

Working as a Data Scientist in Sport



Professional Statisticians' Forum
21 November 2016
Dave Hastie
Director, Sporting Data Science



Introduction

Convincing you to invest the
next hour of your life

- Statistical background
- Business background
- Talk outline

Statistical background



- MSci. Mathematics - University of Bristol - 1996 - 2000
 - 4 year undergraduate degree
 - Large focus on statistics and computing modules in the later years
- Ph.D. Statistics - University of Bristol - 2001 - 2004
 - Reversible jump Markov chain Monte Carlo
 - Methodological, statistical software development
 - No sporting applications
- Postdoctoral researcher - Imperial College London - 2010 - 2012
 - Computationally intensive non-parametric modelling
 - Epidemiology and biostatistics applications

Business background



- Smartodds - 2003 - 2010
 - Head of Quantitative Team
- Onside Analysis - 2011 - 2015
 - Co-Founder, CTO
- Stratagem Technologies - 2014 - 2015
 - Director, CTO
- CricViz - 2015 - present
 - Director (CTO)
- Sporting Data Science - 2015 - present
 - Founder, CEO



- Career highlights by case study
 - Smartodds
 - Onside Analysis
 - Stratagem Technologies
 - CricViz
- Common themes and best practices
- Common challenges
 - Empirical models
 - Model assessment
- My perspective on the future
 - Sporting Data Science



Career Highlights by Case Study

What (I like to tell people) I do

- Smartodds
- Onside Analysis
- Stratagem Technologies
- CricViz



- What the company did [does]
 - Analysis (statistical modelling) of football [sport] using historical data to identify opportunities to place bets with positive expected return

$$\underbrace{E[R]}_{\text{expected return}} = X \underbrace{p(G | \text{data})}_{\text{statistical model}} - \underbrace{1}_{\text{bookmakers' odds}} > 0$$

positive expected return

historical data

binary event on which bet can be placed

Smartodds (cont.)

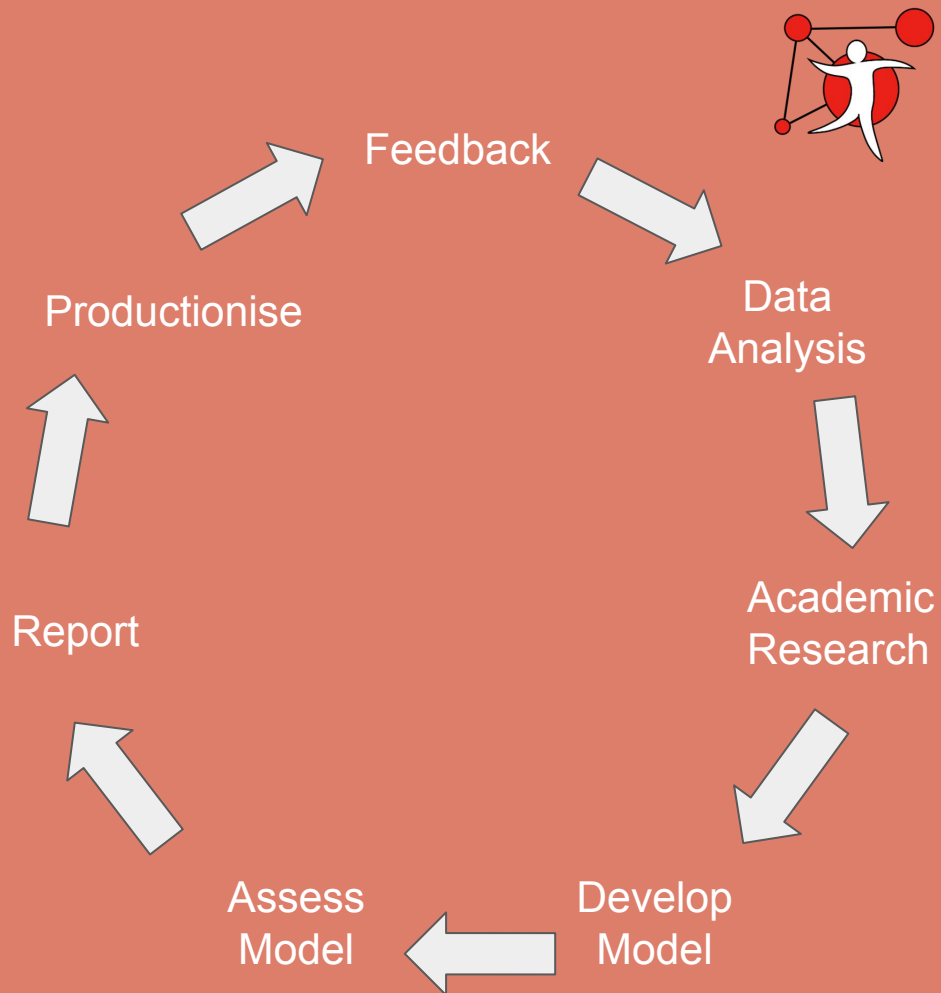
How was data science used?

Data processing

- Multiple sources, often unclean data
- Building pipelines

Model Development Cycle

- Key to reduce time around cycle



Smartodds (cont.)



Example model

$$\log \lambda_H = \mu + \gamma + \alpha_j + \beta_k$$

$$\log \lambda_A = \mu + \alpha_k + \beta_j$$

$$X \sim \text{Poisson}(\lambda_H)$$

$$Y \sim \text{Poisson}(\lambda_H)$$

Considerations:

- Temporal dependence
- Additional count data - shots, corners, chances
- Player data
- Environmental data - pitch, weather, distance

Inside Analysis (2011 - 2015)



- What the company did:
 - Analysis (including statistical modelling) of historical football data to inform and improve future decisions for multiple stakeholders in the football and connected industries.
 - Player assessment
 - Recruitment
 - Salary negotiation
 - Managerial assessment
 - Recruitment
 - Technical analysis
 - Coaching
 - Predictive modelling
 - Betting

Onside Analysis (cont.)



How was data science used?

- Data collection
- Data processing (multiple sources, often unclean)
 - Data pipelines
- Statistical modelling of sporting event data
 - Model development cycle
- Visualisations and reports
- Data delivery
 - API
 - Web dashboards



Example - Evaluating the contribution of players

- Opta data - every touch of the ball is time stamped and coded
- Assessing contribution of player
 - Contribution by scoring goals
 - Contribution by helping team to score goals
 - Contribution by preventing opposing team from scoring goals
- Credibility challenges
 - Disconnect between analysts / C-Level and coaching / playing staff
 - Opta data is an imperfect data source
 - A single implausible model output can be enough for model to be dismissed

Onside Analysis (cont.)



Extract of report

the average leading minutes proportion.¹. In the interests of protecting intellectual property, we omit the full details of this measure, but make the following observations about features of it's designs.

- Goals, times, line-ups and substitutions should be taken into account
- The measure should use the data to draw out increased signal about performance than is possible by simply considering the scoreline
- A player should be rewarded for obtaining positive player supremacy
- The greater the player supremacy, the greater the reward, but not all goals are equal
- Players should be rewarded more for an earlier increase in their player supremacy rather than a later change
- A player making an average contribution to a team should have a score of 0
- A player making a below average contribution should have a negative score
- A player making an above average contribution should have a positive score
- Scores should be comparable across team mates. Comparisons between players in separate teams should be possible, but with extreme caution recognising the limitations
- The measure should be in a form that can be generalised with access to additional high resolution player-action data
- The measure should be straightforward to test and refine as part of a statistical modelling process

The following sections provide some visualisations of OAPI.

Onside Analysis (cont.)



Extract of report

Premier League 2012/2013

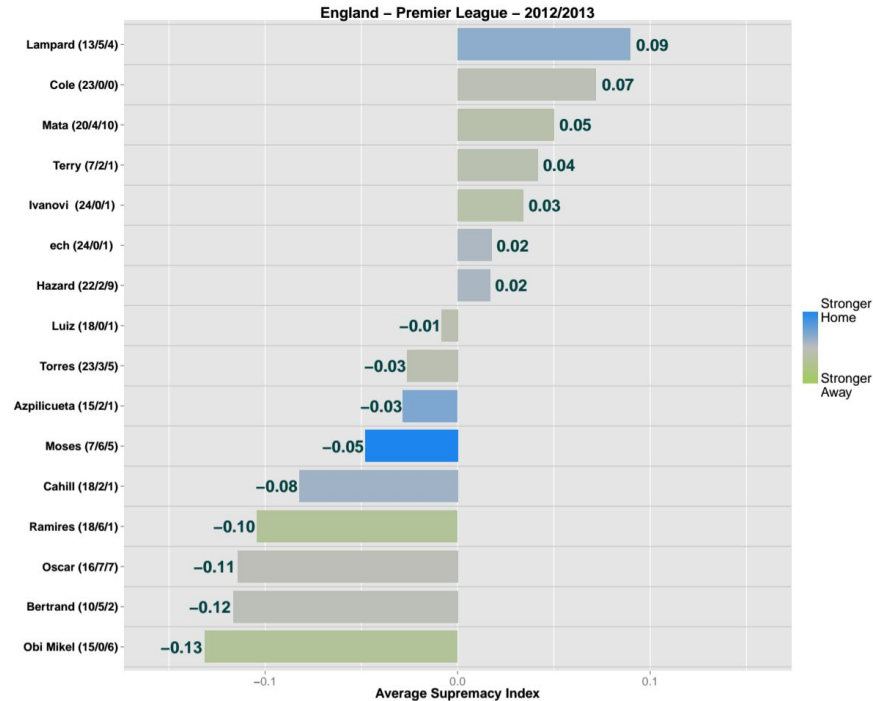


Figure 7: Player index based on average time and supremacy whilst on pitch for 2012/2013 Premier League



- What the company did [does]:

“generating consistent edge in sports prediction trading - a highly liquid asset class with uncorrelated returns consistent edge in sports prediction trading - a highly liquid asset class with uncorrelated returns”

- Similarities to Smartodds
- Taking an approach more inspired by financial trading
 - Portfolio building
 - Trading strategies
 - Risk management
- Trading platform
- Proprietary trading
- Managed fund



How was data science used?

- Data pipelines (collection and processing unclean data)
- Statistical modelling of sporting event data
 - Model development cycle
- Statistical modelling of odds time series
 - Time series methods
 - Reinforcement learning
- Portfolio optimisation
- Visualisations and reports
- Image processing
 - Automated data collection



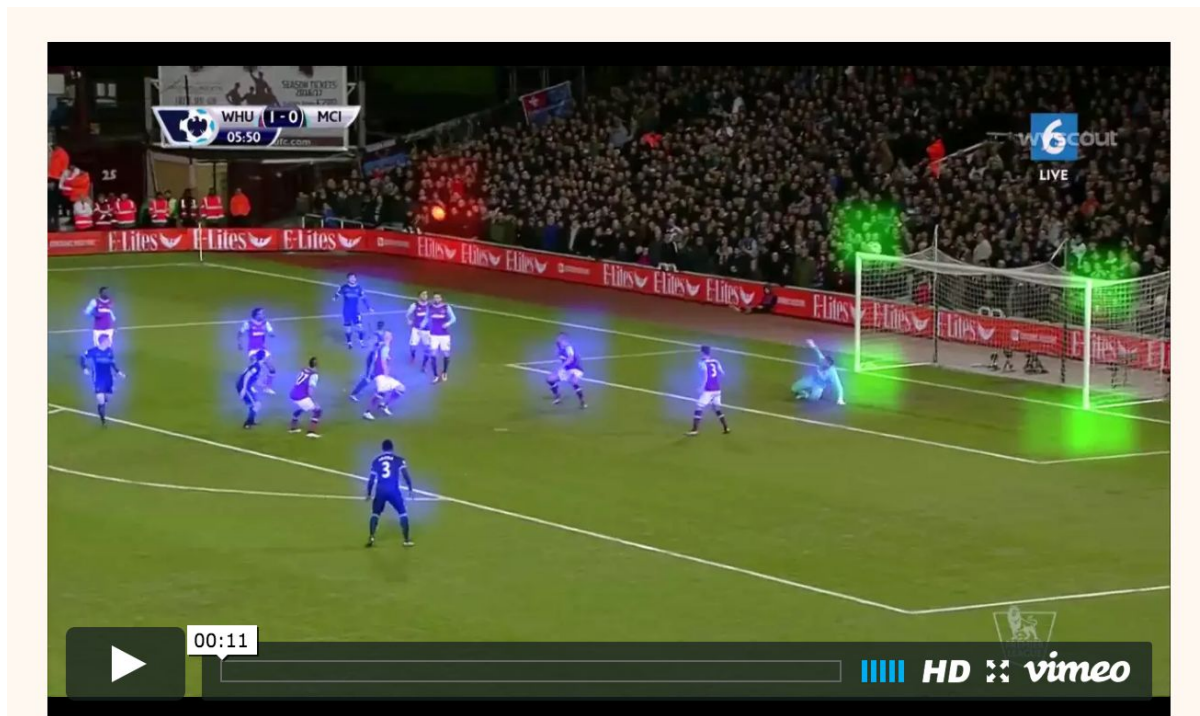
Example application

- Image processing
 - Started out as a Masters student project in partnership with Imperial College London
 - Manual data collection is expensive, can we automate it using image processing
 - Difference categories of “shot” - can we identify them
 - Became a full-time research area

Stratagem Technologies (cont.)



Example application





- What the company does:

“CricViz allows the user to understand cricket in greater detail than ever before. The unique CricViz computer model enables the prediction of match outcome, the interpretation of team and player performance and the anticipation of what is likely to happen next; cricket intelligence at the next level ”

- Very much interested in the modelling and understanding of the sport (cricket)
- Providing services to cricket fans and enthusiasts
 - Building a data driven community
- Providing data and analytics services
 - To broadcasters
 - To the professional cricket teams
 - For the betting industry



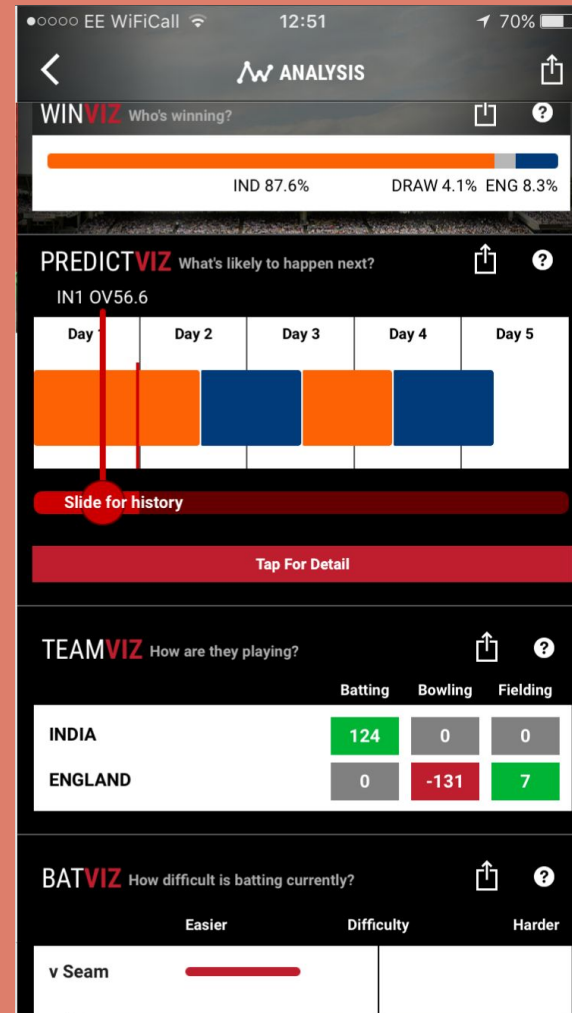
How is data science used?

- Data pipelines (collection and processing unclean data)
 - A large part of business model is built around the company aggregating and utilising multiple data sources that have never before been combined
 - Creating new derived data tailored to customer requirements
- Statistical modelling of sporting event data
 - Modelling outcomes of cricket games
 - Modelling player performance
- Visualisations and reports
- Data dissemination
 - APIs

CricViz (cont.)

Example - mobile app

- Probabilities of match outcome
- Expected length of innings
- Player performance
- Pitch difficulty based on ball tracking
- Demonstrates robust need of data science



CricViz (cont.)



Example - broadcast

PlayViz - fielding
Best fielders this series

	Catches	Drops	Runs impact
Yasir Shah Pak	4	0	58
Alex Hales Eng	5	1	35
Younus Khan Pak	5	3	24
Chris Woakes Eng	3	0	22

sky SPORTS



Common Themes, Best Practices and Challenges

What I try to do and what makes
my job interesting!

- Robust approach
- Pragmatism
- Domain credibility
- Data matching
- Incomplete data
- Model assessment
- Empirical modelling

Theme - Robust practices



Important to adopt robust approaches and utilise best practice

- Use databases
- Invest and take ownership of data processing
 - 50% of work
 - Rubbish in, rubbish out
- Model development framework
 - Repeatable and reproducible
 - Easy and quick to test hypothesis, document and iterate
- Record failed developments as well as successes
- Best software development practices

Theme - Pragmatism



Important to adopt an approach of pragmatic rigour

- Don't discount expert knowledge just because it isn't statistically robust
- Recognise the limitations and assumptions behind a lot of statistical theory
 - Assumptions may not hold in practice
 - Don't discount something just because there is not a mathematical proof
- Be practically rigorous
 - Remember the final goal

Challenge - Domain credibility



- Historically athletes and coaches often typically not highly educated in an academic sense
 - Suspicious of what data can tell them compared to their own expertise
 - Not unique to pro-sports (e.g. betting, or music)
- One shot - unpalatable output is punitive
- In professional sports, sample sizes can be very small
 - Moneyball has helped (and hindered)

Must adopt a professional attitude and obtain gradual buy in

Challenge - Data matching



- Multiple data sources with different identifiers for the same objects
- Linking data sources in real-time has historically been a significant challenge
 - Lots of bespoke algorithms
 - Frustrated by inconsistencies within source
 - Manual matching
- Opportunity for automated machine learning approaches

Challenge - Incomplete data



- Modelling sporting events is often complicated by incomplete data
 - Player / ball positions
 - Defenders non-actions
 - Subjective data
 - Shot quality
- Have to adapt parametric models to take account of these limitations
 - Introduce biases

Challenge Model assessment



It is often hard to decide how to assess if one model is better than another

- Require a robust framework
 - Consistent data set for assessment
 - Reproducible research
 - Part of a wider model development pipeline
 - Document the failures as well as the successes
- How do we measure success?
 - Often lots of competing candidate measures
 - Can we isolate model from the ultimate outcome (e.g. betting)

Challenge - Model assessment (cont.)



Example

- Consider a model for predicting the outcome of a tennis match in-play
 - Simple version is a Markovian based model
 - Can be extended to be non-homogeneous / simulation based
 - Updates after every point
- How do we assess how good the model is
 - At the start?
 - At the end of each set?
 - At the end of each game?
 - Some combination?

Challenge - Model assessment (cont.)



Example (in-play tennis cont.)

- Suppose we can agree on a time frame - what do we measure?
 - RMSPE, not bias - but where is the tradeoff
 - Observed minus expected
 - win probabilities ?
 - set length predictions?
 - What about ultimate outcome?
 - e.g. betting: profit - how do we divorce from strategy?
 - e.g. pro: player rankings - what about outside influences?

Challenge - Model assessment (cont.)



- Rarely a single model that is best
- Must avoid sensitivity to data set selection
- Multiple criteria can lead to multiple models, depending on priority
 - Models can be inconsistent
- We have to be pragmatic
 - Academic theory is hugely useful and where possible should guide us, but often things are messier in practice
 - What is the most important thing that we are trying to measure?
 - Do inconsistent models matter?

Challenge Empirical models



Combining academic statistical modelling theory and domain expertise can be difficult

- Domain expertise is often built up through experience
 - Experts can find it hard to encode
- Where expertise is “encoded” it is often in the form of an empirical model, typically in a spreadsheet
 - Expert inputs
 - Black box constants
 - No “history” of development
 - Very hard to scale and assess
- BUT often very performant and contain valuable insight

Challenge - Empirical models (cont.)



Example - Empirical to Production

- The “expertise”

Challenge - Empirical models (cont.)



Example - Empirical to Production

- Robust data
 - Parsing
 - Processing
 - Wrangling

The screenshot shows the MySQL Workbench interface. The query editor contains the following SQL code:

```
93 SELECT * FROM cricket_models_v1_0.lo_player_viz WHERE game_id = 39623 ORDER BY innings, over, ball;
94 SELECT * FROM cricket_models_v1_0.batting_player_inputs WHERE game_id = 39628;
95
96 SELECT COUNT(b.player_id) AS player_count, SUM(IF(ISNULL(b.impact_or), IF(ISNULL(b.impact), 0, 1), 1)) AS fitted_c
97 MAX(b.recalc) AS recalc
98 FROM cricket_models_v1_0.batting_player_inputs AS b
99 INNER JOIN cricket_project_dev.Game_Player AS gp ON gp.game_id = b.game_id
100 INNER JOIN cricket_project_dev.Player AS p ON p.data_source_id = gp.player_id
101 WHERE b.game_id = 39628 AND b.innings = 0 AND b.batting_order > 0 AND b.player_id = p.id;
102
103 SELECT * FROM Game WHERE id = 39704;
104
105 SELECT * FROM Venue WHERE data_source_id = 262;
106
107 SELECT * FROM cricket_project_probabilities_processing WHERE game_id = 39702;
108 SELECT * FROM cricket_project_probabilities_to_win_new WHERE game_id = 39626;
109 SELECT * FROM cricket_apps_v1_0.v3_by_ball WHERE game_id = 39626;
110
111 SELECT * FROM cricket_project_dev.Game_Player AS gp
112 INNER JOIN cricket_project_dev.Player AS p ON p.data_source_id = gp.player_id
113 WHERE gp.game_id = 39626;
```

The result grid shows the following data:

ngs_id	over_id	ball_no	batsman_id	game_date	team_one_win	team_one_lost	draw	date_refreshed	innings_1_runs	innings_1_wickets	innings_1_overs	inn
0	1	3297	2016-10-20 05:01:45	44.15	34.77	21.08	2016-10-20 05:01:45	451	10	142	362	
0	6	3297	2016-10-20 05:04:20	44.65	34.14	21.21	2016-10-20 05:04:20	451	10	143	362	
1	1	5511	2016-10-20 05:05:54	23.94	25.9	50.16	2016-10-20 05:05:54	461	10	191	361	
1	1	5511	2016-10-20 05:05:54	23.37	28.4	50.23	2016-10-20 05:05:54	456	10	189	362	
1	6	5511	2016-10-20 05:07:53	24.24	25.12	50.64	2016-10-20 05:07:53	462	10	191	361	
2	1	3297	2016-10-20 05:08:51	23.19	25.36	51.45	2016-10-20 05:08:51	464	10	193	362	
2	6	3297	2016-10-20 05:12:03	24.22	25.47	50.31	2016-10-20 05:12:03	464	10	191	362	
3	1	5511	2016-10-20 05:12:53	41.09	19.65	39.26	2016-10-20 05:12:53	463	10	192	322	

The Action Output section shows the following results:

Action	Time	Action	Response	Duration / Fetch Time
1	16:31:21	SELECT * FROM cricket_models_v1_0.t20_resources LIMIT 0, 10000	121 row(s) returned	0.088 sec / 0.00008...
2	17:30:31	Could not connect, server may not be running.	Can't connect to MySQL server on '127.0.0.1' (61)	
3	17:33:03	SELECT * FROM cricket_project_dev.Game_Player AS gp INNER JOIN cricket_...	22 row(s) returned	0.035 sec / 0.00010...
4	17:33:20	SELECT * FROM cricket_project_probabilities_processing WHERE game_id = 3...	960 row(s) returned	0.041 sec / 0.058 sec
5	17:33:30	SELECT * FROM cricket_project_probabilities_to_win_new WHERE game_id = 3...	776 row(s) returned	0.082 sec / 0.061 sec

Challenge - Empirical models (cont.)



Example - Empirical to Production

- Model development framework

```
4 import sqlalchemy as sqlalchemy
5
6 from sdsmodellering.model_runner import ModelRunner
7
8 import cvutils.opta.games as games
9 import cvutils.opta.comps as comps
10
11 import cvmodels.models.limited_over.v1.resources as loresources
12
13 from cvmodels.modelrunners.limited_over.limited_over_model_runner import LimitedOverModelRunner
14 from cvmodels.models.limited_over.v1.win_viz_v1 import LimitedOverWinViz
15
16
17
18 class LimitedOverWinVizModelRunner(ModelRunner, LimitedOverModelRunner):
19     def __init__(self, ds, options):
20         super().__init__(ds, LimitedOverWinViz(), options, self.__get_runs, self.__get_fit_dataset,
21                         self.__get_predict_dataset, self.__get_predict_params,
22                         self.__write_fit, self.__write_predict)
23
24     def __get_runs(self):
25         if 'comp_id' in self._options['dataset']:
26             game_ids = comps.get_competition_games(self._ds.mysql.eng, self._options['dataset']['comp_id'])
27         elif 'game_id' in self._options['dataset']:
28             game_ids = [self._options['dataset']['game_id']]
29         else:
30             game_ids = self.get_covered_games(self._ds.mysql.eng, self._options['dataset']['search_days'],
31                                             [2, 3, 5, 6, 8, 9])
32
33         out = []
34         for gid in game_ids:
35             out = out + self.__get_runs_for_game(gid)
36         return out
```

Challenge - Empirical models (cont.)



Example - Empirical to Production

- An external interface

api.cricviz.com/#!games/get_games

games Show/Hide List Operations Expand Operations

GET /games Get cricket games in a time interval

Implementation Notes

Gets `game` objects, that have been played recently or are being played in near future. If no `days` parameter is passed time interval is 10 days either side of current date.

Optional query param of `days` determines how many days before or after today to look for games. Games can be restricted to a particular competition by passing the optional `comp_id` parameter.

Response Class (Status 200)
Successful response

Model | **Model Schema**

```
  "name": "string"
},
"batting1_team": {
  "id": 0,
  "opta_id": 0,
  "name": "string"
},
"batting2_team": {
  "id": 0,
  "opta_id": 0,
  "name": "string"
}
```

Response Content Type `application/json`

Parameters

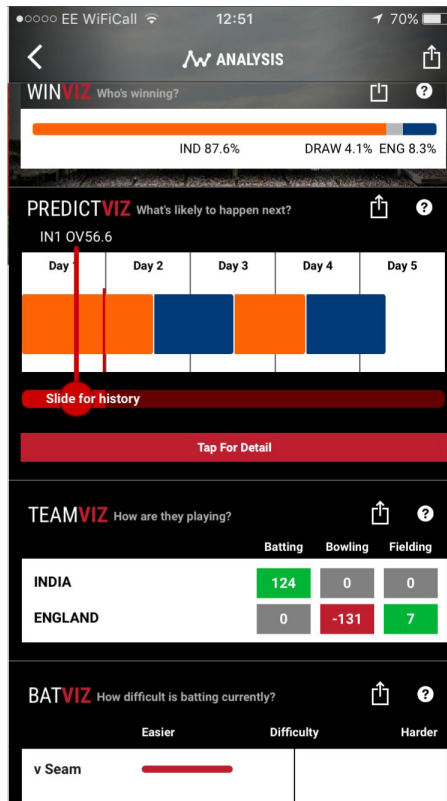
Parameter	Value	Description	Parameter Type	Data Type
days	<input type="text" value="10"/>	Number of days to search for games before and after today. If not supplied, default is 10 days.	query	integer

Challenge - Empirical models (cont.)



Example - Empirical to Production

- The customer view



Challenge - Empirical models (cont.)



Assessing empirical models is even more complicated

- All the usual challenges
- How do you account for the manually inputted domain expertise?
 - Cannot be populated post-hoc
 - Hard to build up a sample that is sufficiently large to robustly assess

Approach similarly to credibility challenges



The future

What we can look forward to

- Sporting Data Science
- Prevalence of data
- Predictive modelling
- Broadthening of scope



My latest project building on career experience to date

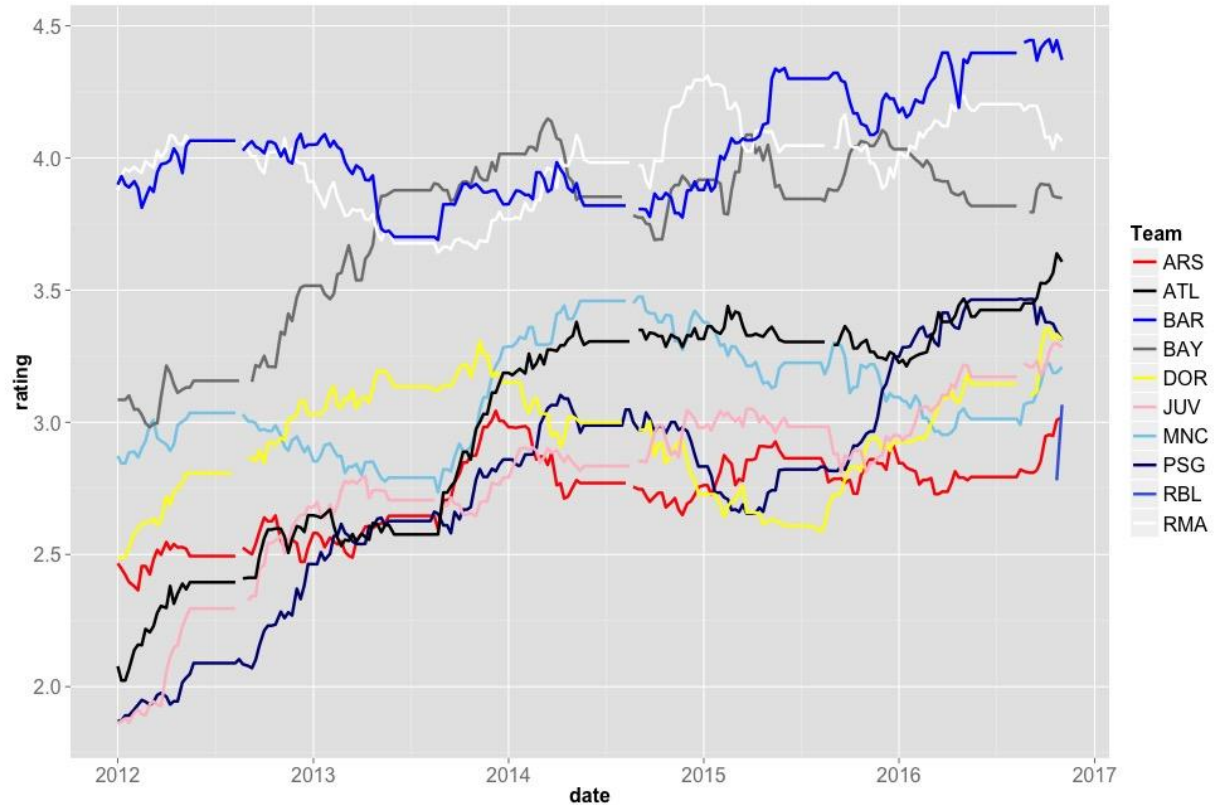
- Consulting to a number of companies
- Traditional predictive modelling for sporting events
- Helping clients with data assets generate a revenue stream from those assets
- Supporting general trend towards data engagement

Sporting Data Science



Example

Team ratings for fan engagement and brands



Prevalence of data



Data is now ubiquitous in sport, as throughout society

- Previously applications were limited by availability of data
- In sport, data has not traditionally been “big data”
 - This has changed with high frequency odds
 - Tracking (ball, player) data
- Big data sets lend themselves to different methods of analysis
- General increase in societal data literacy

Predictive modelling in sport



Role of predictive modelling is changing

- Typically parametric models have been used
 - Interpretable (critically important when data contains a lot of noise)
- Marginal gains in a lot of areas are much diminished but useful in other domains
 - e.g. Betting vs Broadcast
- Machine learning / artificial intelligence approaches can deal with correlations and data dimension better
 - Computer scientist vs statisticians
 - Challenge is to retain interpretability

Breadthening of scope



Changing nature of data means different applications

- Player and ball tracking data
 - Visualisations
 - Coaching applications
- InCrowd Sports
 - Fan engagement platform
 - Data regarding persons movements around match events
- Sportr
 - Natural language processing - summarising multiple news sources



Any questions?

dave.hastie@sportingds.com

- RSS Statistics in Sport section
- Thank you
 - Dixon and Coles, 1997
 - Smartodds
 - Stratagem Technologies
 - CricViz