

STAT 33B Homework 1

Yuanrui Zhu (3034615728)

Sep 10, 2020

This homework is due **Sep 10, 2020** by 11:59pm PT.

Homeworks are graded for correctness.

As you work, write your answers in this notebook. Answer questions with complete sentences, and put code in code chunks. You can make as many new code chunks as you like.

Please do not delete the exercises already in this notebook, because it may interfere with our grading tools.

You need to submit your work in two places:

- Submit this Rmd file with your edits on bCourses.
- Knit and submit the generated PDF file on Gradescope.

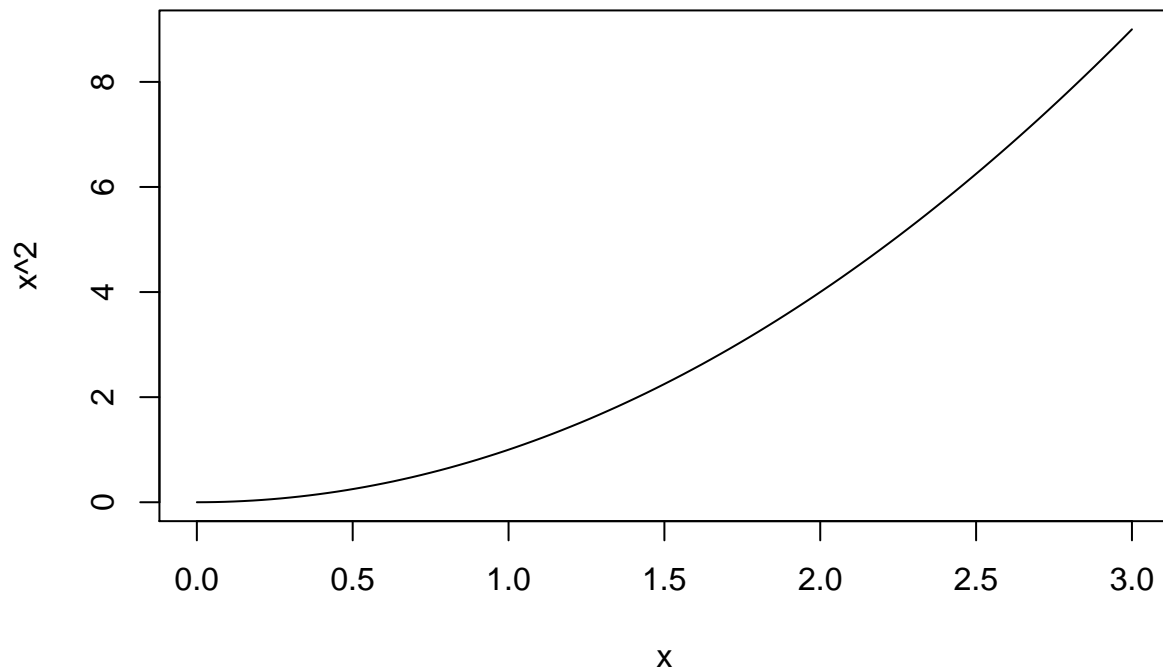
Exercise 1

Note: In R Markdown formatting, a dollar sign $\$$ marks the beginning of a (LaTeX) mathematical expression. For example, to format a linear equation nicely, you can write $ax + b = 0$. You don't need to know LaTeX for this class, so if you don't understand something written inside $\$$, knit the PDF and read that instead.

The `curve()` function plots a curve based on an expression. The basic syntax is `curve(expr, from, to)` where `from` and `to` are the limits of the x-axis.

For example, to plot x^2 between 0 and 3:

```
curve(x^2, 0, 3)
```



The `curve()` function makes it possible to use R as a graphing calculator.

For each expression below:

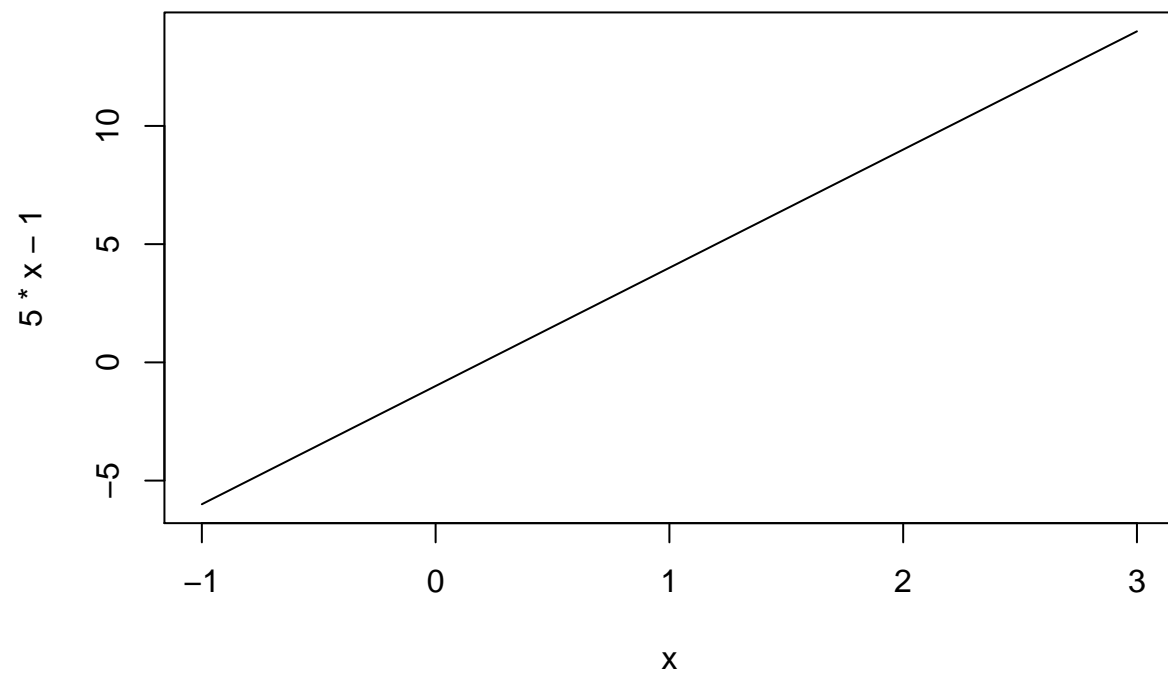
1. Plot the expression so that all points where it is 0 are visible. Experiment to find appropriate limits for the x-axis.
2. At which point(s) is the expression zero? Estimate (within 1 unit) based on your plot rather than computing these mathematically.

Here are the expressions:

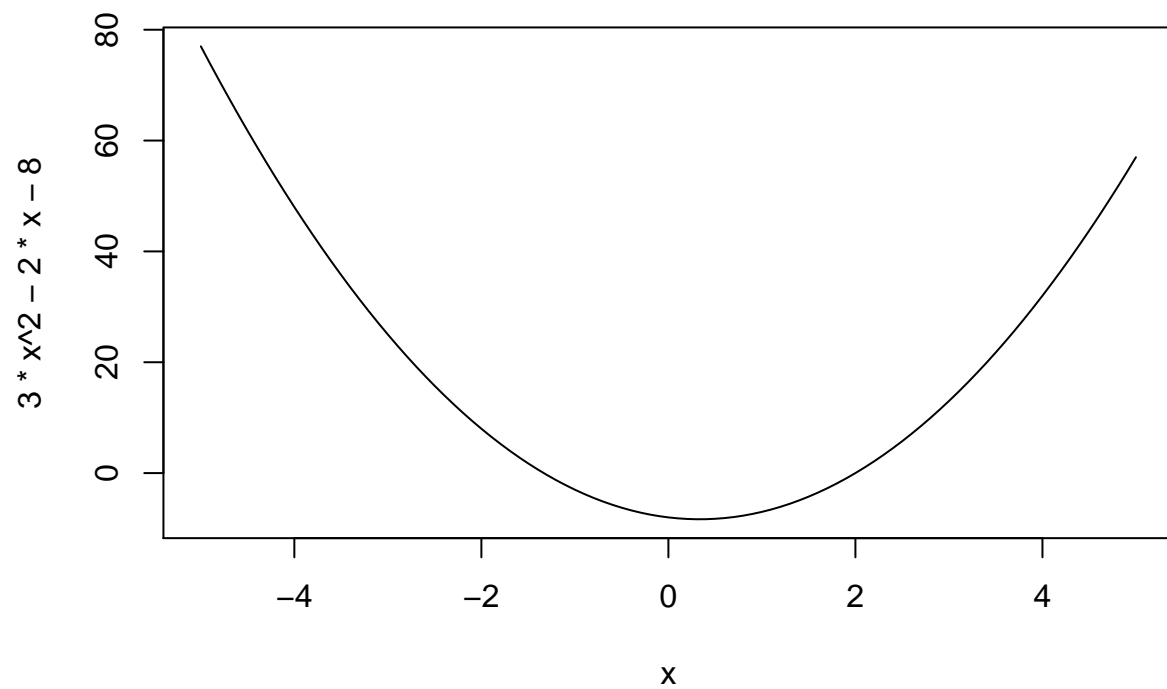
1. $5x - 1$
2. $3x^2 - 2x - 8$
3. $\sqrt{3}x - 4$

YOUR ANSWER GOES HERE:

```
curve(5*x-1, -1, 3) #the expression is zero around .2
```

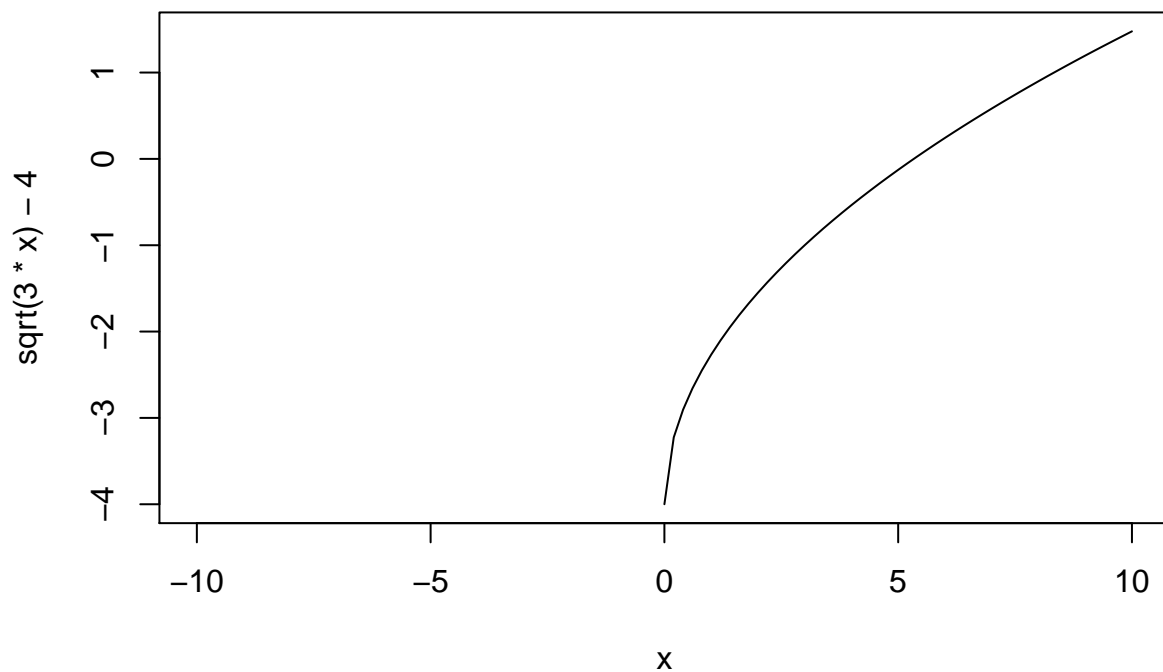


```
curve(3*x^2 - 2*x - 8, -5, 5) # the expression is zero around -1 and 2
```



```
curve(sqrt(3*x) - 4, -10, 10) # the expression is zero around 5
```

```
## Warning in sqrt(3 * x): NaNs produced
```



Exercise 2

A **discrete probability distribution** is a table of mutually exclusive outcomes and their associated probabilities.

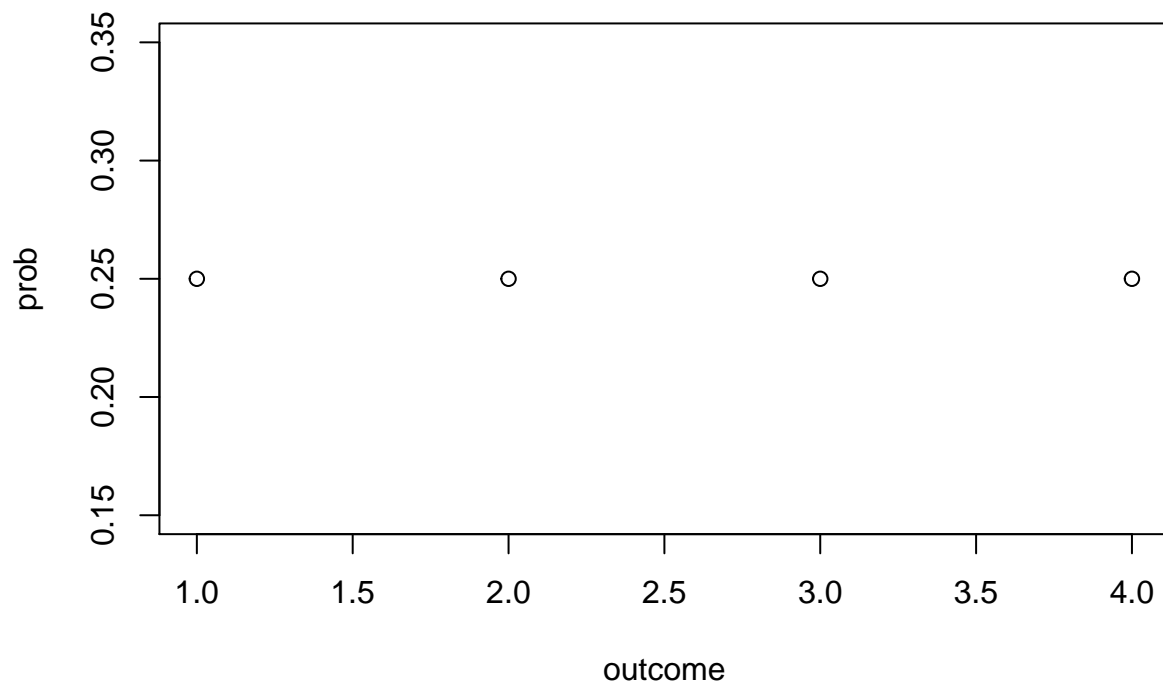
For example, the **discrete uniform distribution** assigns equal probability to each outcome. So the discrete uniform distribution on the integers 1-4 is:

Outcome	Probability
1	0.25
2	0.25
3	0.25
4	0.25

A fair coin toss is another instance of a discrete uniform distribution, where each of the two outcomes has probability 0.5.

You can plot the distribution above with the code:

```
outcome = seq(1, 4)
prob = rep(0.25, 4)
plot(outcome, prob)
```



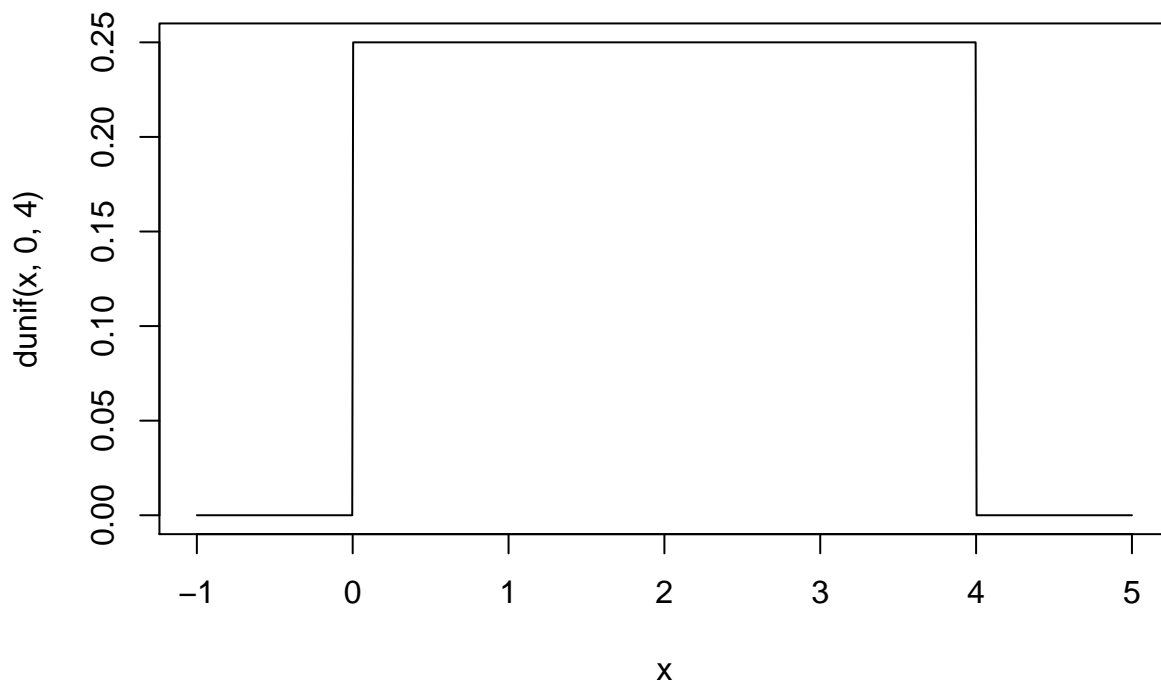
Because the probabilities are all equal, a plot of a uniform distribution is always flat, regardless of the number of outcomes.

The idea of a uniform distribution can also be extended to continuous values. For example, suppose we want to select a decimal value between 0 and 4 (inclusive), with equal probability for any value in the interval. This is a **continuous uniform distribution**.

As with any uniform distribution, a plot of a continuous uniform distribution is flat. However, instead of discrete points, the distribution is a curve along an interval.

You can plot the continuous uniform distribution from 0 to 4 with the code:

```
curve(dunif(x, 0, 4), -1, 5, n = 1000)
```



A plot of a continuous distribution is called a **density plot**.

In a density plot, the y-axis no longer represents probability. Instead, the probability of points in any subinterval is the total area between the curve and the line $y = 0$ over the subinterval. For instance, in the distribution above, the probability of the outcome being a point between 1 and 2 (inclusive) is 0.25.

R provides functions to compute the density of various distributions. The code above uses the `dunif()` function, the **density** plot function for the **uniform** distribution. The function has parameters to control the interval for the distribution.

R also provides functions to sample values randomly from various distributions. For example, the `runif()` function samples **randomly** from a continuous **uniform** distribution. You can sample 10 values from the continuous uniform distribution from -1 to 1 with the code:

```
runif(10, -1, 1)
```

```
## [1] 0.1675466 0.2652415 0.3003270 -0.3940006 0.9811358 -0.5559056
## [7] 0.5439855 0.2751140 0.1171430 0.8984252
```

Since these values are randomly sampled, they will be different each time you run the code.

The **Gaussian distribution** or normal distribution is another continuous probability distribution. The distribution plays an important role in mathematics, statistics, and the natural sciences. The density plot for this distribution is often described as a “bell-shaped curve”. The Gaussian distribution is typically parameterized by its mean (center) and standard deviation (how spread out it is), rather than the interval it covers.

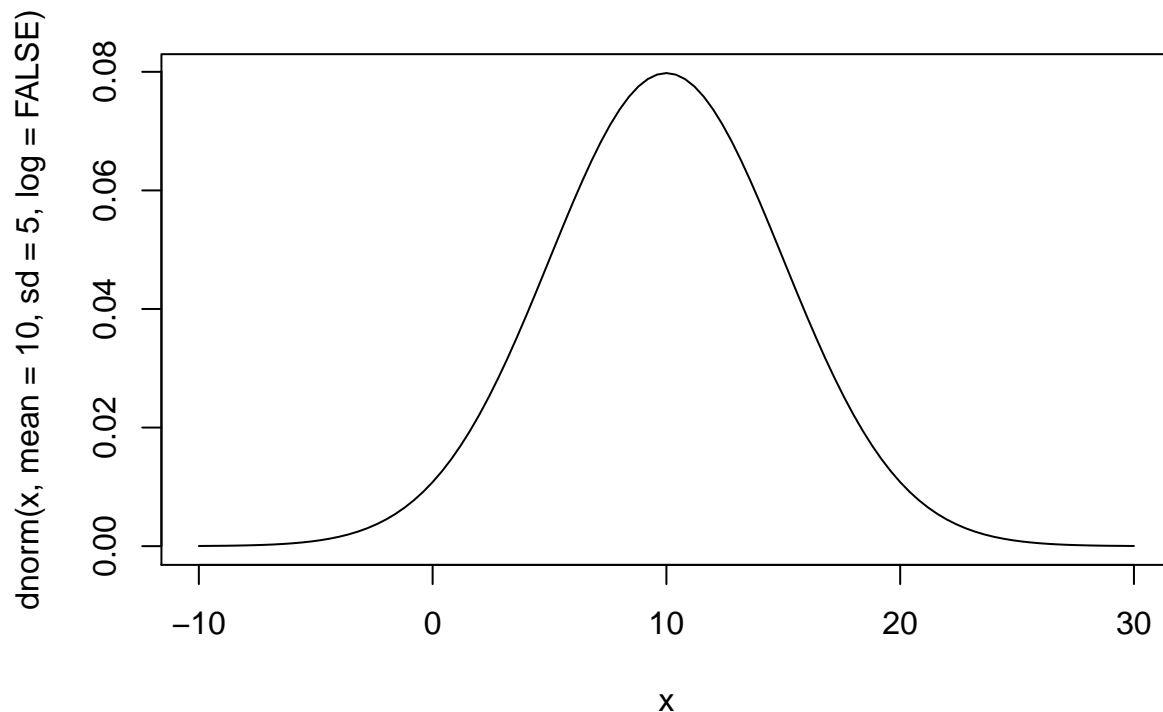
R provides functions `dnorm()` and `rnorm()` for the Gaussian distribution that are analogous to the `dunif()` and `runif()` functions described above.

1. Use the `curve()` function to make a density plot of the Gaussian distribution with mean 10 and standard deviation 5. Set the x-axis to span from -10 to 30.
2. Skim the documentation for `rnorm()`. Then create a variable called `samp` that contains 10 random values from a normal distribution with mean 10 and standard deviation 5.
3. Compute the mean and standard deviation of `samp`. Comment briefly on how much these differ from the true mean 10 and true standard deviation 5.

What happens if you run your code a second time? Are the sample mean and standard deviation the same? Are the sampled points the same?

YOUR ANSWER GOES HERE:

```
curve(dnorm(x, mean = 10, sd = 5, log = FALSE), -10, 30)
```



```
samp = rnorm(10, mean = 10, sd = 5)
abs(mean(samp) - 10)
```

```
## [1] 0.7364551
```

```
abs(sd(samp) - 5)
```

```
## [1] 0.948679
```


#After I run the code the second time, I discovered that the sample mean and standard deviation change

Exercise 3

R's functions for random sampling use a pseudo-random number generator (PRNG). The numbers produced by a PRNG are not truly random, but satisfy conditions that make them close enough to random for most scientific computing tasks.

An advantage of PRNGs is that they are deterministic. Given the same initial parameter, called a **seed**, a PRNG will generate the same sequence of “random” numbers every time. When you start R, the seed is automatically assigned value from your computer that varies (e.g., the system time in milliseconds). This ensures that you get a different sequence of random numbers in each R session.

You can the PRNG seed with the `set.seed()` function. Setting the seed makes your results reproducible. That is, other people can run your code and get the same results even if the code generates random numbers.

Generally, you should only set the seed once at the beginning of an R script or notebook, because it affects all subsequent calls that use the PRNG.

1. Set the seed to 93.
2. Create a variable called **samp** that contains 1000 random values from a normal distribution with mean 10 and standard deviation 5.
3. Compute the mean and standard deviation of **samp**.

What happens if you run your code a second time? Are the sample mean and standard deviation the same? Are the sampled points the same? Explain how step 1 affects this.

4. Compute the number of values in **samp** that are greater than 8. *Hint: there is a fast, simple way to count TRUE values using implicit coercion.*

YOUR ANSWER GOES HERE:

```
set.seed(93)
samp = rnorm(1000, mean = 10, sd = 5)
mean(samp)
```

```
## [1] 10.1505
```

```
sd(samp) # they do not change, as all data are fixed
```

```
## [1] 5.127733
```

```
sum(samp > 8) # there are 668
```

```
## [1] 668
```

Exercise 4

A **Monte Carlo algorithm** is an algorithm that uses random number generation to approximate a non-random result.

One application of Monte Carlo algorithms is to approximate the area of a shape. In this exercise, you'll use a Monte Carlo algorithm to approximate the area of a circle, and thereby approximate pi.

Consider the circle centered at 0 with radius 1. Points on the circle satisfy $x^2 + y^2 = 1$. This is the “target shape” for which you'll approximate the area. In practice, the target shape is usually one whose area is difficult to compute by hand.

In order to approximate the area of the target shape, start by finding a square or rectangle that completely encloses the target shape. For instance, we can use the square with corners (-1, -1) and (1, 1) to enclose the circle.

Next, randomly sample x coordinates and y coordinates uniformly and independently over the square.

Now the proportion of sampled points inside the circle corresponds to the ratio between the area of the circle and the area of the square. Since you know the area of the square, you can use this relationship to compute the area of the circle.

1. Use R and the described algorithm to approximate the value of pi. Use a variable called `n` to control the sample size.

Do not use if-statements or loops in your solution; they are not necessary. Try to take advantage of vectorization and implicit coercion. A typical solution will be 2-6 lines.

2. Test your code for different values of `n` from 10 to 1,000,000. How does changing `n` affect the accuracy of the approximation? How does it affect the time it takes to compute the approximation?

Hint 1: R supports exponent notation, so you can write 1,000,000 in R as `1e6`.

Hint 2: The `microbenchmark` package provides a simple function called `microbenchmark()` for timing expressions. You can use it to time an entire block of code by surrounding the code with curly braces `{ }`.

YOUR ANSWER GOES HERE:

```
simulation <- function(n) {  
  x_cor = runif(n, -1, 1)  
  y_cor = runif(n, -1, 1)  
  return ( 4 * sum((x_cor ** 2 + y_cor ** 2) < 1) / n )  
}
```

```
simulation(10)
```

```
## [1] 2.8
```

```
simulation(100)
```

```
## [1] 3.16
```

```
simulation(1000)
```

```
## [1] 3.152
```

```
simulation(10000)
```

```
## [1] 3.1504
```

```
simulation(100000)
```

```
## [1] 3.13224
```

```
simulation(1000000)
```

```
## [1] 3.142272
```