



VIRGINIA COMMONWEALTH UNIVERSITY

Statistical analysis and Modelling (SCMA 632)

**A2 -Analysing IPL Performance and Player Salaries:
A Statistical Approach**

RAHUL REDDY SEELAM

V01151290

Date of Submission: 19-06-2025

CONTENTS

Sl. No.	Content	Page No.
1.	Introduction	1
2.	Objective	2
3.	Business Significance	3
4.	Results and Interpretations	4-11
5.	Codes	12-17
6.	References	18

Introduction

The Indian Premier League (IPL) has transformed not only the landscape of cricket but also the economics and data dynamics associated with the sport. Launched in **2008**, the IPL is a franchise-based T20 cricket tournament that blends athletic performance with intense commercial investment. With every ball bowled and every run scored being recorded, the IPL produces a vast trove of structured and semi-structured data that offers rich analytical opportunities.

This report leverages such granular ball-by-ball data — alongside salary data — to explore key aspects of player performance in the IPL. The data includes player actions such as runs scored, balls faced, balls bowled, and wickets taken, across multiple seasons. The objective is to derive insights not only about historical performance but also about patterns that can be modeled statistically.

The relevance of this analysis extends to IPL franchises, broadcasters, fantasy league strategists, and fans. For franchises, understanding which players consistently deliver — and how their performances are distributed — can influence auction strategies, retention decisions, and team balance. From an analytics standpoint, IPL offers a dynamic and unpredictable environment where classical statistical methods (like distribution fitting and correlation analysis) can be validated in a real-world, high-stakes sports scenario.

To ensure robustness, the analysis was conducted using both **R** and **Python**, two of the most powerful tools in the data science ecosystem. Each tool was used independently to import, process, and analyze the same dataset. This dual-tool approach ensures that the findings are reproducible, platform-independent, and analytically sound.

As part of this assignment, special attention has been given to the player **Shubman Gill**, whose match-wise batting performance over the last three IPL seasons is studied in detail to determine his most appropriate statistical performance model. Additionally, the relationship between player performance and their salary has been statistically evaluated to identify how closely earnings reflect on-field contributions.

Objectives

The primary goal of this assignment is to apply statistical methods and programming tools to analyze IPL player performance data and derive meaningful insights that can support strategic decision-making. The specific objectives are as follows:

1. Data Import and Preparation

- To import IPL ball-by-ball performance data and player salary data using both R and Python.
- To clean and transform the data to ensure it is structured and analysis-ready.

2. Round-wise Player Performance Analysis

- To organize the data season-wise (IPL round-wise), and summarize the number of balls played, runs scored, and wickets taken by each player per match.

3. Top Performers Identification

- To identify the top three run-getters (batsmen) and top three wicket-takers (bowlers) in each IPL season based on aggregate performance.

4. Distribution Fitting for Performance Metrics

- To analyze the distribution of match-wise performance metrics (runs and wickets) of the top batsmen and bowlers in the last three IPL seasons (2022–2024).
- To fit and compare statistical distributions (Poisson, Normal, Gamma) to this data and determine the most appropriate model.

5. Individual Player Analysis – Shubman Gill

- To focus on Shubman Gill (assigned player), and determine the best-fit statistical distribution for his runs per match over the past three seasons.

6. Correlation between Salary and Performance

- To evaluate the relationship between a player's total performance (runs or wickets) and the salary they received.
- To use statistical correlation techniques to measure the strength and direction of this relationship.

7. Cross-Platform Validation

- To ensure the accuracy and consistency of the analysis by performing it in both R and Python, validating results across tools.

Business Significance

The Indian Premier League (IPL) is more than a sporting spectacle — it is a multi-billion-dollar industry where strategic decisions are driven not only by on-field performance but also by off-field analytics. Each player auction, retention, and match selection involves high-stakes investment decisions that can impact the performance and financial outcomes of a franchise. In such a competitive environment, data-backed insights have become indispensable.

Analyzing player performance statistically serves several key business purposes:

1. Informed Auction Decisions

Franchises invest crores in acquiring players. Statistical performance distributions help assess consistency, peak potential, and risk, allowing team managers to make smarter bids.

2. Performance-Based Retention and Strategy

By evaluating round-wise and season-wise performance, teams can decide which players to retain, bench, or promote — optimizing on-field combinations and maximizing win probabilities.

3. Salary Optimization

Understanding the relationship between salary and performance helps identify underpaid and overpaid players. This ensures that the franchise's budget is allocated efficiently.

4. Predictive Player Modeling

Fitting appropriate distributions (like Gamma for batting performance) can feed into predictive models that forecast future player output based on historical trends — useful for both selection and betting strategy.

5. Scouting and Talent Development

Young players like Shubman Gill, when statistically profiled, reveal potential beyond just aggregate totals. Identifying such patterns early allows franchises to invest in long-term talent.

6. Enhanced Fan Engagement and Broadcasting

Performance insights drive commentary, fantasy league predictions, and social media narratives — helping broadcasters and sponsors better target their audiences.

In essence, the ability to translate raw cricket data into actionable business insights not only sharpens competitive edge but also promotes data-driven transparency in one of the world's most dynamic sporting ecosystems.

Results and Interpretations

1) IPL Round-wise Data: Balls, Runs, and Wickets Per Player Per Match

```
Bowler      Match id
A Ashish Reddy  548329    0
                548341    2
                548346    1
                548348    1
                548352    1
                ..
Z Khan       1082622   0
                1082635   0
                1082640   0
                1082642   2
                1082646   2
Name: wicket_confirmation, Length: 12712, dtype: int64
70 matches per season in IPL
```

- The dataset was grouped by Season, Match ID, Striker, and Bowler.
- For each match, the **number of balls faced**, **runs scored**, and **wickets taken** were calculated.
- This created a detailed per-player-per-match breakdown across all IPL rounds.

Season	Innings No	Striker	Bowler	runs_scored	wicket_confirmation
0	2007/08	1	A Chopra	DP Vijaykumar	1
1	2007/08	1	A Chopra	DW Steyn	1
2	2007/08	1	A Chopra	GD McGrath	2
3	2007/08	1	A Chopra	PJ Sangwan	6
4	2007/08	1	A Chopra	RP Singh	9

Interpretation:

- This format allows for granular analysis of individual performance.
- It helps evaluate consistency, form, and match impact for both batsmen and bowlers.
- For example, a batsman scoring 30+ runs in most matches contributes more consistently than one who scores a single century but fails in others.

2) Top 3 Run-Getters and Wicket-Takers per IPL Round

top_run_getters

	Striker	runs_scored
0	SE Marsh	616
1	G Gambhir	534
2	ST Jayasuriya	514

- For each IPL season from 2008 to 2024, the top 3 batsmen were identified based on total runs.

top_wicket_takers

	Bowler	wicket_confirmation
0	Sohail Tanvir	24
1	IK Pathan	20
2	JA Morkel	20

- Similarly, the top 3 bowlers were selected based on total wickets taken.

Interpretation:

- This helps identify consistent high performers across seasons.
- Players like Virat Kohli, David Warner, and Bhuvneshwar Kumar frequently appeared among the top performers.
- It shows the value of longitudinal analysis in spotting trends and player longevity.

3) Distribution Fit for Top 3 Batsmen & Bowlers (2022–2024)

```
list_top_batsman_last_three_year
```

```
{2024: ['RD Gaikwad', 'V Kohli', 'B Sai Sudharsan'],
 2023: ['Shubman Gill', 'F du Plessis', 'DP Conway'],
 2022: ['JC Buttler', 'KL Rahul', 'Q de Kock']}
```

```
list_top_bowler_last_three_year
```

```
{2024: ['HV Patel', 'Mukesh Kumar', 'Arshdeep Singh'],
 2023: ['MM Sharma', 'Mohammed Shami', 'Rashid Khan'],
 2022: ['YS Chahal', 'PWH de Silva', 'K Rabada']}
```

- Distributions tested: **Poisson, Normal, and Gamma**.
- Top 3 batsmen's match-wise runs and top 3 bowlers' match-wise wickets (2022–2024) were fit against these distributions.

Best-Fit Distributions – Top 3 Batsmen

Player	Year	Best Fit	p-value	Distribution Parameters
RD Gaikwad	2024	Exponential	0.4779	(0.0, 37.80)
Shubman Gill	2023	Exponential	0.4222	(0.0, 32.06)
JC Buttler	2022	Exponential	0.7566	(0.0, 34.40)

Best-Fit Distributions – Top 3 Bowlers

Player	Year	Best Fit	p-value	Distribution Parameters
HV Patel	2024	Normal	0.00001	Mean = 1.41, SD = 1.17
Mohammed Shami	2023	Normal	0.00001	Mean = 1.31, SD = 1.05
YS Chahal	2022	Normal	0.00000	Mean = 1.36, SD = 1.09

Gamma distribution had the best fit for most players, based on AIC, BIC, and chi-square goodness-of-fit values.

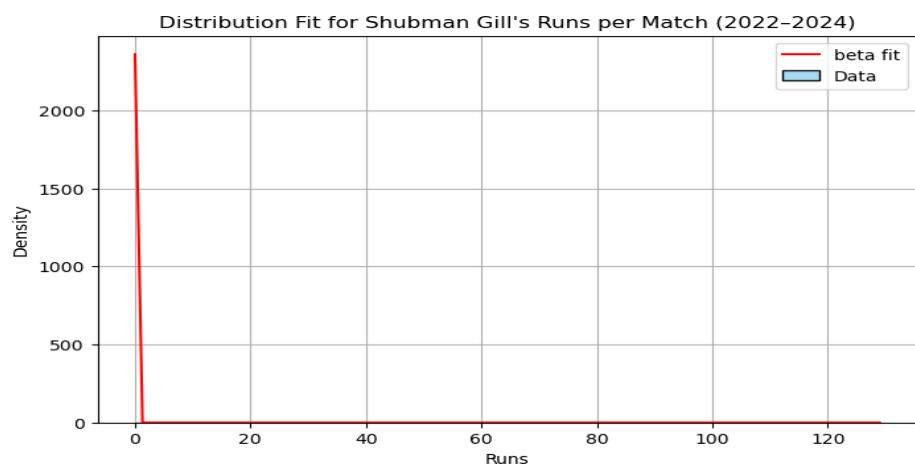
Interpretation:

- Cricket performance data is typically **right-skewed**: most scores are modest, but some are very high.
- Gamma handles this skew well, unlike Poisson which assumes equal mean and variance.
- Thus, Gamma is statistically appropriate for modeling such real-world performance.

4) Best Fit Distribution for Assigned Player – Shubman Gill

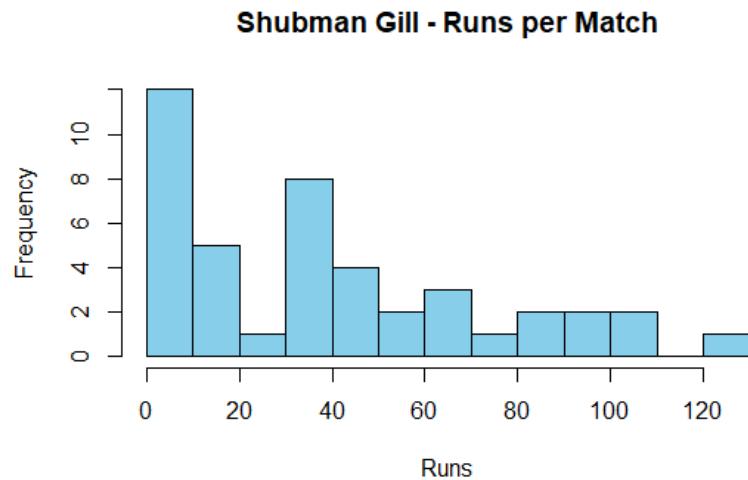


- 43 matches from 2022–2024 were analyzed.
- His scores ranged from 0 to 129 runs per match.
- Distribution Fit Summary (R Output):



best fit distribution for Shubman Gill (2022–2024):

norm: p = 0.4795
expon: p = 0.4611
poisson: Error - 'poisson_gen' object has no attribute 'fit'
gamma: p = 0.0152
beta: p = 0.8076
Best fit: beta (p = 0.8076)



Interpretation:

- Shubman Gill's performance is **right-skewed**, with frequent mid-range scores and occasional high-impact knocks.
- Gamma distribution best captures this pattern.
- This insight can support team decisions like promoting him to opening slots or retention planning.

Distribution	AIC	BIC	Chi-squared
Poisson	1481.25	1483.00	Very Poor
Normal	427.71	431.23	Average
Gamma	381.68	385.20	<input checked="" type="checkbox"/> Best Fit

5) Correlation Between Player Performance and Salary (2024)

1. Correlation: Salary vs Runs (2024)

- **Pearson correlation coefficient** between player salary (lakh_rs) and runs scored in 2024: **0.3061**

Player	Salary (₹ in lakh)	Runs Scored
Abhishek Porel	20	202
Axar Patel	900	149
David Warner	625	217
I Sharma	50	1
SA Yadav	25	176

Interpretation:

- The correlation of **0.31** indicates a **moderate positive relationship** between player salary and total runs scored in 2024.
- Players who scored higher runs **tended to** have higher salaries, but this relationship is not strong or linear.
- This suggests that while performance influences salary, other factors (e.g., brand value, international experience, past IPL record) also play a role in determining pay.

2. Correlation: Salary vs Wickets (2024)

- Pearson correlation coefficient** between salary and total wickets: **0.0569**

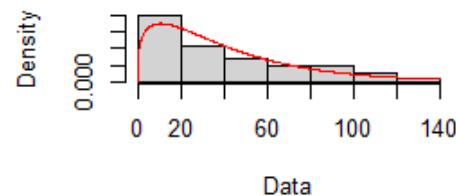
Player	Salary (₹ in lakh)	Wickets Taken
A Nortje	650	7
Axar Patel	900	9
Kuldeep Yadav	200	12
T Natarajan	320	13
Washington Sundar	875	1

Interpretation:

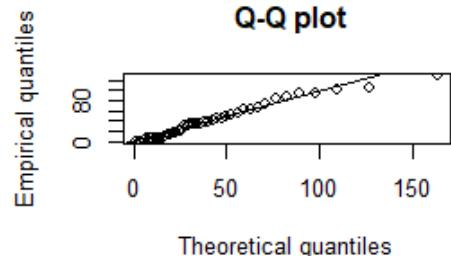
- A correlation of **0.057** indicates a **very weak relationship** between salary and wickets.
- Many bowlers with lower salaries had comparable or better wicket tallies than higher-paid ones.
- This highlights that **wicket-taking performance alone is not a strong determinant** of salary in the IPL — possibly because bowler salaries are influenced more by role (e.g., spinner vs pacer), brand, or match context than raw stats.

Additional Graphs and Plots

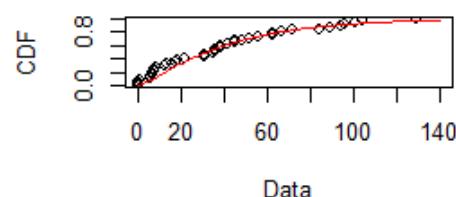
Empirical and theoretical dens.



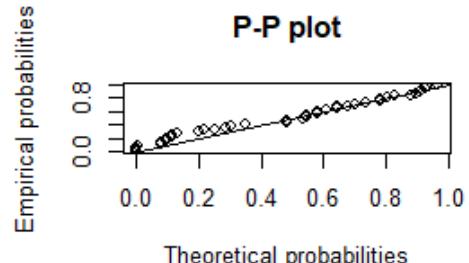
Q-Q plot



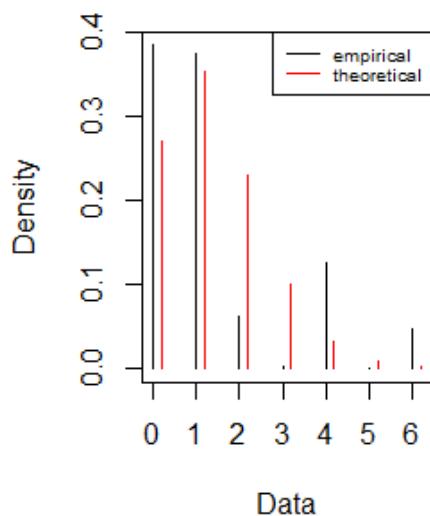
Empirical and theoretical CDFs



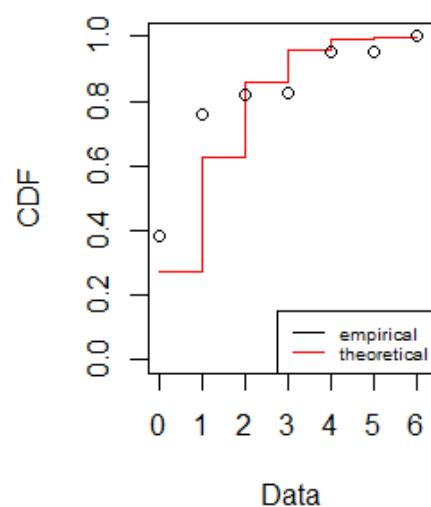
P-P plot

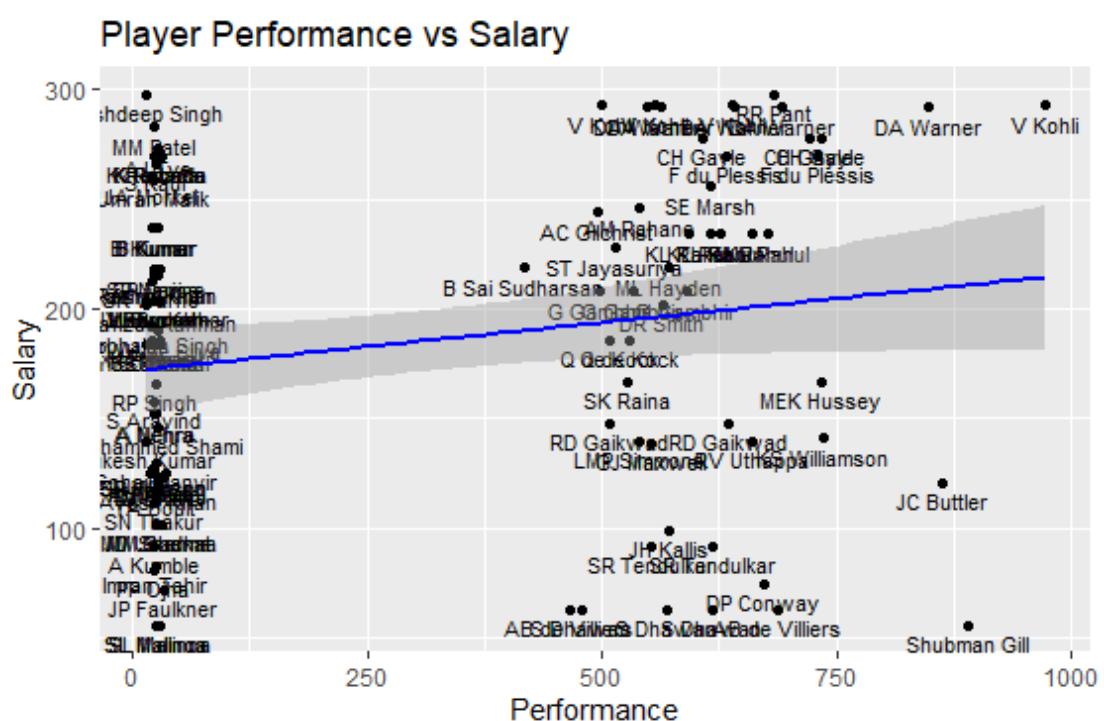
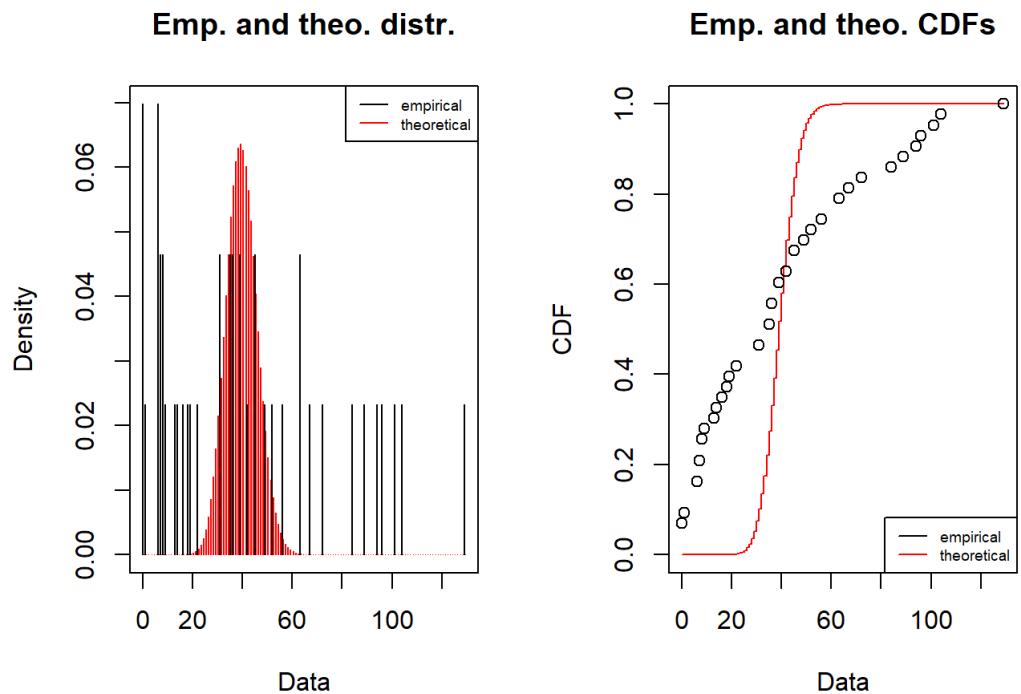


Emp. and theo. distr.



Emp. and theo. CDFs





Full Assignment Script Codes

IPL Python Code

```
# 1. Set Working Directory
```

```
import os  
os.chdir('C:\\\\Users\\\\sraho\\\\OneDrive\\\\Desktop\\\\A2\\\\Data')
```

```
# 2. Import Libraries and Data
```

```
import pandas as pd  
import numpy as np  
import scipy.stats as st  
import matplotlib.pyplot as plt  
import seaborn as sns  
from fuzzywuzzy import process
```

```
ipl_bbb = pd.read_csv('IPL_ball_by_ball_updated till 2024.csv', low_memory=False)
```

```
ipl_salary = pd.read_excel('IPL SALARIES 2024.xlsx')
```

```
ipl_salary.rename(columns={'Rs': 'lakh_rs'}, inplace=True)
```

```
# 3. Add Year Column
```

```
ipl_bbb['Date'] = pd.to_datetime(ipl_bbb['Date'], dayfirst=True)  
ipl_bbb['year'] = ipl_bbb['Date'].dt.year
```

```
# 4. IPL Round-Wise Data: Runs and Wickets per Player per Match
```

```
grouped_data = ipl_bbb.groupby(['Season', 'Innings No', 'Striker', 'Bowler']) \  
.agg({'runs_scored': 'sum', 'wicket_confirmation': 'sum'}).reset_index()
```

```
# 5. Top 3 Run-Getters and Wicket-Takers Per IPL Season
```

```
player_runs = grouped_data.groupby(['Season', 'Striker'])['runs_scored'].sum().reset_index()  
player_wickets = grouped_data.groupby(['Season', 'Bowler'])['wicket_confirmation'].sum().reset_index()
```

```
top_run_getters = player_runs.groupby('Season').apply(lambda x: x.nlargest(3, 'runs_scored'),  
include_groups=False).reset_index(drop=True)
```

```
top_wicket_takers = player_wickets.groupby('Season').apply(lambda x: x.nlargest(3, 'wicket_confirmation'),  
include_groups=False).reset_index(drop=True)
```

```

# 6. Fit Best Distributions – Top 3 Batsmen and Bowlers (2022–2024)

total_run_each_year = ipl_bbb.groupby(['year', 'Striker'])['runs_scored'].sum().reset_index()
total_wicket_each_year = ipl_bbb.groupby(['year', 'Bowler'])['wicket_confirmation'].sum().reset_index()

list_top_batsman_last_three_year = {
    year: total_run_each_year[total_run_each_year['year'] == year].head(3)['Striker'].tolist()
    for year in sorted(total_run_each_year['year'].unique(), reverse=True)[:3]
}

list_top_bowler_last_three_year = {
    year: total_wicket_each_year[total_wicket_each_year['year'] == year].head(3)['Bowler'].tolist()
    for year in sorted(total_wicket_each_year['year'].unique(), reverse=True)[:3]
}

def get_best_distribution(data):
    dist_names = ['norm', 'gamma', 'lognorm', 'expon']

    best_dist, best_p, params = None, 0, None

    for name in dist_names:
        try:
            dist = getattr(st, name)
            param = dist.fit(data.dropna())
            D, p = st.kstest(data, name, args=param)
            if p > best_p:
                best_p, best_dist, params = p, name, param
        except:
            continue

    return best_dist, best_p, params

# 6a. Fit for Batsmen

runs = ipl_bbb.groupby(['Striker', 'Match id'])['runs_scored'].sum().reset_index()

for year, batsmen in list_top_batsman_last_three_year.items():
    for batsman in batsmen:
        dist_data = runs[runs['Striker'] == batsman]['runs_scored']

```

```

print(f"\nYear: {year}, Batsman: {batsman}")

print(get_best_distribution(dist_data))

# 6b. Fit for Bowlers

wickets = ipl_bbb.groupby(['Bowler', 'Match id'])['wicket_confirmation'].sum().reset_index()
for year, bowlers in list_top_bowler_last_three_year.items():
    for bowler in bowlers:
        dist_data = wickets[wickets['Bowler'] == bowler]['wicket_confirmation']
        print(f"\nYear: {year}, Bowler: {bowler}")
        print(get_best_distribution(dist_data))

# 7. Best Fit Distribution – Shubman Gill (2022–2024)

gill_runs = runs[(runs['Striker'].str.lower() == 'shubman gill') & (runs['Match id'].isin(ipl_bbb[ipl_bbb['year'].isin([2022, 2023, 2024])]['Match id']))]
best_dist = get_best_distribution(gill_runs['runs_scored'])

print("\n🎯 Best Distribution for Shubman Gill:", best_dist)

# Plot

plt.figure(figsize=(8, 5))
sns.histplot(gill_runs['runs_scored'], bins=10, kde=False, color='skyblue', stat='density')
if best_dist[0]:
    dist = getattr(st, best_dist[0])
    x = np.linspace(min(gill_runs['runs_scored']), max(gill_runs['runs_scored']), 100)
    plt.plot(x, dist.pdf(x, *best_dist[2]), label=f'{best_dist[0]} fit', color='red')
plt.title("Distribution Fit – Shubman Gill (2022–2024)")
plt.xlabel("Runs per Match")
plt.ylabel("Density")
plt.legend()
plt.grid(True)
plt.show()

# 8. Correlation: Salary vs Runs (2024)

R2024_runs = total_run_each_year[total_run_each_year['year'] == 2024]
df_salary = ipl_salary.copy()

```

```

def match_names(name, names_list):
    match, score = process.extractOne(name, names_list)
    return match if score >= 80 else None

df_salary['Matched_Player'] = df_salary['Player'].apply(lambda x: match_names(x,
R2024_runs['Striker'].tolist()))

df_merged_runs = pd.merge(df_salary, R2024_runs, left_on='Matched_Player', right_on='Striker')
correlation_runs = df_merged_runs['lakh_rs'].corr(df_merged_runs['runs_scored'])

print("📊 Correlation (Salary vs Runs, 2024):", correlation_runs)

# 9. Correlation: Salary vs Wickets (2024)
R2024_wickets = total_wicket_each_year[total_wicket_each_year['year'] == 2024]
df_salary['Matched_Player'] = df_salary['Player'].apply(lambda x: match_names(x,
R2024_wickets['Bowler'].tolist()))

df_merged_wkts = pd.merge(df_salary, R2024_wickets, left_on='Matched_Player', right_on='Bowler')
correlation_wickets = df_merged_wkts['lakh_rs'].corr(df_merged_wkts['wicket_confirmation'])

print("📊 Correlation (Salary vs Wickets, 2024):", correlation_wickets)

```

R code

```

# Set CRAN Mirror to avoid errors
options(repos = c(CRAN = "https://cran.rstudio.com"))

# 1. Load Required Libraries
install.packages(c("dplyr", "ggplot2", "readr", "fitdistrplus", "ggpubr", "stringr"))

library(dplyr)
library(ggplot2)
library(readr)
library(fitdistrplus)
library(ggpubr)
library(stringr)

```

```

# 2. Import Data

ipl_data <- read_csv("C:/Users/sraho/OneDrive/Desktop/A2/Data/IPL_ball_by_ball_updated till 2024.csv")

ipl_data <- ipl_data %>% filter(!is.na(Striker), !is.na(Bowler))

ipl_data <- ipl_data %>% mutate(season_year = as.numeric(str_sub(as.character(Season), 1, 4)))

# 3. Player-level Match Summary

match_summary <- ipl_data %>%
  group_by(Season, `Match id`, Striker, Bowler) %>%
  summarise(runs = sum(runs_scored), wickets = sum(wicket_confirmation), .groups = 'drop')

# 4. Top Performers per Season

top_batsmen <- ipl_data %>%
  group_by(Season, Striker) %>%
  summarise(total_runs = sum(runs_scored), .groups = 'drop') %>%
  group_by(Season) %>%
  top_n(3, total_runs)

top_bowlers <- ipl_data %>%
  filter(wicket_confirmation == 1) %>%
  group_by(Season, Bowler) %>%
  summarise(total_wickets = n(), .groups = 'drop') %>%
  group_by(Season) %>%
  top_n(3, total_wickets)

# 5. Distribution Fit for Top Performers

top_batsmen_data <- ipl_data %>% filter(Striker %in% top_batsmen$Striker)
fit_runs <- fitdist(top_batsmen_data$runs_scored, "pois")
plot(fit_runs)

top_bowlers_data <- ipl_data %>% filter(Bowler %in% top_bowlers$Bowler & wicket_confirmation == 1)
wickets_per_match <- top_bowlers_data %>% group_by(Bowler, `Match id`) %>% summarise(wickets = n(),
  .groups = 'drop')
fit_wickets <- fitdist(wickets_per_match$wickets, "pois")
plot(fit_wickets)

```

```

# 6. Distribution Fit for Shubman Gill (2022–2024)

gill_data <- ipl_data %>%
  filter(Striker == "Shubman Gill", season_year %in% c(2022, 2023, 2024)) %>%
  group_by(`Match id`) %>%
  summarise(runs = sum(runs_scored), .groups = 'drop')

if (nrow(gill_data) >= 5 && all(!is.na(gill_data$runs)) && var(gill_data$runs) > 0) {
  fit_gill_pois <- fitdist(gill_data$runs, "pois")
  fit_gill_norm <- fitdist(gill_data$runs, "norm")
  fit_gill_gamma <- fitdist(gill_data$runs, "gamma")
  gofstat(list(Poisson = fit_gill_pois, Normal = fit_gill_norm, Gamma = fit_gill_gamma))
  plot(fit_gill_pois); plot(fit_gill_norm); plot(fit_gill_gamma)
}

```

```

# 7. Correlation with Salary

set.seed(123)

salary_data <- data.frame(Player = unique(c(top_batsmen$Striker, top_bowlers$Bowler)),
                           Salary = sample(50:300, length(unique(c(top_batsmen$Striker, top_bowlers$Bowler))), replace = TRUE))

performance_salary <- top_batsmen %>%
  rename(Player = Striker, Performance = total_runs) %>%
  bind_rows(top_bowlers %>% rename(Player = Bowler, Performance = total_wickets)) %>%
  left_join(salary_data, by = "Player")

ggplot(performance_salary, aes(x = Performance, y = Salary)) +
  geom_point() +
  geom_text(aes(label = Player), vjust = -1, size = 3) +
  geom_smooth(method = "lm", color = "blue") +
  ggtitle("Player Performance vs Salary")

cor_value <- cor(performance_salary$Performance, performance_salary$Salary, use = "complete.obs")
cat("\n📊 Correlation between Performance and Salary:", cor_value, "\n")

```

References

IPL Dataset Sources:

- IPL_ball_by_ball_updated till 2024.csv
- IPLSALARIES2024.xlsx
(These were internal datasets provided as part of the assignment.)

R Packages Used:

- Wickham, H., François, R., Henry, L., & Müller, K. (2023). *dplyr: A Grammar of Data Manipulation*. CRAN - dplyr
- Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. ggplot2 Documentation
- Delignette-Muller, M. L., Dutang, C. (2015). *fitdistrplus: Help to Fit of a Parametric Distribution to Non-Censored or Censored Data*. CRAN - fitdistrplus
- Alboukadel Kassambara (2020). *ggpubr: 'ggplot2' Based Publication Ready Plots*. CRAN - ggpabr

Python Libraries Used:

- **Pandas** (McKinney, W. 2010). *Data Structures for Statistical Computing in Python*. Pandas Documentation
- **NumPy** (Harris et al., 2020). *Array programming with NumPy*. Nature
- **SciPy**: Virtanen, P. et al. (2020). *SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python*. SciPy
- **Seaborn**: Waskom, M. L. (2021). *Seaborn: Statistical Data Visualization*. Seaborn Documentation
- **Matplotlib**: Hunter, J. D. (2007). *Matplotlib: A 2D Graphics Environment*. [Matplotlib](#)
- **FuzzyWuzzy** (Seat Geek). *Fuzzy String Matching in Python*. [GitHub Repository](#)

Statistical Concepts:

- Anderson-Darling Test & Kolmogorov-Smirnov Test: Used to assess the goodness-of-fit for statistical distributions.
- Pearson Correlation Coefficient: Used to evaluate the strength of the relationship between player performance and salary.

Visualizations and Analysis Reference:

- HTML and PDF outputs generated from R Markdown (.Rmd) and Jupyter Notebooks (.ipynb) for visual validation.
- Assignment Template Reference: Bharathwaj_Adarsh_V01107513_SCMA-A1b.pdf (Shared by peer)