# Cricket Data Deep Dive: Project Report & Documentation

## 1. Executive Summary

This project documents a comprehensive statistical analysis of global cricket data sourced from the Cricsheet database, covering over 10,000 matches across Test, ODI, T20, and IPL formats. The primary goal was to move beyond conventional totals-based statistics to identify **structural trends, measure player efficiency, and determine the actual strategic separators** for success in the modern game.

The analysis confirms that the rise of T20 cricket is the single most defining factor, forcing a shift in player valuation from raw career totals to rate-based efficiency metrics (like strike rate and economy rate). It also reveals that while traditional factors like the toss are overrated, **team discipline and consistency** are critical in deciding the majority of games won by narrow margins.

## 2. Methodology and Technology Stack

The project employed an end-to-end Extract, Transform, Load (ETL) and analysis pipeline using Python and SQL.

| Component | Technology | Role |
| --- | --- | --- |
| **Data Source** | Cricsheet JSON files | Over 10,000 match files across four formats (Test, ODI, T20, IPL). |
| **ETL & Data Modeling** | Python (`transform_load.py`), Pandas, SQLAlchemy | Extracts complex nested JSON data into a normalized SQLite database. |
| **Analysis** | SQL (`analysis_queries.sql`) | Executes complex aggregations, grouping, and filtering to generate 20 specific analytical results. |

| Visualization | Python (`data_visuals.py`), Matplotlib, Seaborn | Generates 10 presentation-ready visualizations based on SQL results. |
| --- | --- | --- |
| Execution | Python (`run_analysis.py`) | Orchestrates the extraction, query execution, and export of results to CSV. |

## 2.1 ETL Process (`transform_load.py`)

The ETL process handles the complex, nested JSON structure of Cricsheet data and transforms it into a relational model within a single SQLite database (`cricket_data.db`).

1. **Extraction:** JSON files are read from format-specific folders (e.g., `data/raw_json/t20`).
2. **Transformation:**
   - **Match Metadata:** The `parse_match_metadata` function extracts high-level details (winner, date, venue, toss decision) into format-specific `[format]_matches` tables. A unique `match_id` is generated for linkage.
   - **Delivery Data:** The `parse_deliveries` function processes the ball-by-ball data, including batter, bowler, runs scored (batter/extras/total), and detailed wicket information (kind, player out). This forms the massive `[format]_deliveries` tables.
3. **Loading:** Pandas DataFrames are loaded into the SQLite database using `to_sql()` with an `if_exists='replace'` strategy, ensuring a clean refresh upon every run.

# 3. Analytical Framework (SQL Queries)

The analysis is driven by 20 named SQL queries (`analysis_queries.sql`) executed directly against the SQLite database. These queries are categorized into two main areas:

## 3.1 Player Performance Metrics

These queries focus on efficiency and specialization, moving beyond career totals:

- **T20 Efficiency:** Calculating `AvgRunsPerMatch` and `HighestT20StrikeRate` (Queries 3, 8) for batsmen, filtering by a minimum number of runs/balls to ensure relevance.
- **Bowling Balance:** Determining the trade-off between `TotalWickets` and `EconomyRate` (Queries 2, 4) in various formats.
- **Power Metrics:** Quantifying the aggressive nature of modern batting by calculating `MostCareerSixes` (Query 5).

- **Influence:** Identifying the true impact players via `MostPlayerOfMatchAwards` (Query 9).

## 3.2 Team and Match Structure Insights

These queries analyze macro-level trends and strategic effectiveness:

- **T20 Impact:** Tracking `YearlyMatchTrend` (Query 18) to quantify the structural growth.
- **Strategic Value:** Quantifying the correlation between winning the toss and winning the match (`TossWinMatchWinPct`, Query 14), which found minimal influence.
- **Team Discipline:** Measuring a team's reliance on opponent mistakes through `TeamExtrasReliancePct` (Query 17).
- **Competitiveness:** Analyzing the distribution of `NarrowestVictoryByRuns` and `NarrowestVictoryByWickets` (Queries 12, 13) to highlight tight margins.

# 4. Key Findings Summary

The visualizations and analytical results presented in the Canvas support the following three major conclusions about the state of modern cricket:

## The T20 Effect

The steep rise in match volume (V1) due to T20 and franchise cricket has forced structural changes. This environment of high-risk, aggressive batting (V2) directly correlates with the dominance of the **'Caught' dismissal (V3)**, making fielding excellence and error-inducing bowling variation the most critical defensive assets.

## Efficiency Over Totals

In a high-volume era, player value is defined by immediate impact. Metrics such as **Average Runs Per Match** (Query 3) and high **Player of the Match** awards (Query 9) are better indicators of current influence than legacy career run totals (V7). For bowlers, the sweet spot is maintaining a high wicket count while simultaneously achieving an **elite economy rate**.

## Strategic Insights

The analysis provides counter-intuitive strategic insights:

- **Toss is Overrated (V4):** Winning the toss offers a negligible statistical advantage (less than 51% win rate regardless of bat/field decision).
- **Narrow Margins (V8):** The majority of matches are decided by narrow margins, confirming competitive balance and highlighting the importance of execution.
- **Discipline is Key (V9):** Teams that score a high percentage of their runs via Extras (opponent mistakes) may mask underlying batting depth issues, emphasizing that **disciplined bowling and batting consistency** are the true separators.

# 5. Conclusion

This project successfully leveraged data analytics to validate intuitive observations and challenge conventional wisdom in cricket. The future of strategic planning and player scouting must be rooted in **rate-based efficiency, specialized skill sets, and rigorous team discipline**.