

# Assignment 8 for Statistical Computing and Empirical Methods: Continuous random variables, independence and laws of large numbers

Henry Reeve

## Introduction

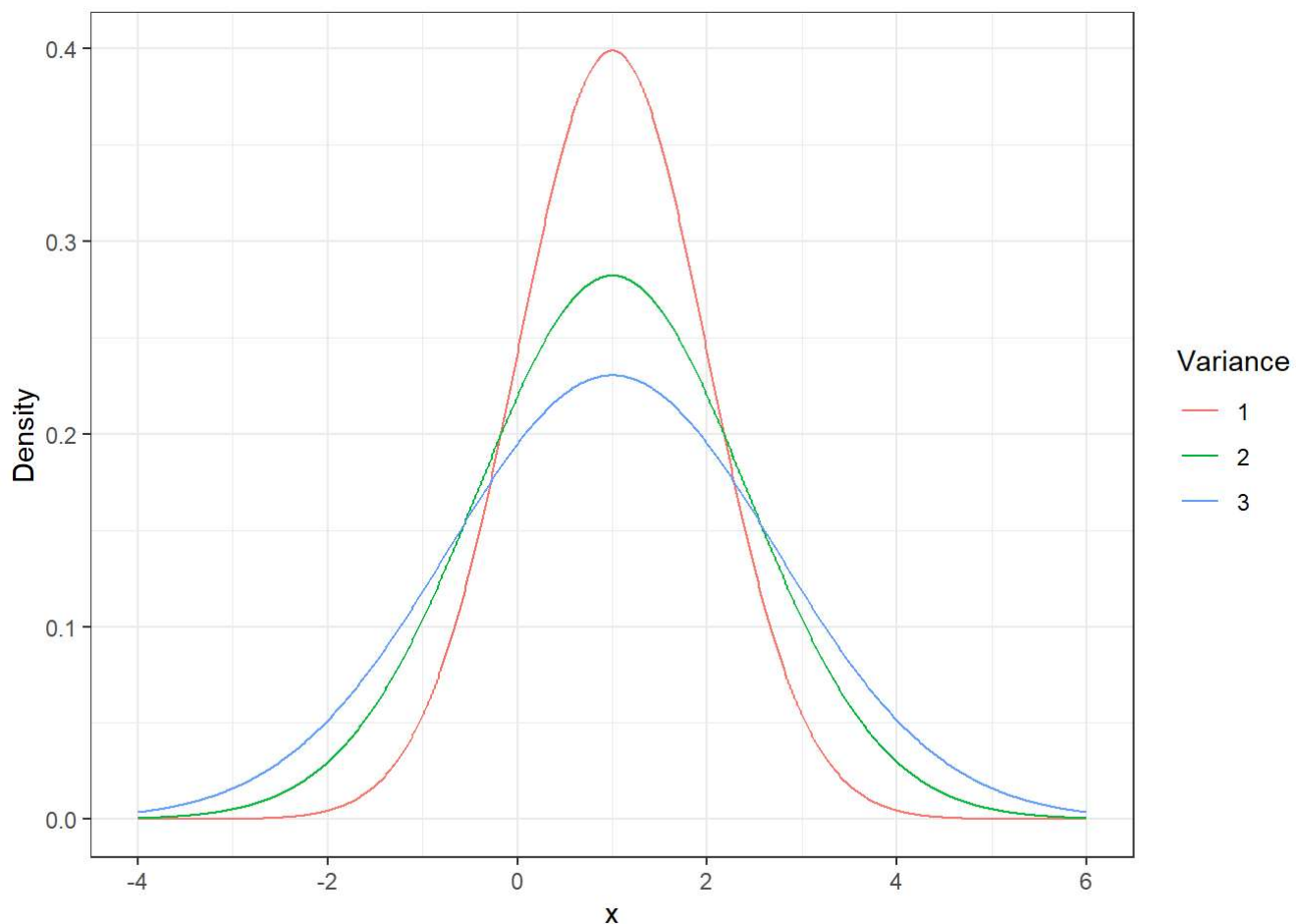
This document describes your eighth assignment for Statistical Computing and Empirical Methods (Unit EMATM0061) on the MSc in Data Science. Before starting the assignment it is recommended that you first watch video lecture 8 entitled “Continuous random variables, independence and laws of large numbers”.

## 1 The Gaussian distribution

Write out the probability density function of a Gaussian random variable with mean  $\mu$  and standard deviation  $\sigma > 0$ .

Use the help function to look up the following four functions: **dnorm()**, **pnorm()**, **qnorm()** and **rnorm()**.

Generate a plot which displays the probability density function for three Gaussian distributions  $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ ,  $X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$  and  $X_3 \sim \mathcal{N}(\mu_3, \sigma_3^2)$  with  $\mu_1 = \mu_2 = \mu_3 = 1$  and variances  $\sigma_1^2 = 1$ ,  $\sigma_2^2 = 2$  and  $\sigma_3^2 = 3$ . Your plot should look something like this:



Generate a similar plot for the cumulative distribution function for three Gaussian distributions  $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ ,  $X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$  and  $X_3 \sim \mathcal{N}(\mu_3, \sigma_3^2)$  with  $\mu_1 = \mu_2 = \mu_3 = 1$  and variances  $\sigma_1^2 = 1$ ,  $\sigma_2^2 = 2$  and  $\sigma_3^2 = 3$ .

Next generate a plot for the quantile function for the same three Gaussian distributions. Describe the relationship between the quantile function and the cumulative distribution function.

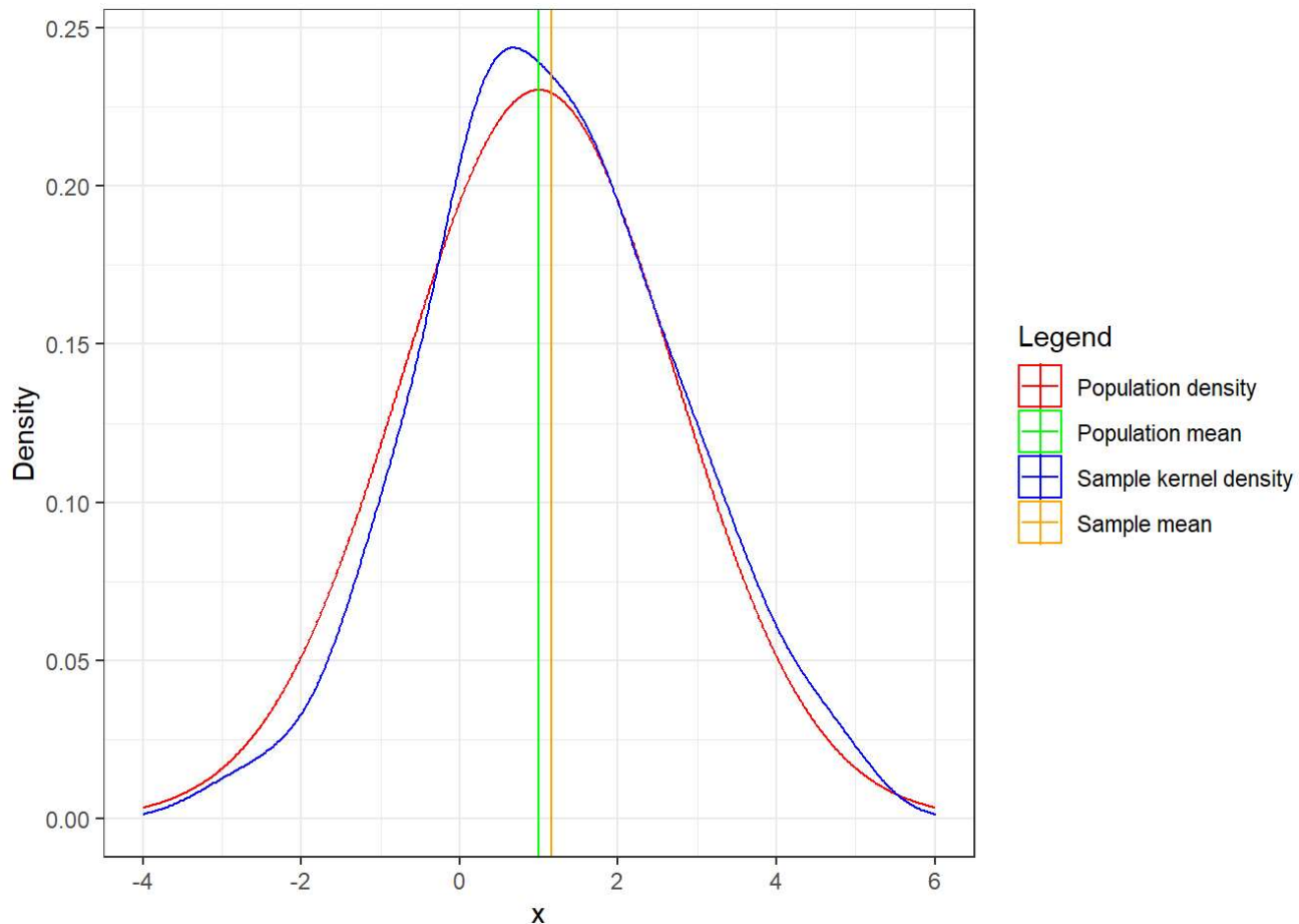
Now use **rnorm()** generate a random independent and identically distributed sequence  $Z_1, \dots, Z_n \sim \mathcal{N}(0, 1)$  so that each  $Z_i \sim \mathcal{N}(0, 1)$  has standard Gaussian distribution with  $n = 100$ . Make sure your code is reproducible by using the **set.seed()** function. Store your random sample in a vector called “standardGaussianSample”.

Without calling the **rnorm()** function again, use your existing sample stored in “standardGaussianSample” to generate a sample of size  $n$  of the form  $Y_1, \dots, Y_n \sim \mathcal{N}(1, 3)$  with mean  $\mu = 1$  and variance  $\sigma^2 = 3$ . Store your second sample in a vector called “mean1Var3GaussianSampleA”.

Reset the random seed to the same value as before using the **set.seed()** function and generate an i.i.d. sample of the form  $Y_1, \dots, Y_n \sim \mathcal{N}(1, 3)$  using the **rnorm()** function by setting the mean and standard deviation. Store this sample in a vector called

“mean1Var3GaussianSampleB”. Compare the vectors mean1Var3GaussianSampleA and mean1Var3GaussianSampleB.

Now generate a graph which includes both a kernel density plot for your sample mean1Var3GaussianSampleA and the population density (the probability density function) generated using **dnorm()**. You can also include two vertical lines which display both the population mean and the sample mean. Your plot should something like the following:



## 2 Bayes theorem

Suppose that there is a rare medical condition and an associated test for that condition. Let  $X$  and  $Y$  be a pair of binary random variables corresponding to a randomly selected member of the population. More precisely both  $X$  and  $Y$  have outcome space  $\{0, 1\}$ , with  $X$  defined by

$$X = \begin{cases} 1 & \text{if the person has the condition} \\ 0 & \text{if the person doesn't have the condition.} \end{cases}$$

Similarly,  $Y$  is defined by

$$Y = \begin{cases} 1 & \text{if the test is positive} \\ 0 & \text{if the test is negative.} \end{cases}$$

Suppose that the probability that a random person within the population has the condition is  $\mathbb{P}(X = 1) = 0.001$ . Suppose further that a person who has the condition will have a positive test result with probability  $\mathbb{P}(T = 1|X = 1) = 0.95$ . Similarly, a person who doesn't have the condition will have a negative test result with probability  $\mathbb{P}(T = 0|X = 0) = 0.95$ .

Compute the probability that a person who has a positive test result, actually has the condition. That is, compute  $\mathbb{P}(X = 1|T = 1)$ .

### 3 The exponential distribution

Let  $\lambda > 0$  be a positive real number. An exponential random variable  $X$  with parameter  $\lambda$  is a continuous random variable with density  $p_\lambda : \mathbb{R} \rightarrow (0, \infty)$  defined by

$$p_\lambda(x) := \begin{cases} 0 & \text{if } x < 0 \\ \lambda e^{-\lambda x} & \text{if } x \geq 0. \end{cases}$$

First prove that  $p_\lambda$  is a well-defined probability density function.

Compute the population mean and variance of an exponential random variable  $X$  with parameter  $\lambda$ .

Compute the cumulative distribution function and the quantile function for exponential random variables with parameter  $\lambda$ .

### 4 Transformations of continuous random variables

Suppose that  $X$  is a continuous real-valued random variable with density  $p_X : \mathbb{R} \rightarrow (0, \infty)$ . Suppose that  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  is a function and define  $Y = \varphi(X)$  by

$$\mathbb{P}(\varphi(X) \in A) = \mathbb{P}(X \in \varphi^{-1}(A)),$$

for  $A \subseteq \mathbb{R}$ . Is  $Y = \varphi(X)$  a discrete or continuous random variable?

Suppose that  $\varphi$  is a strictly increasing differentiable function with derivative  $\partial\varphi/\partial x > 0$ . Show that  $\varphi(X)$  is a continuous random variable and compute its density.

What is the probability density function of the random variable  $\varphi(X)$  when  $\varphi$  is a strictly decreasing differentiable function with  $\partial\varphi/\partial x < 0$ ?

Suppose now that  $\alpha \in \mathbb{R} \setminus \{0\}$ ,  $\beta \in \mathbb{R}$ . Define  $\varphi_{\alpha,\beta} : \mathbb{R} \rightarrow \mathbb{R}$  by  $\varphi_{\alpha,\beta}(z) = \alpha \cdot z + \beta$ . Compute the probability density function for the random variable  $Y_{\alpha,\beta} = \varphi_{\alpha,\beta}(X)$ .

When  $\varphi = \varphi_{\alpha,\beta}$  we have  $\mathbb{E}[\varphi_{\alpha,\beta}(X)] = \varphi_{\alpha,\beta}(\mathbb{E}[X])$ . Is this case for other functions  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  in place of  $\varphi_{\alpha,\beta}$ ?

As an optional extra look up Jensen's inequality.

## 5 The Binomial distribution and the central limit theorem

Two important discrete distributions are the Bernoulli distribution and the Binomial distribution. We say that a random variable  $X$  has Bernoulli distribution with parameter  $q \in [0, 1]$  if  $X$  has outcome space  $\{0, 1\}$  with  $\mathbb{P}(X = 1) = q$ . This is often abbreviated as  $X \sim \mathcal{B}(q)$ . For more details on the Bernoulli distribution we refer to the lecture.

Given  $n \in \mathbb{N}$  and  $q \in [0, 1]$ , we say that  $Z$  is Binomially distributed random variable with parameters  $n$  and  $q$  if  $Z = X_1 + \dots + X_k$  where  $X_i \sim \mathcal{B}(q)$  and  $X_1, \dots, X_k$  are independent and identically distributed. This is often abbreviated as  $Z \sim \text{BINOM}(n, q)$ .

Compute the expectation and variance of  $Z \sim \text{BINOM}(k, q)$ . You can use the following two useful facts:

1. Given any sequence of random variables  $W_1, \dots, W_k$  we have

$$\mathbb{E} \left[ \sum_{i=1}^k W_i \right] = \sum_{i=1}^k \mathbb{E} [W_i].$$

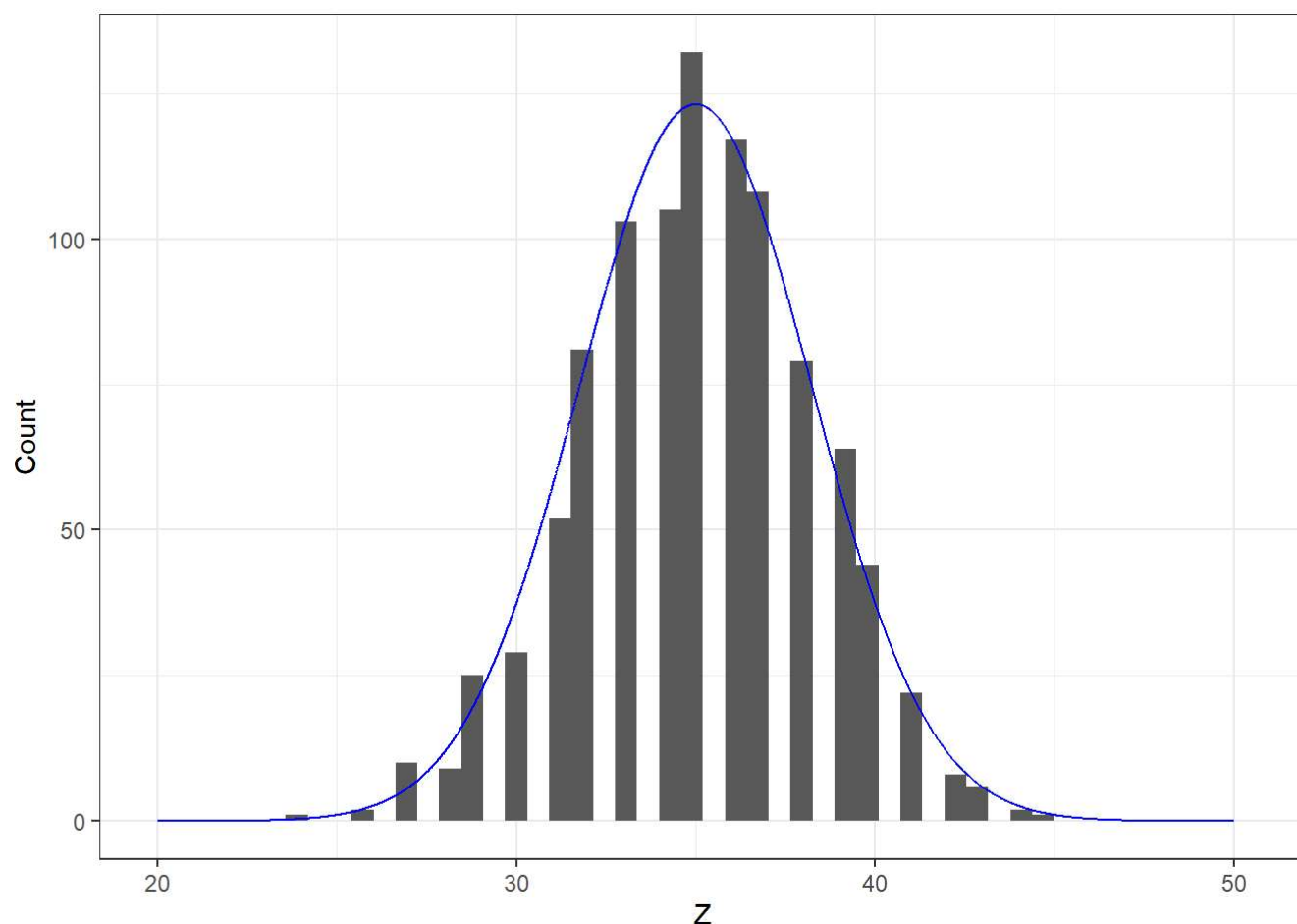
2. Given a sequence of **independent** random variables  $W_1, \dots, W_k$  we have

$$\text{Var} \left( \sum_{i=1}^k W_i \right) = \sum_{i=1}^k \text{Var} (W_i).$$

Is it always true that  $\text{Var} \left( \sum_{i=1}^k W_i \right) = \sum_{i=1}^k \text{Var} (W_i)$ , even if  $W_1, \dots, W_k$  are not independent?

Use the **rbinom()** function to generate a sample of size  $n$  consisting of independent Binomial random variables  $Z_1, \dots, Z_n \sim \text{BINOM}(n, q)$ , with  $n = 1000$ ,  $k = 50$  and  $q = 0.7$ .

Generate a histogram for your sample. In addition, plot a rescaled density for the normal distribution with mean  $\mu = k \cdot q = 700$  and variance  $\sigma^2 = k \cdot q \cdot (1 - q) = 210$ . You should rescale by  $n = 50000$ . Your plot should something like this:



As a challenging optional extra, try to explain this behavior using the central limit theorem?

## 6 Why do we rescale the median absolute deviation?

In our lecture on exploratory data analysis we introduced the Median Absolute Deviation. This can be computed within R using the **mad()** function. Use the **help()** function to investigate the arguments for this function. Note that there is an optional “constant” argument that defaults to 1.4826. In this question we will look at why this is the case.

The following code generates the sample standard deviation as a function of the sample size for some randomly generated Gaussian data.

```
total_sample_size<-5000 # set the total sample size

num_trials<-8 # set the number of trials

set.seed(123) # set the random seed

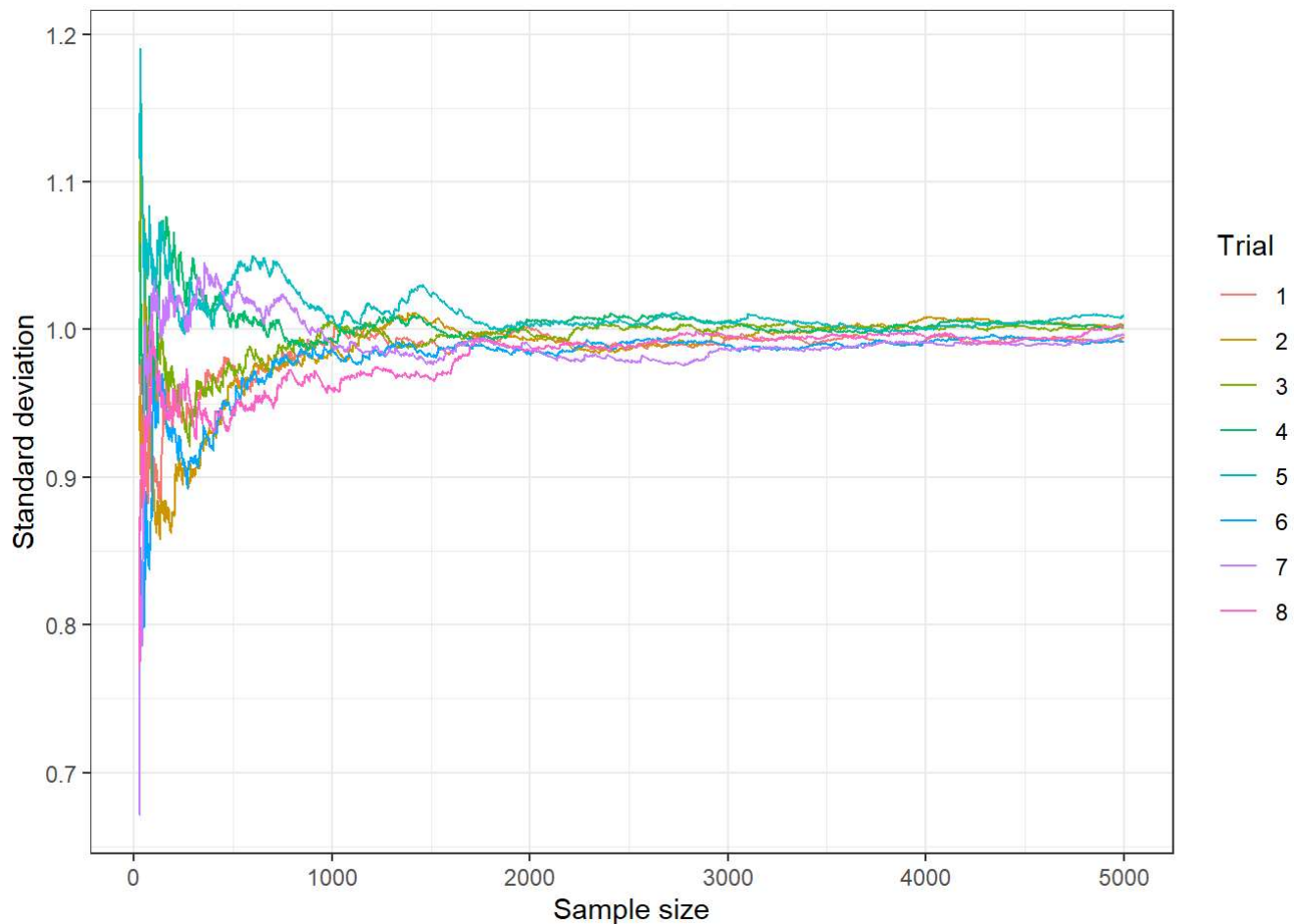
# generate a sample standard deviations as function of the sample size for some randomly
generated Gaussian data
gaussian_sample_sd_by_sample_size<-function(){

  return(
    data.frame(sample_size=seq(total_sample_size),X=rnorm(total_sample_size))%>% #generate
normal data
    mutate(sd=map_dbl(row_number(),~sd(X[1:.x]))) # compute sd of the initial segment
  )
}

df<-data.frame(trial=seq(num_trials))%>%

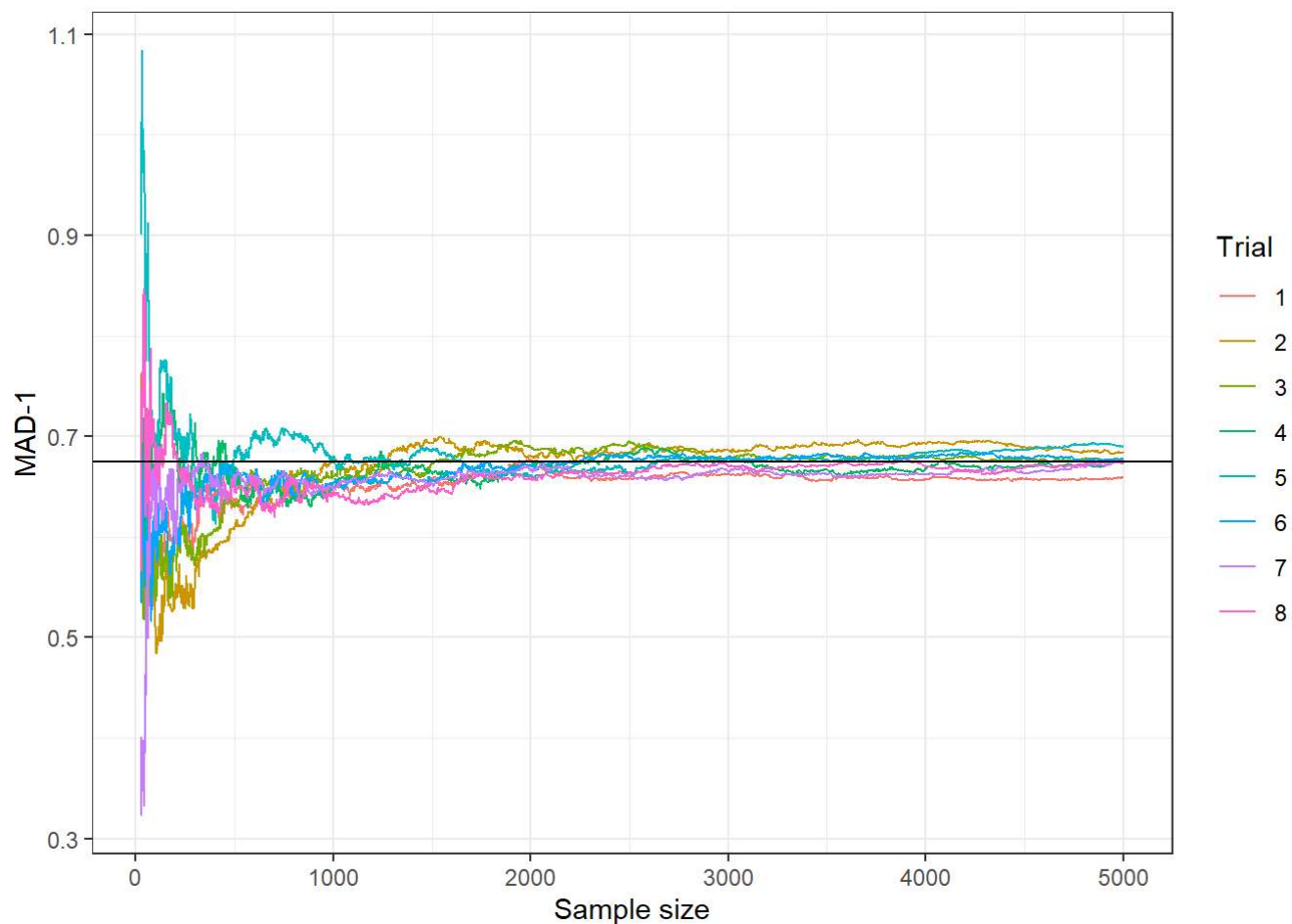
  mutate(data=map(trial, ~gaussian_sample_sd_by_sample_size()))%>% # apply simulation for each trial
  unnest(cols=data) # unnest over the different trials

ggplot(df)%>%
  filter(sample_size>25),aes(x=sample_size,y=sd,color=as.character(trial)))+
  geom_line()+theme_bw()+labs(color="Trial",x="Sample size", y="Standard deviation") # Plot results
```



By default **mad()** includes a scale factor of 1.4826. However, we can set the scale factor to 1 by applying `constant=1`. We shall denote this form of the median absolute deviation by “MAD-1”. Copy and modify the above code to generate a data frame which contains a column for the median absolute deviation with a scale factor of one (MAD-1). Plot your data and include a horizontal line at  $0.6744908 = 1/1.4826$ . Your plot should look something like this:





## 7 Chebyshev's law of large numbers

This is an entirely optional and much more advanced exercise. It is intended only for those who already have significant experience of probability. You can safely leave this out if you have insufficient time.

Firstly, give a proof of Chebyshev's law of large numbers.

Secondly, look up Hoeffding's inequality and explain the advantage of Hoeffding's inequality over the law of large numbers?