

Fitting distributions with R

The why and how of distribution fitting in R

Sean Manzi¹

¹PenARC(NIHR Applied Research Collaboration South West Peninsula)
University of Exeter

NHS-R conference, 2021

Workshop overview

- 1 What is a distribution?
- 2 When would I use a distribution?
- 3 Sampling data from a distribution
- 4 Fitting data to a distribution
- 5 A distribution fitting app template

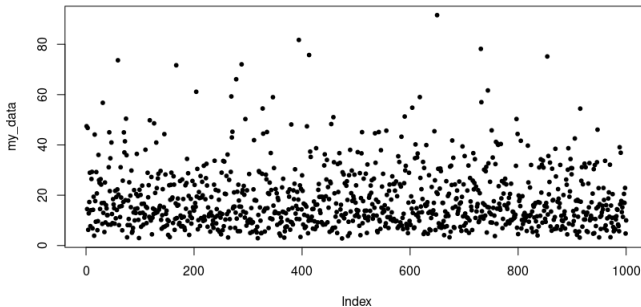
What is a distribution?

A distribution is:

- A description of the shape of some data
- It describes the probability of a value occurring in the data
- A concise mathematical description of a curve that describes the spread of potential values

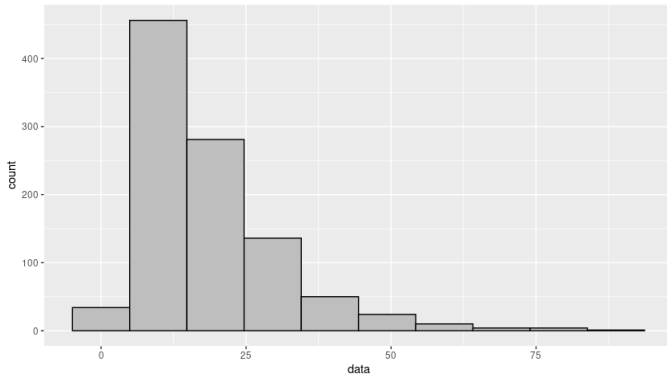
What makes a distribution - raw data

A scatter plot of some raw data



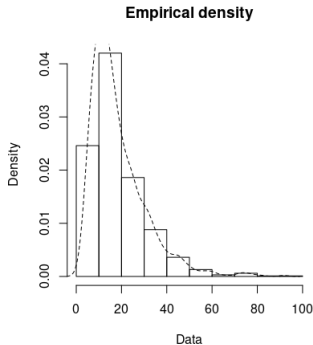
What makes a distribution - histogram

A histogram of the raw data



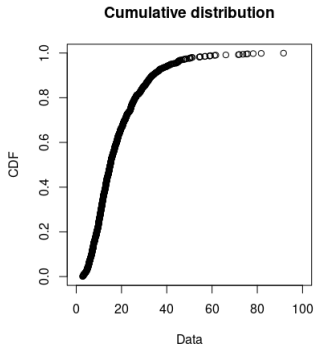
What makes a distribution - PDF

The probability density function of the data distribution



What makes a distribution - CDF

The cumulative density function of the data distribution



Types of distributions

There are two main types or groups of distributions:

- Discrete distributions
- Continuous distributions

Discrete named distributions

The most common discrete distributions are binomial, uniform discrete, poisson and geometric

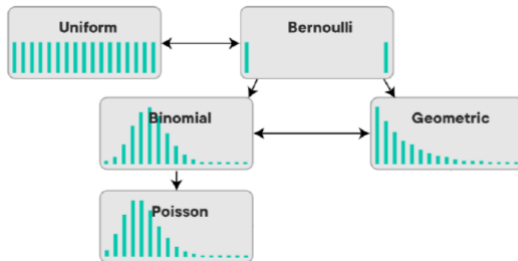


Figure: Image credit:

<https://bernard-mlab.com/post/probability-distribution/>

Continuous named distributions

The most common continuous distributions are normal, log-normal, exponential, weibull, beta and gamma

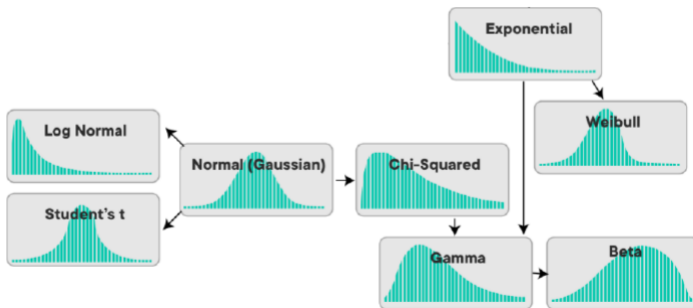


Figure: Image credit:

<https://bernard-mlab.com/post/probability-distribution/>

Generalised uses for distributions

Distributions are used for:

- Inferential statistics
 - Hypothesis testing
 - Determining uncertainty
- Predictive modelling
 - Parameterising a model
- Data engineering
 - Imputing missing data

Real world applications

Probability distributions are used in many areas such as:

- Manufacturing process control
- Work flow planning
- Robotics
- Insurance
- Health service modelling

Where I use distributions in health service modelling

- Determining inter-arrival times (IAT's)
- Determining how long an activity will take
- Attributing characteristics to individuals e.g. male/female, age etc.
- Determining the probability of an event occurring
- Determining the probability of a particular treatment or process outcome

When I might need to sample from a distribution?

When we have a known distribution which has been derived from a fitting process, we need to be able to sample data from that distribution for use in our model or algorithm

This can be achieved using sampling functions for named distributions

Sampling approaches

For most named distributions there are built in functions in R to generate random numbers from a defined distribution

There are two main approaches for sampling from a known distribution

- Directly sampling one or more values from a distribution
- Generating a set of two or more values from a distribution and sampling from these

The first approach is used when you want to be able to sample from the entire distribution with replacement

The second approach is used when you want to be able to pre-determine the number or range of values sampled from the distribution and sample these with or without replacement

Sampling functions

Distribution random sampling functions tend to have the general form of

function name(number of values to be sampled, kwargs for distribution parameters)

For example, when sampling from the uniform distribution the function name is 'runif', we define the number of values to be sampled (e.g. 100), the minimum value 'min=' and the maximum value 'max='

```
uniform <- runif(100, min=0, max=90)
```


Other distribution sampling functions

Examples of other functions for sampling from named distributions:

```
normal <- rnorm(100, mean=5, sd=1.5)
exponential <- rexp(100, rate=1.6)
poisson <- rpois(1000, lambda=4)
lognormal <- rlnorm(1000, meanlog=2.7, sdlog=0.6)
```

If we want to sample from the distribution values that we have created we can use the **sample** function. The args for this function are: a vector of values from which the sample is to be drawn, the number of values to be sampled, whether a value can be sampled more than once (with replacement). An example of this using the data we previously produced is:

```
our_sample <- sample(normal, 10, replace=TRUE)
```

Task

Try creating different distributions using the functions that we have just looked at. Change the argument inputs to examine the impact these have on the shape of the data by plotting it as a histogram – `hist()`.

Tip: Try changing the number of bins in your histogram to get different levels of resolution on the sampled data.

A basic distribution fitting process

- Look at the shape of your data
- Fit your data to likely distributions
- Check the fit of the data

Stop IMPORTANT!

Why can't I just use my real data, why use a distribution at all?

- Your real data is only a sample of all possible values; the population
- The use of real data in a stochastic model causes over fitting
- We need variation aside from that contained in the data to better approximate the population

Stop IMPORTANT!

What do I do when multiple distributions or no distributions fit my data very well?

- Distribution fitting is as much an art as it is a science
- The process relies on interpretation, experience and perseverance
- Real world data is messy, you will have to try multiple distributions and tweak them until they work

Stop IMPORTANT!

I have 1,000,000 data points, should I use them all?

- Definitely not!
- Using too many data points will again cause over fitting
- A rule of thumb is to use between 2,000 and 10,000 data points
- The greater the variation in the data the more data that will be needed to estimate that variation

Plot your data

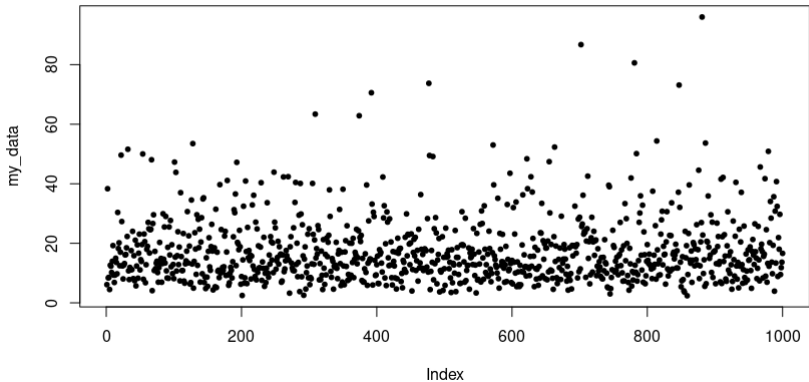
Make a simple plot of your data and look at the shape of it

```
library(fitdistrplus)
library(ggplot2)
```

```
set.seed(12)
#Import data and plot
my_data <- rlnorm(1000, meanlog=2.7, sdlog=0.6)
plot(my_data, pch=20)
```


Plot your data

A simple scatter plot of the data



Plot your data

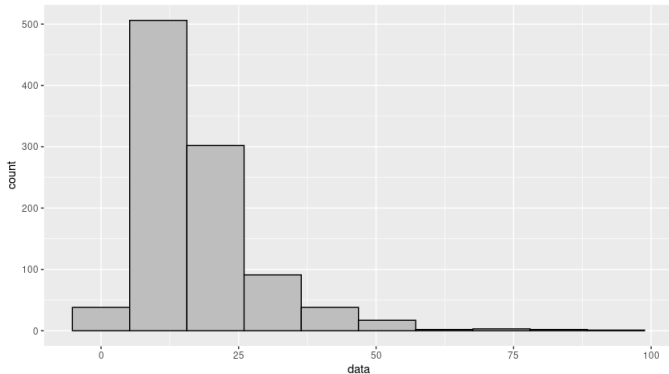
A histogram is a more useful way to view your data

```
df = data.frame("data" = my_data)
ggplot(data=df) +
  geom_histogram(mapping=aes(x=data),bins=10,
                 col="black",
                 fill="grey")
```

Here we use `ggplot2` to plot the data or you could use the base R `hist()` function

Plot your data

The histogram of the data shows us the frequency of the values within the data

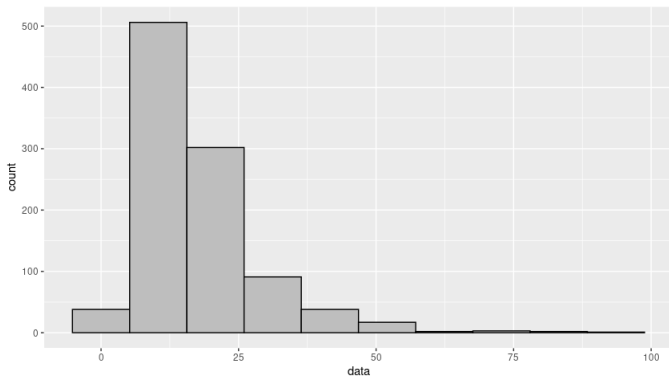


What you are looking for in the shape of your data:

- Is it symmetrical? Likely normal/gaussian distribution
- Is it positively or negatively skewed?
- Is there one or more modal values? Mode = most common
- Is the mode the lowest or highest value? likely an exponential distribution

Plot your data

What characteristics can you see in the data? Symmetry? Skew?
Mode number and location?



Empirical density and cumulative distribution

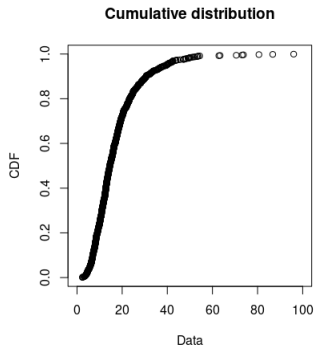
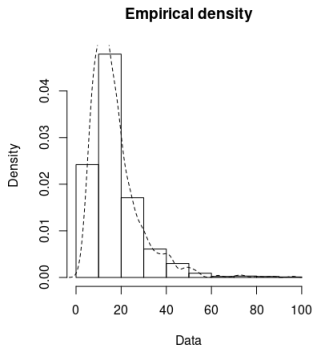
You can use the **plotdist** function to see the Probability Density Function (PDF) and Cumulative Density Function (CDF) of your data

- Empirical density (PDF) - equivalent to histogram giving probability density of observations
- Cumulative distribution (CDF) - Adds up probability density of observations

```
plotdist(my_data, histo = TRUE, demp = TRUE)
```

Empirical density and cumulative distribution

The PDF (left) and CDF (right) of the data



Cullen and Frey graph

This type of graph is used to assess the potential fit of the data in terms of skewness (+ve or -ve skew) and kurtosis (sharpness of the peak of the curve)

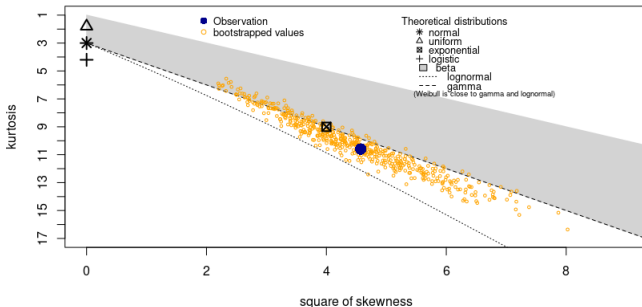
```
descdist(my_data, discrete=FALSE, boot=500)
```

If fitting a discrete distribution change the discrete argument value to TRUE

Cullen and Frey graph

In the Cullen and Frey graph note how the observation and bootstrap values are more closely aligned to the gamma distribution line. This is misleading and we will see why in a minute.

Cullen and Frey graph



Fitting your data

Now that you have an indication of which distributions might fit your data best you can start to test them

The process for this is as follows:

- Fit the data to likely distributions
- Visually assess the fit using:
 - Theoretical PDF
 - Theoretical CDF
 - Q-Q plot: comparison of quantiles, also most sensitive
 - P-P plot: comparison of CDF's
- Calculate goodness of fit and uncertainty estimates
- Select the best fitting distribution and extract the distribution parameters
- Or retest using different distributions

Fitting your data

You use the **fitdist** function to fit a particular distribution to your data and obtain parameters for a theoretical fitted curve

```
fit_w <- fitdist(my_data, "weibull")
summary(fit_w)
```

The function takes the vector of your sample data as the first argument and the name of the distribution to fit as the second argument

By default the maximum likelihood estimate (mle) method is used but there are other methods that can be used. However, mle is a good generic estimation method for distribution fitting. For more information see <https://www.rdocumentation.org/packages/fitdistrplus/versions/1.1-6/topics/fitdist>

Fitting your data

The output from the `fitdist` function is viewed using the `summary` function and takes the fitting output as its input. Below is an example of the fitting summary output

Fitting of the distribution ' weibull ' by maximum likelihood Parameters :

```

            estimate Std. Error
shape  1.721185  0.03836032
scale 19.545755  0.38103448
Loglikelihood: -3637.931  AIC:  7279.862  BIC:  7289.678
Correlation matrix:

```

	shape	scale
shape	1.0000000	0.3344119
scale	0.3344119	1.0000000

Fitting your data

To expedite the fitting process you can fit more than one distribution at a time using lists and for loops as can be seen below

```
dists <- c("gamma","lnorm","weibull")
fit <- list()
for (i in 1:length(dists))
  fit[[i]] <- fitdist(my_data, dists[i])
```

```
for (i in 1:length(dists))
  print(summary(fit[[i]]))
```

Fitting your data

Fitting multiple
distributions
enables you to print
the summaries the
multiple summaries

```
Fitting of the distribution ' gamma ' by maximum likelihood
Parameters :
      estimate Std. Error
shape 3.1561369 0.134341294
rate  0.1824111 0.008415521
Loglikelihood: -3580.725   AIC: 7165.45   BIC: 7175.266
Correlation matrix:
      shape      rate
shape 1.0000000 0.9225761
rate  0.9225761 1.0000000
```

```
Fitting of the distribution ' lnorm ' by maximum likelihood
Parameters :
      estimate Std. Error
meanlog 2.6841381 0.01819330
sdlog   0.5753227 0.01286443
Loglikelihood: -3550.253   AIC: 7104.505   BIC: 7114.321
Correlation matrix:
      meanlog sdlog
meanlog      1      0
sdlog        0      1
```

```
Fitting of the distribution ' weibull ' by maximum likelihood
Parameters :
      estimate Std. Error
shape 1.721185 0.03836032
scale 19.545755 0.38103448
Loglikelihood: -3637.931   AIC: 7279.862   BIC: 7289.678
Correlation matrix:
      shape      scale
shape 1.0000000 0.3344119
scale 0.3344119 1.0000000
```

Fitting your data

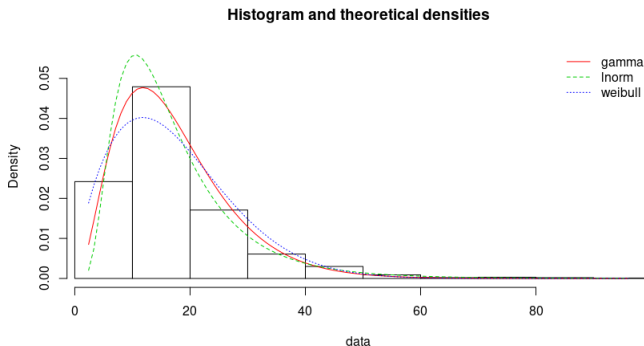
You then use the output from the `distfit` function to produce graphs of the PDF's, CDF's, Q-Q plots and P-P plots. Note the creation of a legend using the distribution names in a list format.

```
par(mfrow=c(2,2))
plot.legend <- dists
denscomp(fit, legendtext = plot.legend)
cdfcomp (fit, legendtext = plot.legend)
qqcomp (fit, legendtext = plot.legend)
ppcomp (fit, legendtext = plot.legend)
```

Note: The par function is used to set the plot parameters to produce a 2 x 2 grid of the four plots. This can be omitted and the plots produced individually

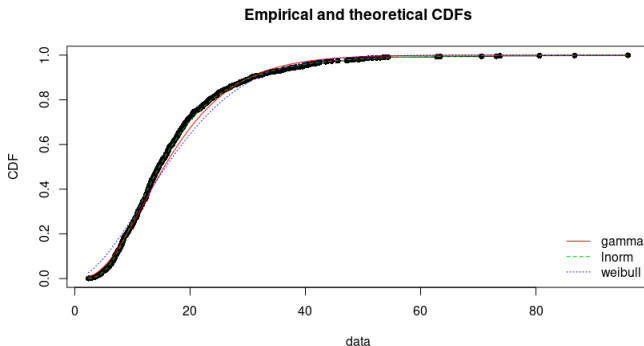
Visually assessing the fit

On the PDF plot you are trying to judge how closely the PDF curve of the theoretical distributions approximates the data



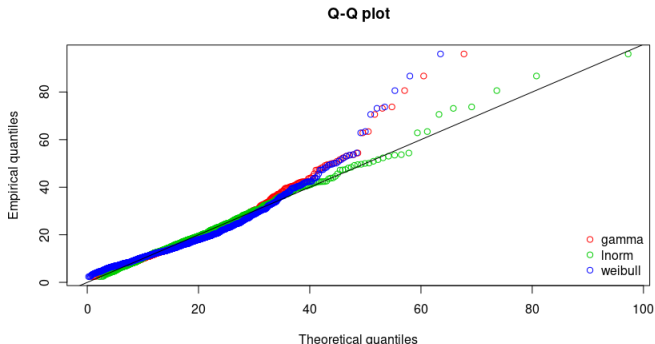
Visually assessing the fit

On the CDF plot you are trying to judge how closely the CDF curve of the theoretical distributions approximates the data



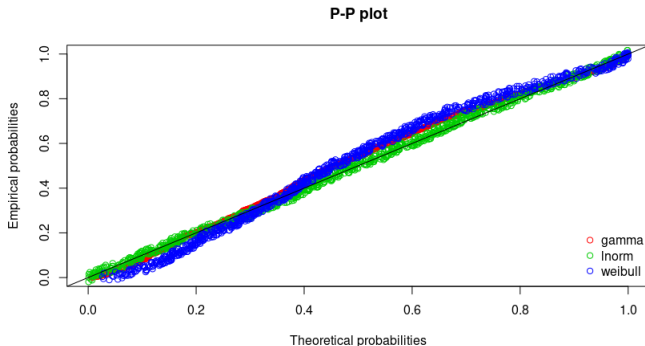
Visually assessing the fit

On the Q-Q plot you are trying to judge how much the theoretical distribution quantiles deviate from the straight line through the plot. This is similar to the least squares method where you seek to minimise the residuals



Visually assessing the fit

On the P-P plot you are trying to judge how much the theoretical CDF deviates from the straight line through the plot. This again is similar to the least squares method where you seek to minimise the residuals



Goodness of fit

To calculate the goodness of fit statistics for our fitted distributions we use the **gofstat** function passing in our list of fitted data and the names of the distributions as a list

```
f <- gofstat(fit, fitnames=c("gamma","lnorm","weibull"))
```

○

When we view the output we receive three different statistics each calculating goodness of fit in slightly different ways. These are all calculating the deviation of the theoretical distributions from the real data. We are aiming to minimise these statistical values.

Goodness-of-fit statistics

	gamma	lnorm	weibull
Kolmogorov-Smirnov statistic	0.05344531	0.02079443	0.07977389
Cramer-von Mises statistic	0.83505652	0.06739170	2.17626274
Anderson-Darling statistic	4.93931231	0.38442499	14.24261892

Goodness-of-fit criteria

	gamma	lnorm	weibull
Akaike's Information Criterion	7165.450	7104.505	7279.862
Bayesian Information Criterion	7175.266	7114.321	7289.678

Goodness of fit

If we look at the full output of the gofstat calculations we can see more information. In the example below we see that statistically the fit of these distributions have been rejected. Do not be alarmed this often happens due to the difficulty of fitting data. This is why we need to use all of the previous steps rather than relying only on statistical tests of fit.

Name	Type	Value
f	list [15] (S3: gofstat.fitdist, fi	List of length 15
chisq	double [3]	64.5 28.9 137.0
chisqbreaks	double [27]	5.05 6.43 7.17 7.97 8.51 9.38 ...
chisqpvalue	double [3]	2.42e-05 2.68e-01 2.03e-17
chisqdf	double [3]	25 25 25
chisqtable	double [28 x 4]	36.0 36.0 36.0 36.0 36.0 36.0 53.2 42.0 26.5 ..
cvm	double [3]	0.8351 0.0674 2.1763
cvmtest	character [3]	'rejected' 'not computed' 'rejected'
ad	double [3]	4.939 0.384 14.243
adtest	character [3]	'rejected' 'not computed' 'rejected'
ks	double [3]	0.0534 0.0208 0.0798
kstest	character [3]	'rejected' 'not rejected' 'rejected'
aic	double [3]	7165 7105 7280

Parameter uncertainty

The fitted distribution parameters are only estimates of the 'true' distribution. The **bootdist** function undertakes a bootstrapping simulation to estimate the confidence interval for the parameters. This can be useful for bounding any tweaking of the distribution parameters when implementing them in a model.

```
for (i in 1:length(fit))
  ests <- bootdist(fit[[i]], niter = 1e3)
  print(paste0("****",dists[i],"****"))
  print(summary(ests))
```

Parameter uncertainty

Here is the output from the bootdist function. This provides the median, upper bound and lower bound values of the parameter estimates based on a bootstrap simulation

```
[1] "****gamma****"
Parametric bootstrap medians and 95% percentile CI
      Median      2.5%      97.5%
shape 3.1538563 2.9151636 3.4313436
rate  0.1826134 0.1673977 0.1993214
[1] "****lnorm****"
Parametric bootstrap medians and 95% percentile CI
      Median      2.5%      97.5%
meanlog 2.6851781 2.6488932 2.7188573
sdlog   0.5750318 0.5493428 0.5988646
[1] "****weibull****"
Parametric bootstrap medians and 95% percentile CI
      Median      2.5%      97.5%
shape 1.717855 1.638245 1.805857
scale 19.531898 18.842982 20.309847
```


The fitted distribution parameters

The aim of this whole distribution fitting process has been to gain the parameters for a distribution that has been fitted to your real world data. These parameters are found in the `fitdist` function output under `estimate`. You can then refine the distribution when applying it in practice, if necessary, using the upper and lower bounds of the uncertainty estimates. Original distribution parameters: $\text{meanlog} = 2.7$, $\text{sdlog} = 0.6$

```
Fitting of the distribution 'lnorm' by maximum likelihood
      estimate Std. Error
meanlog 2.6841381 0.01819330
sdlog    0.5753227 0.01286443
```

```
Parametric bootstrap medians and 95% percentile CI
      Median      2.5%      97.5%
meanlog 2.6851781 2.6488932 2.7188573
sdlog    0.5750318 0.5493428 0.5988646
```

Distribution fitting exercise

In the workspace there are two datasets named 'fitting_task_1.csv' and 'fitting_task_2.csv'. Your task is to use the fitting process that we have just discussed to fit a distribution to each of these datasets and estimate the parameters.

There is also an R script call 'fitting_task_code.R' that contains the basic code needed for this task. You can choose to use this or write your own script.

More distributions with actuar

The `fitdistrplus` package contains a limited number of named (common) distributions that you can fit to your data. The **actuar** package contains more named distributions that can be used for distribution fitting in combination with the `fitdistrplus` package. To find out more about the `actuar` package see <https://www.rdocumentation.org/packages/actuar/versions/2.3-3>

A distribution fitting app template

Once you start using distributions you will find they are really useful! Distribution fitting however, is not an exact science and is a bit of an art. Repeatative testing is required to refine your distribution and its parameters for real world use.

Building a custom app front end for your code can make repetitive tasks easier.

As a starter template the folder 'app_template_files' contains an app which uses the code we have just looked at

Open the folder, run the 'app.R' file and try it out using the 'example_app_data.csv' file

Distribution fitting resources

- fitdistrplus documentation

cran.r-project.org/web/packages/fitdistrplus/fitdistrplus.pdf

- actuar documentation

<https://www.rdocumentation.org/packages/actuar/versions/2.3-3>

Thank you for listening
Any questions?
For more useful R training check out
<https://rforhealthcare.org/>