



## 1 ARTICLE INFORMATION

### 2 Article title

3 NBA odds dataset linking sportsbook implied probabilities to game outcomes (2020–2024)

### 4 Authors

5 Hieu Vu

### 6 Affiliations

7 University of Virginia

### 8 Corresponding author's email address and Twitter handle

9 brr2tu@virginia.edu

### 10 Keywords

11 *sports wagering; basketball; implied probability; calibration; relational database; API integration;*  
12 *time-series snapshots*

### 13 Abstract

14 This dataset was created by collecting and integrating National Basketball Association (NBA) game  
15 results with sportsbook betting odds in order to support reproducible analyses of implied  
16 probabilities, calibration, and betting-return metrics across seasons and bookmakers. Data collection  
17 was performed programmatically in Python using two public web APIs. NBA team metadata and  
18 game-level results were retrieved from the baldontlie API, including game identifiers, UTC start  
19 timestamps, season labels, home/away teams, and final scores. Betting odds were retrieved from  
20 The Odds API using its historical odds endpoint to capture daily snapshots over a multi-year window.  
21 Each snapshot returns event metadata (including commence time and listed teams), bookmaker  
22 metadata, and market outcomes for head-to-head (moneyline) markets with American odds.

23

24 During transformation, team names from both sources were standardized via lightweight string  
25 normalization to reduce naming mismatches. American odds were converted into implied  
26 probabilities for each outcome, and the odds data were stored in a relational schema (SQLite)  
27 separating events, bookmakers, and quotes. A linking table was produced to connect odds events to  
28 completed NBA games by matching normalized home/away team names and selecting the closest  
29 event by absolute time difference within a configurable threshold, with an additional pass that  
30 considers swapped home/away orientation. The final dataset supports reuse for evaluating  
31 sportsbook probability calibration, favorite/underdog performance, bookmaker comparisons, and  
32 strategy backtesting over multiple seasons, and it can be extended with additional markets  
33 (spreads/totals) and more robust entity-resolution rules.

34



## 35 SPECIFICATIONS TABLE

|                              |  |
|------------------------------|--|
| <b>Subject</b>               | <u>Social Sciences</u>   |
| <b>Specific subject area</b> | Sports analytics and betting markets analysis using implied probabilities and game outcomes  |
| <b>Type of data</b>          | <p>Type of data</p> <ul style="list-style-type: none"><li>- Table (relational database tables)</li><li>- Chart</li><li>- Graph</li><li>- Figure</li></ul> <p>Data format</p> <ul style="list-style-type: none"><li>- Raw</li><li>- Processed</li><li>- Analyzed</li></ul>  |
| <b>Data collection</b>       | NBA game results and team metadata were collected via the balldontlie API using scripted HTTP requests in Python. Historical sportsbook odds were collected daily using The Odds API historical endpoint, capturing head-to-head markets with American odds across multiple bookmakers. Data were normalized by standardizing team names, converting timestamps to UTC, and transforming American odds into implied probabilities. Odds events were linked to completed games by matching normalized team names and minimizing time differences within a fixed threshold. Data ingestion, cleaning, and storage were implemented using Python, pandas, SQLAlchemy, and SQLite. |
| <b>Data source location</b>  | Data was collected from publicly available web APIs and stored locally and version-controlled at the University of Virginia, Charlottesville, Virginia, USA.   |
| <b>Data accessibility</b>    | <p>Repository name: NBA_SPORTS_BETTING_DB</p> <p>Data identification number: N/A</p> <p>Direct URL to data: <a href="https://github.com/CSHieuV/NBA_SPORTS_BETTING_DB">https://github.com/CSHieuV/NBA_SPORTS_BETTING_DB</a></p> <p>Instructions for accessing these data: Clone the repository using Git, then open the SQLite database (sports_pipeline.db) locally. Data tables can be queried using standard SQL tools or accessed programmatically via Python using pandas and SQLAlchemy. The repository also includes scripts for reproducing data extraction, transformation, and visualization.</p>  |



|                          |      |
|--------------------------|------|
| Related research article | None |
|--------------------------|------|

36

## 37 VALUE OF THE DATA

38

39 Why are these data valuable?

- 40 - This data provides a link between historical NBA game outcomes and sportsbook betting  
41 odds, allowing researchers to work with both realized results and implied probabilities in a  
42 single, normalized dataset. The inclusion of standardized team names, bookmaker identifiers,  
43 odds, and relational joins reduces the technical overhead typically required to combine  
44 sports results with betting-market data from multiple sources. As such, a typical user will not  
45 have to do any of the data cleaning required before using this data.

46 How can these data be reused by other researchers?

- 47 - The dataset can be reused to study probability calibration, forecasting accuracy, and  
48 uncertainty in prediction markets by comparing implied probabilities against observed  
49 outcomes. Because both raw American odds and derived implied probabilities are included,  
50 researchers may apply alternative probability conversions or statistical models without  
51 needing to recollect the underlying data.

52 What types of analyses do these data support?

- 53 - This data supports descriptive, comparative, and methodological analyses, including  
54 bookmaker-level comparisons, season-by-season trends, favorite versus underdog  
55 performance, and return-on-investment simulations. The relational structure also enables  
56 time-based filtering and aggregation for other relevant studies.

57 How can the dataset be extended or combined with other data sources?

- 58 - The data can be extended by incorporating additional betting markets, other sports leagues,  
59 or external covariates such as team statistics or player-level metrics. The existing extraction  
60 and transformation scripts provide a reusable framework for integrating new sources while  
61 preserving the original schema.

## 62 BACKGROUND

63 This dataset was compiled to support structured analysis of sports betting markets in the context of  
64 professional basketball, with a focus on linking sportsbook odds to realized game outcomes in a  
65 reproducible manner. In sports analytics, betting odds are commonly interpreted as expressions of  
66 implied probabilities, yet these data are often dispersed across proprietary platforms, inconsistent in  
67 format, and difficult to align with official game results over extended time periods.

68 Methodologically, the dataset was generated within a data engineering framework that emphasizes  
69 reproducible extraction, transformation, and loading (ETL) pipelines. Public web APIs were used to

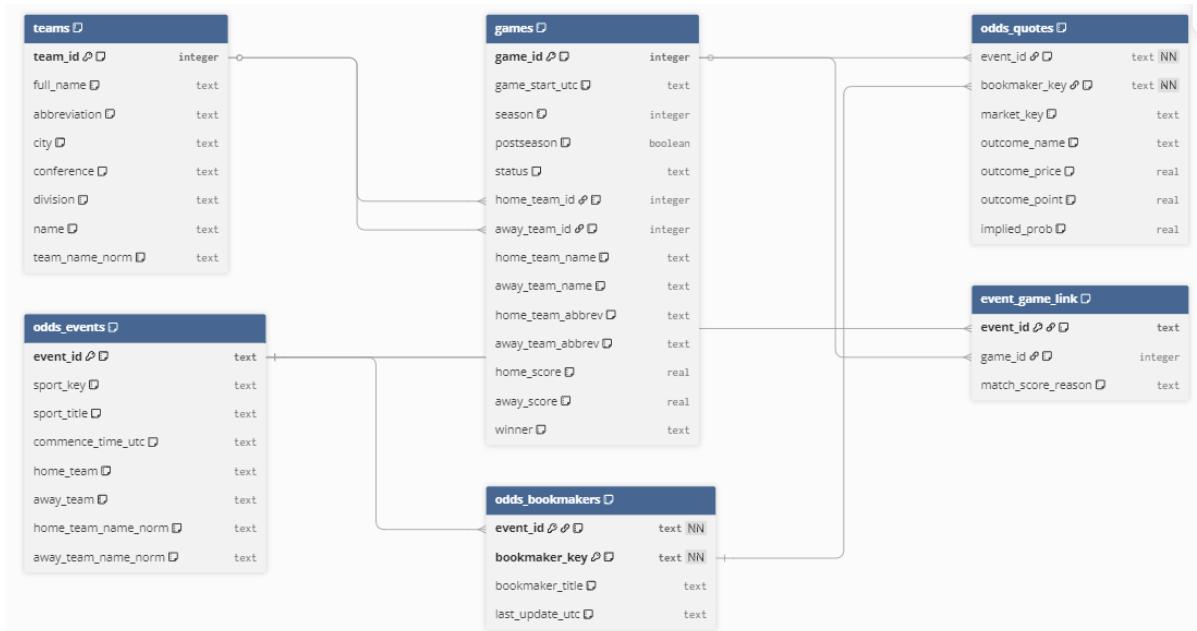
70 collect NBA game metadata and results alongside historical sportsbook odds. The data was organized  
71 into a relational schema to preserve the hierarchical structure of events, bookmakers, and market  
72 outcomes, while minimizing redundancy.

73 The dataset was developed in the context of coursework and independent research in data science  
74 and computer science, where there is a need for realistic, multi-source datasets that demonstrate  
75 challenges such as entity resolution, time alignment, and probability transformation. By providing  
76 both raw odds and derived implied probabilities, the dataset enables reuse for methodological  
77 demonstrations without requiring direct access to the original APIs.

78

## 79 DATA DESCRIPTION

80



81

82 Link to Database Documentation: <https://dbdocs.io/brr2tu/DS-6600-Data-Science-Final-Project>

83 The dataset is organized as a single SQLite relational database (`sports_pipeline.db`) that integrates  
84 NBA game results with historical sportsbook betting odds. The database schema is designed to  
85 preserve source provenance, support reproducible joins across APIs, and enable querying at the  
86 team, game, event, bookmaker, and market-outcome levels. The following tables comprise the  
87 dataset.

88

89 Teams table:

90 The teams table is a reference (dimension) table containing National Basketball Association team  
91 metadata retrieved from the BallDontLie API. Each row corresponds to one NBA team and includes a  
92 stable team identifier, full team name, abbreviation, city, conference, and division. A normalized



93 team name field is included to support joins with external data sources that use different naming  
94 conventions.

95 Games table:

96 The games table contains NBA game-level data collected from the BallDontLie API, with one row per  
97 game. Fields include the game identifier, season label, postseason indicator, game status, start time  
98 in UTC, home and away team identifiers, team names and abbreviations at ingestion time, final  
99 scores, and a computed winner label. Foreign keys link home and away teams to the teams table.

100 Odds events table:

101 The odds\_events table stores event-level metadata retrieved from The Odds API historical snapshot  
102 endpoint. Each row represents a single betting event and includes the event identifier, sport  
103 metadata, commence time in UTC, and raw home and away team names as provided by the odds  
104 source. Normalized team name fields are included to facilitate matching with NBA games.

105 Odds bookmakers table:

106 The odds\_bookmakers table contains bookmaker-specific metadata for each odds event. Rows are  
107 uniquely identified by a composite primary key consisting of the event identifier and bookmaker key.  
108 This table records the bookmaker's stable identifier, display name, and last update timestamp for the  
109 event.

110 Odds quotes table:

111 The odds\_quotes table stores individual market outcome quotes associated with each event and  
112 bookmaker. Fields include the market type (e.g., head-to-head), outcome label, American odds price,  
113 optional point values, and computed implied probabilities. This table may contain multiple rows per  
114 event and bookmaker due to repeated historical snapshots.

115 Event game link table:

116 The event\_game\_link table connects betting events from The Odds API to completed NBA games  
117 from BallDontLie. Each row maps a single odds event to at most one NBA game using normalized  
118 team names and time proximity heuristics. A descriptive field records the matching rationale.

119

## 120 EXPERIMENTAL DESIGN, MATERIALS AND METHODS

121

122 The dataset was generated using a scripted data engineering pipeline implemented in Python to  
123 reproducibly collect, normalize, and store NBA game data and sportsbook betting odds from two  
124 independent public web APIs. NBA team metadata and game-level results were retrieved from the  
125 BallDontLie API via authenticated HTTP requests, with team data collected as a reference table and  
126 game data collected for specified NBA seasons. Retrieved fields included stable team and game  
127 identifiers, season labels, postseason indicators, game status, team names and abbreviations, game  
128 start timestamps, and final scores when available. All timestamps were converted to UTC and stored  
129 as ISO 8601 strings, and a winner label was computed from final scores when both team scores were



130 present. Historical sportsbook odds were collected from The Odds API using its historical snapshot  
131 endpoint, restricted to the basketball\_nba sport key and head-to-head (moneyline) markets, with  
132 data queried on a per-date basis to capture repeated snapshots. Each API response included event-  
133 level metadata, bookmaker metadata, and outcome-level pricing data expressed as American odds,  
134 which were stored in separate relational tables to preserve source hierarchy.

135

136 Data normalization and transformation were applied uniformly across sources. Team names from  
137 both APIs were normalized, with both raw and normalized values retained. American odds were  
138 converted to implied probabilities using standard formulas without vigorish removal, and boolean  
139 values were stored using SQLite compatible integer representations. Odds events were linked to NBA  
140 games using a heuristic matching procedure based on normalized team names and time proximity,  
141 considering both original and swapped home-away orientations and selecting the closest match  
142 within a predefined time threshold; each event was linked to at most one game, with a descriptive  
143 matching rationale recorded. All data was stored in a single SQLite database (sports\_pipeline.db)  
144 using a relational schema comprising team, game, event, bookmaker, quote, and linkage tables.  
145 Database creation and population were managed using SQLAlchemy, with pandas used for  
146 intermediate data manipulation and validation. API credentials were handled via environment  
147 variables stored in .env files, and all scripts used for data acquisition, transformation, and storage are  
148 included in the associated public GitHub repository.

## 149 LIMITATIONS

150 One limitation of the dataset is that the volume and temporal coverage of sportsbook odds data  
151 were constrained by the number of API requests permitted due to budget constraints on the  
152 creator's end for The Odds API. As a result, historical odds were collected at discrete snapshot  
153 intervals rather than continuously, and some games may have fewer recorded odds updates than  
154 others. This limitation may affect the density of temporal snapshots available for certain events or  
155 seasons but does not impact the structure or integrity of the stored data. No other significant  
156 limitations were encountered during data collection or curation.

## 157 ETHICS STATEMENT

158 The authors have read and complied with the ethical requirements for publication in Data in Brief.  
159 This work does not involve human subjects, animal experiments, or data collected from social media  
160 platforms. All data were obtained from publicly accessible web APIs and do not contain personal,  
161 sensitive, or identifiable information.

## 162 CRediT AUTHOR STATEMENT

163 Hieu Vu: Conceptualization; Data curation; Software; Methodology; Validation; Formal analysis;  
164 Investigation; Resources; Writing - original draft; Writing - review & editing; Visualization;  
165 Supervision; Project administration.

166



## 167 ACKNOWLEDGEMENTS

168 There are no contributors to acknowledge who do not meet the criteria for authorship. This research  
169 did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit  
170 sectors.

171

## 172 DECLARATION OF COMPETING INTERESTS

173 The authors declare that they have no known competing financial interests or personal relationships  
174 that could have appeared to influence the work reported in this paper

## 175 REFERENCES

- 176 [1] S. D. Levitt, "Why are gambling markets organised so differently from  
177 financial markets?" *The Economic Journal*, vol. 114, no. 495, pp.  
178 223–246, 2004. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1468-0297.2004.00207.x>
- 179 [2] ——, "Why are gambling markets organised so differently from  
180 financial markets?" *University of Chicago Price Theory*, Tech. Rep.,  
181 2004, working paper. [Online]. Available: <https://pricetheory.uchicago.edu/levitt/Papers/LevittWhyAreGamblingMarkets2004.pdf>
- 182 [3] E. Snowberg and J. Wolfers, "Explaining the favorite-  
183 longshot bias: Is it risk-love or misperceptions?" *IZA*  
184 Institute of Labor Economics, Tech. Rep. DP4884,  
185 2010. [Online]. Available: <https://www.iza.org/publications/dp/4884/explaining-the-favorite-longshot-bias-is-it-risk-love-or-misperceptions>
- 186 [4] M. Ottaviani and P. Sørensen, "The favorite-longshot bias: An  
187 overview of the main explanations," *NBER*, Tech. Rep. w15923, 2010.  
188 [Online]. Available: [https://www.nber.org/system/files/working\\_papers/w15923/w15923.pdf](https://www.nber.org/system/files/working_papers/w15923/w15923.pdf)
- 189 [5] FiveThirtyEight, "Fivethirtyeight data repository," <https:////data.fivethirtyeight.com/>, 2017, cC-BY-4.0. Accessed Oct. 23,  
190 2025.
- 191 [6] ——, "Nba elo dataset," <https://www.kaggle.com/datasets/>



- 197    [fivethirtyeight/fivethirtyeight-nba-elo-dataset](https://github.com/fivethirtyeight/nba-elo-dataset), 2019, accessed Oct.  
198    23, 2025.
- 199    [7] ——, “Nfl historical scores and elo (forecasting game data),” <https://github.com/fivethirtyeight/nfl-elo-game>, 2018, accessed Oct. 23, 2025.
- 200    [8] The Odds API, “Sports odds api,” <https://the-odds-api.com/>, 2025, JSON  
201    feed for bookmaker odds. Accessed Oct. 23, 2025.
- 202    [9] ——, “Odds api documentation (v4),” <https://the-odds-api.com/liveapi/guides/v4/>, 2025, endpoints for moneyline, spreads, totals. Accessed Oct.  
203    23, 2025.
- 204    [10] Odds API IO, “Odds api: Real-time sports betting odds (250+ book-  
205    makers),” <https://odds-api.io/>, 2025, accessed Oct. 23, 2025.
- 206    [11] football-data.org, “Api reference,” <https://www.football-data.org/documentation/api>, 2025, rEST endpoints for fixtures and results.  
207    Accessed Oct. 23, 2025.
- 208    [12] ——, “Api quickstart (v4),” <https://www.football-data.org/documentation/quickstart>, 2022, accessed Oct. 23, 2025.
- 209    [13] BALLDONTLIE, “Balldontlie sports apis,” <https://balldontlie.io/>, 2025,  
210    multi-league coverage. Accessed Oct. 23, 2025.
- 211    [14] ——, “Nba api (documentation),” <https://nba.balldontlie.io/>, 2025, his-  
212    torical and current NBA data. Accessed Oct. 23, 2025.
- 213    [15] CollegeFootballData.com, “College football data api,” <https://api.collegefootballdata.com/>, 2025, games, play-by-play, advanced stats.  
214    Accessed Oct. 23, 2025.
- 215    [16] ——, “College football data api: Documentation,” <https://www.postman.com/api-evangelist/college-football-data/documentation/el6ukhf/college-football-data-api>, 2025, postman docs. Accessed Oct.  
216    23, 2025.
- 217    [17] S. Gilani and SportsDataVerse, “cfbfastr: College football data in  
218    r (sportsdataverse),” <https://cfbfastr.sportsdataverse.org/reference/index>.

- 226 html, 2025, open-source wrappers around CollegeFootballData API.
- 227 Accessed Oct. 23, 2025.
- 228 [18] H. Vu, “NBA sports betting database,” GitHub repository, 2025.
- 229 [Online]. Available: [https://github.com/CSHieuV/NBA\\_SPORTS\\_BETTING\\_DB](https://github.com/CSHieuV/NBA_SPORTS_BETTING_DB).
- 230 Accessed Oct. 23, 2025]