# March Madness Basketball Tournament



357 Teams in Division 1

20-30 games in the regular season

68 teems seeded for NCAA March Madness
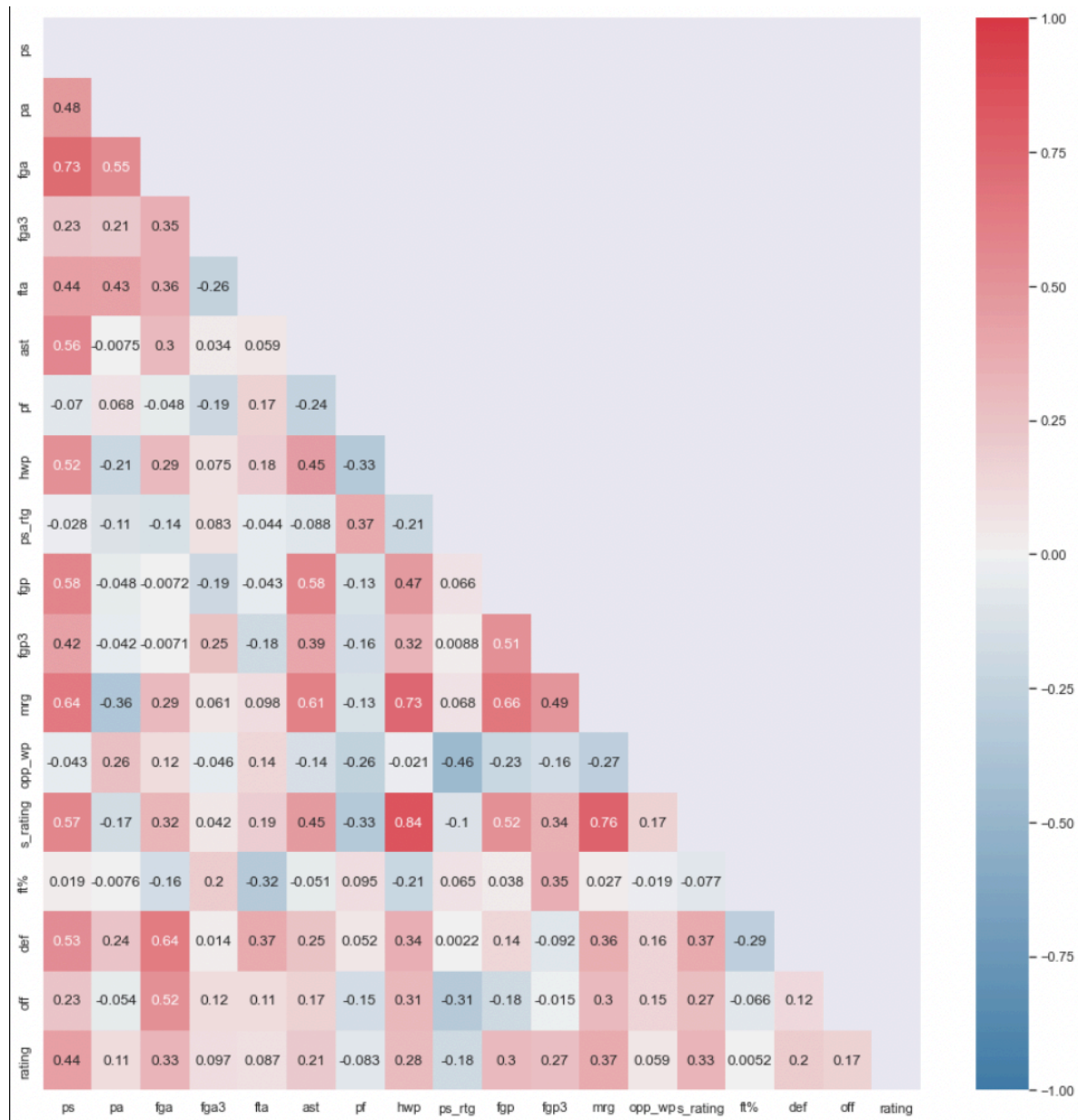
Q : Given regular season data, who is going to win NCAA?

.

**FIRST FOUR**

16 Norfolk St. (16-7) 54
16 App St. (17-11) 53  **W**

11 Wichita St. (16-5) 53
11 Drake (25-4) 53  **W**

MARCH 18

12 Mt. St. Mary's (12-10) 16
Texas So. (16-8) 16  **E**

Michigan St. (15-12) 11
UCLA (17-9) 11  **E**

**WEST**

1 Gonzaga (26-0) 98
16 Norfolk St. 55
— 1 Gonzaga 87
8 Oklahoma (15-10) 72
9 Missouri (16-9) 68
— 8 Oklahoma 71
— 1 Gonzaga 87
5 Creighton (20-8) 63
12 UCSB (22-4) 62
— 5 Creighton 72
4 Virginia (18-6) 58
13 Ohio (16-7) 62
— 13 Ohio 58
— 5 Creighton 65
— 1 Gonzaga 85
6 USC (22-7) 72
11 Drake 56
— 6 USC 85
3 Kansas (20-8) 93
14 Eastern Wash. (16-7) 84
— 3 Kansas 51
— 6 USC 82
7 Oregon (20-6) 95
10 VCU (19-7)
— 7 Oregon 95
2 Iowa (21-8) 86
15 Grand Canyon (17-6) 74
— 2 Iowa 80
— 7 Oregon 68
— 6 USC 68
— 1 Gonzaga 93

**EAST**

1 Michigan (20-4) 82
16 Texas Southern 66
— 1 Michigan 86
8 LSU (18-9) 76
9 St. Bonaventure (16-4) 61
— 8 LSU 78
— 1 Michigan 76
5 Colorado (22-8) 96
12 Georgetown (13-12) 73
— 5 Colorado 53
4 Florida St. (16-6) 64
13 UNC Greensboro (21-8) 54
— 4 Florida St. 71
— 4 Florida St. 58
— 1 Michigan 49
6 BYU (20-6)
11 UCLA 73
— 11 UCLA 67
3 Texas (19-7) 52
— 11 UCLA 80
— 11 UCLA 90

**SOUTH**

Baylor (22-2) 1 79
Hartford (15-8) 16 55
— Baylor 1 78
North Carolina (18-10) 8 62
Wisconsin (17-12) 9 85
— Wisconsin 9 63
— Baylor 1 62
Villanova (16-6) 5 73
Winthrop (23-1) 12 63
— Villanova 5 84
Purdue (18-9) 4 69
North Texas (17-9) 13 78
— North Texas 13 81
— Villanova 5 51
— Baylor 1 81
Texas Tech (17-10) 6 65
Utah St. (20-8) 11 53
— Texas Tech 6 66
Arkansas (22-6) 3 85
Colgate (14-1) 14 68
— Arkansas 3 68
— Arkansas 3 72
Florida (14-9) 7 75
Virginia Tech (15-6) 10 70
— Florida 7 75
Ohio St. (21-9) 2 72
Oral Roberts (16-10) 15 75
— Oral Roberts 15 81
— ORU 15 70
— Arkansas 3 72
— Baylor 1 78

**MIDWEST**

Illinois (23-6) 1 78
Drexel (12-7) 16 49
— Illinois 1 58
Loyola Chicago (24-4) 8 71
Georgia Tech (17-8) 9 60
— Loyola Chi. 8 71
— Loyola Chi. 8 58
Tennessee (18-8) 5 56
Oregon St. (17-12) 12 70
— Oregon St. 12 80
Oklahoma St. (20-8) 4 69
Liberty (23-5) 13 60
— Oklahoma St. 4 70
— Oregon St. 12 65
— Oregon St. 12 81
San Diego St. (23-4) 6 62
Syracuse (16-9) 11 78
— Syracuse 11 78
West Virginia (18-9) 3 84
— Syracuse 11 75
— Syracuse 11 46
— Houston 2 59

**FINAL FOUR** — Indianapolis

**NATIONAL CHAMPIONSHIP**
APRIL 5

SEMIFINALS
1 Gonzaga 70

SEMIFINALS
Baylor 1 86

**Baylor**

#MARCHMADNESS
Watch the tournament on these networks
or online at NCAA.COM/MARCHMADNESS

# Challanges to start with

- We had structured, but vast data (20 CSV files)

- Each team plays with a different set of teams.

- Post season tournements are indicative but small in sample size

- How to account for recency and fixture difficulty

# Challanges to start with

| | Win% | PS Rtg | Opponent % | Rating |
|---|---|---|---|---|
| 0 | 82.608696 | 3 | 39.870094 | 22.478789 |
| 1 | 20.000000 | 0 | 52.375424 | -27.624576 |
| 2 | 61.904762 | 1 | 49.282581 | 11.187343 |
| 3 | 80.000000 | 4 | 56.272717 | 36.272717 |
| 4 | 40.000000 | 1 | 43.971455 | -16.028545 |
| 5 | 22.222222 | 1 | 43.768685 | -34.009093 |
| 6 | 43.750000 | 1 | 49.216030 | -7.033970 |
| 7 | 31.578947 | 0 | 49.454457 | -18.966596 |
| 8 | 40.000000 | 0 | 57.754902 | -2.245098 |
| 9 | 54.166667 | 5 | 49.165630 | 3.332296 |

# Regression with almost all variables

```
==================================================================================
Dep. Variable:                    rating   R-squared:                       0.320
Model:                               OLS   Adj. R-squared:                  0.102
Method:                    Least Squares   F-statistic:                     1.470
Date:                   Fri, 25 Feb 2022   Prob (F-statistic):              0.149
Time:                           02:36:52   Log-Likelihood:                 -307.17
No. Observations:                     67   AIC:                             648.3
Df Residuals:                         50   BIC:                             685.8
Df Model:                             16
Covariance Type:               nonrobust
==================================================================================
                 coef    std err          t      P>|t|      [0.025      0.975]
----------------------------------------------------------------------------------
const       1237.7168   1259.060      0.983      0.330   -1291.179    3766.613
ps            14.0759     10.914      1.290      0.203      -7.845      35.997
pa             5.5790      5.529      1.009      0.318      -5.527      16.685
fga          -13.9223     15.485     -0.899      0.373     -45.025      17.180
fga3          -7.0952      6.159     -1.152      0.255     -19.465       5.275
fta          -14.8414     12.167     -1.220      0.228     -39.279       9.596
ast           -3.8815      2.791     -1.391      0.170      -9.487       1.724
pf             0.6271      3.033      0.207      0.837      -5.466       6.720
hwp           -0.3558      0.623     -0.571      0.571      -1.608       0.896
ps_rtg        -3.7474      2.968     -1.262      0.213      -9.709       2.215
fgp          -21.5417     19.728     -1.092      0.280     -61.167      18.083
fgp3          -3.0010      4.194     -0.716      0.478     -11.425       5.422
mrg            8.4968      5.583      1.522      0.134      -2.718      19.711
opp_wp         0.6815      1.164      0.586      0.561      -1.656       3.019
s_rating       0.1676      0.965      0.174      0.863      -1.771       2.106
ft%           -3.9242      2.937     -1.336      0.188      -9.822       1.974
def           -2.3734      2.817     -0.843      0.403      -8.031       3.284
```
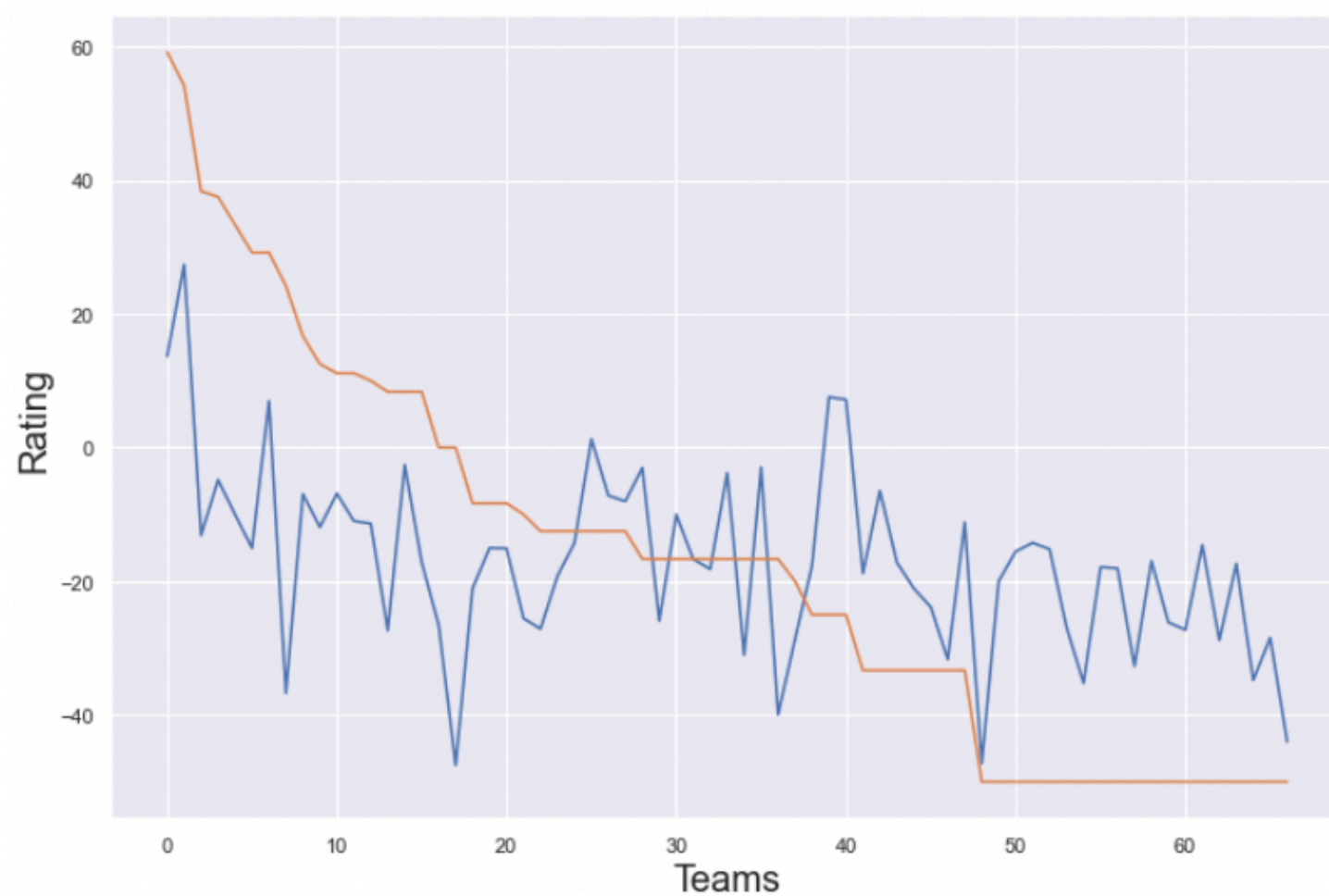
# Fit of the Model

# Lowest AIC Regression

```
                         OLS Regression Results
==============================================================================
Dep. Variable:                 rating   R-squared:                       0.223
Model:                            OLS   Adj. R-squared:                  0.199
Method:                 Least Squares   F-statistic:                     9.208
Date:                Fri, 25 Feb 2022   Prob (F-statistic):           0.000306
Time:                        02:40:40   Log-Likelihood:                -311.61
No. Observations:                  67   AIC:                             629.2
Df Residuals:                      64   BIC:                             635.8
Df Model:                           2
Covariance Type:            nonrobust
==============================================================================
                 coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
const       -194.7093     47.217     -4.124      0.000    -289.036    -100.382
ps             2.4840      0.628      3.957      0.000       1.230       3.738
ps_rtg        -3.0799      1.989     -1.548      0.126      -7.054       0.894
==============================================================================
Omnibus:                        4.185   Durbin-Watson:                   0.414
Prob(Omnibus):                  0.123   Jarque-Bera (JB):                3.731
Skew:                           0.494   Prob(JB):                        0.155
Kurtosis:                       2.400   Cond. No.                     1.11e+03
==============================================================================
```

Fit of the Model (Lowest AIC)

# How did it go / what did we find:

Turns out that College Basketball is not so easily predictable:

- NCAA Tournament results are hard to interpret

- Teams play only against a small portion of all teams, so having (un)favorable match-ups skew the perceived strength during the „Season"

Conclusions:

- The fact that we didn't find any obvious correlation doesn't mean there can't be some – actually there is groups that have been using models based on Machine Learning tools in the past, which predicted the results quite well (be it thanks to luck, or model quality)

- Would be cool to look further into this, but currently outside our skill range