

PUBG Data Royale

CSCI 4502 Final Project

Matthew Davenport

Federico Aragon

University of Colorado
Boulder, Colorado, USA

Thomas Alder

Haotian Dong

PROBLEM STATEMENT/MOTIVATION

In this project we are looking to analyze data from the results of a player's matches they played in a Battle Royale video game named PlayerUnknown's Battlegrounds. Using this dataset, we hope to answer questions regarding the game's mechanisms through an analysis of player damage and deaths and how this statistic could be improved. We would also like to see what the likelihood of a player being a hacker is through understanding the inflicted damage they cause and whether there are significant outliers which may suggest an unfair advantage. Lastly, we would like to find the best strategy to use when playing the game in terms of determining the most dangerous places to avoid, places which may be beneficial, and the general strategy of players that tend to win more.

LITERATURE SURVEY

1 Final Circle Heatmap

Kevin Pei. 2018. Final Circle Heatmap: PUBG Match Deaths and Statistics.

<https://www.kaggle.com/skihikingkevin/final-circle-heatmap>

This previous work explored the "final circle" of games played in PlayerUnknown Battleground's to visualize patterns in the ultimate positions of the ever-decreasing circle of survival. Throughout a match in PUBG, players are restricted to a circled playing area which continuously has periods in which it decreases in circumference and times

where it stagnates. Using python and Jupyter notebooks, Kevin Pei uses the deaths of players who placed first and second place who did not die outside of the circle to determine the location of each final circle. He then plots the results on a heatmap which shows the distribution of final circle locations demonstrating areas of congestion in which the final circle ended up in a concentrated area. Visualizing both maps shows that the edges of the map never see the final circle at their location while numerous pockets of high density are dotted in areas of more diverse end locations.

2 Survival Analysis Using Kaplan-Meier Estimator

Robert Kwiatowski. 2018. PUBG - Survival Analysis (Kaplan-Meier): PUBG Match Deaths and Statistics.

<https://www.kaggle.com/datark1/pubg-survival-analysis-kaplan-meier>

This previous work shows the damage and die rate data for each player outside the play zone before PC 1.0 Patch #15 of PlayerUnknown Battleground's to visualize the survival situations. The data set also has the data of the party situation including SOLO, DUO, and SQUAD, the use of vehicle for each player, and the numbers of players for each game. He plots the results for different groups with those conditions to illustrate the survival rate for players. Then, Kwiatowski uses a Kaplan-Meier estimator to use non-parametric analysis using the number of events that occur

during a time frame and the survived units from the previous time frame. Each time frame represents the amount of time the circle is a stagnant play zone. Using this, he creates three plots which shows the survival probability based on the team members a player has which are solo (no team members), duo (one team member), and a squad (2-3 team members). Here it is seen that the highest chance of survival is for squads. The same process is done for two different scenarios in which a player uses a vehicle versus when they do not which yields results that show using a vehicle has a greater probability of survival. A final stratification is done based on the number of players in a match between small (< 33 players), medium (33-66 players), and large (> 66 players). Here the results show that the smaller number of players creates the larger survival chances.

3 Data Visualizations

Peter Ott. 2018. PUBG - You're not as bad as you think: PUBG Match Deaths and Statistics. Various data visualizations.

<https://www.kaggle.com/ottpeterr/pubg-you-re-not-as-bad-as-you-think>

In this previous work, Peter Ott creates a variation of data visualizations mainly histograms which plots different relationships which can be found in the data sets. First, he plots the relationship between the rank of team placement and player kills before splitting the dataset into quartiles based on the final ranking of teams following a match. Here he finds that making it into the top five isn't a task that requires a lot of kills which many might think with a "King of the Hill" style of game. Through his next visualizations, the same theme occurs. He finds that a lot of kills is necessary to be a winner otherwise known as being part of the "Chicken Dinner Club" with most winners falling into the three to six kill range.

PROPOSED WORK

1 Data Cleaning

- Shorter games in the dataset which have missing data must be removed.
- Longer games that have missing values will be dropped as this is minimal throughout the dataset.
- If the number of players or teams in a game is too small, then remove this data as well as it might skew the information we want to gain.

2 Data Preprocessing

- Reducing the number of attributes through correlation analysis.
- Removing columns which have no association with match statistics like player_name and team_id.
- Transforming data when necessary, for example if percentages are stored as integers or floats greater than 1, then convert them to their corresponding decimal values.
- Normalization of attributes.
- Removal of games that are not played on the map of Erangel seen in figure 1.

3 Data Integration

- Merging all the separate csv files for the deaths and aggregates before an integration of the final deaths and aggregates all of which are provided in the PUBG Match Deaths and Statistics Kaggle page.

4 Data Warehousing

- Storing the historical data with a Python Cubes data warehouse which then can be queried for the data mining process.

5 Data Mining

- Analyze the data to gain useful knowledge and answer our initial questions.
- Determine significant attributes which contribute best to a winning strategy.

6 Data Visualization

- Presenting the information found in a manner that is understandable by players of the game and non-players alike.



Figure 1: Map of Erangel which will be used with geographic player locations to produce a heatmap of player deaths.

Our proposed work looks to be different than the previous literature survey. While these previous works have created data visualizations and used Kaplan-Meier to show trends and estimate survival respectively, we hope to make more fleshed-out data visualizations which incorporate the Erangel map more as well as creating our own survival prediction score using the attributes provided in the dataset.

DATA SET

For this project, we are going to be using the dataset PUBG Match Deaths and Statistics which can be found on Kaggle¹. This dataset was supplied by Kaggle user skihikingkevin or Kevin Pei who

compiled the results of over 720,000 matches extracting the data from pubg.op.gg which tracks games of PUBG. Currently, the dataset is on Thomas's and Matthew's machine.

EVALUATION METHODS

Ideally with the correct results we would be able to predict the placements of players given certain attributes about a match like their total damage given or taken, time spent walking and driving, or the length of game. We would also like to predict match results given certain scenarios and environmental factors within a game such as landing locations, total number of teams, equipment picked up, etc. These factors can be displayed on heatmap data visualizations which portray "hot spots" around the two different maps of the game, scores can be given to different areas regarding their "dangerous" levels.

TOOLS

Throughout this project we will be using a variety of tools much of which will revolve around Python as we will use it for our main programming language. With Python we can clean and integrate the data, create data visualizations using libraries like Numpy, Pandas, and SciPy for machine learning as well as using Python Cubes to create a data warehouse for our historical data with the potential to use it for further cleaning and integration. For project reports and documentation, we will be using Google Drive and Docs and for the centralized project repository and milestones we will be using GitHub.

MILESTONES

By the progress report, we plan to have our data set cleaned, integrated, and hosted on a Python Cubes OLAP server. Once we have accomplished that, we can mine the model and data warehouse created with Python cubes using drill-down and aggregation techniques. By the final report, we will have the final project completed with the data mined for the interesting questions, correlations, and analysis we are looking to find as well as creating data visualizations which we hope will

¹ <https://www.kaggle.com/skihikingkevin/pubg-match-deaths/home>

provide a more statistical insight of the game, its mechanics and scenarios in which players can learn or understand how to better play PUBG.