

# MT5763\_2\_220021614

Nico Herrig

2022-10-15

## Problem 1

Description: Consider the following independent random variables:

$X \sim N(\mu = 4, \sigma^2 = 10)$

$Y \sim U(a = 2, b = 8)$ .

Compute  $Pr(X > Y)$

Use bootstrapping to derive the sampling distribution for your estimate of  $Pr(X > Y)$

Show how the sample variance of this sampling distribution changes as a function of the number of Monte Carlo simulations.

---

For the underlying problem, a sample containing 100,000 random deviates from  $X \sim N(\mu = 4, \sigma^2 = 10)$  and  $Y \sim U(a = 2, b = 8)$  is used. To simulate “real-world conditions”, the solution is obtained from only the below given vectors for X and Y.

```
probability_calculator <- function(n) {  
  
  X <- rnorm(n, mean = 4, sd = sqrt(10))  
  Y <- runif(n, min = 2, max = 8)  
  
  # Calculating Pr(X>Y)  
  Pr_hat <- sum(X > Y) / n  
  
  output <- list(X = X,  
                 Y = Y,  
                 Pr_hat = Pr_hat)  
  
  return(output)  
}  
  
results <- probability_calculator(100000)  
  
print(results[3])
```

```
## $Pr_hat  
## [1] 0.39243
```

Calculating  $\widehat{Pr}(X > Y)$  from the initial sample without any further methods, we derive a value of 0.39243 .

To derive the distribution of  $\widehat{Pr}(X > Y)$ , we use a non-parametric bootstrap. One benefit of this technique is that it does not rely on the assumption of normally distributed data. The following chunk of code defines a function for such bootstrap.

```
# Function for bootstrap procedure, using parallel computing technique for
# speeding up the computation.
bootstrap_multicore_problem1 <- function(n_straps, vec1 = X, vec2 = Y) {
  prob_vector <- unlist(mclapply(1:n_straps, function(n = n, vec1 = X, vec2 = Y) {
    # sampling both vectors X & Y
    resX <- vec1[sample(1 : length(vec1), length(vec1), replace = TRUE)]
    resY <- vec2[sample(1 : length(vec2), length(vec2), replace = TRUE)]

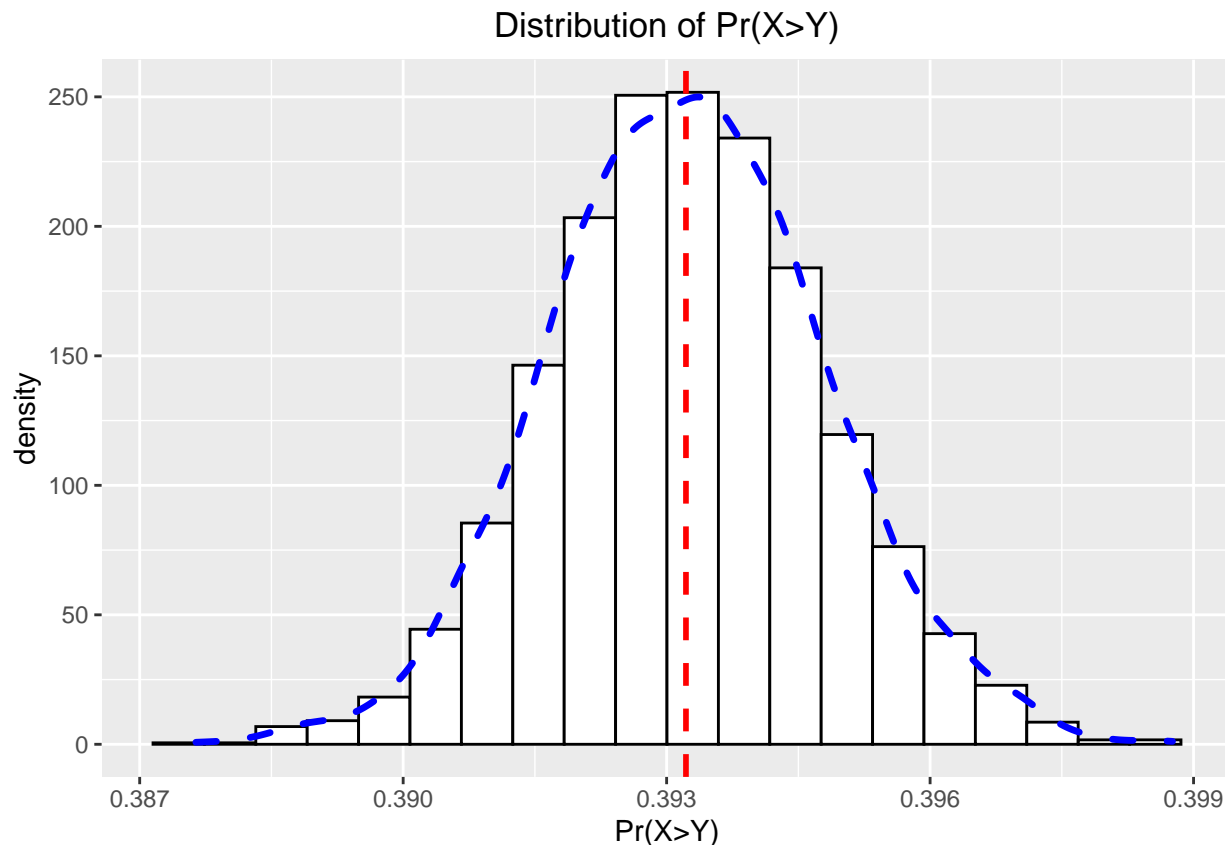
    Prob <- sum(resX > resY) / length(resX) # calculating Pr(X>Y) of each
    return(Prob) # bootstrap sample
  }, mc.cores = 6))
  return(prob_vector)
}
```

The bootstrap algorithm above re-samples the vectors  $X$  and  $Y$  *with replacement*, calculates the resulting  $\widehat{Pr}(X > Y)$ , and repeats this procedure  $n$  times. The algorithm generates a vector with  $n$  probabilities. Using the bootstrap algorithm, the distribution of  $\widehat{Pr}(X > Y)$  can now be evaluated.

1. Point estimates (Quantile):

```
##          2.5%          50%          97.5%
## 0.3902497 0.3932200 0.3963615
```

2. Visualization (histogram):



From the observed data we can infer that, after bootstrapping the original sample, the data is normally distributed around the central value (median) of 0.39322

Lastly, it is of interest how the sample variance of the sampling distribution changes with dependence on the simulations run. This part of the experiment is based on the *law of large numbers* theorem. As an observer, we can assume that with increasing number of simulations (increasing number of random deviates used), our *simulated*  $\widehat{\Pr}(X > Y)$  approximates the *theoretical* value of  $\Pr(X > Y)$  (Dekking et al. (2005)).

To test this assumption, we analyse the behavior of  $\widehat{\Pr}(X > Y)$  for  $n$  simulations, where  $n$  is a vector from 5 to 300,000 by steps of 100 simulations.

```
rm(.Random.seed)

# vector for the number of simulations/deviates, as a sequence from 5 to 300k
# by steps of 100.
n_deviates <- seq(5, 300000, by = 100)

# Generating n deviates, corresponding to the vector defined above, and
# calculating  $\Pr(X>Y)$  from the used deviates.
prob <- unlist(mclapply(n_deviates, function(i) {

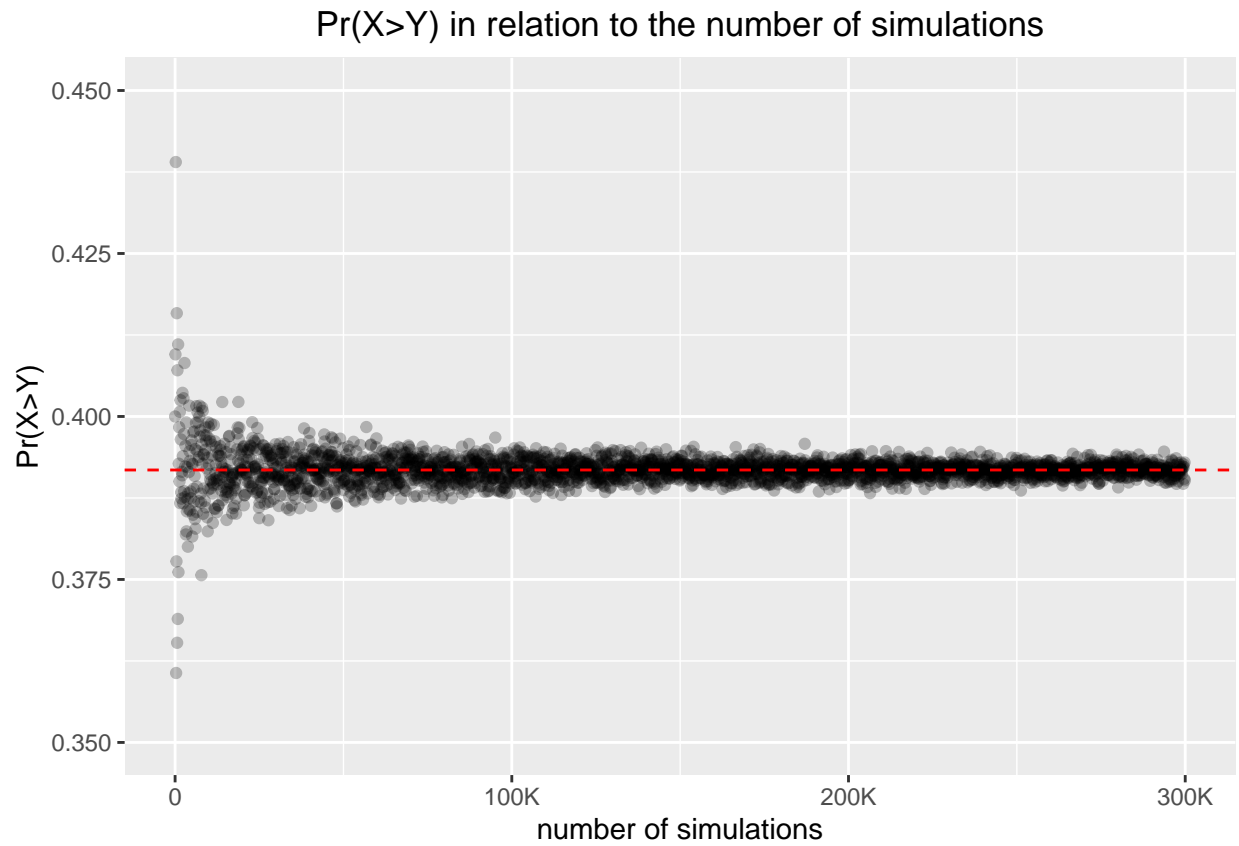
  x <- rnorm(i, mean = 4, sd = sqrt(10))
  y <- runif(i, min = 2, max = 8)

  prob <- sum(x > y) / length(x)
  deviates <- i
```

```

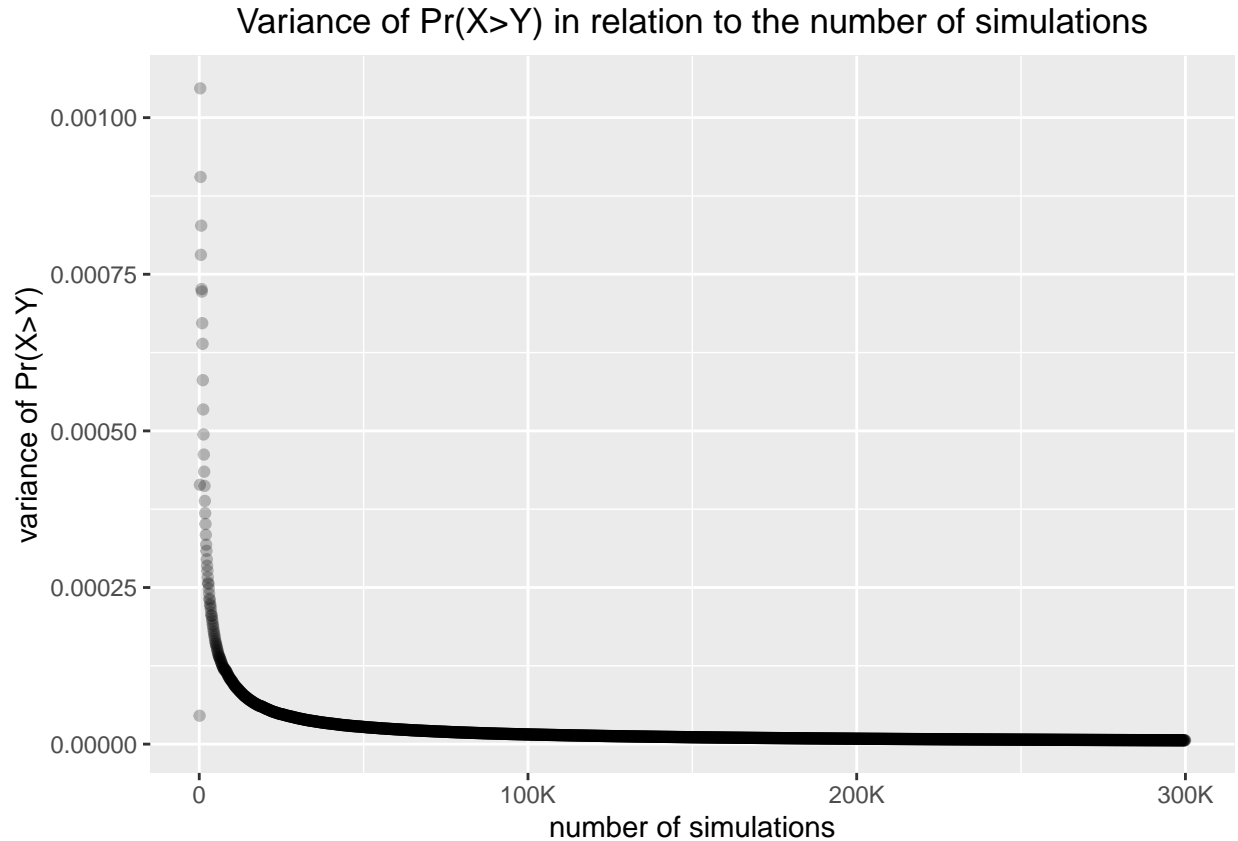
    return(prob)
}, mc.cores = 6, mc.set.seed = TRUE)
)

```



The plot above shows that with an increasing number of simulations used, the observed  $\widehat{Pr}(X > Y)$  converges to its theoretical value (Britannica (2020)).

It can be assumed that for our simulation,  $|\widehat{Pr}(X > Y) - Pr(X > Y)| \rightarrow 0$  as  $n \rightarrow \infty$ .



A look at the total sample variance with dependence on the number of deviates used for the calculation shows that the variance approaches zero arbitrarily closely when more than 100,000 simulations are run, which underlines the assumption made above. It can be deduced that more than 100,000 simulations should be run to calculate the empirical probability  $\widehat{Pr}(X > Y)$ .

## Problem 2

Description: Consider the following football tournament format: a team keeps playing until they accrue 7 wins or 3 losses (whichever comes first - no draws allowed). Assume a fixed win rate  $P \in [0, 1]$  across all rounds (they are paired at random).

Plot how the total number of matches played (i.e. wins + losses) varies as a function of  $p$ .

Comment on the observed win rate relative to the assumed win rate  $p$  (i.e. if a team obtains 2 wins - 3 losses, the maximum likelihood point estimate for their win rate is 40%). Specifically, focus on the effect driven by the format of this tournament.

First, an algorithm simulating the above described tournament is needed. The code below shows a function for simulating the process.

```
# function for simulating the tournament
tournament_sim <- function(p_win) {
  p_loss <- 1 - p_win #loss rate
  outcome <- c("win", "loss")
  results_storage <- c() #empty vector for storing match results
```

```

# as you never play more than 9 games
# (6 wins + 3 losses (=9) or 2 losses and 7 wins (=9))
for (i in 1 : 9) {

  # simulating matches with sampling from c("win", "loss") with given
  # probabilities and replacement.
  # Adding the result to storage vector.
  results_storage <- c(sample(outcome, size = 1,
                             replace = TRUE,
                             prob = c(p_win, p_loss)), results_storage)

  # Conditions for winning or loosing the tournament
  if (length(results_storage[results_storage == "win"]) == 7) {
    break
  }

  if (length(results_storage[results_storage == "loss"]) == 3) {
    break
  }
}

true_winrate <- (sum(str_count(results_storage, pattern = "win")) /
                length(results_storage))

#binding output together as a list
output <- list(played_matches = length(results_storage),
              overview = table(results_storage),
              true_winrate = true_winrate)

return(output)
}

```

To solve the problem, a fixed theoretical win rate  $p$  has to be defined. The following assumes a fixed win rate per round of 75%, so  $p = 0.75$ .

```

# declaring p
p <- 0.75

#simulations one single tournament
results <- tournament_sim(p_win = p)

```

Playing one tournament, the observed results are as following:

1. Overall results:

```

## results_storage
## win
## 7

```

2. Number of matched played:

```
## [1] 7
```

3. True win rate

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

It now should be taken in consideration how the *number of total matches* changes if we alter the winning rate  $p$ . As the conditions for ending a tournament are either 7 wins or 3 losses, the number of played matches can be defined as  $x \in \mathbb{N}[3, 9]$ , where  $x$  is the number of matches played. Win rates from  $p = 0.1$  to  $p = 0.9$  are taken into consideration, in steps by 0.1.

```
# Generating a vector containing the win rates of interest  
win_rates <- seq(from = 0.1, to = 0.90, by = 0.1)
```

The algorithm described below takes a vector of win rates (*rates*) and the number of simulations per fixed win rate  $n$  as input. It then simulates  $n$  tournaments per win rate and stores the number of played matches per tournament in a corresponding storing matrix.

```
tournament_stats_ngames <- function(n, rates = win_rates) {  
  
  # generating a matrix for storing the results  
  store <- matrix(nrow = length(rates), ncol=n)  
  
  # iterations i are corresponding the the number of rates the algorithm  
  # shall take into consideration  
  for (i in 1 : length(rates)) {  
    prob <- rates[i]  
  
    # simulating n tournaments per win rate and storing the number of matches  
    # played into the matrix "store"  
    store[i,] <- unlist(mclapply(1 : n, function(i){  
      as.integer(tournament_sim(p_win = prob)[1])  
    }, mc.cores = 6))  
  }  
  return(store)  
}
```

As a next step, 10,000 matches *for each*  $p$  are simulated and stored in a matrix (*table\_n\_matches*).

```
rm(.Random.seed)  
  
# simulating 10,000 tournaments per fixed win rate (variable win_rates [0.1, 0.9] is  
# set as default)  
table_n_matches <- tournament_stats_ngames(n = 10000)
```

As exact results are of interest, the matrix containing the number of matches per tournament is put into a non-parametric bootstrap algorithm. The algorithm, using a matrix as an input as well as the number of bootstraps  $n$  (*per winning rate*), generates the mean of the strapped sample in a first step (i.e.,  $n$  means per fixed winning rate). In a second step, the algorithm then calculates the mean of the calculated means from

the bootstrapping in step 1.

The aim of this algorithm is to compute the most accurate estimator for the numbers of matches played in a tournament, corresponding to the underlying fixed winning rate  $p$ .

```
# Generating a bootstrap function with number of straps and matrix of results
# as Input. Generates a Data.frame as Output, containing the average value of
# the sample's medians and the corresponding win rate.

bootstrap_problem_two <- function(n_straps, store){

  # generating a store matrix with 9 rows (one per win rate) and one column
  # per bootstrap, storing the median of the bootstrap sample
  avg_storer <- matrix(nrow = 9, ncol = n_straps)

  for (j in 1 : 9){ # 9 iterations (as we use 9 win rates)

    #bootstrap and calculate the mean of each bootstrap sample
    avg_storer[j,] <- unlist(mclapply(1 : n_straps, function(i){
      mean(store[j, sample(1 : dim(store)[2], size = dim(store)[2],
        replace = TRUE)])
    }, mc.cores = 8))
  }

  # generating a vector to store the median of each row
  substore <- rep(NA, 9)

  # calculating the median of each row of avg_storer
  for (i in 1 : 9) {
    substore[i] <- mean(avg_storer[i,])
  }

  # binding together the win rate and its corresponding average matches per
  # tournament
  output <- data.frame(Prob = seq(from = 0.1, to = 0.90, by = 0.1),
    average_value = substore)

  return(output)
}
```

Using this algorithm on our matrix from above (*table\_n\_matches*) with  $n = 1000$  bootstraps, we derive the following results:

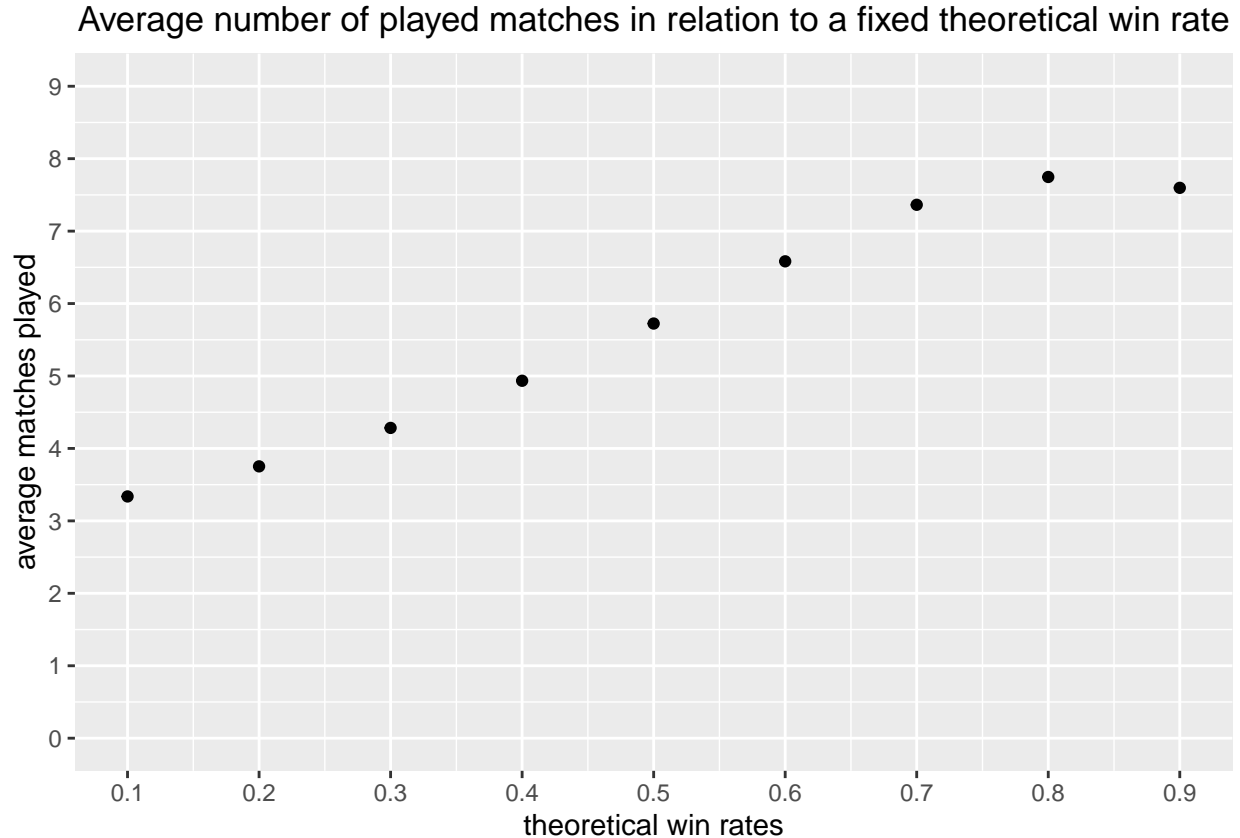
```
## Warning in rm(.Random.seed): object '.Random.seed' not found
```

```
##   probability_theoretical matches_played      delta
## 1                0.1         3.337501 0.0000000
## 2                0.2         3.752637 0.4151362
## 3                0.3         4.283493 0.5308561
## 4                0.4         4.934104 0.6506105
## 5                0.5         5.724208 0.7901044
## 6                0.6         6.583107 0.8588990
## 7                0.7         7.363277 0.7801695
```

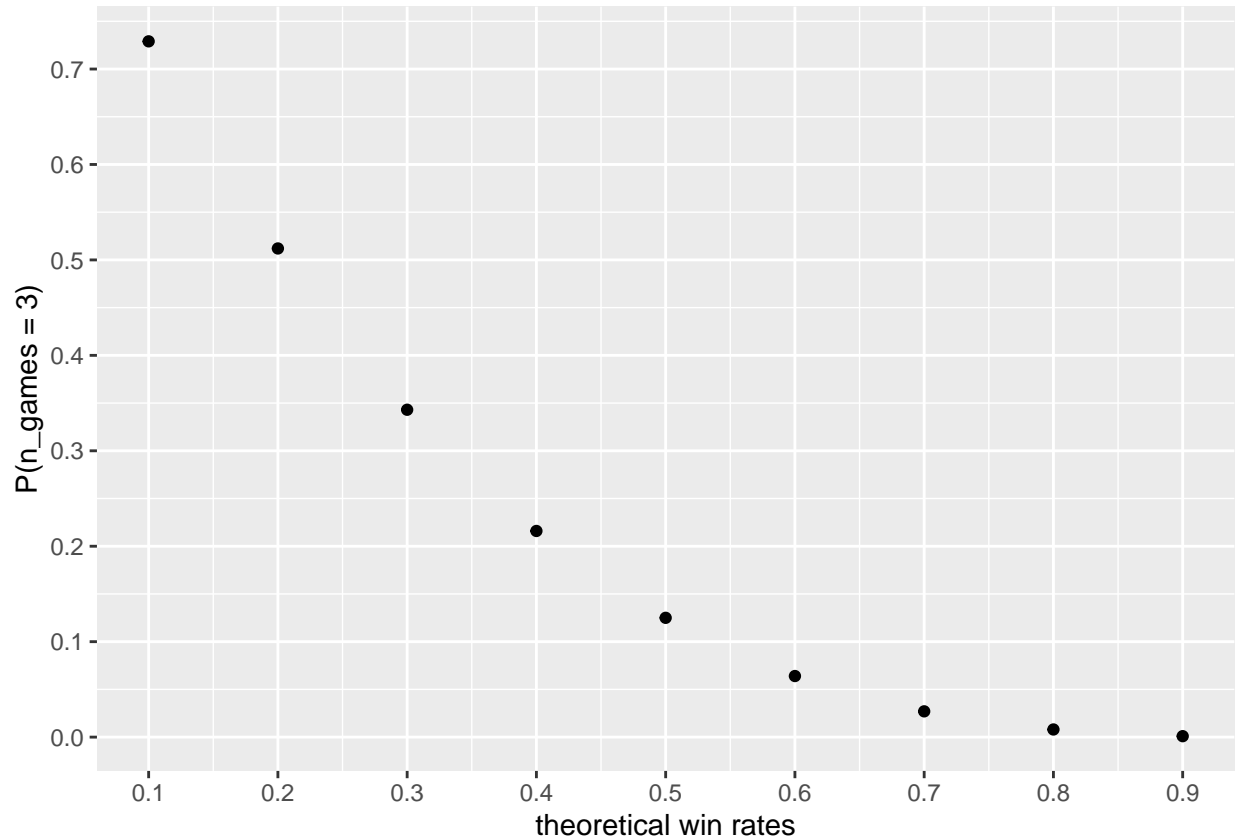


```
## 8          0.8          7.747312  0.3840351
## 9          0.9          7.597410 -0.1499016
```

```
## Warning in rm(.Random.seed): object '.Random.seed' not found
```



Analyzing the average number of matches played within a tournament as a function of  $p$ , we can observe an increase in the number of matches played from  $p = 0.1$  until its peak at  $p=0.8$  and afterwards decreasing. The reason lies in the format of the tournament itself. Due to the conditions of 7 wins or 3 losses, the tournament implies specific characteristics in terms of observable probabilities. In example, the probability for just playing the minimum amount of games, implying a team loses all of its first three games, can be expressed as  $p(n_{games} = 3) = (1 - p)^3$ , while winning *at least* one game can be expressed as  $1 - p(n_{games} = 3)$ . Assuming  $p = 0.1$ ,  $p(n_{games} = 3)$  is  $(1 - 0.1)^3 = 0.729$ , implying that in 73 out of 100 tournaments observed, no more than 3 games are played if  $p = 10\%$ . Thinking of  $p(n_{games} = 3)$  as a *bottleneck*, the probability of losing the first three games decreases exponentially with increasing  $p$ .



This explains the growing  $\delta$  values observed, which are peaking at  $p = 0.6$ .

Another interesting factor is the peak of average played matches at  $p=0.8$  and a negative  $\delta$  at 0.9. Thinking about the problem as a maximization problem of matches played per tournament, we want to know for which win rate  $p$  the probability of playing 9 matches in a tournament is the highest. We are looking at two functions: One for 7 wins and 2 losses ( $f_1(p)$ ) and the other for 6 wins and 3 losses ( $f_2(p)$ ):

$$f_1(p) = \binom{9}{7} p^7 (1-p)^{9-7}$$

$$f_2(p) = \binom{9}{2} (1-p)^2 (1 - (1-p))^{9-2}$$

```
# Using simple grid optimization
p <- seq(0.1, 0.95, 0.001) #vector with fixed win rates

y <- unlist(lapply(p, function(i){ # f(p)_1
  (choose(9,7) * i^7 * (1 - i)^2)
}))

index_max1 <- which.max(y) # index for maximum of f(p)_1
maximum_func1 <- p[index_max1]

z <- unlist(lapply(p, function(i){ # f(p)_2
  (choose(9,3) * (1-i)^3 * (1 - (1-i))^6)
}))

index_max2 <- which.max(z) # index for maximum of f(p)_2
maximum_func2 <- p[index_max2]

print(maximum_func1)
```

```
## [1] 0.778
```

```
print(maximum_func2)
```

```
## [1] 0.667
```

Using a simple grid search algorithm for seeking the maximum of both  $f_1(p)$  and  $f_2(p)$ , both functions peek at a  $p$  of 0.778 and 0.667, respectively. In conclusion, the probability of playing the maximum amount of matches in a tournament decreases with a  $p > 0.78$ . Therefore, we can observe a negative  $\delta$  at  $p = 0.9$

The last part of the problem asks for the true win rates (observed win rate) of a team in comparison to the underlying theoretical win rate of the simulation. To solve this problem with computation, an appropriate algorithm is needed.

The algorithm below is an altered version of the *tournament\_stats\_ngames* algorithm. The only difference is that it now uses the third element of the output list of the original underlying algorithm *tournament\_sim*, which is the proportion of matches won in a tournament.

```
tournament_stats_truewinrate <- function(n, rates = win_rates) {
  store <- matrix(nrow = length(rates), ncol=n)

  for (i in 1 : length(rates)) {
    prob <- rates[i]

    store[i,] <- unlist(mclapply(1:n, function(i){
      as.numeric(tournament_sim(p_win = prob)[3]) # output of tournament_sim is
    }, mc.cores = 6))                          # now true_winrate, as numeric
  }
  return(store)
}
```

We again simulate 10,000 tournaments and use the bootstrap algorithm (1000 bootstrap-samples per theoretical win rate) for deriving the most accurate estimates.

```
true_w_rates <- tournament_stats_truewinrate(10000)

df_true_rates <- bootstrap_problem_two(n_straps = 1000, store = true_w_rates)

df_true_rates$diff <- df_true_rates$average_value - df_true_rates$Prob

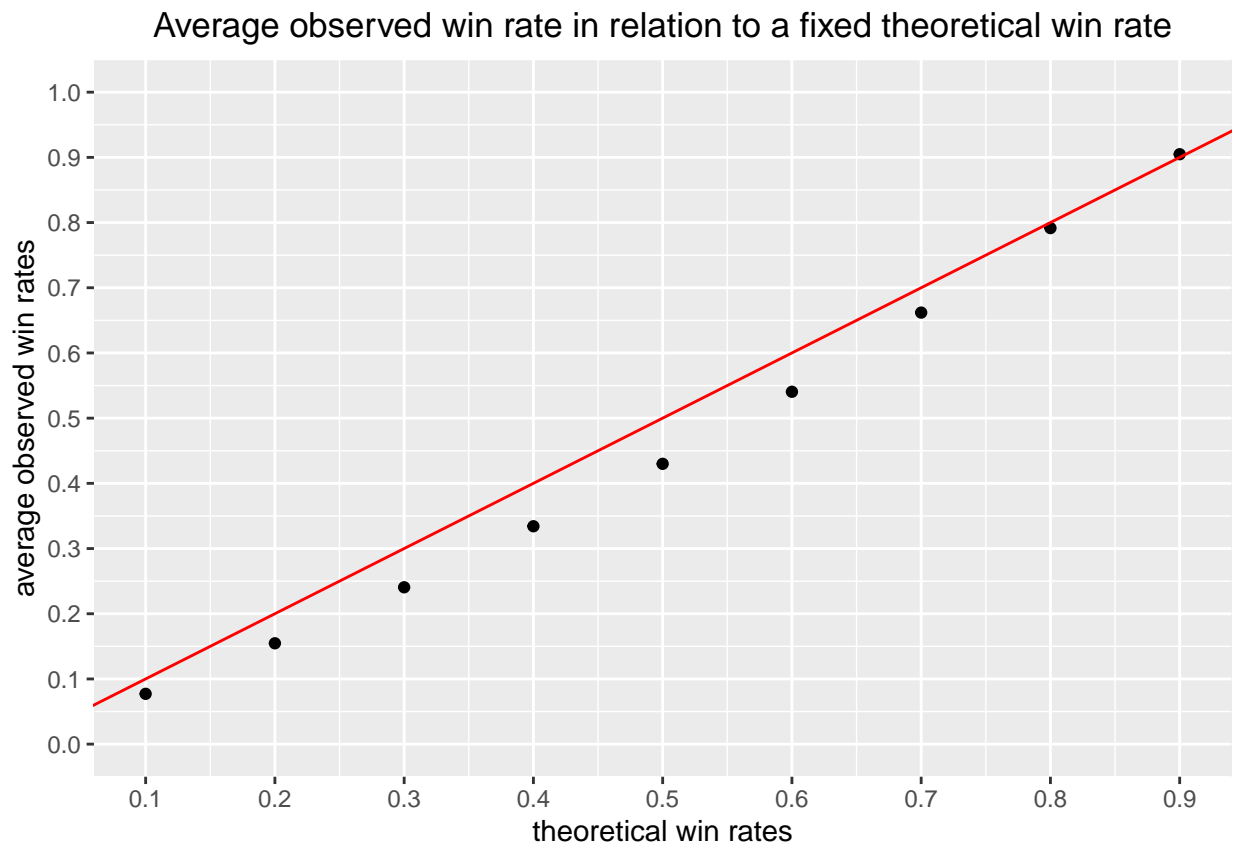
df_true_rates <- df_true_rates%>%
  rename(probability_theoretical = 1,
         win_rate_observed = 2,
         difference = 3)

print(df_true_rates)
```

```
## probability_theoretical win_rate_observed difference
## 1                      0.1              0.07716387 -0.022836134
```

```
## 2          0.2      0.15469365 -0.045306352
## 3          0.3      0.24057159 -0.059428410
## 4          0.4      0.33419858 -0.065801423
## 5          0.5      0.42995967 -0.070040328
## 6          0.6      0.54046562 -0.059534384
## 7          0.7      0.66191646 -0.038083542
## 8          0.8      0.79154954 -0.008450461
## 9          0.9      0.90482995  0.004829948
```

```
df_true_rates %>%
  ggplot(aes(x = probability_theoretical, y = win_rate_observed))+
  geom_point()+
  scale_x_continuous(name="theoretical win rates", limits=c(0.1, 0.9),
                     n.breaks = 9)+
  scale_y_continuous(name="average observed win rates", limits=c(0.0, 1.0),
                     n.breaks = 11)+
  ggtitle("Average observed win rate in relation to a fixed theoretical win rate")+
  theme(plot.title = element_text(hjust= 0.5))+
  geom_abline(slope = 1, colour = "red")
```



The graphic above shows the average (*mean*) observed win rate per match in relation to its underlying fixed win rate. Although simulating 10,000 tournaments and using bootstrapping to obtain the most accurate estimate, we can observe a certain difference between the theoretical win rate per match and the actual observed counterpart.

This discrepancy is again due to the format of the tournament, i.e. the condition of 7 wins or 3 losses ending the tournament explained above. As  $p(n_{\text{games}} = 3) = 0.729$  when  $p = 0.1$ , we can observe a 0% win rate at 73 out of 100 tournaments with such  $p$ , implying a difference between the theoretical and the observed win

rate of -0.1 each time.

The greatest difference between the theoretical and the average observed win rate can be seen at a win rate of 40%, with a difference of -18% from the theoretical winning rate. If we apply the logic from the example above, the probability of no wins at all is  $(1 - 0.4)^3 = 0.216$ , implying a difference of -0.4 in 22 out of 100 tournament compared to the theoretical winning rate.

#### References:

- Britannica, The Editors of Encyclopaedia. 2020. "Convergence." *Encyclopedia Britannica*.
- Dekking, Frederik Michel, Cornelis Kraaikamp, Hendrik Paul Lopuhaä, and Ludolf Erwin Meester. 2005. *A Modern Introduction to Probability and Statistics: Understanding Why and How*. Vol. 488. Springer.