

(Much) better than lotto tickets: Analytics and NCAA tournament win probabilities

Michael J. Lopez

Assistant Statistics Professor
Skidmore College
@StatsbyLopez

March 2, 2016

Outline

Introduction

- Kaggle contest

- What we did

- Lucky or good?

Discussion and Advice

- Kaggle

- Traditional pools

- Sports analytics

Background

- ▶ Single elimination tournament with 64-68 teams.
- ▶ \$2.5 billion wagered on the tournament in 2012 (Boudway 2014; Tsu 2014)
- ▶ Reminder: gambling is still illegal in most places

Scoring formats:

1. Traditional 1:2:4:8:16:32 (Yahoo, ESPN, etc)
 - ▶ Points for picking winners only, done before tournament
2. Kaggle “March Machine Learning Mania”
 - ▶ All games judged on predictive probabilities

Kaggle description

- ▶ ~ 400 entries from 248 teams in 2014 (500 teams in 2015)
- ▶ Predict win probability for every *possible* game (2278 contests)
- ▶ Only 63 games actually played used in scoring
- ▶ Scoring function:

$$\text{LogLoss} = -y \times \log(\hat{y}) + (1 - y) \times \log(1 - \hat{y})$$

where \hat{y} is the predicted probability of a win and y is the actual outcome (0 or 1).

- ▶ We* won this contest in 2014.
- ▶ I'll address a few main questions:
 - ▶ What did we do?
 - ▶ How lucky did we get?
 - ▶ Traditional NCAA pools
 - ▶ Lessons transferable to sports analytics

*Jointly with Gregory Matthews (Loyola-Chicago)

- ▶ Two sources of data:
 - ▶ Las Vegas point spread data (Model M_1)
 - ▶ Ken Pomeroy efficiency (Pomeroy 2012) ratings (Model M_2)
 - ▶ Why efficiency?
- ▶ Logistic regression: outcome variable of 1 for a win and 0 for a loss.
 - ▶ **Model M_1** : 1 predictor
 - ▶ Spread
 - ▶ **Model M_2** : 5 predictors
 - ▶ Offensive efficiency (home, away)
 - ▶ Defensive efficiency (home, away)
 - ▶ Neutral indicator

- ▶ Find w to minimize LogLoss of $w \times \hat{y}_{M_1} + (1 - w) \times (1 - \hat{y}_{M_2})$
- ▶ In-sample versus out-of-sample testing
- ▶ Our submissions:
 - ▶ $S_1 = 0.75\hat{y}_{m_1} + 0.25\hat{y}_{m_2}$
 - ▶ $S_2 = 0.25\hat{y}_{m_1} + 0.75\hat{y}_{m_2}$ (Winning entry)

How lucky were we?

- ▶ We simulated a the tournament 10,000 times with differing “true” underlying win probabilities:
 - ▶ S_1, S_2
 - ▶ Mean of top 10 entries.
 - ▶ Mean of all entries.
 - ▶ All games 0.5
- ▶ In each simulated tournament we scored all entries and counted how often we won.

Results, Kaggle tournament

- ▶ Given our probabilities as true probabilities:
 - ▶ Each entry \sim 15% of winning
 - ▶ Each entry \sim 50% of top-10 finish
- ▶ We finished 4th in 2015

Outline

Introduction

- Kaggle contest

- What we did

- Lucky or good?

Discussion and Advice

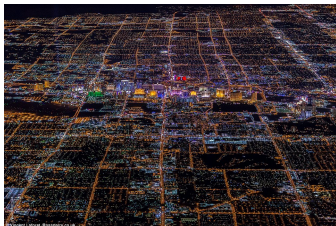
- Kaggle

- Traditional pools

- Sports analytics

Lessons, Kaggle tournament

1. You always need luck
2. 2 prediction models combined together outperform either alone
3. Better data \geq complex models



Lessons, traditional pools

1. Find upset pools - people ~~are idiots~~ pick too passively

Lessons, traditional pools

1. Find upset pools - people ~~are idiots~~ pick too passively
2. Game Theory
 - ▶ Opponents picks known in expectation
 - ▶ Find value: Duke 2010 (✓), Kentucky 2015 (X)
 - ▶ Consider n

Lessons, traditional pools

1. Find upset pools - people ~~are idiots~~ pick too passively
2. Game Theory
 - ▶ Opponents picks known in expectation
 - ▶ Find value: Duke 2010 (✓), Kentucky 2015 (X)
 - ▶ Consider n
3. Tools
 - ▶ kenpom.com, fivethirtyeight.com
 - ▶ <http://www2.isye.gatech.edu/~jsokol/lrmc/>
 - ▶ 'Who picked whom'

Lessons, traditional pools

1. Find upset pools - people ~~are idiots~~ pick too passively
2. Game Theory
 - ▶ Opponents picks known in expectation
 - ▶ Find value: Duke 2010 (✓), Kentucky 2015 (X)
 - ▶ Consider n
3. Tools
 - ▶ kenpom.com, fivethirtyeight.com
 - ▶ <http://www2.isye.gatech.edu/~jsokol/lrmc/>
 - ▶ 'Who picked whom'
4. $P(\text{Loss}) > P(\text{Win})$

Lessons, sports analytics

1. Research
2. Visualize
3. Practice
4. Share
5. Improve
6. Generalize

Sidebar: R-statistical software

Citations

- ▶ Lopez, M and Matthews, G.J. “Building an NCAA men’s basketball predictive model and quantifying its success.” *Journal of Quantitative Analysis in Sports*, 11:1 (2015): 5-12.
- ▶ Carlin, B. P. 1996. ‘Improved NCAA Basketball Tournament Modeling Via Point Spread and Team Strength Information.’ *The American Statistician* 50:3943.
- ▶ Pomeroy, K. 2012. Ratings Glossary. URL <http://bit.ly/1LGb79q> (accessed June 1, 2014).