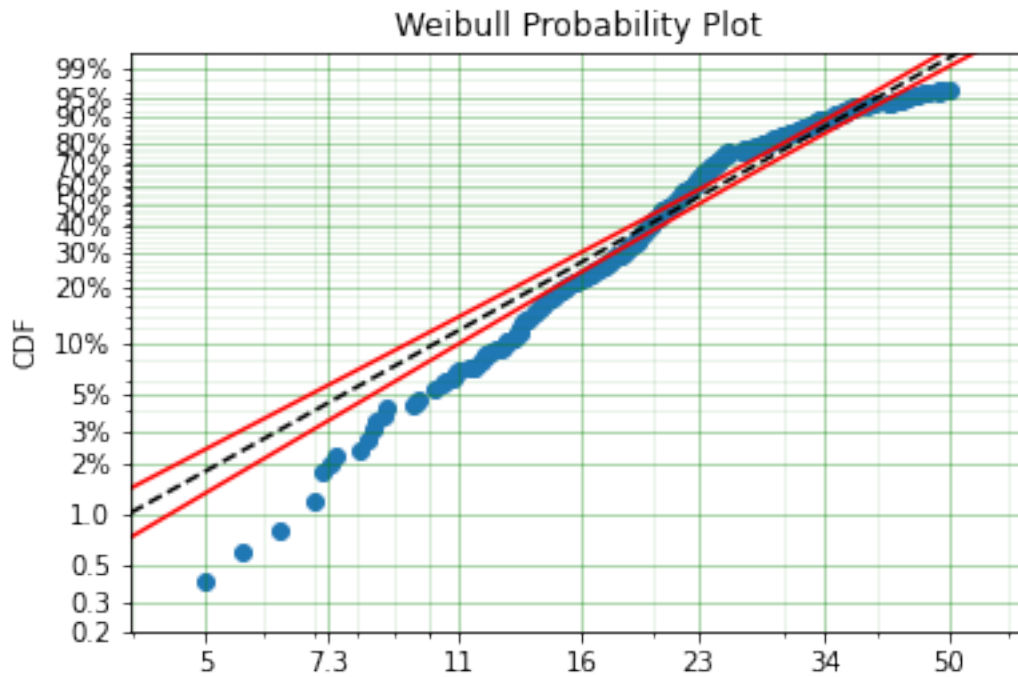
model = surv.Weibull.fit(x, c, n)
print(model)
model.plot()

```

```

Parametric SurPyval Model
=====
Distribution      : Weibull
Fitted by        : MLE
Parameters       :
    alpha: 25.53669601993307
    beta:  2.469159446459548

```



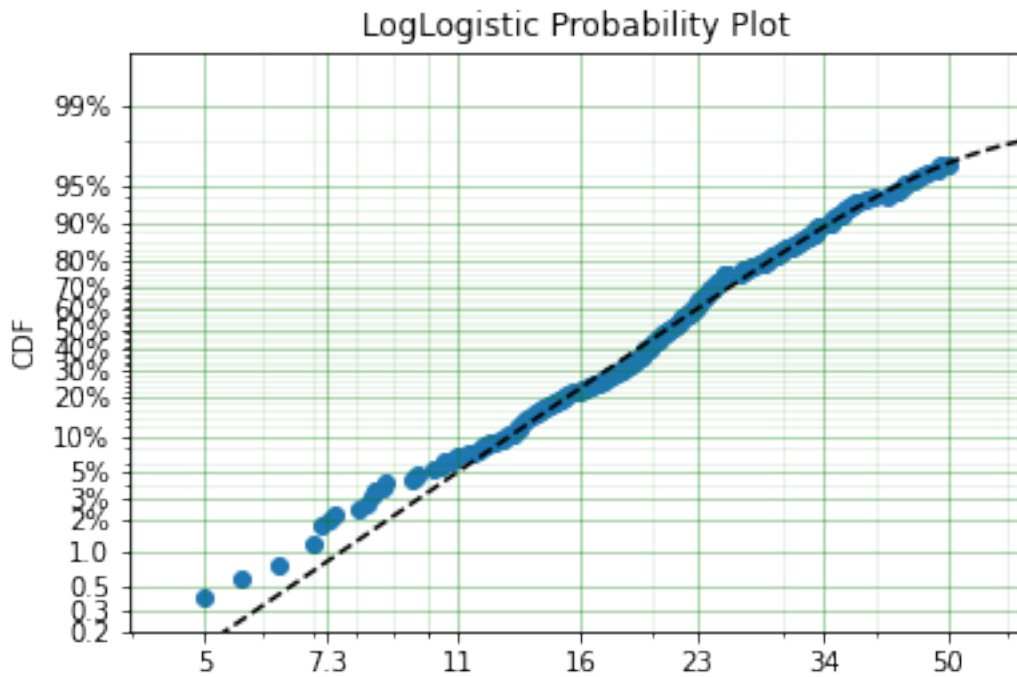
We can see that the model has changed slightly, however, there appears to be a ‘disconnect’ near 24 that makes the model a poor fit above and below the value. Let’s see if a different distribution will improve the fit:

```
import surpyval as surv
x = Boston.df['medv'].values
x, c, n = surv.xcn_handler(x)
# Right censor the highest value
c[-1] = 1

model = surv.LogLogistic.fit(x=x, c=c, n=n, lfp=True)
print(model)
model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : LogLogistic
Fitted by        : MLE
Max Proportion (p) : 0.9861133787129936
Parameters       :
    alpha: 20.804405478058186
    beta: 4.56190516414644
```





This appears to be a much better fit, however, there is still quite a bit of difference between the data and the model in the middle of the distribution. Lets create a custom spline to see if we can perfect the fit.

```
import surpyval as surv
x = Boston.df['medv'].values
x, c, n = surv.xcn_handler(x)
# Right censor the highest value
c[-1] = 1

def Hf(x, *params):
    x = np.array(x)
    Hf = np.zeros_like(x)
    knot = params[0]
    params = params[1:]
    dist1 = surv.Weibull
    dist2 = surv.LogLogistic
    Hf = np.where(x < knot, dist1.Hf(x, *params[0:2]), Hf)
    Hf = np.where(x >= knot, (dist1.Hf(knot, *params[0:2])
                             + dist2.Hf(x, *params[2::])), Hf)
    return Hf

bounds = ((0, 50), (0, None), (0, None), (0, None), (0, None),)
param_names = ['knot', 'alpha_w', 'beta_w', 'alpha_ll', 'beta_ll']
name = 'WeibullLogLogisticSpline'
support = (0, np.inf)

WeibullLogLogisticSpline = surv.Distribution(name, Hf, param_names, bounds, support)

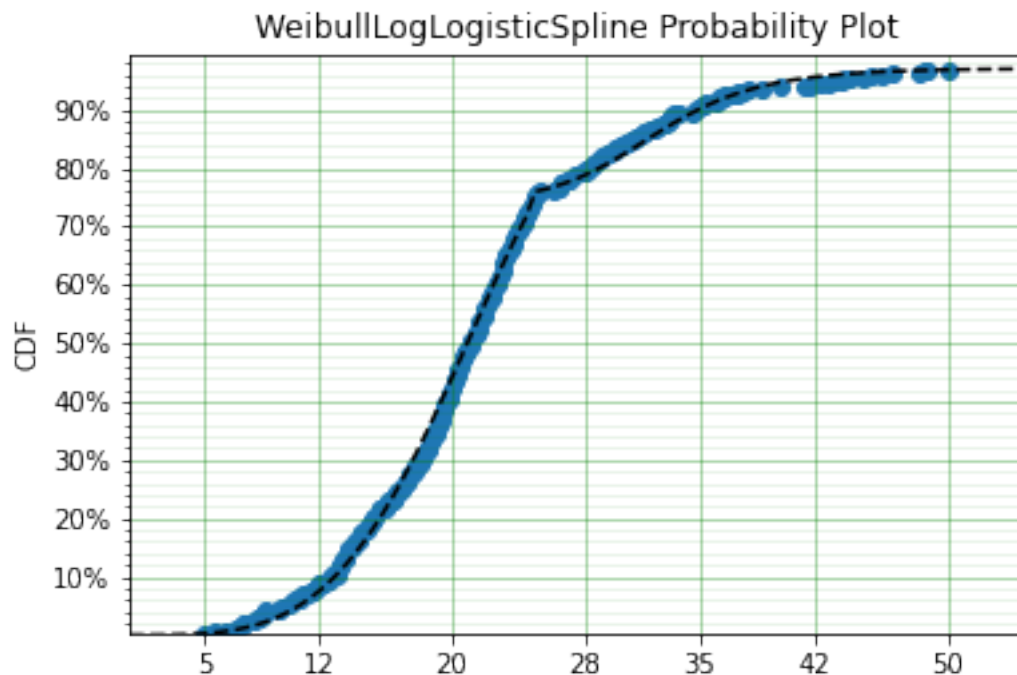
model = WeibullLogLogisticSpline.fit(x=x, c=c, n=n, lfp=True)

print(model)
model.plot()
```

```

Parametric SurPyval Model
=====
Distribution      : WeibullLogLogisticSpline
Fitted by        : MLE
Max Proportion (p) : 0.9711459340639835
Parameters       :
    knot: 25.0000103742294
    alpha_w: 22.735658691657452
    beta_w: 3.926996942307611
    alpha_ll: 32.2716411336919
    beta_ll: 10.120540049344006

```



Much better!

It must be said that this is a bit ‘hacky’. There is no theory that we are using to guide the choice of the spline model, we are simply finding the best fit to the data. For example, this model would not be able to be used for extrapolation too far beyond \$50,000, this is because the model is limited to 97.1% of houses. A separate spline would be needed to model those data. However, the example shows the importance of censoring and the power of the surpyval API!

### 1.13.2 Applied Reliability Engineering

In reliability engineering we might be interested in the proportion of a population that will experience a particular failure mode. We do not want to ship the items that will fail so that our customers do not have a poor experience. But, we will want to determine the minimum duration of a test that can establish whether a component will fail. This is because a test that is too long we will waste time and money in testing and if a test is too short we will ship too many items that will fail in the field. We need to optimise this interval to minimize the cost of testing but also the number of items at risk in the field.

Using data from the paper that introduced the Limited Failure Population model (also known as the Defective Sub-population) to the reliability engineering world [Meeker] we can show how surpyval can be used in part to calculate

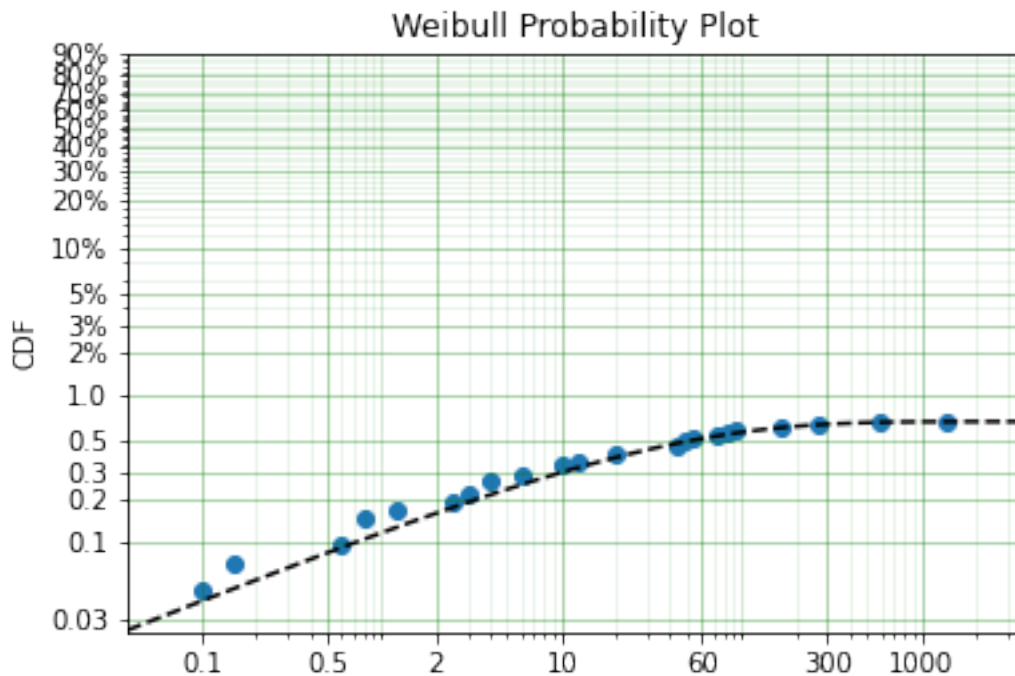
an optimal ‘burn-in’ test duration.

```
import surpyval as surv

f = [.1, .1, .15, .6, .8, .8, 1.2, 2.5, 3., 4., 4., 6., 10., 10.,
     12.5, 20., 20., 43., 43., 48., 48., 54., 74., 84., 94., 168., 263., 593.]
s = [1370.] * 4128

x, c, n = surv.fs_to_xcn(f, s)
model = surv.Weibull.fit(x, c, n, lfp=True)
print(model)
model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : Weibull
Fitted by        : MLE
Max Proportion (p) : 0.006744450944727198
Parameters       :
    alpha: 28.367193779799038
    beta: 0.4959762140288241
```



We can see from these results that at maximum we will have approximately 0.67% fail. If the company accepts a 0.1% probability of their products failing in the field then we can calculate the interval at which the difference between the total population and the proportion failed in the test is 0.1%.

```
from scipy.optimize import minimize
fun = lambda x : (0.001 - np.abs(model.p - model.ff(x)))**2

res = minimize(fun, 10, tol=1e-50)
print(res.x)
```

```
[104.43741352]
```

Therefore we should do a burn in test up to approximately 104.4 to make sure we minimize the number of items shipped that are defective while also minimizing the duration of the test. We can simply change the value of 0.001 in the above code to any value we may wish to use.

### 1.13.3 Demographics / Actuarial

In demographics and actuarial studies, the distribution of the life of a population is of interest. For the demographer, it is necessary to understand how a population might change, in particular, how the expected lifespan is changing over time. The same applies to an actuary, an actuary is interested in lifetimes to understand the risk of payouts among those who own a life insurance policy.

The [Gompertz-Makeham](#) is a distribution used in demography and actuarial studies to estimate the lifetime of a population. This can be implemented in `surpyval` with relative ease.

```
import surpyval as surv
from autograd import numpy as np
from matplotlib import pyplot as plt
from scipy.special import lambertw

bounds = ((0, None), (0, None), (0, None),)
support = (0, np.inf)
param_names = ['lambda', 'alpha', 'beta']
def Hf(x, *params):
    Hf = params[0] * x + (params[1]/params[2])*(np.exp(params[2]*x))
    return Hf

GompertzMakeham = surv.Distribution('GompertzMakeham', Hf, param_names, bounds,
↪support)
```

We now have a GM distribution object that can be used to fit data. But we need some data:

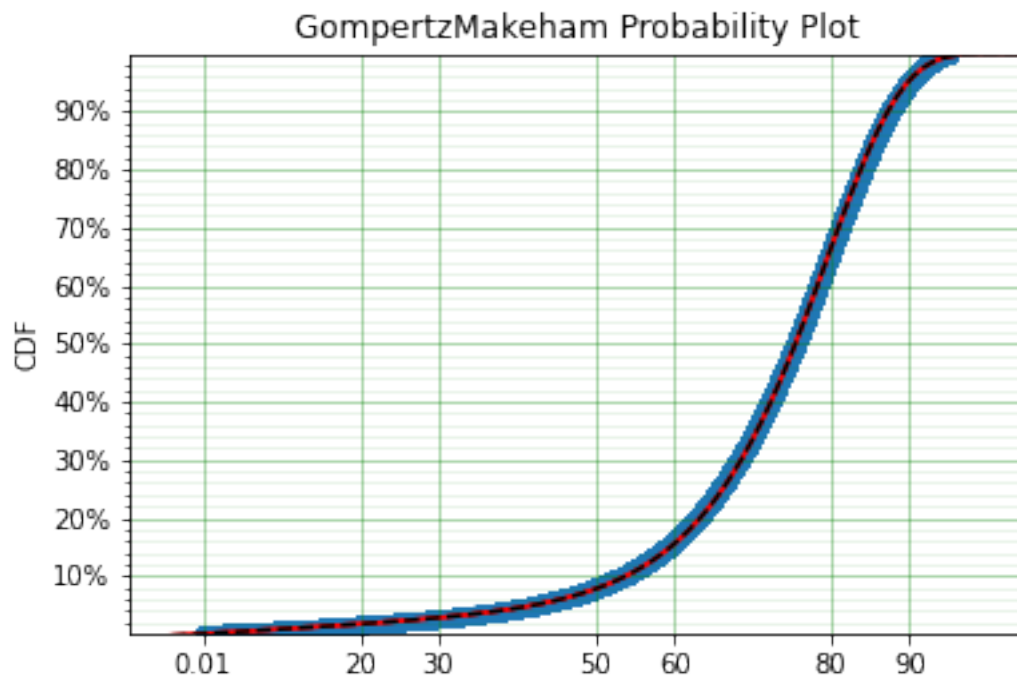
```
# GM qf()
def qf(p, params):
    lambda_ = params[0]
    alpha = params[1]
    beta = params[2]
    return (alpha/(lambda_ * beta) - (1./lambda_)*np.log(1 - p)
           - (1./beta)*lambertw((alpha*np.exp(alpha/lambda_)*(1 - p)**(-(beta/lambda_
↪))))/(lambda_))).real

np.random.seed(1)
params = np.array([.68, 28.7e-3, 102.3])/1000
x = qf(np.random.uniform(0, 1, 100_000), params)
# Filter out some numeric overflows.
x = x[np.isfinite(x)]
```

The parameters for the distribution come from [\[Gavrilov\]](#), specifically the parameters for the lifespans of the 1974-1978 data. So in this case we have (simulated) data on the lifespans of 100,000 thousand people and we need to determine the GM parameters. This can be compared to the historic parameters to see if the age related mortality has changed or has remained roughly constant. To do so, all we need do with `surpyval` is to put the data to the `fit()` method.

```
model = GompertzMakeham.fit(x)
model.plot(alpha_ci=0.99, heuristic='Nelson-Aalen')
model
```

```
Parametric SurPyval Model
=====
Distribution      : GompertzMakeham
Fitted by        : MLE
Parameters       :
  lambda: 0.0007827108147066848
  alpha: 2.1199267751549727e-05
  beta: 0.10690878152126947
```



You can see that the model is a good fit to the data. Using the model we can determine the probability of death in a given term for a random individual from the population. This is useful to price the premium of a life insurance policy. For example, if a 60 year old was to take out a two year policy, what premium should we charge them for the policy. First, we need to determine the probability of death:

```
p_death = model.fff(62) - model.fff(60)
policy_payout = 100_000
expected_loss = policy_payout * p_death
print(p_death, expected_loss)
```

```
0.025337351289907883 2533.7351289907883
```

From the results above, you can see that the probability of death over the two year interval is approximately 2.5%. Given the contract is to payout \$100,000 in this event, the expected loss is therefore \$2,533.74. Therefore, to make a profit, the policy will need to cost more than \$2,533.74. So say the company has a strategy of making 10% from each policy, the policy cost to the individual would therefore be \$2,787.11. If we divide this payment scheme into a per month basis over the two years we get a monthly payment of \$116.13 for two years (in the case of death the amount

owing can be subtracted from the payout).

Although this is a basic example, as insurance companies would have much more sophisticated models, it shows the basics of how demographic and actuarial data can be used. This shows the application of surpyval to actuarial and demographic studies.

## 1.13.4 Applied Reliability Engineering - 2

In reliability engineering you can come across the case where a new product has been built that is similar in design to a previous, but has better materials, geometry, seals.. etc. You have data from the tests of the old product and new results for the same test on the new product. The only problem, the new product only had one failure in the test! What will you do?

Given the similarities, it is common to use the same shape parameter, the  $\beta$  value, from a similar product as an initial estimate. In this case, we may need to know the reliability of the item in the field. We can create a model of this new product, but first the old product:

```
import surpyval as surv

x_old = [ 5.2, 10.7, 16.3, 22. , 32.9, 38.6, 42.1, 58.7, 92.8, 93.8]
old = surv.Weibull.fit(x_old)
print(old)
```

```
Parametric SurPyval Model
=====
Distribution      : Weibull
Fitted by        : MLE
Parameters       :
    alpha: 45.27418484669478
    beta:  1.377623372184365
```

We can use the above value of beta with the new data:

```
x_new = [87, 100]
c_new = [0, 1]
n_new = [1, 9]

surv.Weibull.fit(x_new, c_new, n_new, fixed={'beta' : 1.3776}, init=[100.])
```

```
Parametric SurPyval Model
=====
Distribution      : Weibull
Fitted by        : MLE
Parameters       :
    alpha: 525.1398140084557
    beta:  1.3776
```

The characteristic life of the new bearing is over 10 times higher! Quite an improved new design. This new model can be used as part of the sales of the new product (10x more life!) and to provide recommendations for maintenance.

## 1.13.5 Social Science / Criminology

Another application of surpyval is when encountering extreme values. The Weibull distribution is one of the limiting cases of the Generalized Extreme Value distribution. In other words, the Weibull distribution is the distribution that can model the strength of a chain because it can model the extreme value, in this case the minimum, of a collection of

distributions. A chain is as only as strong as it's weakest link. If there are many many links in a chain (which is a fair assumption) then links of which follow a known strength distribution, then the strength of the chain will follow a Weibull distribution. It is for this reason that the Weibull distribution is so widely used.

Another extreme value is the maximum. The maximum extreme value distribution is the Frechet distribution. But, if you simply inverse a minimum, you can get a maximum. Therefore, if we know our data is following a process of finding a maximum, then we can use the Weibull distribution to model the phenomena.

**Warning:** This may be a distressing topic for some readers.

Social scientists and criminologists are interested in understanding the phenomena of mass shootings in an effort to eliminate the scourge from society. A mass shooting is an extreme event, and an extreme event can be modelled to understand the risks of future occurrence, and with that understanding, the effect of interventions can also be understood.

Using the gun violence data from [Kaggle](https://www.kaggle.com/jameslko/gun-violence-data) we can model the process. That is, if we take the maximum number of deaths in a given month over several years, we have data that can be used to estimate the probability of something even worse occurring. This data covers the period from 2013 to 2018, see Kaggle for more details.

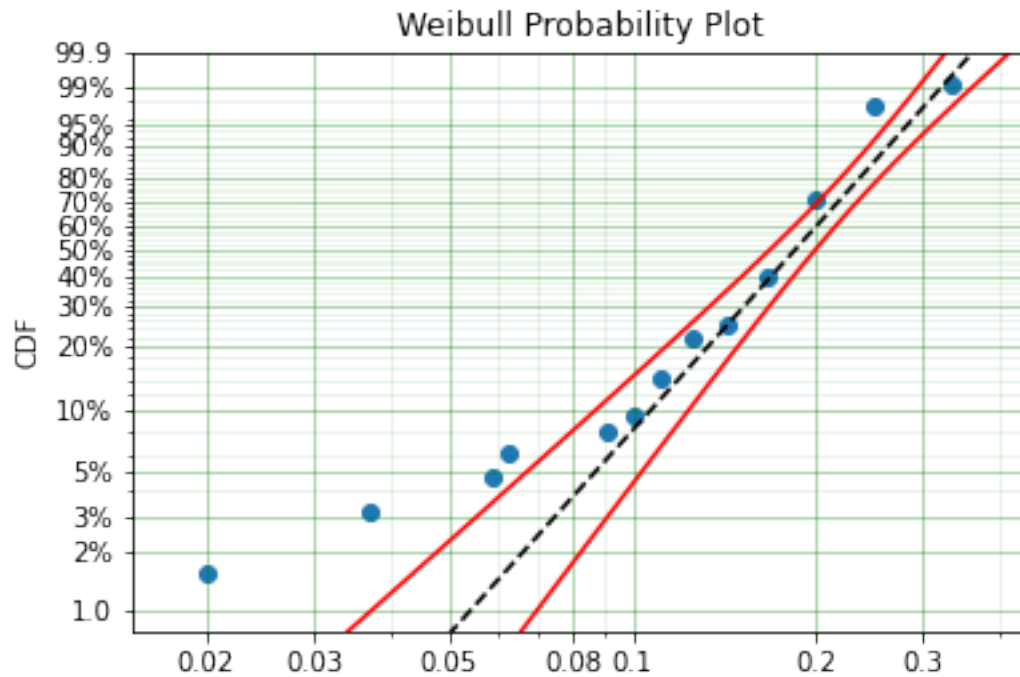
```
import surpyval as surv
import pandas as pd

# Data not in surpyval, available at https://www.kaggle.com/jameslko/gun-violence-data
gun_violence_df = pd.read_csv('../gun-violence-data_01-2013_03-2018.csv', parse_
    ↪dates=['date'])

# Find the maximum number of people killed each month
gun_violence_df = gun_violence_df.groupby(pd.Grouper(key='date', freq='M')).agg({'n_
    ↪killed' : 'max'})

x = df['n_killed'].values

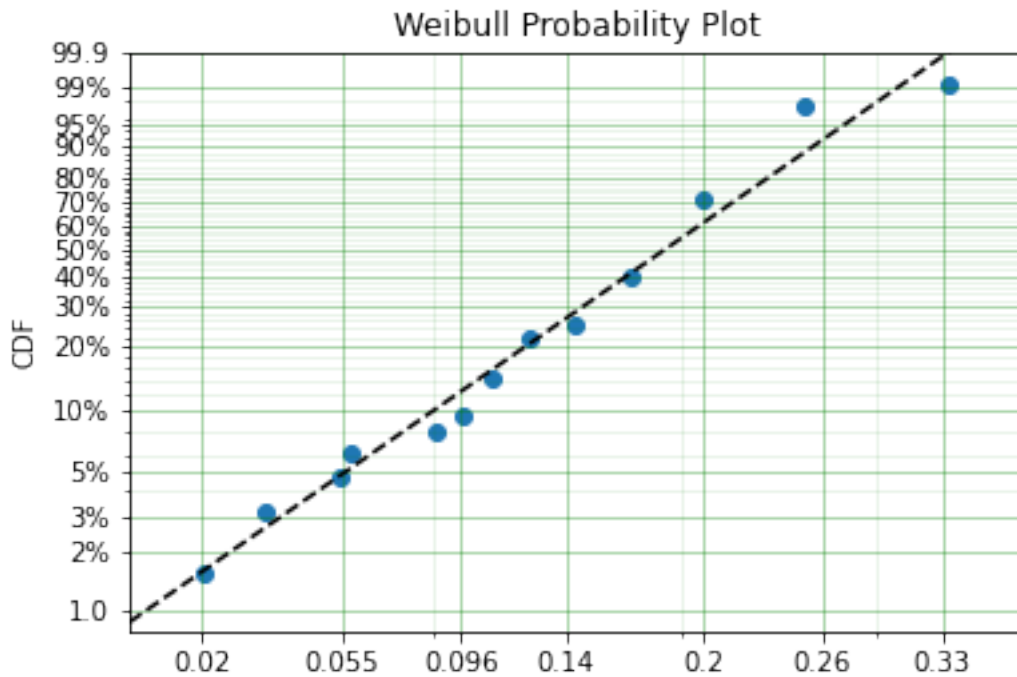
# Inverse the data to get the maximum
model = surv.Weibull.fit(1./x)
model.plot()
```



It is worth reminding that since we have taken the inverse, it is the lower values that represent more victims. And it is the extremes that we are trying to capture. You can see from the above plot that the model does not fit the data from 0.02 to 0.1 very well. We can try using a different approach

```
# Inverse the data to get the maximum
mpp_model = surv.Weibull.fit(1./x, how='MPP', offset=True)
mpp_model.plot()
```





You can see that this model is a much better description of the data. However, the problem is that it cannot have a real interpretation. Because the offset is negative, that means there is a non-zero probability of 0, which because the data was inversed, means that there is a non-zero probability of having a shooting with infinite victims. This model is therefore not a good option for such extreme extrapolations. The model can however, be used to estimate the probability of having a shooting as bad or worse than the most extreme event up to 2040.

```
p_happening = mpp_model.ff(1./50)
p_not_happening = 1 - p_happening
# Months from 2022 to 2040
months = 12 * (2040 - 2022)
p_not_happening_before_2040 = (p_not_happening)**(months)
(1 - p_not_happening_before_2040)*100
```

```
(1.6077640040390584, 96.98325003600236)
```

The model estimates that there is an approximately 1.6% chance of an event killing 50 or more people in a given month, which may seem low, however, because there are 216 months between 2022 and 2040 the chances of not having as extreme an event over that time period becomes horrifyingly small. The model suggests that the probability of having a month in which an event with more than 50 people will be killed, has a 97.0% chance of happening from 2022 to 2040. Chilling.

This is a bit higher than other reports of the same prediction, see [Duwe] who report at 35% probability, which is some, but not even close to complete, relief.

### 1.13.6 Economics

Economists are interested in the times between recessions. This information helps them formulate policy proscriptions that may (or may not) reduce the duration of a recession, or the time between recessions. Using data from Tadeu Cristino et al. [TC] we can use real data to estimate the probability of a recession.

```
import pandas as pd
import surpyval as surv

start = [np.nan, "June 1857", "October 1860", "April 1865", "June 1869",
         "October 1873", "March 1882", "March 1887", "July 1890", "January 1893",
         "December 1895", "June 1899", "September 1902", "May 1907", "January 1910",
         "January 1913", "August 1918", "January 1920", "May 1923", "October 1926",
         "August 1929", "May 1937", "February 1945", "November 1948", "July 1953",
         "August 1957", "April 1960", "December 1969", "November 1973",
         "January 1980", "July 1981", "July 1990", "March 2001", "December 2007"]

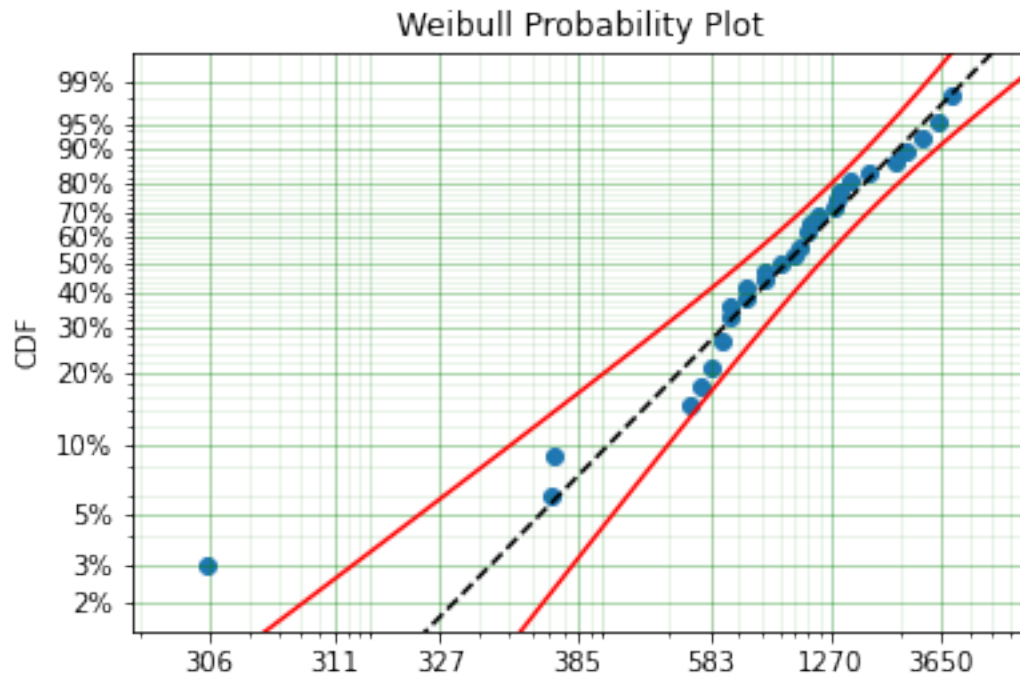
end = [
    "December 1854", "December 1858", "June 1861", "December 1867", "December 1870",
    "March 1879", "May 1885", "April 1888", "May 1891", "June 1894", "June 1897",
    "December 1900", "August 1904", "June 1908", "January 1912", "December 1914",
    "March 1919", "July 1921", "July 1924", "November 1927", "March 1933",
    "June 1938", "October 1945", "October 1949", "May 1954", "April 1958",
    "February 1961", "November 1970", "March 1975", "July 1980", "November 1982",
    "March 1991", "November 2001", "June 2009"
]

df = pd.DataFrame({'start' : pd.to_datetime(start),
                  'end' : pd.to_datetime(end)})

# Compute time from end of last recession to peak of next.
x = (df.start - df.end.shift(1)).dropna().dt.days.values

model = surv.Weibull.fit(x, offset=True)
print(model)
model.plot()
```

```
Parametric SurPyval Model
=====
Distribution      : Weibull
Fitted by        : MLE
Offset (gamma)    : 304.0659320899125
Parameters       :
    alpha: 895.3220718605215
    beta: 1.0629492868473804
```



You can see from the above the data is a good fit to the model! Great. So now what?

We can communicate what the expected time between recessions is:

```
model.mean()
```

```
1178.2499033086633
```

Therefore the average growth period is 1,178 days, or about 3.2 years between recessions.

## 1.13.7 References

## 1.14 API

### 1.14.1 Non-Parametric

### 1.14.2 Parametric

#### Distribution Classes

#### Exponential

#### Gamma

#### Gumbel

Logistic

LogLogistic

LogNormal

Normal

Uniform

Weibull

Exponentiated Weibull

Parametric Class

Parametric Fitter

Parametric Mixture Model

### 1.14.3 Datasets

## 1.15 Changelog

### 1.15.1 v0.11.0 (planned)

- General ALT fitter full release
- General PH fitter full release
- Formulas
- Add more than [Breslow](#) to the CoxPH methods.
- Parameter confidence bound
- Document the rationale behind using Fleming-Harrington as the default.
- Docs on how to integrate with Pandas
- Docs for CoxPH
- Docs for Accelerated Life fitters
- Create a `RegressionFitter` class. I keep copying code across the three fitters.
- Allow truncation with zi and lfp models.
- Allow truncation with regression

### 1.15.2 v0.10.1.0 (25 Mar 2022)

- Changed plot methods to now take 'Axis' object. This allows a user to pass in an existing axis.

- plot functions now return an `Axis` object instead of the `Lines2D` object. Allows for easy user update after plotting.
- Added `fs_to_xcn` as it was dropped in 10.0.1.
- Changed all imports for numpy to be done from the `surpyval` module. This will allow for easy maintenance in future in the event of deprecated `autograd`.

### 1.15.3 v0.10.0.1 (22 Nov 2021)

- Removed `fsl_to_xcn` function and replaced with `fsli_to_xcn` function that is able to take any combination of `fsli`.

### 1.15.4 v0.10.0 (9 Aug 2021)

- Version snapshot for JOSS review

### 1.15.5 v0.9.0 (5 Aug 2021)

- Better initial estimates in the `_parameter_initialiser` for the `lfp` data (use max F from `nonp` estimate...)
- [issue #13](#) - Better failures when insufficient data provided.
- [issue #12](#) - Created `fsli_to_xcn` helper function.
- Fixed bug in confidence bounds implementation for offset distributions. CBs were not using the offset and were therefore way out. Now fixed.
- Created a `NonParametric.cb()` method to match `Parametric` API for confidence bounds.
- Cleaned up `NonParametric` code (removed some technical debt and duplicated code).
- Changed the `__repr__` function in `NonParametric` to be aligned to `Parametric`
- Updated the docstring for `fit()` for `NonParametric`
- Fixed bug in `NonParametric` that required the `x` input to be in order for the functions (e.g. `df` etc.).
- `CoxPH` released.
- General AL fitter in beta
- General PH fitter in beta
- Created `Linear`, `Power`, `InversePower`, `Exponential`, `InverseExponential`, `Eyring`, `InverseEyring`, `DualPower`, `PowerExponential`, `DualExponential` life models.
- Created `GeneralLogLinear` life model for variable stress count input.
- For each combination of a `SurPyval` distribution and life model, there is an instance to use `fit()`. For example there are `WeibullDualExponential`, `LogNormalPower`, `ExponentialExponential` etc.
- **Docs Updates:**
  - **Add application examples to docs:**
    - \* Reliability Engineering
    - \* Actuary / Demography
    - \* [Social Science/Criminology](#)
    - \* Boston Housing

- \* Medical science
- \* [Economics](#)
- \* Biology - Ware, J.H., Demets, D.L.: Reanalysis of some baboon descent data. *Biometrics* 459–463 (1976).

### 1.15.6 v0.8.0 (27 July 2021)

- Made backwards incompatible changes to LFP models, these are now created with the `lfp=True` keyword in the `fit()` method
- Created ability to fit zero-inflated models. Simply pass the `zi=True` option to the `fit()` method.
- Changes to `utils.xcnt_handler` to ensure `x`, `xl`, and `xr` are handled consistently.
- changed the way `__repr__` displays a `Parametric` object.
- Changed the default for plotting to be Fleming-Harrington. This was a result of seeing how poorly the Nelson-Aalen method fits zero inflated models. FH therefore offers the best performance of a Non-Parametric estimate at the low values of the survival function (as KM reaches 0 for fully observed data) and at high values (KM is good but NA is poor).
- Added a Fleming-Harrington method to the `Turnbull` class.
- Improved stability with dedicated `log_sf`, `log_ff`, and `log_df` functions. Less chance of overflows and therefore better convergence.
- Changed interpolation method of `NonParametric`. Allows for use of cubic interpolation
- Changed `from_params` to accept `lfp` and `zi` (or any combo)
- Changed `random()` in `Parametric` so that `lfp` or `zi` models can be simulated!
- Improved the way `surpyval` fails
- Substantial docs updates.

### 1.15.7 v0.7.0 (19 July 2021)

- Major changes to the confidence bounds for `Parametric` models. Now use the `cb()` method for every bound.
- Removed the `OffsetParametric` class and made `Parametric` class now work with (or without) an offset.
- Minor doc updates.

## 1.16 Support

If you need help with survival analysis, please ask a question on [stats.stackexchange](https://stats.stackexchange.com).

If you've searched the `surpyval` documentation for what you've been looking for and can't find it, please add as suggestion for a feature on GitHub. `SurPyval` is a growing tool. Or, if you need help with `surpyval` feel free to email Derryn at [derryn.knife@gmail.com](mailto:derryn.knife@gmail.com).

## 1.17 Contributing

If you want to contribute to SurPyval, please do! Please review the current open [feature requests](#) to see if your desired feature is in the requests. If not, please raise a new one to notify the community. We can assign you feature for you to branch and develop.

SurPyval is in the process of complying with the PEP8 standard so please make all contributions as per that standard.

## 1.18 Acknowledgements

The Cox Proportional Hazards, Competing Risks, and Fine and Gray sections were developed for Cartiga LLC. A big thanks to Cartiga for allowing the code to be open sourced.



The core of the surpyval package continues to be maintained with the support of Reliafy.



## 1.19 Installation

*surpyval* can be installed easily with the pip command:

```
$ pip install surpyval
```





## CHAPTER 2

---

### Indices and tables

---

- `genindex`
- `modindex`
- `search`



---

## Bibliography

---

- [Bagdonavicius] Bagdonavicius, V., & Nikulin, M. (2001). Accelerated life models: modeling and statistical analysis. CRC press.
- [KM] Kaplan, E. L., & Meier, P. (1958). Nonparametric estimation from incomplete observations. *Journal of the American statistical association*, 53(282), 457-481.
- [NA] Nelson, Wayne (1969). Hazard plotting for incomplete failure data. *Journal of Quality Technology*, 1(1), 27-52.
- [FH] Fleming, Thomas R and Harrington, David P (1984). Nonparametric estimation of the survival distribution in censored data. *Communications in Statistics-Theory and Methods*, 13(20), 2469-2486.
- [TB] Turnbull, Bruce W (1976). The empirical distribution function with arbitrarily grouped, censored and truncated data. *Journal of the Royal Statistical Society: Series B (Methodological)*, 38(3), 290-295.
- [TC] Tadeu Cristino, C., Żebrowski, P., & Wildemeersch, M. (2020). Assessing the time intervals between economic recessions. *PloS one*, 15(5), e0232615.
- [Cole] Cole SR, Hudgens MG. Survival analysis in infectious disease research: describing events in time. *AIDS*. 2010;24(16):2423-31.
- [Duwe] Duwe, G., Sanders, N. E., Rocque, M., & Fox, J. A. (2021). Forecasting the Severity of Mass Public Shootings in the United States. *Journal of Quantitative Criminology*, 1-39.
- [Gavrilov] Gavrilov, L. A., Gavrilova, N. S., & Nosov, V. N. (1983). Human life span stopped increasing: why?. *Gerontology*, 29(3), 176-180.
- [Meeker] William Q. Meeker (1987) Limited Failure Population Life Tests: Application to Integrated Circuit Reliability, *Technometrics*, 29(1), 51-65