

Analytics in the Octagon: Modeling UFC Fight Prediction

Group 4 Project Proposal Review/Progress Report

Evan Embry

CSPB 4502

University of Colorado-Boulder

Boulder, CO , USA

evem6983@colorado.edu

ABSTRACT

Mixed Martial Arts (MMA), specifically within the Ultimate Fighting Championship (UFC) organization, has gained massive global popularity, making it a fascinating domain for sports analytics. This project will use historical UFC fight data to develop and compare machine learning models that predict bout outcomes. By examining fighter attributes (e.g., height, reach, stance, experience) and contextual factors (such as whether the fight is a title bout or not), I aim to identify the most significant predictors of victory and compare our models against real-world metrics like bookmaker odds.

1. PROBLEM STATEMENT/ MOTIVATION

Sports analytics has traditionally been dominated by team sports like baseball and basketball, but MMA presents a unique context for data-driven exploration. Unlike team-based competitions, MMA focuses on individual matchups, offering a more direct way to study performance factors. My core objective is to determine if machine learning methods can accurately forecast the winner of a UFC fight, considering both fighter-specific attributes and event-level context (e.g., venue, date, title match). I believe that identifying these key predictors can offer practical insights for coaches, fighters, fans, and even sports bettors.

2. LITERATURE SURVEY

Prior research on MMA has ranged from descriptive overviews of fighter statistics to advanced modeling efforts. Some Kaggle contributors and academic authors have shown that tree-based models like

Random Forest or XGBoost can capture complex interactions among fighting styles, physical attributes, and performance metrics. Other work has highlighted the influence of grappling skills and stance on match outcomes, suggesting that these factors significantly shape fight dynamics. Meanwhile, several sports economics studies have examined how accurately betting odds reflect actual probabilities. It appears that while bookmakers generally offer efficient markets, certain data-driven strategies can still outperform the implied probabilities in niche sports. My project distinguishes itself by integrating multiple derived features (e.g., reach difference, height difference, experience gap) with contextual details, then directly comparing our model outputs to betting odds to gauge real-world utility.

3. PROPOSED WORK

3.1 Data Collection

I will use the publicly available “Ultimate UFC Dataset” on Kaggle, provided by “mdabbert.” This dataset covers thousands of UFC fights with numerous features, including fighter profiles, fight outcomes, and odds data.

3.2 Preprocessing

I plan to clean and prepare the data by:

- Imputing missing physical stats (height, reach) where possible, often grouped by light class.
- Removing or consolidating duplicates.
- Creating features that highlight differences between Red and Blue corners (e.g., reach difference, height difference) and encoding categorical variables like stance.

- Addressing potential class imbalance by using techniques such as SMOTE or adjusting model class weights if one corner systematically wins more often.

3.3 Feature Engineering

While raw data on strikes, takedowns, and other stats are helpful, I will also explore derived metrics. For example, I may add variables related to “days since last fight,” “average fight duration,” or “win/loss streaks,” which could shed light on momentum and ring readiness.

3.4 Modeling

I plan to train multiple classification models:

1. **Logistic Regression** for a clear, interpretable baseline.
2. **Random Forest** for strong performance on tabular data and to access feature importance insights.
3. **XGBoost** for advanced boosting capabilities, often recognized for high performance in Kaggle competitions.

These results will be compared to naive baselines (like always picking the bookmaker favorite) and to each other. I will tune hyperparameters to achieve the best possible predictive accuracy.

3.5 Evaluation

I intend to measure Accuracy, Precision, Recall, F1 Score, and AUC-ROC to gauge predictive performance. Since betting odds serve as a practical benchmark, I will convert odds into implied probabilities and compare them to our model outputs, potentially examining calibration metrics such as log loss or Brier Score.

4. DATA SET

I have confirmed access to the “Ultimate UFC Dataset” from Kaggle:

<https://www.kaggle.com/datasets/mdabbert/ultimate-ufc-dataset>

This resource offers:

- **Fighter Attributes:** Name, age, height, reach, stance, and so on.
- **Fight-Specific Info:** Date, location, title-bout status, and weight class.

- **Outcome Details:** Winner, method, round, and betting odds.

Any inconsistencies or missing values will be carefully addressed before modeling.

5. EVALUATION METHODS

5.1 Metrics

To ensure robust model comparisons, I will examine:

- **Accuracy:** The proportion of correct predictions across all fights.
- **Precision & Recall:** Especially relevant if I discover an imbalance where one corner wins significantly more.
- **F1 Score:** Offers a single measure that balances both precision and recall.
- **AUC-ROC:** Reflects how well the model ranks potential winners across different classification thresholds.

5.2 Comparison to Odds

Sports bettors and oddsmakers base their wagers on implied probabilities derived from betting lines. I will similarly translate the odds in our dataset to implied probabilities, then compare these to our model’s predicted probabilities. This step helps us assess real-world viability and may reveal whether a data-driven approach can outperform market expectations in certain scenarios.

6. TOOLS

Our primary tools and libraries will be:

- **Python (with libraries like pandas and NumPy)** for data cleaning and manipulation.
- **scikit-learn and XGBoost** for building and tuning our classification models.
- **Jupyter Notebook** for exploratory analysis and organized experimentation.
- **GitHub** for version control, enabling the team to collaborate effectively and keep track of all progress.

7. MILESTONES

- **Weeks 1–2:** Complete an initial data audit, clean the dataset, and perform exploratory data analysis (EDA).
- **Weeks 3–4:** Implement and evaluate baseline models (Logistic Regression, Random Forest), handle any class imbalance strategies.
- **Week 5:** Introduce XGBoost, tune hyperparameters, and compare predictive outcomes to implied betting probabilities.
- **Week 6:** Finalize all models, generate visual aids (e.g., ROC curves, feature importance plots), and assemble the final report.

MILESTONES COMPLETED

- **Data Collection & Cleaning:** I successfully gathered and preprocessed the Ultimate UFC Dataset from Kaggle. Missing values for fighter attributes were handled using imputation techniques, and new features such as reach difference and height difference were engineered.
- **Exploratory Data Analysis (EDA):** Initial visualizations confirmed certain expected correlations (e.g., height and reach advantages) and highlighted class imbalance in fight outcomes.
- **Baseline Model Implementation:** Logistic Regression and Random Forest models were trained, establishing benchmarks for prediction accuracy.
- **Initial Model Evaluation:** I computed Accuracy, Precision, Recall, F1 Score, and AUC-ROC for these baseline models, confirming that they outperform random guessing but require refinement.

MILESTONES TO DO

- **Advanced Model Development:** I need to finalize the implementation of XGBoost, tune hyperparameters, and refine feature selection.

- **Comparison to Betting Odds:** I plan to assess how well our model predictions align with implied probabilities from betting lines.
- **Further Feature Engineering:** Additional variables such as 'days since last fight' and 'win/loss streaks' will be integrated for potential performance improvements.
- **Final Evaluation & Report Compilation:** Once modeling is complete, I will generate visual aids (e.g., ROC curves, feature importance plots) and analyze findings for the final project submission.

RESULTS SO FAR

- **Performance Metrics:** Initial models show moderate predictive accuracy, with Random Forest outperforming Logistic Regression.
- **Early Correlations:** Certain fighter attributes, such as reach and recent fight frequency, appear to have strong predictive power.
- **Data Challenges:** Some inconsistencies in betting odds data require additional validation before final integration.

ACKNOWLEDGMENTS

I sincerely appreciate the Kaggle user, “mdabbert,” for providing and maintaining the UFC dataset, as well as public UFC data resources that offer additional references on fighter records and historical events.

REFERENCES

- [1] K. Williams. 2021. “Predicting MMA Fight Outcomes Using Machine Learning.” *Proceedings of 2021 Sports Analytics Conference*.
- [2] R. Hernandez. 2020. “Exploring Reach, Stance, and Takedown Defense in UFC Fighters.” *International Journal of Combat Sports Analytics*.
- [3] A. Smith, B. Johnson, and C. Lee. 2019. “Betting Markets and Probabilistic Forecasts:

Evidence from the UFC.” *Journal of Sports Economics*.

[4] David Kosiur. 2001. *Understanding Policy-Based Networking* (2nd. ed.). Wiley, New York, NY.