# Motivation

- Gain deeper understanding of NBA

- How did the game developed over the years?

- Distribution of player performance

- Why/How are some teams/players better than others?

- Is it possible to predict the outcome a game?

# Outline

- Data format and loading pipeline
- Exploration of data on different levels:
  - Player
  - Seasons
  - Game
- Prediction tasks:
  - Game result
  - All-NBA

# The raw data

Raw Play-By-Play base data set

- 19 seasons ( 2000-01 – 2018-19), data split by season
- 10 389 755 plays/events (shot, foul, turnover, …)
- 35 columns
- nominal, continuous, time series

- Main workhorse
  - Python 3 + jupyter notebooks

- Data acquisition
  - beautiful soup, selenium, tdqm, webdriver-manager

- Data analysis and prediction
  - numpy, pandas, plotly, sklearn, pickle

- Version control
  - GitHub

# Data loader pipeline

Base data set of 1.6 GB

- long time to load
- especially with extra processing

Solution:

- multiple data loaders created
- special parameter for loading
- intense use of pandas masks
- storing computed results in intermediate pickle files
- ~1650 lines of code
- data loading optimized from minutes to seconds or milliseconds

# Player data

When searching for player data no complete dataset were found

1. We created our own data scraper
2. Used selenium and chrome driver to simulate browser usage since static loading failed
3. Now it can retrieve all player data available on the NBA website

# Exploring player characteristics

We managed to learn some common characteristics of NBA players

**Average statistics about players through all seasons**

|  | Age (years) | Height (cm) | Weight (kg) |
|---|---|---|---|
| Mean | 27 | 200 | 100 |
| Standard deviation | 0.35 | 0.26 | 0.94 |

Age and weights of the players for season 2018-19 (with random noise)



Distribution of player stays in NBA



Distribution of club changes

# Combining age with extracted performance

Answering questions that our raw data cannot directly answer:

Where is the peak point of players performance during their careers?

# Diving into individual players

Impact of variables on player performance

- club changed negatively impacted player performance (points)
- problems because of incomplete data



Brian Skinner statistics through the seasons, with marked club changes



LeBron James statistics through the seasons, with marked club changes

# Star players dominate the statistics

Few players played many games or scored a lot of points

The median number of points scored is 742

The median number of games played is 142

# Distribution of shot distances per season

How game changed over years or what has changed in the game?



Shoting attempts per season

# Shift in play-style confirmed by data

Goodbye to mid-range shots!!



Shot Attempts by Type (% of all shots)

# Analysis of games

- Large corpus of 22,965 games (every game in 19 consecutive regular seasons)
- Extract the essence of an NBA game -> Contrast with domain knowledge
- Heartbeat of a game:

number of field goal makes/misses by minute:

# Using data to answer questions

- Why do players shoot worse at the end of close games?
  - Defense vs. Pressure/Fatigue
- Use shot types as indicator:

Field goal accuracy

Free throw accuracy



Shows power of combining data analysis with domain knowledge

# Final game score/result

- Most important/interesting feature of a game
- Product of an artificial competition
- Unique distribution:



- What impacts winning?
- Can we predict a game's result?

# Finding features correlated with wins

Some do  (Field-goal accuracy)

Some don't (Field goal volume)

# Putting games into temporal context

- Order games by date in time -> Series of games
- Past performance is indicative of future success:

Wins by season record prior to game:



Legend: ■ Home team wins  ■ Visitor team wins

Chart 1 (Home team has better record): Home team wins 74.5%, Visitor team wins 25.5%

Chart 2 (Teams have even record): Home team wins 61%, Visitor team wins 39%

Chart 3 (Visitor team has better record): Home team wins 45.2%, Visitor team wins 54.8%

➔ Central for game result prediction

# Predictive mining

1. **Winner of the game**

   - game level

   - sports betting?

2. **All-NBA team**

   - season level

   - justifying the journalists choices / who deserved the reward?

   - find minimal production stats for achieving All-NBA

- Importance of data loader pipeline

   - easy to add new features

# Winner of the game - overview

- Statistical baseline
  - 60 % games won by home team
- Window size (rolling averages)
- Unbalanced data
- Scaled features
- Labeled team IDs
- Hyperparameter search
  - regularization, kernel, max depth
- Feature selection
  - manual selection, RFE, select k best, fastener

Mean accuracy of LogisticRegression models for different window sizes of recent stats

# Winner of the game - comparison of models

### Random forest

| Accuracy | Precision | Recall | f1-score |
|---|---|---|---|
| 0.65 | 0.64 | 0.62 | 0.61 |

### SVM

| Accuracy | Precision | Recall | f1-score |
|---|---|---|---|
| 0.66 | 0.64 | 0.64 | 0.64 |

### SVM hyperparameters

| Parameter | Value |
|---|---|
| C | 24.8 |
| kernel | "linear" |
| decision_function_shape | "ovo" |
| random_state | 0 |

### Gradient boosting classifier

| Accuracy | Precision | Recall | f1-score |
|---|---|---|---|
| 0.65 | 0.64 | 0.63 | 0.64 |

### Logistic regression

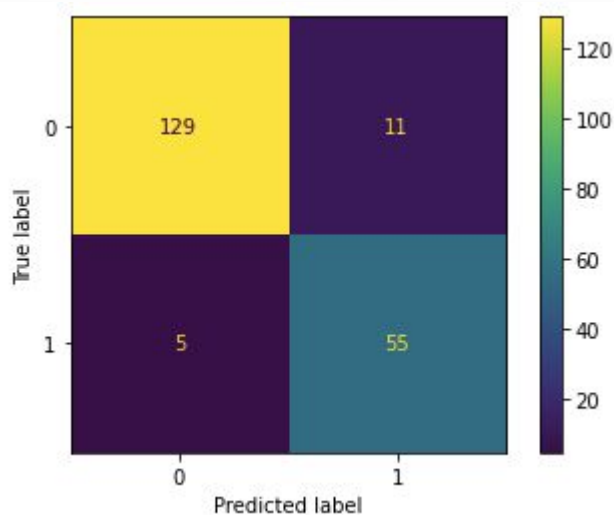| Accuracy | Precision | Recall | f1-score |
|---|---|---|---|
| 0.65 | 0.65 | 0.65 | 0.64 |

# All-NBA team - overview

- 15 players selected by the journalists
- Closely related to previously extracted features/analysis
- Scraped additional data
- Filtering by number of games played (145 players per season)
- Very unbalanced –> needed balancing
- Logistic regression

1. First 15 seasons as training set, predict last 4 seasons
2. Dataset is shuffled, seasons are not relevant

# All-NBA team - evaluation of model

### Chronological data



| Accuracy | Precision | Recall | f1-score |
|----------|-----------|--------|----------|
| 0.92 | 0.83 | 0.92 | 0.87 |

### Shuffled data



| Accuracy | Precision | Recall | f1-score |
|----------|-----------|--------|----------|
| 0.90 | 0.87 | 0.78 | 0.82 |

# Conclusion

- Play-By-Play data allows for multiple views on data
  - Extracted many features on different levels
  - Per game, player, season …
- Strong data-loading pipeline crucial for efficient working in a team
- Combining previous domain knowledge with data analysis allowed us to answer interesting questions in a data-driven manner
- Putting games into order reveals temporal dependencies
- Prediction of game result / All-NBA is possible