

Assignment 5 for Statistical Computing and Empirical Methods

Dr. Henry WJ Reeve

Teaching block 1 2021

Introduction

This document describes your fifth 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 lectures 11 and 12.

You are encouraged to discuss these questions with your colleagues.

Begin by creating an Rmarkdown document with html output. You are not expected to hand in this piece of work, but it is a good idea to get used to using Rmarkdown.

Several **optional** extra questions have also been included. These are marked by stars (“*”). It is recommended that you first complete all unstarred questions before proceeding through the starred questions. **Do not** be concerned if you do not have time to complete the starred questions!

1 Expectation and variance of a discrete random variable

Suppose that $\alpha, \beta \in [0, 1]$ with $\alpha + \beta \leq 1$ and let X be a discrete random variable with distribution supported on $\{0, 1, 5\}$. Suppose that $\mathbb{P}(X = 1) = \alpha$ and $\mathbb{P}(X = 5) = \beta$ and $\mathbb{P}(X \notin \{0, 1, 5\}) = 0$.

What is the probability mass function $p_X : \mathcal{S} \rightarrow [0, 1]$ for X ?

What is the expectation of X ?

What is the variance of X ?

2 Simulating data with the uniform distribution

We shall now use the uniform distribution to simulate data from the discrete random variable discussed in the previous question. A uniformly distributed random variable U is a continuous random variable with probability density function

$$p_U(x) = \begin{cases} 1 & \text{if } x \in [0, 1] \\ 0 & \text{otherwise.} \end{cases}$$

Show that for any region pair of numbers $a, b \in \mathbb{R}$ with $0 \leq a \leq b \leq 1$ we have $\mathbb{P}(U \in [a, b]) = b - a$.

Now let's return to the discrete random variable discussed in the previous question in which $\mathbb{P}(X = 1) = \alpha$ and $\mathbb{P}(X = 5) = \beta$ and $\mathbb{P}(X = 0) = 1 - \alpha - \beta$. First consider the case in which $\alpha = \beta = 0.25$. You can generate a sequence of i.i.d. copies X_1, \dots, X_n of X as follows:

```
set.seed(0)

n<-1000

sample_X<-data.frame(U=runif(n))%>%
  mutate(X=case_when(
    (0<=U)&(U<0.25)~1,
    (0.25<=U)&(U<0.5)~5,
    (0.5<=U)&(U<=1)~0))%>%
  pull(X)
```

Why does this `sample_X` correspond to a sequence of i.i.d. copies X_1, \dots, X_n of X where $\mathbb{P}(X = 1) = \alpha$ and $\mathbb{P}(X = 5) = \beta$ and $\mathbb{P}(X = 0) = 1 - \alpha - \beta$ with $\alpha = \beta = 0.25$?

Now create a function called `sample_X_015()` which takes as inputs α , β and n and outputs a sample X_1, \dots, X_n of independent copies of X where $\mathbb{P}(X = 1) = \alpha$ and $\mathbb{P}(X = 5) = \beta$ and $\mathbb{P}(X = 0) = 1 - \alpha - \beta$.

Next take $\alpha = 1/2$ and $\beta = 1/10$, and use your function `sample_X_015()` to create a sample of size $n = 10000$, of the form X_1, \dots, X_n consisting of independent copies X for each value of β . What is the sample average of X_1, \dots, X_n ? How does this compare with $\mathbb{E}(X)$? Use your understanding of the law of large numbers to explain this behaviour.

In addition, compute the sample variance of X_1, \dots, X_n and compare with $\text{Var}(X)$.

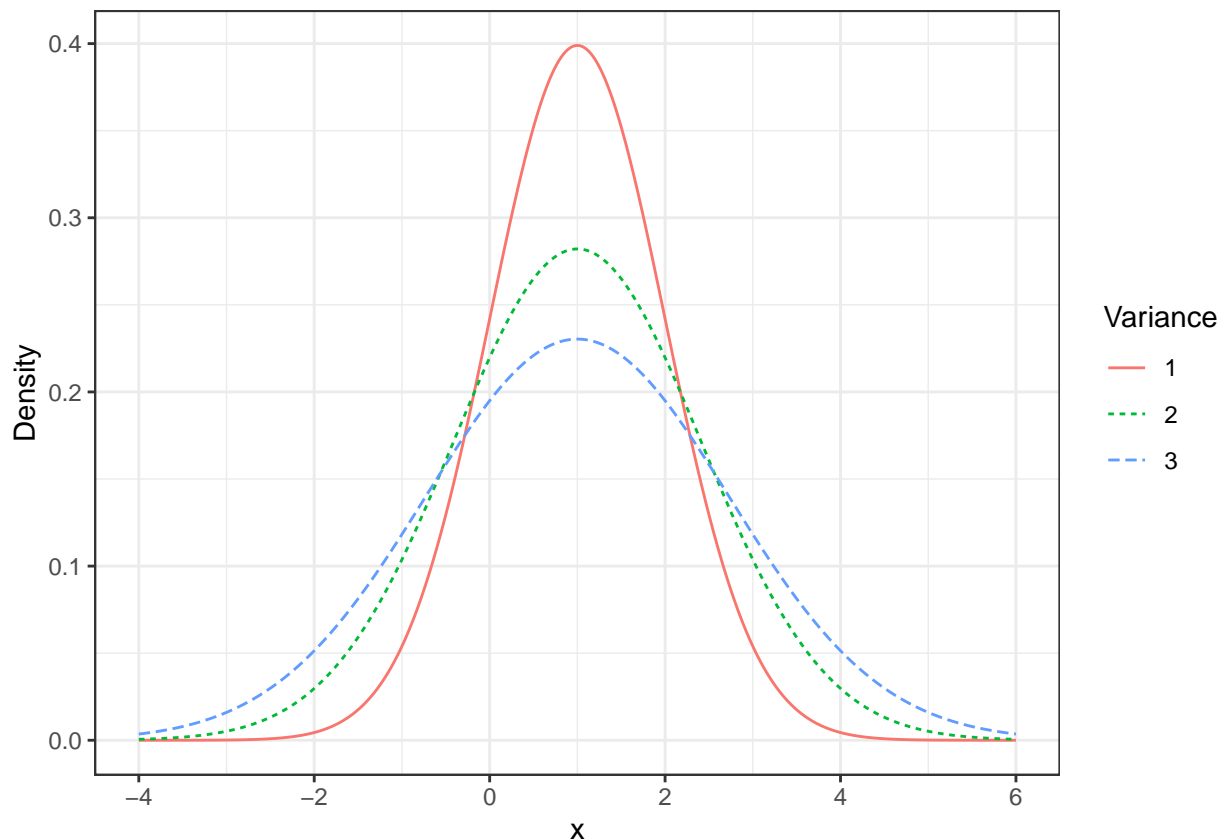
Now take $\alpha = 1/10$ and vary β in increments of 0.01 from 0 to 9/10, using your function `sample_X_015()` to create a sample of size $n = 100$, X_1, \dots, X_n consisting of independent copies X for each value of β . Create a plot of the sample averages as a function of β .

3 The Gaussian distribution

Write out the probability density function of a Gaussian random variable with mean μ 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 corresponding 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.

(*) Recall that for a random variable $X : \Omega \rightarrow \mathbb{R}$ is said to be Gaussian with expectation μ and variance σ^2 ($X \sim \mathcal{N}(\mu, \sigma^2)$) if for any $a, b \in \mathbb{R}$ we have

$$\mathbb{P}(a \leq X \leq b) = \int_a^b \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{z-\mu}{\sigma}\right)^2} dz.$$

Suppose $Z \sim \mathcal{N}(0, 1)$ is a Gaussian random variable. Take $\alpha, \beta \in \mathbb{R}$ and let $W : \Omega \rightarrow \mathbb{R}$ be the random variable given by $W = \alpha Z + \beta$. Apply a change of variables to show that W is a Gaussian random variable with expectation β and variance α^2 .

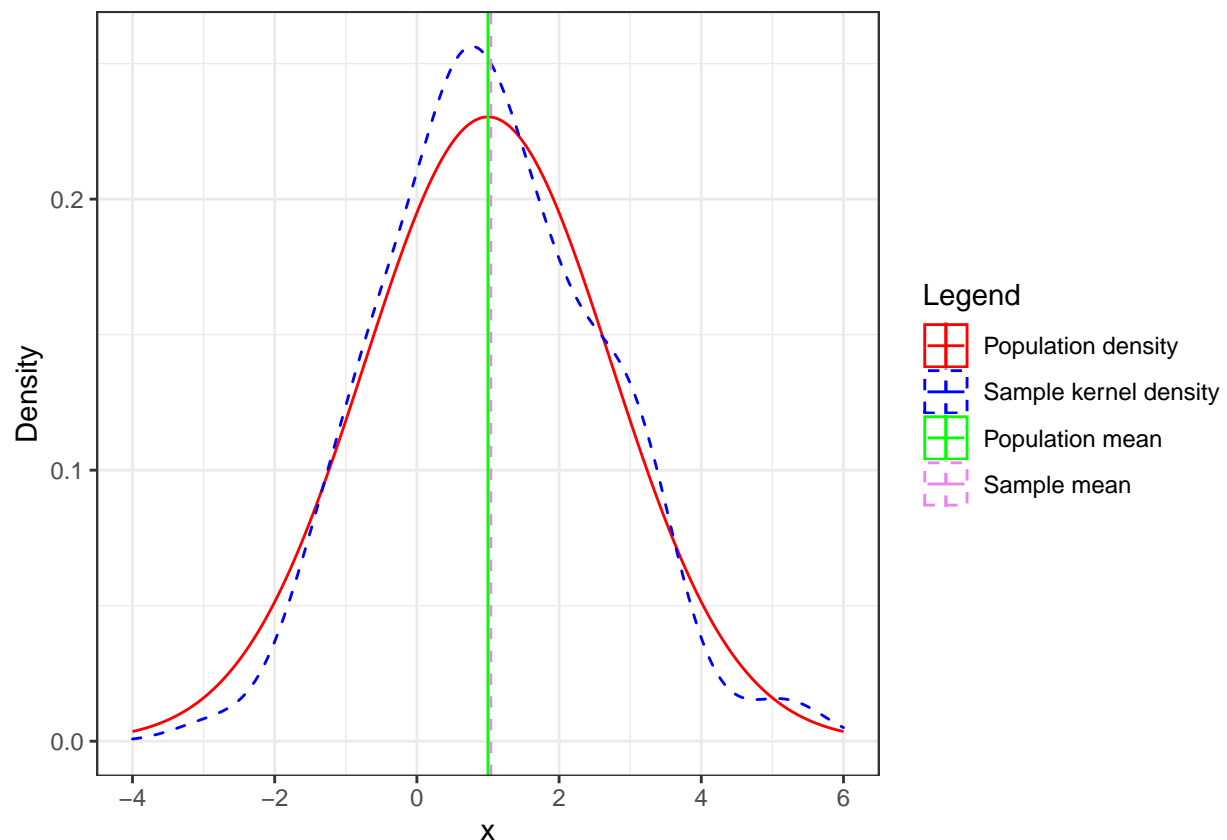
Note: If you are short on time you can simply take this fact as given and move onto the next part of the question.

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".

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 expectation $\mu = 1$ and population variance $\sigma^2 = 3$. Store your second sample in a vector called `mean1Var3GaussianSampleA`. The i -th observation in the sample `mean1Var3GaussianSampleA` should be of the form $Y_i = \alpha \cdot Z_i + \beta$, for appropriately chosen $\alpha, \beta \in \mathbb{R}$, where Z_i is the i -th observation in the sample `standardGaussianSample`.

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. 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. You may want to use the `geom_density()` and `geom_vline()` functions. Your plot should something like the following:



4 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 $p \in [0, 1]$ if $\mathbb{P}(X = 1) = p$ and $\mathbb{P}(X = 0) = 1 - p$. This is often abbreviated as $X \sim \mathcal{B}(p)$.

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

(Q) Compute the expectation and variance of $Z \sim \text{Binom}(n, p)$. You may want to make use of 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 **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?

The function `dbinom()` in R allows us to compute the probability mass function of a Binomial random variable $Z \sim \text{Binom}(n, p)$. By taking $x \in \{0, 1, \dots, n\}$ `size=n` and `prob=p` as arguments, the function `dbinom(x,size=n,prob=p)` will return the probability mass function $p_Z(x) = \mathbb{P}(Z = x)$ evaluated at x . You can run `?dbinom` in the R console to find out more.

Consider the case where $n = 50$ and $p = 7/10$. Use the `dbinom()` to generate a dataframe called `binom_df` with two columns - `x` and `pmf`. The first column contains the numbers $\{0, 1, \dots, 50\}$ inclusive. The second column gives the corresponding value of the probability mass function $p_Z(x) = \mathbb{P}(Z = x)$ with $Z \sim \text{Binom}(50, 7/10)$. Use the `head()` function to observe the first 3 rows as your data frame. The result should look as follows:

```
##      x      pmf
## 1  0 7.178980e-27
## 2  1 8.375477e-25
## 3  2 4.787981e-23
```

The function `dnorm()` in R allows us to compute the probability density function of a Gaussian random variable $W \sim \mathcal{N}(\mu, \sigma^2)$ with expectation μ and variance σ^2 . By taking $x \in \mathbb{R}$, `mean=mu` and `sd=sigma` as arguments, the function `dnorm(x,mean=mu,sd=sigma)` will return the probability density function $f_W(x)$ for $W \sim \mathcal{N}(\mu, \sigma^2)$, evaluated at x . You can run `?rnorm` in the R console to find out more.

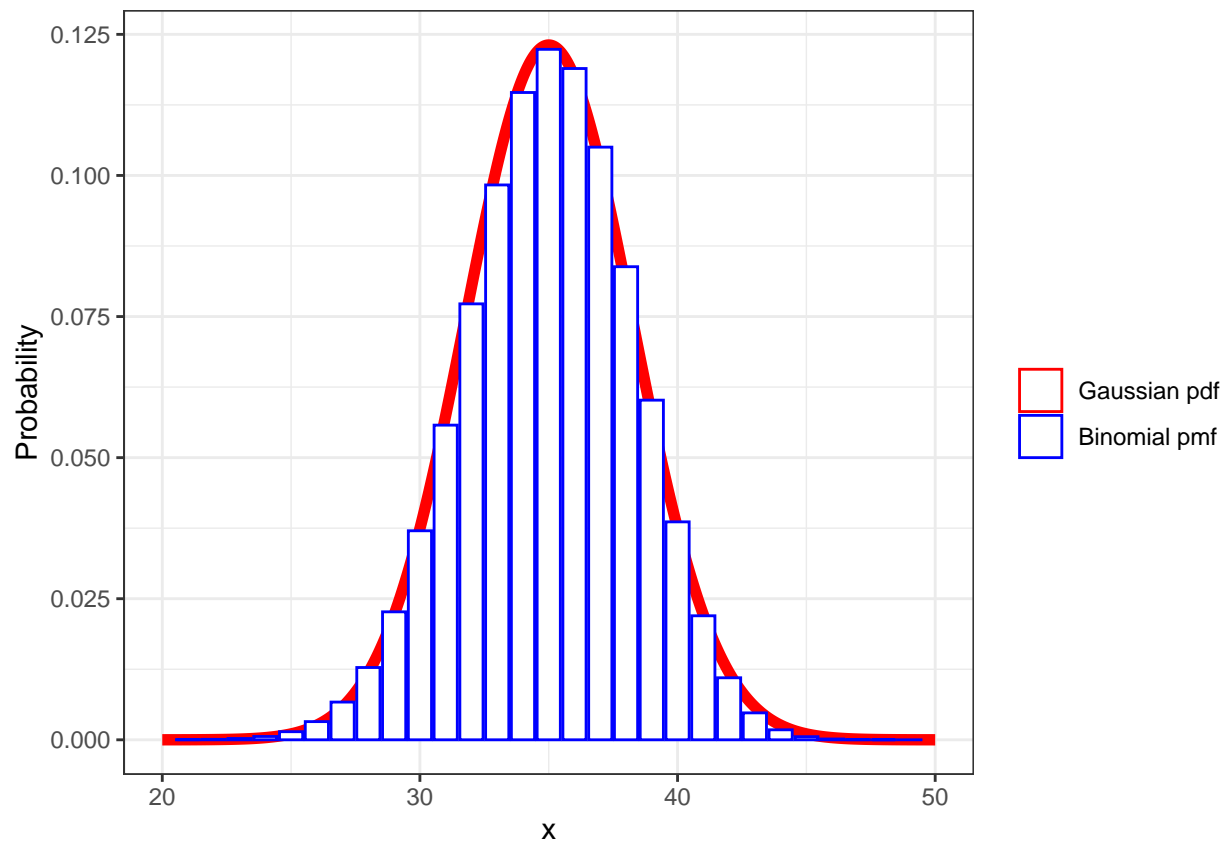
We shall consider a case where $\mu = 50 \cdot 0.7$ and $\sigma = \sqrt{50 \cdot 0.7 \cdot (1 - 0.7)}$. Use the `rnorm()` to generate a dataframe called `norm_df` with two columns - `x` and `pdf`. The first column contains the numbers $\{0, 0.01, 0.02, 0.03, \dots, 49.99, 50\}$. The second column gives the corresponding value of the probability density function $f_W(x)$ with $W \sim \mathcal{N}(5, 7/10)$. Use the `head()` function to observe the first 3 rows as your data frame. Your result should look as follows:

```
##      x      pdf
## 1 0.00 5.707825e-27
## 2 0.01 5.901264e-27
## 3 0.02 6.101201e-27
```

Next, use the following code to create a plot which compares the probability density for your Gaussian distribution $W \sim \mathcal{N}(\mu, \sigma^2)$ where $\mu = n \cdot p$ and $\sigma = \sqrt{n \cdot p(1-p)}$ and the probability mass function for your Binomial distribution $Z \sim \text{Binom}(n, p)$.

```
colors<-c("Gaussian pdf"="red", "Binomial pmf"="blue")
fill<-c("Gaussian pdf"="white", "Binomial pmf"="white")

ggplot()+labs(x="x",y="Probability")+theme_bw()+
  geom_line(data=gaussian_df,
    aes(x,y=pdf,color="Gaussian pdf"),size=2)+
  # create plot of Gaussian density
  geom_col(data=binom_df,
    aes(x=x,y=pmf,color="Binomial pmf",fill="Binomial pmf"))+
  scale_color_manual(name = "", values=colors)+
  scale_fill_manual(name = "", values=fill)+
  xlim(c(20,50))
```



(*) Now use the central limit theorem to explain the results you observe.

5 Exponential distribution

Let $\lambda > 0$ be a positive real number. An exponential random variable X with rate parameter λ 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_λ is a well-defined probability density function.

Compute the population mean and variance of an exponential random variable X with parameter λ .

Compute the cumulative distribution function and the quantile function for exponential random variables with parameter λ .

Now implement a function called `my_cdf_exp()`. The function `my_cdf_exp()` should take as input two numbers $x \in \mathbb{R}$ and $\lambda > 0$ and output the value of the cumulative distribution function $F_X(x)$ where X is an exponential random variable with rate parameter λ .

Check your function `my_cdf_exp()` gives rise to the following output:

```
lambda<-1/2  
  
map_dbl(.x=seq(-1,4),.f=~my_cdf_exp(x=.x,lambda=lambda))
```

```
## [1] 0.0000000 0.0000000 0.3934693 0.6321206 0.7768698 0.8646647
```

Type `?pexp` into your R console to learn more about R's inbuilt cumulative distribution function for the exponential distribution. We can now confirm that our when $\lambda = 1/2$ as follows:

```
test_inputs<-seq(-1,10,0.1)  
my_cdf_output<-map_dbl(.x=test_inputs,.f=~my_cdf_exp(x=.x,lambda=lambda))  
inbuilt_cdf_output<-map_dbl(.x=test_inputs,.f=~pexp(q=.x,rate=lambda))  
  
all.equal(my_cdf_output,inbuilt_cdf_output)
```

```
## [1] TRUE
```

Next implement a function called `my_quantile_exp()`. The function `my_quantile_exp()` should take as input two arguments $p \in [0, 1]$ and $\lambda > 0$ and output the value of the quantile function $F_X^{-1}(p)$ where X is an exponential random variable with rate parameter λ .

Once you have implemented your function compare with R's inbuilt `qexp` function using the same procedure as we used above for the cumulative distribution function for inputs $\lambda = 1/2$ and $p \in \{0.01, 0.02, 0.03, \dots, 0.99\}$. Note that you don't need to consider inputs $p \leq 0$ or $p \geq 1$ here.

6 Poisson distribution

Many discrete random variables have distributions supported on a finite set (eg. Bernoulli, Binomial). Poisson random variables are a family of discrete random variables with distributions supported on $\mathbb{N}_0 := \{0, 1, 2, 3, \dots\}$. Poisson random variables are frequently used to model the number of events which occur at a constant rate in situations where the occurrence of individual events are independent. For example, we might use the Poisson distribution to model the number of mutations of a given strand of DNA per time unit, or the number of customers who arrive at store over the course of a day. Hence, like Binomial random variable, Poisson random variables can be used to model count data. The key difference is that Poisson random variables apply more readily to situations where there is no natural upper bound on the total count.

Take $\lambda > 0$. The Poisson random variable X with parameter λ has probability mass function $p_\lambda : \mathbb{R} \rightarrow (0, \infty)$ defined by

$$p_\lambda(x) = \begin{cases} \frac{\lambda^x e^{-\lambda}}{x!} & \text{for } x \in \mathbb{N}_0 \\ 0 & \text{for } x \notin \mathbb{N}_0. \end{cases}$$

Show that p_λ is a well-defined probability mass function. More precisely:

1. $p_\lambda(x) \geq 0$ for all $x \in \mathbb{N}_0$
2. $\sum_{x \in \mathbb{R}} p_\lambda(x) = 1$.

Compute both the expectation and the variance of a Poisson random variable X with probability mass function p_λ .

()** Both Binomial and Poisson random variables can be used to model count data. Whilst Binomial random variables $Z \sim \text{Binom}(n, p)$ apply to situations where there is an upper bound n on the number of successes, Poisson random variables apply to situations where there is no natural upper bound on the total count. In fact the Poisson random variable X can be viewed as an approximation to Binomial random variable with very large n and $np \approx \lambda$. This approximation is often used for computational reasons in situations where we want to model a Binomial random variable with a very large n and a small value of p .

As an optional extra, show the following: Suppose we fix $\lambda \in \mathbb{R}$ and a Poisson random variable X with expectation λ , and take probabilities $p_n = \frac{\lambda}{n}$ for each $n \in \mathbb{N}$, and Binomial random variables $Z \sim \text{Binom}(n, p_n)$. Then for each $k \in \mathbb{N}_0 = \{0, 1, 2, 3, \dots\}$ we have

$$\lim_{n \rightarrow \infty} \mathbb{P}(Z_n = k) = \mathbb{P}(X = \lambda) = \frac{\lambda^k e^{-\lambda}}{k!}.$$

You may want to use the fact that for any real number $t \in \mathbb{R}$ we have $\lim_{n \rightarrow \infty} \left(1 + \frac{t}{n}\right)^n = e^t$.

7 (**) The law of large numbers and Hoeffding's inequality

Prove the following version of the weak law of large numbers.

Theorem (A law of large numbers). *Let $X : \Omega \rightarrow \mathbb{R}$ be a random variable with a well-defined expectation $\mu := \mathbb{E}(X)$ and variance $\sigma^2 := \text{Var}(X)$. Let $X_1, \dots, X_n : \Omega \rightarrow \mathbb{R}$ be a sequence of independent copies of X . Then for all $\epsilon > 0$,*

$$\lim_{n \rightarrow \infty} \mathbb{P} \left(\left| \frac{1}{n} \sum_{i=1}^n X_i - \mu \right| \geq \epsilon \right) = 0.$$

You may want to begin by looking up Chebyshev's inequality.

Now investigate Hoeffding's for sample averages of bounded random variables. How does this compare to the law of large numbers?