# Bracketology

Bracketology is the process of predicting the field of college basketball participants in the NCAA Basketball Tournament, named as such because it is commonly used to fill in tournament brackets for the postseason. It incorporates some method of predicting what the NCAA Selection Committee will use as its Ratings Percentage Index in order to determine at-large (non-conference winning) teams to complete the field of 64 teams, and, to seed the field by ranking all teams from first through sixty-eighth. Bracketology also encompasses the process of predicting the winners of each of the brackets. In recent years the concept of bracketology has been applied to areas other than basketball.

Various methods are used to predict the winners in a bracket. While some use math and statistics, others make selections based on team mascots or colors. President Barack Obama became famous for his bracket predictions. After entering office, he presented his projected winners annually on ESPN in a segment called Barack-etology.

Bracketology as a discipline has spread beyond a focus on basketball, into other sports, as well as pop culture, history, nature, and other topics where a loose application of binary opposition may be profitable for study or enjoyment, albeit without the label of "bracketology" itself.

## Below are the steps of the Bracketology for NCAA Men's Basketball.

**Part 1:** Create a machine learning model to predict the winner based on seed and team name.

Choosing a model type

We have chosen Logistic Regression. We need a model that generates a probability for each possible discrete label value, which in our case is either a 'win' or a 'loss'. Logistic regression is a good model type to start with for this purpose. The good news is that the ML model will do all the math and optimization during model training.
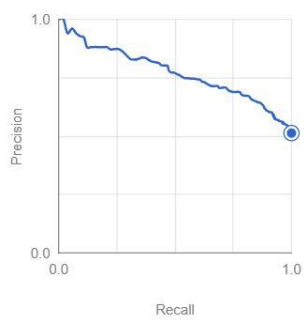
**Step 2**: -

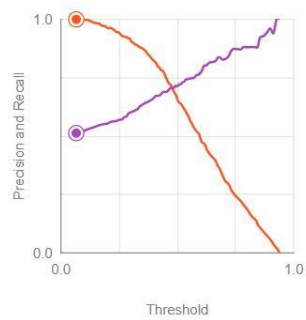View model training details and stats.

# Learn rate
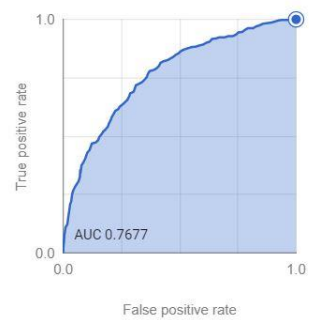


## Precision-Recall curve



## Precision and Recall vs Threshold



## ROC curve



AUC 0.7677

## ncaa_model

Details   **Training**   Evaluation   Schema

View as ⚪ Graphs ⬤ Table

| Iteration | Training Data Loss | Evaluation Data Loss | Learn Rate | Duration (seconds) |
|---|---|---|---|---|
| 3 | 0.5044 | 0.5778 | 1.6000 | 4.38 |
| 2 | 0.5508 | 0.5927 | 0.8000 | 4.27 |
| 1 | 0.6090 | 0.6297 | 0.4000 | 5.45 |
| 0 | 0.6595 | 0.6670 | 0.2000 | 6.03 |

**Step 3**: -

See what the model learned about our features.

To learn this, we can execute the following query in bigquery.

SELECT category, weight FROM UNNEST(( SELECT category_weights FROM ML.WEIGHTS(MODEL `bracketology.ncaa_model`) WHERE processed_input = 'seed')) # try other features like 'school_ncaa' ORDER BY weight DESC

Query complete (0.8 sec elapsed, 24.5 KB processed)

Job information   **Results**   JSON   Execution details

| Row | category | weight |
|---|---|---|
| 1 | 01 | 0.584699252031549 |
| 2 | 02 | 0.3812926215523544 |
| 3 | 03 | 0.25577817883745174 |
| 4 | 04 | 0.08085337884663339 |
| 5 | 06 | 0.04796437984873807 |
| 6 | 05 | -0.0053866903047856875 |
| 7 | 07 | -0.01976989150187833 |
| 8 | 08 | -0.1518538285462349 |
| 9 | 11 | -0.17889419395575332 |
| 10 | 10 | -0.20913129101187983 |

As we can see, if the seed of a team is very low (1,2,3) or very high (14,15,16) the model gives it a significant weight (max is 1.0) in determining the win loss outcome. Intuitively this makes sense as we expect very low seed teams to perform well in the tournament. The real magic of machine learning is that we didn't create a ton of hardcoded IF THEN statements in SQL telling the model IF the seed is 1 THEN give the team a 80% more chance of winning. Machine learning does away with hardcoded rules and logic and learns these relationships for itself.

**Step 4**: -

Evaluate model performance

| Row | precision | recall | accuracy | f1_score | log_loss | roc_auc |
|-----|-----------|--------|----------|----------|----------|---------|
| 1 | 0.71717171717171 | 0.662004662004662 | 0.6925837320574163 | 0.6884848484848485 | 0.5777782337963289 | 0.767683 |

The value will be around 69% accurate. While it's better than a coin flip, there is room for improvement.

**Step 5**: -

Making predictions

We can make predictions is as simple as calling ML.PREDICT on a trained model and passing through the dataset we want to predict on.

```
1  CREATE OR REPLACE TABLE `Bracketology.predictions` AS (
2
3  SELECT * FROM ML.PREDICT(MODEL `Bracketology.ncaa_model`,
4
5  # predicting for 2018 tournament games (2017 season)
6  (SELECT * FROM `data-to-insights.ncaa.2018_tournament_results`)
7  )
8  )
```

Valid.

Query results

Query complete (1.7 sec elapsed, 42.8 KB processed)

This statement created a new table named windy-lyceum-233722:Bracketology.predictions.

Job information    **Results**    JSON    Execution details

| Row | predicted_label | confidence | correct_label | game_date | round | seed | school_ncaa | points | opponent_seed | opponent_school_ncaa | opponent_points |
|-----|------|-------|------|------------|----|----|-------------|----|----|--------------|----|
| 1 | win | 0.864736635598085 | loss | 2018-03-15 | 64 | 04 | Arizona | 68 | 13 | Buffalo | 89 |
| 2 | win | 0.8981839119836009 | loss | 2018-03-16 | 64 | 01 | Virginia | 54 | 16 | UMBC | 74 |
| 3 | win | 0.8528053004359512 | loss | 2018-03-16 | 64 | 04 | Wichita St. | 75 | 13 | Marshall | 81 |
| 4 | loss | 0.8895132514511246 | win | 2018-03-15 | 64 | 13 | Buffalo | 89 | 04 | Arizona | 68 |
| 5 | loss | 0.874429817399404 | win | 2018-03-16 | 64 | 16 | UMBC | 74 | 01 | Virginia | 54 |
| 6 | loss | 0.8701799722994595 | win | 2018-03-16 | 64 | 13 | Marshall | 81 | 04 | Wichita St. | 75 |
| 7 | loss | 0.8058240185589365 | win | 2018-03-18 | 32 | 05 | Clemson | 84 | 04 | Auburn | 53 |

**Step 6**: -

Now, we can check how our model performed on predictions.

How many did our model get right for the 2018 NCAA tournament?

SELECT * FROM `bracketology.predictions`

WHERE predicted_label <> label

Out of 134 predictions (67 March tournament games), our model got it wrong 38 times.

70% overall for the 2018 tournament matchup.

**Part 2:** Using skillful ML model features

**Step 1**: -

Create a new ML dataset with these skillful features. Below are some of the features

season, win, win_seed, , win_school_ncaa, lose_seed, lose_school_ncaa

Now, we will preview the new features

training_new_features

<span style="float:right">Q QUERY TABLE    COPY TABLE    DELETE TABLE    EXPORT ▼</span>

Schema   Details   Preview

| Row | season | label | seed | school_ncaa | pace_rank | poss_40min | pace_rating | efficiency_rank | pts_100poss | efficiency_rating | opponent_seed | opponent_school_ncaa | opp_pace_rank | opp_poss_ |
|-----|--------|-------|------|-------------|-----------|------------|-------------|-----------------|-------------|-------------------|---------------|----------------------|---------------|-----------|
| 1 | 2015 | win | 01 | Duke | 188.0 | 68.714 | 43.785 | 13.0 | 23.354 | 97.032 | 05 | Utah | 268.0 | 6 |
| 2 | 2015 | win | 01 | Duke | 188.0 | 68.714 | 43.785 | 13.0 | 23.354 | 97.032 | 02 | Gonzaga | 227.0 | 6 |
| 3 | 2015 | win | 01 | Duke | 188.0 | 68.714 | 43.785 | 13.0 | 23.354 | 97.032 | 07 | Michigan St. | 247.0 | 6 |
| 4 | 2015 | win | 01 | Duke | 188.0 | 68.714 | 43.785 | 13.0 | 23.354 | 97.032 | 08 | San Diego St. | 327.0 | 6 |
| 5 | 2015 | win | 01 | Duke | 188.0 | 68.714 | 43.785 | 13.0 | 23.354 | 97.032 | 16 | Robert Morris | 166.0 | 6 |
| 6 | 2015 | win | 01 | Duke | 188.0 | 68.714 | 43.785 | 13.0 | 23.354 | 97.032 | 01 | Wisconsin | 342.0 | 6 |
| 7 | 2014 | loss | 01 | Wichita St. | 282.0 | 62.907 | 21.29 | 16.0 | 22.5 | 96.478 | 08 | Kentucky | 256.0 | 6 |
| 8 | 2014 | win | 01 | Wichita St. | 282.0 | 62.907 | 21.29 | 16.0 | 22.5 | 96.478 | 16 | Cal Poly | 346.0 | 5 |
| 9 | 2015 | loss | 01 | Kentucky | 214.0 | 68.338 | 38.396 | 9.0 | 25.902 | 98.174 | 01 | Wisconsin | 342.0 | 6 |
| 10 | 2015 | win | 01 | Kentucky | 214.0 | 68.338 | 38.396 | 9.0 | 25.902 | 98.174 | 08 | Cincinnati | 311.0 | 6 |
| 11 | 2015 | win | 01 | Kentucky | 214.0 | 68.338 | 38.396 | 9.0 | 25.902 | 98.174 | 03 | Notre Dame | 326.0 | 6 |

**Step 2**: -

Train the new model.

Query editor

<span style="float:right">⊡ HIDE EDITOR</span>

```
1  CREATE OR REPLACE MODEL
2    `Bracketology.ncaa_model_updated`
3  OPTIONS
4    ( model_type='logistic_reg') AS
5
6  SELECT
7    # this time, dont train the model on school name or seed
8    season,
9    label,
10
```

✓ Valid.

▶ Run ▾   ⤓ Save query   ⦂⦂⦂ Save view   ⊙ Schedule query ▾   ⚙ More ▾      This query will process 101.4 KB (ML) when run. ✓

Query results

Query complete (37.0 sec elapsed, 0 B (ML) processed)

Job information   **Results**   JSON

✓ This statement created a new model named windy-lyceum-233722:Bracketology.ncaa_model_updated.      Go to model

**Step 3**: -

Evaluate the new model's performance.

SELECT * FROM ML.EVALUATE(MODEL `bracketology.ncaa_model_updated`)

**Step 4**: -

Inspect what the model learned.



**Step 5**: -

Now, we will inspect what the model has learned. Which features does the model weigh the most in win / loss outcome? Here, We've taken the absolute value of the weights in our ordering so the most impactful (for a win or a loss) are listed first.

**Step 6**: -

Now, we will make prediction. For that we will use 2018 season.

**Step 7**: -

Prediction analysis: -

**Step 8**: -

Where were the upsets in March 2018?



The major upset was the same which we have found in the previous model : UMBC vs Virginia.

**Step 9**: -

Comparing model performance.



The model predicted a Florida St. (09) upset of Xavier (01) and they did.

The upset was correctly predicted by the new model (even when the seed ranking said otherwise) based on new skillful features like pace and shooting efficiency.

**Step 10**: -

Predicting for the 2019 March Madness tournament.



**Step 11**: -

Create a matrix of all possible games.

Since we don't know which teams will play each other as the tournament progresses, we'll simply have them all face each other.

In SQL, an easy way to have a single team play every other team in a table is with a CROSS JOIN.

Add in 2018 team stats (pace, efficiency).



Prepare 2019 data for prediction.

Make predictions



Here we filtered the model results to see all of Duke's possible games and we can see that model has predicted that Virginia has the highest chance of winning and that is what happened. Virginia won 2019 NCAA Basketball men's final.