# Can March Madness be Tamed with Math?

An exploration of constructing a probabilistic NCAA tournament bracket.



Baltimore-python
4/15/2015
Daniel O'Neill

# Overview

- Background on Competition
- Benchmarks and Early Models
- Refining the Models
- Tracking the Competition
- Results
- Take-aways/Future Work/Recommendations
- Questions

# NCAA College Tournament

| Second Round MARCH 19-20 | Third Round MARCH 21-22 | Regional Semifinals MARCH 26-27 | Regional Finals MARCH 28-29 | National Semifinals APRIL 4 | | National Semifinals APRIL 4 | Regional Finals MARCH 28-29 | Regional Semifinals MARCH 26-27 | Third Round MARCH 21-22 | Second Round MARCH 19-20 |

FILL OUT YOUR BRACKET
NATIONAL BRACKET DAY
MARCH 16

**First Four®**

| 16 Manhattan (19-13) 64 | 11 Ole Miss (20-12) 94 | First Round* | 77 N. Florida (23-11) 16 | 55 Boise St. (25-8) 11 |
| 16 Hampton (16-17) 74 | 11 BYU (25-9) 90 | DAYTON MARCH 17-18 | R. Morris (19-14) 16 | Dayton (25-8) 11 |

MW   W   S   E

Watch On truTV

**MIDWEST** CLEVELAND Mar 28

1 Kentucky (34-0) 79
16 Hampton 56
8 Cincinnati (22-10) 66
9 Purdue (21-12) 65
5 West Virginia (23-9) 68
12 Buffalo (23-9) 62
4 Maryland (27-6) 65
13 Valparaiso (28-5) 62
6 Butler (22-10) 56
11 Texas (20-13) 48
3 Notre Dame (29-5) 69
14 Northeastern (23-11) 65
7 Wichita St. (28-4) 81
10 Indiana (20-13) 76
2 Kansas (26-8) 75
15 New Mexico St. (23-10) 56

1 Kentucky 64
8 Cincinnati 51
5 WVU 69
4 Maryland 59
6 Butler 64
3 Notre Dame 67
7 Wichita St. 78
2 Kansas 65

1 Kentucky 78 (Louisville Mar 21)
5 WVU 39 (Columbus Mar 22)
3 Notre Dame 81 (Pittsburgh Mar 21)
7 Wichita St. 70 (Omaha Mar 22)

1 Kentucky 68 (Cleveland Mar 26)
3 Notre Dame 66 (Cleveland Mar 26)

1 Kentucky 64
April 4

**Final Four®** INDIANAPOLIS APRIL 4 AND 6

**National Championship** APRIL 6

| 1 Wisconsin 63 | Duke | 68 Duke 1 |

**EAST** SYRACUSE Mar 29

61 Mich. St. 7
April 4

70 Louisville 4 (Syracuse Mar 27)
76 Mich. St. 7 (Syracuse Mar 27)

65 NC St. 8
75 Louisville 4
58 Oklahoma 3
62 Mich. St. 7

Villanova 1 (Pittsburgh Mar 21) 68
NC St. 8 71
UNI 5 53
Louisville 4 66
Dayton 11 (Columbus Mar 22) 66
Oklahoma 3 72
Mich. St. 7 60
Virginia 2 (Charlotte Mar 22) 79

93 Villanova (32-2) 1
52 Lafayette (20-12) 16
66 N. Carolina St. (20-13) 8
LSU (22-10) 9
71 UNI (30-3) 5
Wyoming (25-9) 12
57 Louisville (24-8) 4
55 UC Irvine (21-12) 13
Providence (22-11) 6
Dayton 11
Oklahoma (22-10) 3
60 Albany (24-8) 14
70 Michigan St. (23-11) 7
63 Georgia (21-11) 10
Virginia (29-3) 2
67 Belmont (22-10) 15

**WEST** LOS ANGELES Mar 28

1 Wisconsin (31-3) 86
16 Coastal Caro. (24-9) 72
8 Oregon (25-9) 79
9 Oklahoma St. (18-13) 73
5 Arkansas (26-8) 56
12 Wofford (28-6) 53
4 North Carolina (24-11) 67
13 Harvard (22-7) 65
6 Xavier (21-13) 76
11 Ole Miss 57
3 Baylor (24-9) 56
14 Georgia St. (24-9) 57
7 VCU (26-9) 72
10 Ohio St. (23-10) 75
2 Arizona (31-3) 93
15 Texas Southern (22-12) 72

1 Wisconsin 72
8 Oregon 65
5 Arkansas 78
4 UNC 72
6 Xavier 60
14 Georgia St. 67
10 Ohio St. 75
2 Arizona 66

1 Wisconsin 79 (Omaha Mar 22)
1 Wisconsin 85 (Los Angeles Mar 26)
4 UNC 87 (Jacksonville Mar 21)
6 Xavier 60 (Jacksonville Mar 21)
2 Arizona 73 (Los Angeles Mar 26, Portland Mar 21)

1 Wisconsin 71 (Los Angeles Mar 28)

**SOUTH** HOUSTON Mar 29

61 Duke 1
April 4

Duke 1 (Charlotte Mar 22)
Houston Mar 27

66 Duke 1
57 Utah 5

Duke 1 66 (Houston Mar 27)
Gonzaga 2

85 Duke (29-4) 1
56 Robert Morris 16
76 San Diego St. (26-8) 8
64 St. John's (21-11) 9
57 Utah (24-8) 5
50 S.F. Austin (29-4) 12
84 Georgetown (21-10) 4
74 Eastern Wash. (26-8) 13
59 SMU (27-6) 6
UCLA (20-13) 11
59 Iowa St. (25-8) 3
60 UAB (19-15) 14
65 Iowa (21-11) 7
Davidson (24-7) 10
86 Gonzaga (32-2) 2
76 N. Dakota St. (23-9) 15

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

tbs   CBS   TNT   truTV

**ALL TIMES EASTERN**

The NCAA opposes all forms of sports wagering.

# Machine Learning Mania 2015
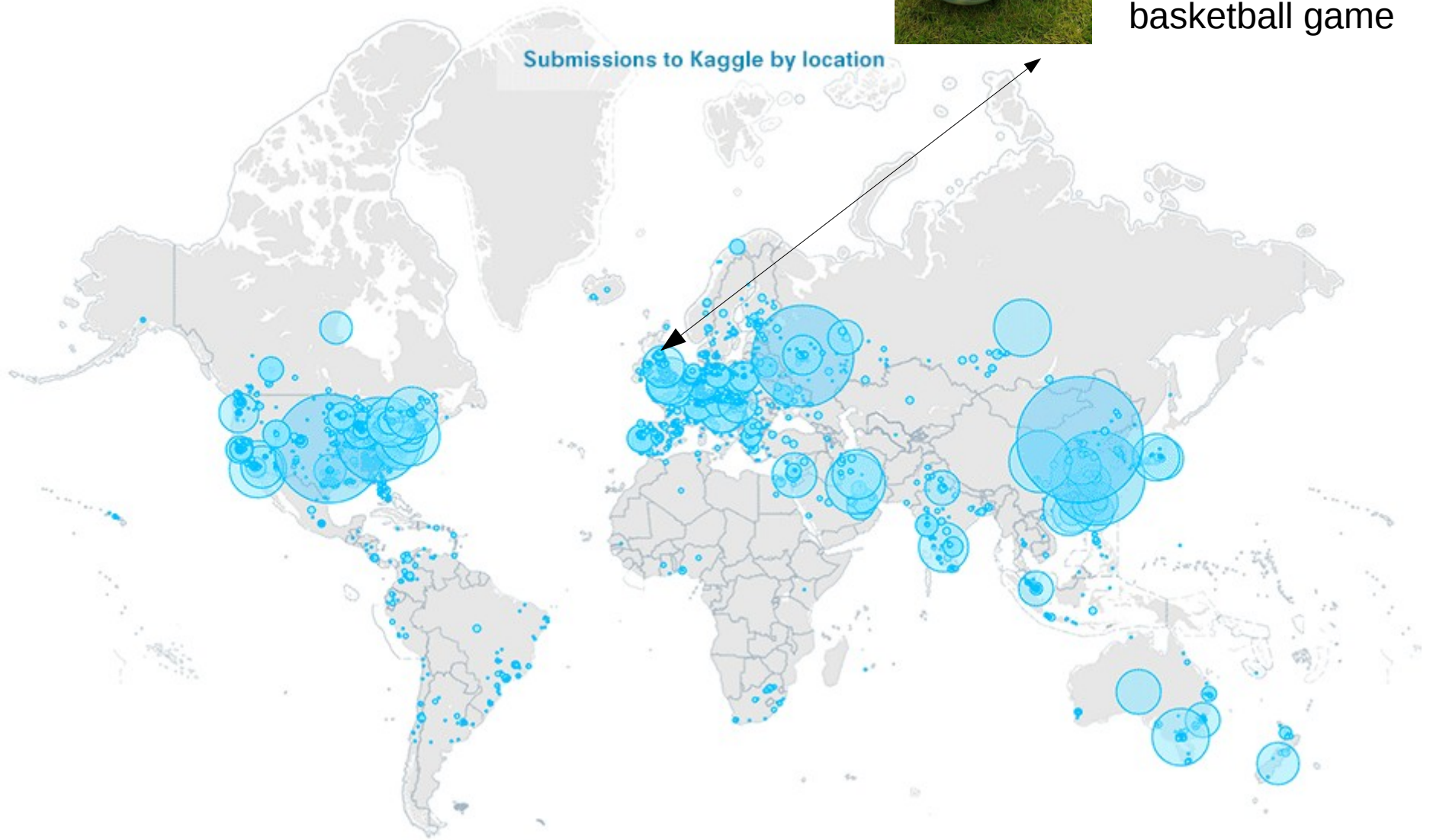


- www.kaggle.com
- Founded in 2010
- Over 200,000 kagglers
- Devoted to a community of data scientists and learners solving data science problems
- Competition and Learning

Submissions to Kaggle by location

Bluefool:
Kaggle Master that had never seen a basketball game

# Rules of the Game

- Provide Probabilities of Team A beating Team B for every possible combination (68*67/2=2278)

- Scored on Games that Actually Happen (63)

- All games weighted equally

- Lowest LogLoss wins!

- Correctly predict a game with 80% :0.22

- Incorrectly predict a game with 80%:1.61

$$\mathrm{LogLoss} = -\frac{1}{n}\sum_{i=1}^{n}\left[y_i\log(\hat{y}_i) + (1-y_i)\log(1-\hat{y}_i)\right],$$

# Data Provided

- Stage 1 (Feb 2-Selection Sunday):
  - Team Name File (Mapped team name to number)
  - Regular Season / Tourney results (1985-2014)
  - Regular Season / Tourney detailed results (2003-2014)
  - Tourney Seeds
  - Tourney Slot
  - Evaluated on 2010-2014 tournaments
- Stage 2 (Submissions due Day before Tourney):
  - 2015 Versions of data
  - Follow leader board live
- Additional Data: Massey Ordinals, Vegas Spreads, www.kenpom.com

# My Philosophy on the Competition

- There is an element of luck
  - Sample size is very small for each prediction (1 game) True probabilities are never evaluated
  - Probability of winning better than ESPN tournament
- Submission does not have to be great
  - Has to be reasonable
  - Want it to be slightly different than a standard approach to separate from other competitors
- Wanted to learn something

# Overview

- Background on Competition
- Model Building
- Tracking the Competition
- Results
- Take-aways/Future Work/Recommendations
- Questions

# Pandas

## Python Data Analysis Library

- Pandas is convenient for both live exploratory analysis and for incorporating into scripts

- Uses DataFrame as primary object

- Comparable to R

- Allows for quick selection and sub-selection of data

- Column names add to readability

- Easy to aggregate

- reg_season=pd.read_csv('regular_season_compact_results.csv')

- train=reg_season[(reg_season.season>2000)&(reg_season.season<2010)]

- train['win_pct']=1.*train['wins']/(1.*train['wins']+1.*train['losses'])

- combined=pd.merge(results,seeds,how='inner',left_on=['wteam','season'],right_on=['team','season'])

- Describe()

- Head(), tail()

- Sorting

# Initial Benchmarks

- Data was loaded into a Pandas dataframe

- Submission file was created using Python itertools to find every combination of two teams

- Lower number team placed first and then the probability of Team A beating Team B would be calculated

- Checked with a 0.5 submission for all games

- Moved to a linear model based on seeds

# PageRank

- Weights items based on their relative importance

- Originally used in links for weighting webpages and the probability of a user moving to that page

- Incorporates a level of randomness to ensure that the model does not spiral away from convergence

$$M = (1 - p) \cdot A + p \cdot B \qquad B = \frac{1}{n} \cdot \begin{bmatrix} 1 & 1 & \ldots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \ldots & 1 \end{bmatrix}$$

$M^k z$ converges to the vector v*

# Applying it to Basketball

- All teams are a node in the graph

- Each game is an edge

- The losing team points to the winning team

  – Later iterations of my model used score differential to weight the edges

  – The score differential was adjusted to reflect home or away teams

  – Took the log of the difference to adjust for some scaling problems

# Adjusting the Randomness

- Weight of the matrix that all teams have equal importance was given more probability in increments until a LogLoss on the test data was achieved

- Ended up close to only using roughly 10% of the weighted matrix.  This demonstrated that all teams are fairly equal

- Rated Miami high because it beat Duke by a large differential.  Did not make tournament.

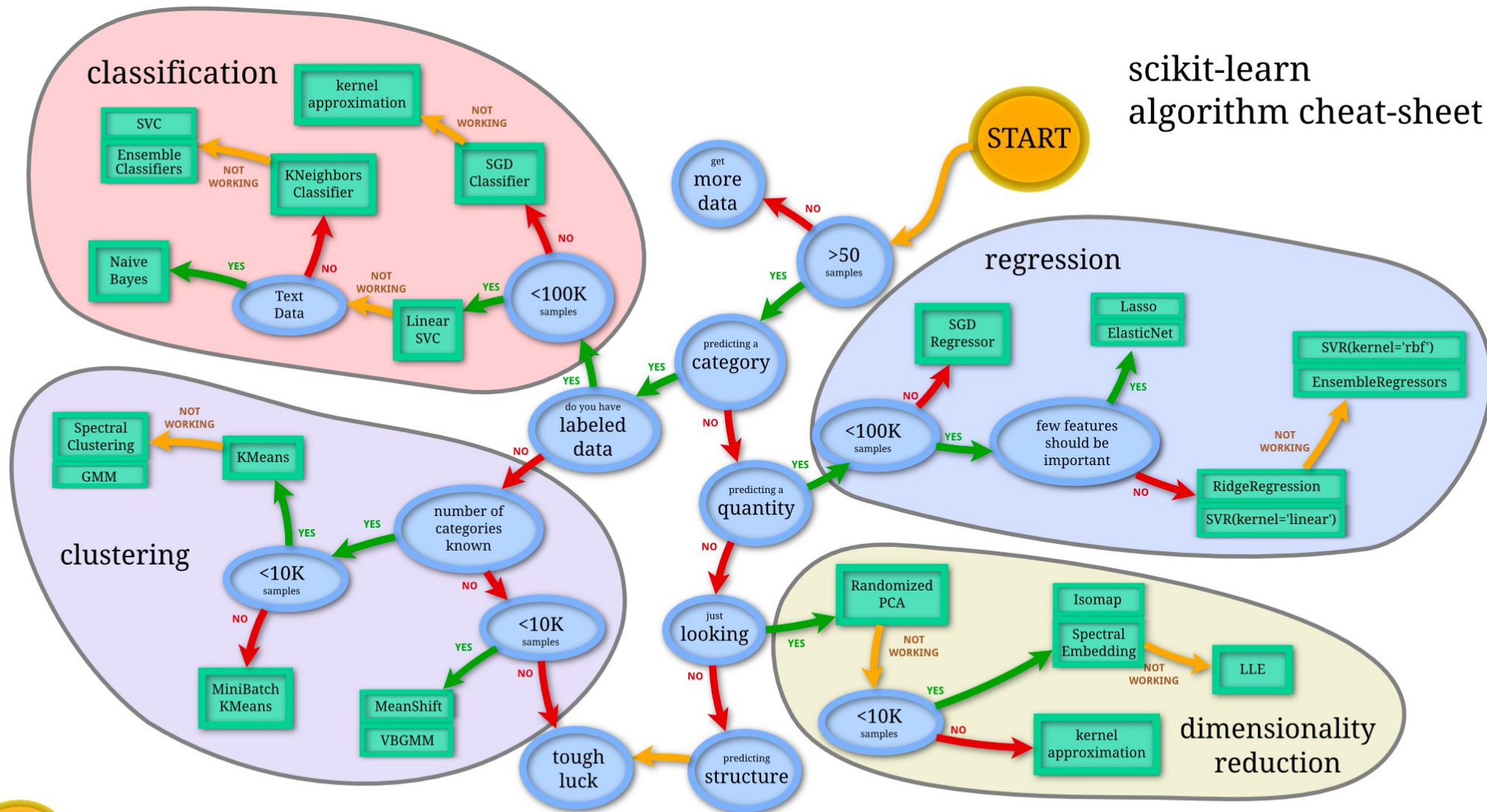| 95.00% | 5.00% |
|---|---|
| Duke | Duke |
| Kentucky | Notre Dame |
| Wisconsin | NC State |
| Virginia | Virginia |
| Gonzaga | North Carolina |
| North Carolina | Kentucky |
| Villanova | Miami FL |
| Arizona | Georgetown |
| Kansas | Villanova |
| VA Commonwealth | Kansas |
| Ohio St | Wisconsin |
| Baylor | Butler |
| Wichita St | Oklahoma |
| Notre Dame | Maryland |
| Northern Iowa | Providence |
| BYU | BYU |
| Arkansas | Baylor |
| Maryland | Syracuse |
| Cincinnati | Iowa St |
| Utah | West Virginia |

# SciKit Learn

## Machine Learning Library built on numpy scipy and matplotlib

- Extensive Machine Learning library for Python
- Tools for Data Splitting
  - cross_validation train_test_split, Kfold, StratifiedKfold
- Data Preprocessing
  - PCA, Data Scaling, BoxCox
- Classification-Nearest neighbors, randomforest,GBM
- Regression-Linear, Logistic, Ridge Lasso
- Model Selection: GridSearch, RandomGridSearch
- Well documented
- All models follow similar structure of fit and predict
  - Ensembling advantages

# scikit-learn algorithm cheat-sheet

**START**

## classification

kernel approximation

SVC

Ensemble Classifiers

NOT WORKING

KNeighbors Classifier

NOT WORKING

SGD Classifier

NOT WORKING

Naive Bayes

YES

Text Data

NO

NOT WORKING

Linear SVC

YES

<100K samples

NO

get more data

NO

>50 samples

YES

predicting a category

YES

do you have labeled data

NO

## regression

SGD Regressor

NO

Lasso ElasticNet

YES

SVR(kernel='rbf')

EnsembleRegressors

NOT WORKING

<100K samples

YES

few features should be important

NO

RidgeRegression

SVR(kernel='linear')

predicting a quantity

YES

NO

## clustering

Spectral Clustering

GMM

NOT WORKING

KMeans

YES

<10K samples

YES

number of categories known

NO

NO

MiniBatch KMeans

<10K samples

YES

MeanShift

VBGMM

NO

just looking

YES

## dimensionality reduction

Randomized PCA

NOT WORKING

<10K samples

YES

Isomap

Spectral Embedding

NOT WORKING

LLE

NO

kernel approximation

NO

tough luck

predicting structure

NO

**Back**

scikit learn

# Constructing Models

- All data provided was given in terms of the winning team

- Binary Classification

- All data was doubled and inversed

  - Seed difference was multiplied by -1

  - All team attributes were flipped

# Primarily Used Linear and Logistic Regression

- Created Several Models Based on a Linear or Logistic Framework

- Averaged them to reduce extremes

- No regularization was used throughout

# Kentucky

## Top Left Region

1 Villanova
16 Lafayette
1 Villanova .68
8 NC State
9 LSU
8 NC State .32
1 Villanova .61
5 Northern Iowa
12 Wyoming
5 Northern Iowa .46
4 Louisville
13 UC Irvine
4 Louisville .54
4 Louisville .39
1 Villanova .38
6 Providence .51
11 Boise St .49
6 Providence .35
3 Oklahoma
14 Albany NY
3 Oklahoma .65
3 Oklahoma .32
7 Michigan St
10 Georgia
7 Michigan St .06
2 Virginia
15 Belmont
2 Virginia .94
2 Virginia .68
2 Virginia .62
2 Virginia .99

## Top Right Region

1 Duke .99
16 North Florida .01
1 Duke .59
8 San Diego St
9 St John's
8 San Diego St .41
1 Duke .65
5 Utah
12 SF Austin
5 Utah .54
4 Georgetown
13 E Washington
4 Georgetown .46
5 Utah .35
1 Duke .65
6 SMU
11 UCLA
11 UCLA .56
3 Iowa St
14 UAB
14 UAB .44
11 UCLA .29
1 Duke .41
7 Iowa
10 Davidson
7 Iowa .54
2 Gonzaga
15 N Dakota St
2 Gonzaga .46
7 Iowa .71
7 Iowa .35

2 Virginia .41

## Bottom Left Region

1 Kentucky .99
16 Manhattan .01
1 Kentucky .89
8 Cincinnati
9 Purdue
8 Cincinnati .11
1 Kentucky .85
5 West Virginia
12 Buffalo
5 West Virginia .44
4 Maryland
13 Valparaiso
4 Maryland .56
4 Maryland .15
1 Kentucky .74
6 Butler
11 Texas
6 Butler .52
3 Notre Dame
14 Northeastern
3 Notre Dame .48
6 Butler .45
7 Wichita St
10 Indiana
7 Wichita St .45
2 Kansas
15 New Mexico St
2 Kansas .55
2 Kansas .55
2 Kansas .26
1 Kentucky .67

1 Kentucky .99

## Bottom Right Region

1 Wisconsin .78
16 Coastal Car
1 Wisconsin .61
8 Oregon
9 Oklahoma St
8 Oregon .22
1 Wisconsin .57
5 Arkansas
12 Wofford
5 Arkansas .43
4 North Carolina
13 Harvard
4 North Carolina .57
4 North Carolina .39
1 Wisconsin .33
6 Xavier
11 Mississippi
6 Xavier .45
3 Baylor
14 Georgia St
14 Georgia St .55
14 Georgia St .24
7 VA Commonwealt
10 Ohio St
10 Ohio St .35
2 Arizona
15 TX Southern
2 Arizona .65
2 Arizona .76
2 Arizona .43

# Overview

- Background on Competition
- Benchmarks and Early Models
- Model Building
- Tracking the Competition
- Results
- Take-aways/Future Work/Recommendations
- Questions

# More Madness!

- NetProphet
  - http://netprophetblog.blogspot.com/
- Released his submission with random noise incorporated
- BlueFool and others incorporated his model
- One contestant copied his submission exactly and then gave Kentucky a 100% chance of victory in all games

# March Machine Learning Mania (Round of 64 Predictions)

# March Machine Learning Mania (Round of 16 Predictions)

# Overview

- Background on Competition
- Benchmarks and Early Models
- Model Building
- Tracking the Competition
- Results
- Take-aways/Future Work/Recommendations
- Questions

# Results

- After Day 2 I was 18$^{th}$ out of 341.  Feeling good

- At the end I was 216$^{th}$. Glad I don't bet

- The winner had some odd predictions

  - 100% chance that #14 Georgia St would beat #3 GA

- BlueFool only trained on tournament data and was able to come in 8$^{th}$

# Overview

- Background on Competition

- Benchmarks and Early Models

- Model Building

- Tracking the Competition

- Results

- Take-aways/Future Work/Recommendations

- Questions

# Take Aways

- There is a level of gamesmanship with such a small sample size

- Favorites do win but all teams are playing at a high level

- Last year's winners ran simulations of how many times there model would win and found almost the entire board had a shot

# Future Models

- Incorporate Coach data
- Perform more feature analysis and possible feature engineering
- Weight the later part of the season
- Focus more on blending models
- Incorporate more ranking systems

# Favorite Resources on Data Science

- Applied Predictive Modeling- Max Kuhn, Kjell Johnson

- The LION Way: Learning plus Intelligent Optimization- Mauro Brunato, Roberto Battiti

- Podcasts: Talking Machines, O'Reilly Data Show

- Coursera

- edX

# Questions?