



# FC Barcelona La Liga Performance: A Markov Model

Ayushi Gupta, Krantik Das, Harsh Mittal, Subhanan Maity

September 2023

## 0.1 Problem Statement

The data of the match result of a particular football club (Barcelona FC), for the season 2020-2021 of La Liga has been collected. Barcelona played a total of 76 matches in the entire season. The win loss and draw has been annotated as W L D. We find it trivial to assume that the result of the  $i$ th match only depends on the result of the  $(i - 1)$ th match. To avoid biasness the result of the past 3 matches have been clubbed together, e.g, if they have a result sequence like Win, Loss, Loss, we have annotated it as WLL. Every possible such combination has been taken. Then, it has been attempted to predict the probabilities of transition from one particular state space to another in the sense that if Barcelona FC has won the past 2 matches and lost the third one (WWD), what is the probability that it will win the fourth match i.e. transition to state space WDW. In order to model it as a markov chain, we have made certain assumptions. Firstly, the match result in the past three matches gives sufficient information to predict the match result of the next match. Secondly, all other variables, both quantifiable and non quantifiable, such as psychological factors, weather conditions, injuries to the players etc. are kept as constant. Thus, according to our assumptions, we are of the opinion that the current problem statement can be modeled as a Markov Chain.

## 0.2 Data

The table below shows the first ten data points. Here, the dates have been indexed as time spaces with integer values. Each match result has been assigned a match point. Win, lose and draw are assigned values 1,-1, and 0 respectively. Total score at each time period is the cumulative sum of the match points till that time period. State is the main variable of concern which annotates the match result of the past three matches for each time period. Thus considering every possible state space, there can exist a total of 27 state spaces. For 76 matches, because the first two time periods will not have a state space, a total of 74 data points remain for state spaces. Specific to our data, only 22 state spaces were observed for Barcelona in the 2020-21 season and the same 22 states have been considered.

## 0.3 Methodology

### 0.3.1 State Spaces

According to the match result of the team, the new variable state was created as shown in the before mentioned table. Every  $i$ th data point of this new variable codes the information of the match results of the  $i$ th,  $(i - 1)$ th, and  $(i - 2)$ nd match. In total, we have 74 such data points because of the presence of 76 matches played by the team. There can be three values  $W$ ,  $L$ ,  $D$  that each

match result can have, and as a result, there are 27 unique combinations for the new variable.

As per the data of Barcelona, all the values of the new variable were found out. There were 22 observed unique values which have been taken as the state spaces. While coding in R, the sequence of the new variable has been stored as a vector called `match_results`. All the unique values of this vector have been taken as the state spaces and coded as a new vector named `state_spaces`.

```

1 # Create a sequence of states (replace with your own data)
2 # Define the possible states
3 match_results <- c(
4   "WDD", "DDD", "DDW", "DWD", "WDW", "DWD", "WDD", "DDW",
5     "DWW", "WWW",
6   "WWD", "WDW", "DWD", "WDW", "DWW", "WWD", "WDD", "DDW",
7     "DWD", "WDL",
8   "DLL", "LLL", "LLW", "LWW", "WWW", "WWW", "WWW", "WWL",
9     "WLW", "LWW",
10  "WWW", "WWL", "WLW", "LWD", "WDW", "DWD", "WDD", "DDD",
11    "DDL", "DLW",
12  "LWW", "WWW", "WWD", "WDD", "DDD", "DDD", "DDW", "DWL",
13    "WLW", "LWW",
14  "WWW", "WWW", "WWL", "WLD", "LDL", "DLW", "LWL", "WLW",
15    "LWD", "WDD",
16  "DDW", "DWD", "WDD", "DDW", "DWW", "WWW", "WWW", "WWW",
17    "WWW", "WWW",
18  "WWD", "WDW", "DWL", "WLW"
19 )
20 # Define the possible state spaces
21 state_spaces <- unique(match_results)

```

Listing 1: State Space Sequence

Date	TimeState	Match Result	Match Point	Total Score	STATE	
27-09-2020	1	W	1	1		
04-10-2020	2	D	0	1		
24-10-2020	3	D	0	1	WDD	
07-11-2020	4	D	0	1	DDD	
29-11-2020	5	W	1	2	DDW	
13-12-2020	6	D	0	2	DWD	
16-12-2020	7	W	1	3	WDW	
19-12-2020	8	D	0	3	DWD	
29-12-2020	9	D	0	3	WDD	
31-01-2021	10	W	1	4	DDW	
13-02-2021	11	W	1	5	DWW	

Table 1: Match Results for Barcelona FC

### 0.3.2 Transition Probability Matrix

Using the *"match\_results"* vector which displays the sequence of the state spaces chronologically, a transition probability matrix was created. The dimension of the matrix has a dimension of 22 rows and 22 columns. Each element of the matrix represents the probability of transitioning from one state space to another. The probability of transition from state A to B is calculated by dividing the frequency of state A going to state B by the total frequency of state A going to all the 22 state spaces.

The same has been coded in R. The transition probability matrix has been calculated and named *"transition\_matrix\_prob"*. The printed *"transition\_matrix\_prob"* is shown as output below. For visual simplicity, only a 4\*4 matrix, that is first 16 entries, have been printed.

```
1
2 # Create an empty transition probability matrix
3 n_states <- length(state_spaces)
4 transition_matrix <- matrix(0, nrow = n_states, ncol =
5     n_states, dimnames = list(state_spaces, state_spaces))
6
7 # Function to update the transition matrix
8 update_transition_matrix <- function(matrix, sequence) {
9     from_state <- sequence[1]
10    to_state <- sequence[2]
11    matrix[from_state, to_state] <- matrix[from_state,
12        to_state] + 1
13    return(matrix)
14 }
15
16 # Loop through the match results and update the transition
17 # matrix
18 for (i in 1:(length(match_results) - 1)) {
19     transition_matrix <-
20         update_transition_matrix(transition_matrix,
21             c(match_results[i], match_results[i + 1]))
22 }
23
24 # Normalize the transition matrix to obtain probabilities
25 row_sums <- rowSums(transition_matrix)
26 transition_matrix_prob <- transition_matrix / row_sums
27
28 # Print the transition probability matrix
29 print(transition_matrix_prob)
30
31 View(transition_matrix_prob)
```

Listing 2: Transition Probability Matrix

The transition probability matrix is shown below:

	WDD	DDD	DDW	WDW
WDD	0.0	0.4285714	0.5714286	0.0
DDD	0.0	0.2500000	0.5000000	0.0
DDW	0.0	0.0000000	0.0000000	0.5
WDW	0.0	0.0000000	0.0000000	0.6

After calculating the transition probability matrix, it can be clearly observed that the sum of each row of the transition probability matrix is equal to 1. Thus, we can effectively conclude that the model represents a Markov Chain.

### 0.3.3 Visualization of Transition Probability Matrix

Since the number of state spaces is huge, it seemed imprudent to construct a transition probability diagram. Instead, the transition probability matrix has been visualized using a heatmap. A heatmap of a transition probability matrix visually represents the probabilities of transitioning from one state to another in a Markov chain or similar stochastic process. Each cell in the heatmap corresponds to the probability of transitioning from the row state to the column state. The color of each cell is typically used to encode the magnitude of the transition probability. The intensity of color in each cell indicates the probability value. Darker colors represent higher probabilities, while lighter colors represent lower probabilities. This allows us to quickly identify which transitions are more likely or less likely.

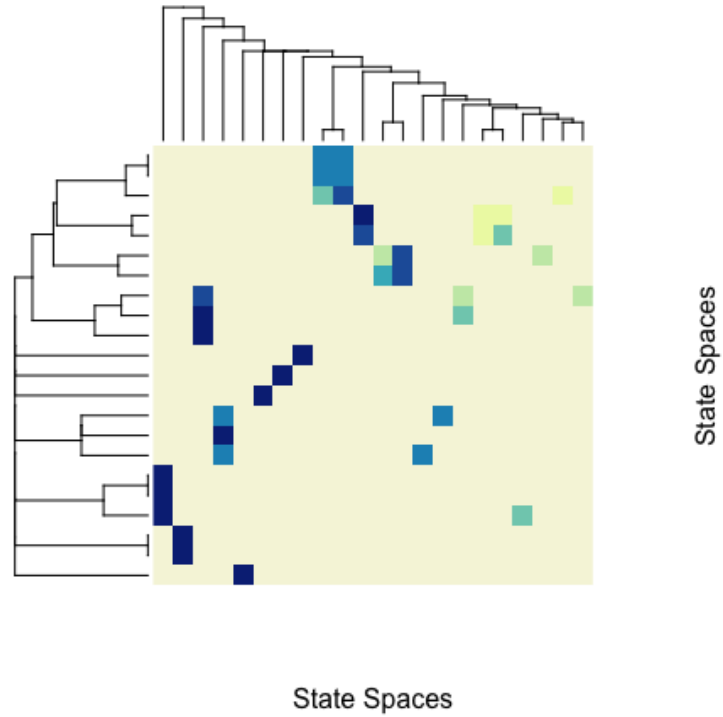
The same has been done in R with the following code. Below is the heatmap obtained:

```

1
2 # Load the required package for color palettes
3 library(RColorBrewer)
4
5
6 # Define the custom color palette
7 custom_palette <- c("beige", brewer.pal(8, "YlGnBu"))
8
9 # Create a heatmap with the custom color palette
10 heatmap(transition_matrix_prob,
11         col = custom_palette,
12         xlab = "State Spaces",
13         ylab = "State Spaces",
14         labRow = NA, # Remove row labels
15         labCol = NA) # Remove column labels

```

Listing 3: Heatmap



### 0.3.4 MCMC Simulation using Gibbs Sampling

The goal of the Markov Chain Monte Carlo simulation is to sample prior and posterior distributions for our Markov chain. While there are many methods in MCMC, here the Gibbs sampling method has been used. The following steps are involved in the process of Gibbs sampling to get the prior and posterior distribution:

1. In order to initialize our variables, empty lists named “prior\_samples” and “posterior\_samples” have been created to store samples from the prior and posterior distributions.
2. The Gibbs sampling loop involves the generation of samples from the prior and posterior distributions. Gibbs sampling is a technique for generating samples from a complicated distribution by iteratively sampling from simpler conditional distributions. In each iteration of the loop, an initial “current\_state” has been selected from the “match\_results”.
3. In the Gibbs sampling process, in order to get the prior distribution, a new state is chosen based on the transition probabilities from the current state to all state spaces. This new state is then included in the “prior\_samples” list.

4. For the posterior distribution, a similar process is followed but based on the transition probabilities from the new state that was added in the “prior\_samples” list.
5. We repeat this process for a specified number of iterations.
6. Then, the indices are converted back into the state spaces for easier interpretation.
7. To get the prior and posterior distributions, we normalize the prior and posterior samples by dividing by the total number of samples to get probabilities.

We have run these codes in R and have gotten the prior and posterior distributions. They are then printed, and for better visualization of the posterior distribution, a bar plot has been constructed for the same.

```

1 # Define the number of iterations for Gibbs sampling
2 num_iterations <- 1000
3
4
5 # Number of burn-in iterations
6 burn_in_iterations <- 100
7
8 # Ensure state_sequence contains unique state sequences
9 match_results <- unique(match_results)
10
11 # Create an index map to match state sequences to matrix
    indices
12 state_indices <- match(match_results,
    rownames(transition_matrix_prob))
13
14 # Normalize the transition probability matrix
15 transition_matrix_prob_normalized <- transition_matrix_prob
    / rowSums(transition_matrix_prob)
16
17 # Initialize an initial state
18 current_state_index <- sample(1:length(match_results), 1)
19
20 # Initialize empty vectors to store prior and posterior
    samples
21 prior_samples <- character(0)
22 posterior_samples <- character(0)
23
24 # Gibbs sampling loop
25 for (iteration in 1:num_iterations) {
26   # Sample the prior distribution
27   prior_sample_index <- sample(1:length(match_results), 1,
    prob =
    transition_matrix_prob_normalized[current_state_index,
    ])

```

```

28 prior_samples <- c(prior_samples,
29                     match_results[prior_sample_index])
30
31 # Update the current state
32 current_state_index <- prior_sample_index
33
34 # Sample the posterior distribution
35 posterior_sample_index <- sample(1:length(match_results),
36                                 1, prob =
37                                 transition_matrix_prob_normalized[current_state_index,
38                             ])
39 posterior_samples <- c(posterior_samples,
40                       match_results[posterior_sample_index])
41 }
42
43 # Optionally, you can convert state indices back to state
44 # sequences for interpretation
45 posterior_samples_sequences <-
46   match_results[match(posterior_samples, match_results)]
47
48 # Calculate the prior and posterior distributions
49 prior_distribution <- table(prior_samples) /
50   length(prior_samples)
51 posterior_distribution <-
52   table(posterior_samples_sequences) /
53   length(posterior_samples_sequences)
54
55 # Print or visualize the prior and posterior distributions
56 print("Prior_Distribution:")
57 ## [1] "Prior Distribution:"
58 print(prior_distribution)
59 ## prior_samples
60 ##   DDD   DDL   DDW   DLL   DLW   DWD   DWL   DWW   LDL
61 ##   LLL   LLW   LWD   LWL
62 ## 0.044 0.013 0.067 0.016 0.032 0.082 0.020 0.035 0.019
63 ## 0.016 0.016 0.030 0.016
64 ##   LWL   WDL   WDW   WLD   WLW   WWD   WWL   WWW
65 ## 0.062 0.081 0.016 0.069 0.019 0.060 0.054 0.043 0.190
66
67 print(posterior_distribution)
68 ## posterior_samples_sequences
69 ##   DDD   DDL   DDW   DLL   DLW   DWD   DWL   DWW   LDL
70 ##   LLL   LLW   LWD   LWL
71 ## 0.050 0.010 0.065 0.016 0.032 0.077 0.027 0.032 0.019
72 ## 0.016 0.016 0.030 0.020
73 ##   LWL   WDL   WDW   WLD   WLW   WWD   WWL   WWW
74 ## 0.058 0.092 0.012 0.062 0.014 0.065 0.057 0.045 0.185
75 # Create a bar plot of the posterior distribution
76 barplot(posterior_distribution, names.arg =

```

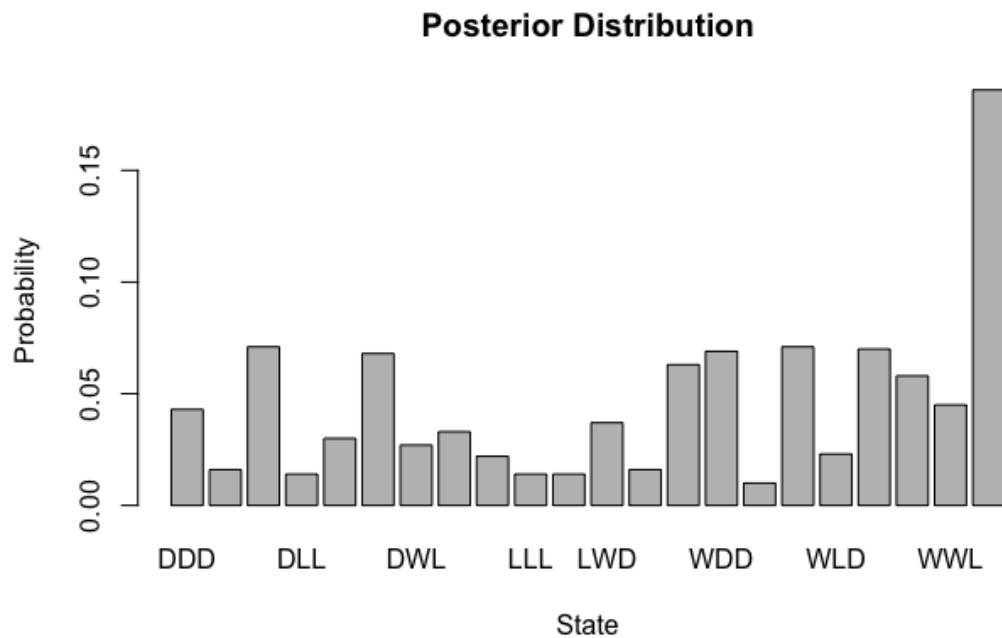


```

names(posterior_distribution),
  xlab = "State", ylab = "Probability", main =
    "Posterior_Distribution")

```

Listing 4: Gibbs Sampling



### 0.3.5 Convergence

There are many techniques to check for convergence of a markov chain. Here, we have used the trace plot to visualize how the sampled states change over iterations. This indicates whether the markov chain has reached a stable distribution or not. The trace plot has been plotted for the posterior distribution using R. The below shown trace plot comes out to be a horizontal line at value 0. This indicates a stable distribution and Markov Chain Monte Carlo (MCMC) algorithm has converged to a state where the sampled values of the parameter are centred around zero.

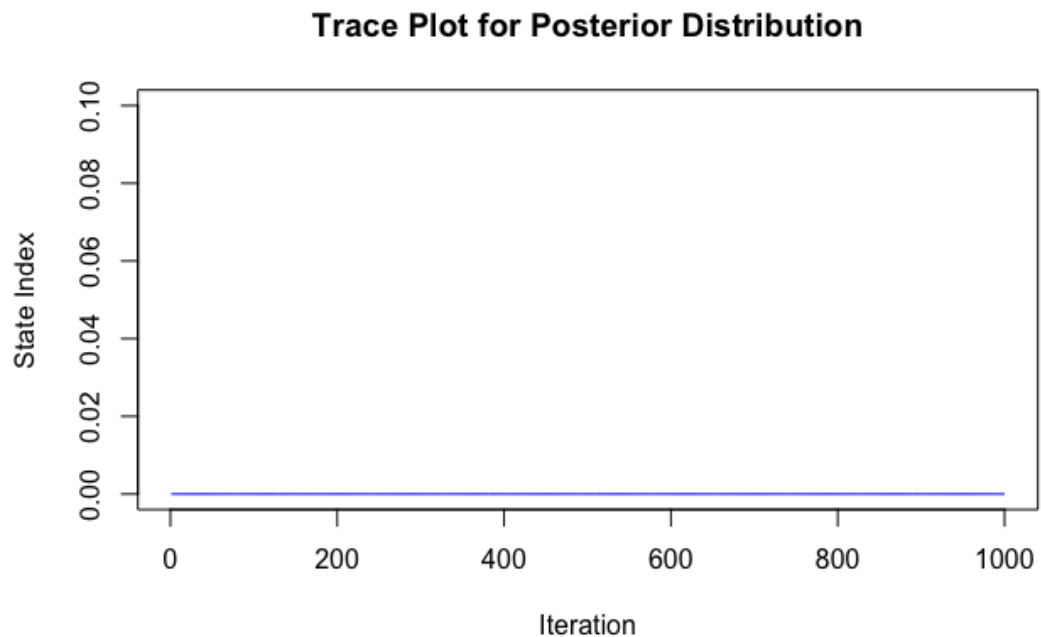
The same has been done in R with the following code. Below is the traceplot obtained:

```

1
2 #CONVERGENCE
3 # Initialize the vector for trace plot
4 trace_plot <- numeric(num_iterations)
5
6 y_axis_limits <- c(0.0, 0.1)
7 # Create a trace plot for the posterior distribution
8 plot(1:num_iterations, trace_plot, type = "l", col =
  "blue", xlab = "Iteration", ylab = "State_Index", main =
  "Trace_Plot_for_Posterior_Distribution",ylim =
  y_axis_limits)

```

Listing 5: Trace Plot



### 0.3.6 Stationarity

#### Stationary Distribution computation through Eigen Values

According to our Markov Chain model, we have finite number of state spaces. The markov chain is also irreducible, because there exists some number of matches, that Barcelona FC can play, such that it can move from any state

space of any state space. Thus, our Markov chain must have a unique stationary probability distribution.

In the context of Markov chains, eigenvalues and eigenvectors are used to find the stationary distribution of the chain, which describes its long-term behaviour. So we have calculated the eigen values and eigen vectors of the transpose of the transition probability matrix. We have then stored all the eigen values with values close to one in a vector to get the stationary distribution. This is because stationary distribution corresponds to the eigenvector associated with the eigen value 1. After normalizing the stationary distributions, we have printed the same normalized stationary distribution.

```
1
2 # Compute the eigenvalues and eigenvectors
3 eigen_result <- eigen(t(transition_matrix_prob)) #
   Transpose the matrix for column-wise eigenvectors
4
5 # Find the index of the eigenvalue closest to 1 (stationary
   distribution)
6 index <- which(abs(eigen_result$values - 1) < 1e-8)
7
8 # Extract the corresponding eigenvector as the stationary
   distribution
9 stationary_distribution <- eigen_result$vectors[, index]
10
11 # Normalize the stationary distribution to sum to 1 (if
   needed)
12 stationary_distribution <- stationary_distribution /
   sum(stationary_distribution)
13
14 # Print the stationary distribution
15 print(stationary_distribution)
```

Listing 6: Stationary Distribution Calculation

Index	Value
1	0.08435003
2	0.04820002
3	0.07230003
4	0.07887276
5	0.07120457
6	0.03834092
7	0.18624882
8	0.05576080
9	0.01314546
10	0.01314546
11	0.01314546
12	0.01314546
13	0.06040038
14	0.04298050
15	0.06813300
16	0.03406650
17	0.01205000
18	0.02637684
19	0.02629092
20	0.01432683
21	0.01432683
22	0.01318842

Table 2: Values as a Table

If a Markov chain reaches stationarity at time period  $N$ , the transition probability matrix at time period  $N$  will be stagnant for all the transition probability matrices at time periods greater than  $N$ , i.e.,  $N + 1, N + 2, N + 3, \dots$ . Thus, we have performed a simulation of matrix multiplication in R using the transition probability matrix. We have created an R loop where the transition probability matrix is multiplied by itself until we get the "current matrix" =  $(\text{transition\_matrix\_prob})^N$ .

In order to find out at what time period the transition probability matrix becomes stagnant, we have calculated this  $N$ , which comes out to be time period 42. This means that the transition probability matrix for time periods later than 42 should be identical to the current matrix. To verify the same, we have directly calculated the 43rd and 45th power of the transition probability matrix. On comparing them, the three matrices—current matrix and transition probability matrix to the power 43 and 45—come out to be equal.

```

1
2 \subsection{Matrix Convergence}
3
4 #MATRIX MULTIPLICATION OVER A NUMBER OF TIME PERIODS
5

```

```

6 # Define the matrix A
7 A <- transition_matrix_prob
8
9 # Define the tolerance level
10 tolerance <- 1e-6
11
12 # Initialize variables
13 N <- 1
14 previous_matrix <- A
15 current_matrix <- A %*% A # A^2
16
17 # Loop until convergence
18 while (sum(abs(current_matrix - previous_matrix)) >
19       tolerance) {
20   N <- N + 1
21   previous_matrix <- current_matrix
22   current_matrix <- current_matrix %*% A # Multiply by A
23   again
24 }
25
26 # N is the smallest power at which convergence occurs
27 cat("Convergence_at_power_N=", N, "\n")
28
29 # The stationary matrix is in 'current_matrix'
30 cat("Stationary_Matrix_(A^N):\n")
31 print(current_matrix)
32 view(current_matrix)
33
34 # Calculate A^43 directly
35 power_43_direct <- expm::'%^%'(A, 43)
36 print(power_43_direct)
37 view(power_43_direct)
38
39 # Calculate A^45 directly
40 power_45_direct <- expm::'%^%'(A, 45)
41 print(power_45_direct)
42 view(power_45_direct)

```

Listing 7: Stationarity

## 0.4 Results

1. After calculating the transition probability matrix, it can be clearly observed that the sum of each row of the transition probability matrix is equal to 1. Thus, we can effectively conclude that the model represents a Markov Chain.
2. The heatmap shows that, there are very few non zero values in the transition probability matrix. This can also be logically verified by taking a

simple instance, that is, if Barcelona won the first three matches (WWW), it cannot go to WLD in the immediate next time period.

3. By performing the Monte Carlo Markov Chain simulation using the Gibbs Sampling algorithm, we were successfully able to arrive at a prior distribution and at a posterior distribution.
4. We have also taken a burn in period to take care of any biasedness that might arrive due to the sampling of the prior distribution.
5. The convergence of the posterior distribution was assessed using the traceplot diagram. The horizontal line at value zero verifies convergence. Additionally, we can conclude that the sampled values of the parameter are centred around zero.
6. The presence of finite state space and irreducibly in our Markov Chain concludes to a unique stationary distribution.
7. The Stationary distribution was calculated by finding out the eigen values of the transposed transition probability metrics. The resultant stationary distribution was summing up to one.
8. The posterior distribution is equal to this stationary distribution with some tolerance indicating a positive sign for convergence.
9. The verification of stationarity in a Markov Chain was performed by matrix multiplication of transition probability matrix for multiple time periods.
10. It was found out that the Markov Chain reaches stationarity at time period 42.

## 0.5 Inference

It can be observed that Barcelona's past 3 performances are sufficient to predict the result of the upcoming match. The probability transition matrix for Barcelona converges at time period 42, that is, 44th match. After playing 44 matches, the overall performance of the team becomes constant and we can predict the next match outcomes using the same set of probability values.

### Team Members:

- Ayushi Gupta (2022EE02) - Modelling, coding, analysis, interpretation, and presentation
- Krantik Das (2022QF12) - Modelling, coding, analysis, interpretation, and presentation
- Harsh Mittal (2022EE06) - Modelling, coding
- Subhanan Maity (2022FE33) - Analysis, interpretation