# Statistics in Sports: Football (Soccer) Overview

ZACHARY BINNEY, PHD MPH

OXFORD COLLEGE OF EMORY UNIVERSITY

FALL 2023

# Disclosure

- Not a soccer analytics guy


- Incomplete, may not even focus on most important current topics

# Roadmap

- 1. Intro and Types of Soccer Data
  - Events vs. Tracking

- 2. Expected Goals (xG) Models

- 3. Beyond Shots

- 4. Player Evaluation
  - Radar Charts

# Intro and Types of Soccer Data

# Intro to Soccer Analytics

- Let's watch a goal...        But in reverse...

- Highlights from Atlanta United-Orlando City FC match, October 2021

- Video in lecture folder

# Intro to Soccer Analytics

- Let's watch Atlanta United's first match and first goal ever (!)

  - [2017 vs. NY Red Bull, start around 26:00](#)

  - How might you break this down into data a computer could analyze?

# Intro to Soccer Analytics

- What if you paused every time something "interesting" or "noteworthy" happened, and logged that?

  - **Events Data**

JUN. 10, 2014, AT 3:58 PM

## The People Tracking Every Touch, Pass And Tackle in the World Cup

By Carl Bialik

moods and preferences. Throughout the year, 350 part-time analysts working in London and a half-dozen other Opta branches in Europe and North and South America record every pass, header and goal while watching live or recorded video of more than 14,000 matches around the world. The

team's match to confirm every goal was attributed correctly. And I watched as Opta's media team processed the raw numbers — 1,600 to 2,000 events per game — into TV-ready factoids, which they heard commentators repeat

But most of the work is logging routine passes. Opta's analysts log each one by dragging and clicking a mouse at the spot where the pass was received, then keying in the player who received it. Their monitors have an image of a

# Intro to Soccer Analytics

- **Events Data** example: Statsbomb (105 variables)

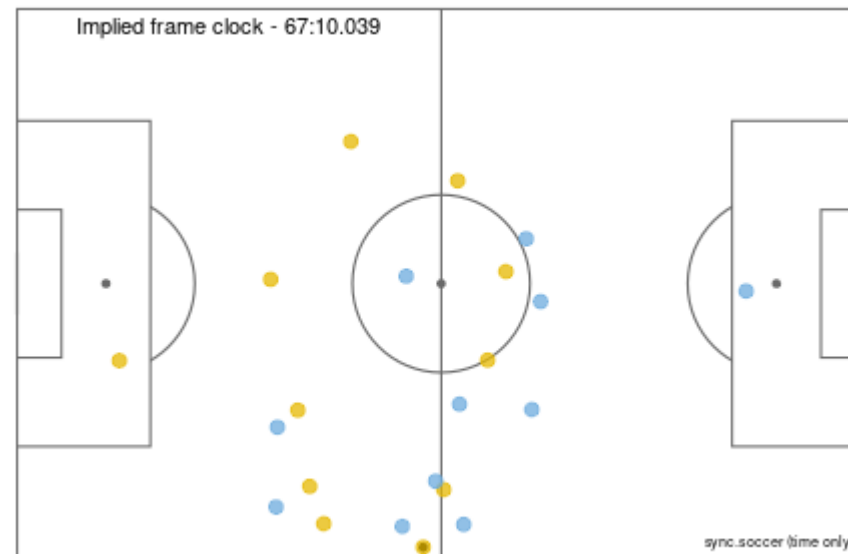| | period | minute | second | possession | possession_team.name | player.name | type.name | position.name | duration | location |
|---|---|---|---|---|---|---|---|---|---|---|
| 4 | 1 | 0 | 0 | 1 | Houston Dash | NA | Half Start | NA | 0.000 | NULL |
| 5 | 1 | 0 | 0 | 2 | Utah Royals | Diana Matheson | Pass | Center Attacking Midfield | 1.204 | c(60, 40) |
| 6 | 1 | 0 | 1 | 2 | Utah Royals | Katrina Gorry | Ball Receipt* | Right Attacking Midfield | NA | c(53, 39) |
| 7 | 1 | 0 | 1 | 2 | Utah Royals | Katrina Gorry | Pass | Right Attacking Midfield | 3.070 | c(93, 18) |
| 8 | 1 | 0 | 4 | 3 | Houston Dash | Amber Brooks | Pass | Right Center Back | 1.372 | c(28, 63) |
| 9 | 1 | 0 | 6 | 3 | Houston Dash | Kealia Ohai | Ball Receipt* | Right Center Midfield | NA | c(49, 71) |
| 10 | 1 | 0 | 6 | 3 | Houston Dash | Kealia Ohai | Carry | Right Center Midfield | 3.720 | c(49, 71) |
| 11 | 1 | 0 | 7 | 3 | Houston Dash | Katrina Gorry | Pressure | Right Attacking Midfield | 2.292 | c(68, 14) |
| 12 | 1 | 0 | 9 | 3 | Houston Dash | Kealia Ohai | Pass | Right Center Midfield | 2.000 | c(63, 74) |
| 13 | 1 | 0 | 11 | 3 | Houston Dash | Nichelle Patrice Prince | Ball Receipt* | Right Center Forward | NA | c(105, 70) |
| 14 | 1 | 0 | 11 | 3 | Houston Dash | Nichelle Patrice Prince | Carry | Right Center Forward | 3.200 | c(105, 70) |
| 15 | 1 | 0 | 15 | 3 | Houston Dash | Nichelle Patrice Prince | Pass | Right Center Forward | 0.804 | c(118, 68) |
| 16 | 1 | 0 | 15 | 3 | Houston Dash | Rachel Daly | Ball Receipt* | Left Center Forward | NA | c(115, 45) |
| 17 | 1 | 0 | 15 | 3 | Houston Dash | Rachel Corsie | Clearance | Left Center Back | 0.000 | c(5, 33) |
| 18 | 1 | 0 | 42 | 4 | Houston Dash | Sofia Huerta | Pass | Center Attacking Midfield | 2.193 | c(120, 80) |
| 19 | 1 | 0 | 44 | 4 | Houston Dash | Diana Matheson | Clearance | Center Attacking Midfield | 0.000 | c(12, 47) |
| 20 | 1 | 0 | 47 | 4 | Houston Dash | Diana Matheson | Pressure | Center Attacking Midfield | 0.649 | c(9, 51) |
| 21 | 1 | 0 | 47 | 4 | Houston Dash | Nichelle Patrice Prince | Pass | Right Center Forward | 1.893 | c(108, 21) |
| 22 | 1 | 0 | 49 | 4 | Houston Dash | Rebecca Elizabeth Sauerbrunn | Clearance | Right Center Back | 0.000 | c(13, 48) |

# Intro to Soccer Analytics

- **Events Data** example: Statsbomb (105 variables)

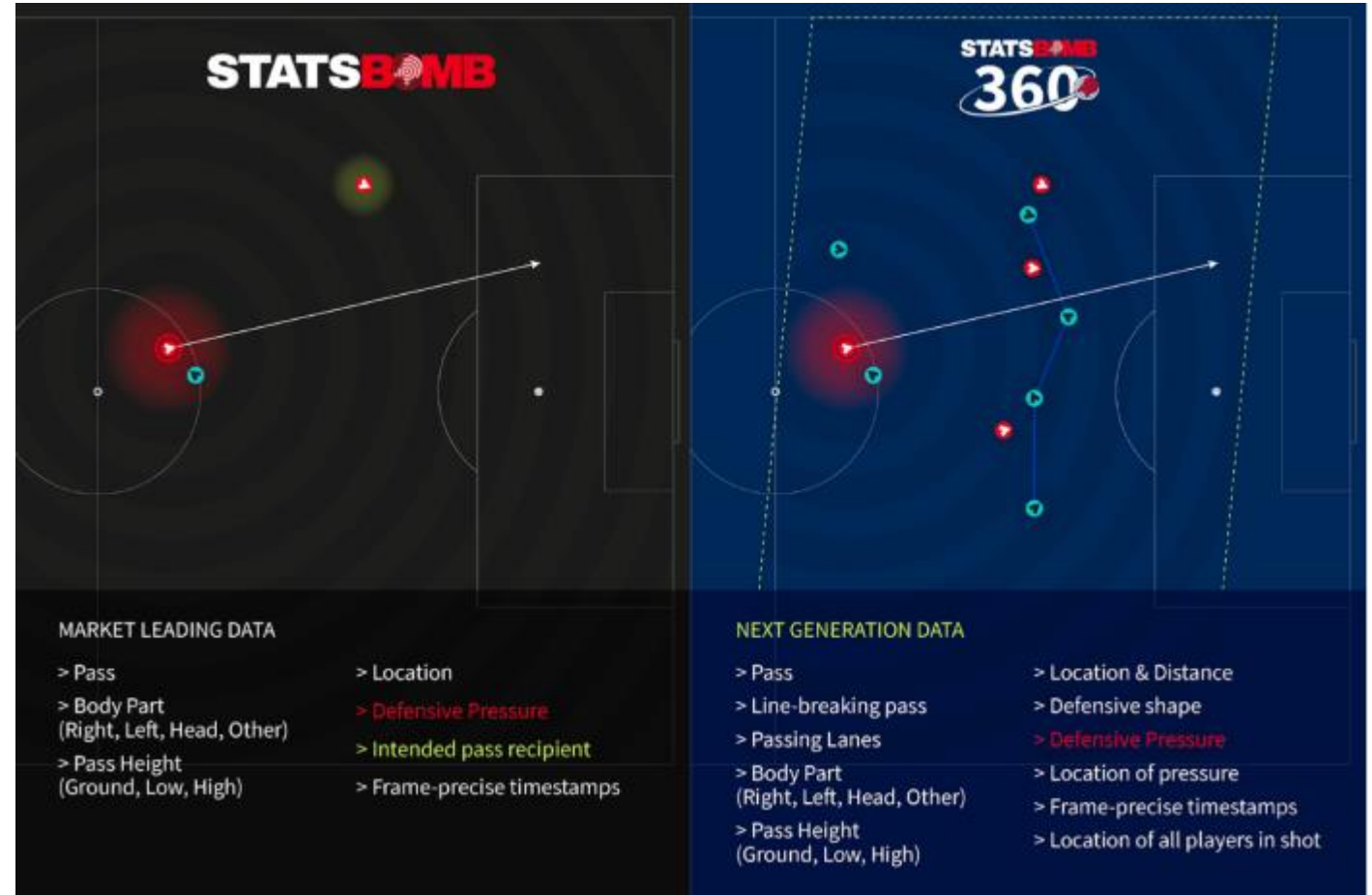| play_pattern.name | team.name | pass.length | pass.angle | pass.end_location | pass.cross | pass.switch | pass.aerial_won | pass.assisted_shot_id |
|---|---|---|---|---|---|---|---|---|
| Regular Play | Houston Dash | NA | NA | NULL | NA | NA | NA | NA |
| From Kick Off | Utah Royals | 7.071068 | -2.99969550 | c(53, 39) | NA | NA | NA | NA |
| From Kick Off | Utah Royals | NA | NA | NULL | NA | NA | NA | NA |
| From Kick Off | Utah Royals | 0.000000 | 0.00000000 | c(93, 18) | NA | NA | NA | NA |
| Regular Play | Houston Dash | 22.472204 | 0.36397895 | c(49, 71) | NA | NA | NA | NA |
| Regular Play | Houston Dash | NA | NA | NULL | NA | NA | NA | NA |
| Regular Play | Houston Dash | NA | NA | NULL | NA | NA | NA | NA |
| Regular Play | Utah Royals | NA | NA | NULL | NA | NA | NA | NA |
| Regular Play | Houston Dash | 42.190044 | -0.09495170 | c(105, 70) | NA | NA | NA | NA |
| Regular Play | Houston Dash | NA | NA | NULL | NA | NA | NA | NA |
| Regular Play | Houston Dash | NA | NA | NULL | NA | NA | NA | NA |
| Regular Play | Houston Dash | 20.099750 | -1.67046500 | c(116, 48) | TRUE | NA | NA | NA |
| Regular Play | Houston Dash | NA | NA | NULL | NA | NA | NA | NA |
| Regular Play | Utah Royals | NA | NA | NULL | NA | NA | NA | NA |
| From Corner | Houston Dash | 47.296936 | -1.80551900 | c(109, 34) | NA | TRUE | NA | NA |
| From Corner | Utah Royals | NA | NA | NULL | NA | NA | NA | NA |
| From Corner | Utah Royals | NA | NA | NULL | NA | NA | NA | NA |
| From Corner | Houston Dash | 12.000000 | 1.57079640 | c(108, 33) | TRUE | NA | NA | NA |
| From Corner | Utah Royals | NA | NA | NULL | NA | NA | NA | NA |

# Events vs. Tracking Data

- **Events Data** is distinct from and simpler than…

  - **Tracking Data**, which is similar to that in (American) football and comes from GPS chips, RFID tags, and/or cameras tracking X-Y location, speed, and direction of every player and ball multiple times per second
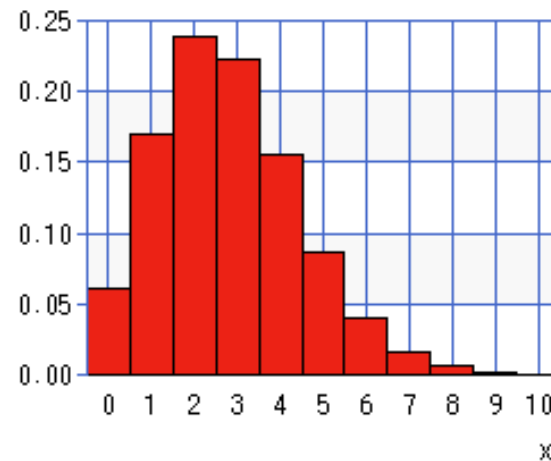
# Middle Ground?

- Context allows for better assessment of decision-making (e.g. was a better passing lane open, did a pass beat the defense's line, etc.)

- Tracking data can grant even more (what happens between events), but...
  - Requires tech + data investments
  - Voluminous and unwieldy
  - Have to figure out value and how to analyze

# Expected Goal (xG) Models

# xG Models

- How do you win a soccer match?

- Problem: (actual) **goals** are rare, noisy
  - Poisson distribution with mean of 2.8 goals (roughly Norwich City, 2021/22):



- To reduce noise, instead of how a team did look at how they *should* have done

# xG Models

- **Expected Goals (xG) Models**. Based on:

**Location of shooter**: How far was it from the goal and at what angle on the pitch?
**Body part**: Was it a header or off the shooter's foot?
**Type of pass**: Was it from a through ball, cross, set piece, etc?
**Type of attack**: Was it from an established possession? Was it off a rebound? Did the defense have time to get in position? Did it follow a dribble?

the location of all players on the pitch at the moment the shot was taken. Was the goalkeeper in position? Was it an open goal or were there a number of defenders between the shooter and the goal? Was the shooter being pressured? Was it a 1v1 situation with the keeper?

- Average result (% of goals that went in) in similar situations → xG for that shot
  - What method do you think is (or could be) used to estimate xG values?

- How do you think each of these affects xG?

# xG Models: Location

- **xG** models and position on pitch

# xG Models: Uses

- **xG** model uses



Sam Kerr, Shot Map
FA Women's Super League, 2020/21

● Head  ▲ Left Foot  ◆ Right Foot

Expected Goals Value

0.0  0.2  0.4  0.6  0.8

STATSBOMB

# xG Models: Uses

- **xG** model uses

# xG Models: Uses

- **xG** model uses

# xG Models: Uses

- Can add in **expected assists (xA)** for more holistic picture



**Expected Goal Contribution**
FA Women's Super League, 2020-21

STATSBOMB

| Player | |
|---|---|
| Samantha May Kerr | |
| Vivianne Miedema | |
| Janine Elizabeth Beckie | |
| Francesca Kirby | |
| Bethany England | |
| Chloe Kelly | |
| Tobin Powell Heath | |
| Caitlin Jade Foord | |
| Lauren Hemp | |
| Jordan Nobbs | |
| Pernille Mosegaard Harder | |
| Christen Annemarie Press | |
| Samantha June Mewis | |
| Jill Roord | |
| Martha Thomas | |

Per 90

■ NPxG  ■ xG Assisted

Minimum 600 minutes
NPxG = Value of shots taken (no penalties)
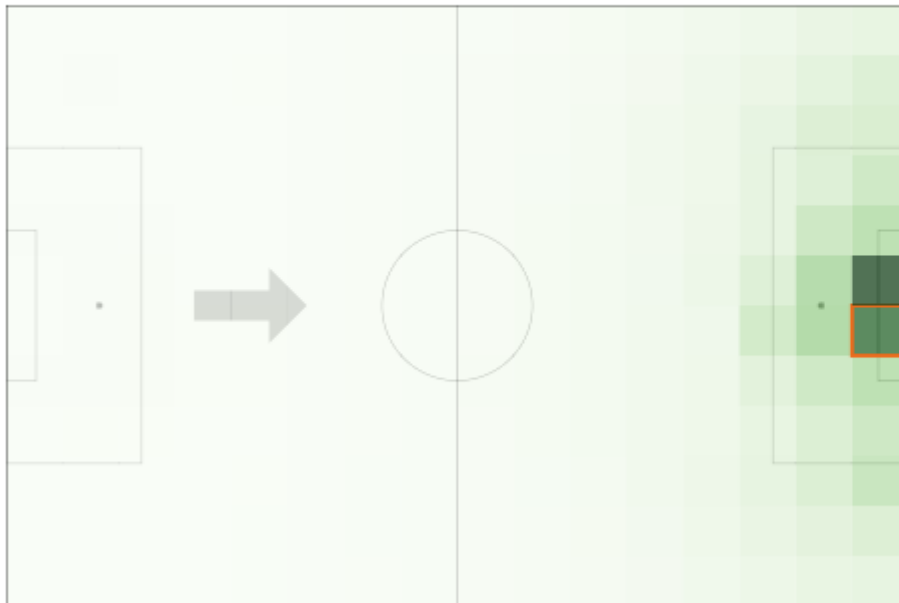xG assisted = Value of shots assisted

# xG Models: Extensions

- Related model: **expected thread (xT)**



Expected Threat (xT) = 0.371

i.e. when the team has the ball in the highlighted zone, they will score in the next **5** actions **37.1%** of the time.
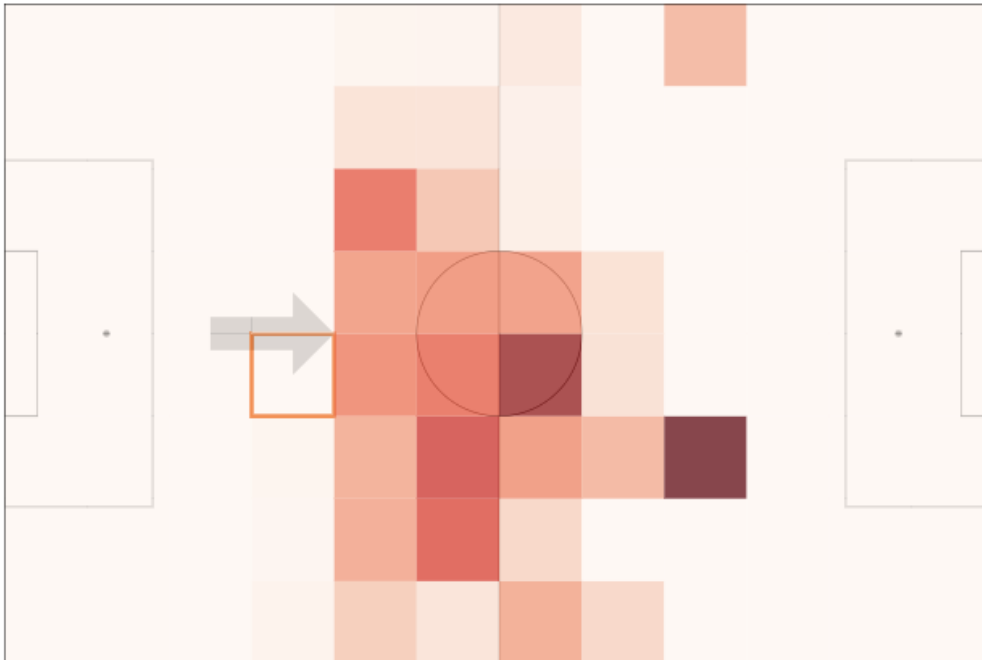
# xG Models: Extensions

- Related model: **expected thread (xT)**. Take a step back...

# xG Models: Extensions

- Ultimate goal: **possession value** to be able to estimate value of individual actions on the field regardless of position

- Can you link this quest with any concepts we've seen in other sports?

**Soccer possession models are gaining steam**

Key soccer possession models by publication year, with type of model and possession information
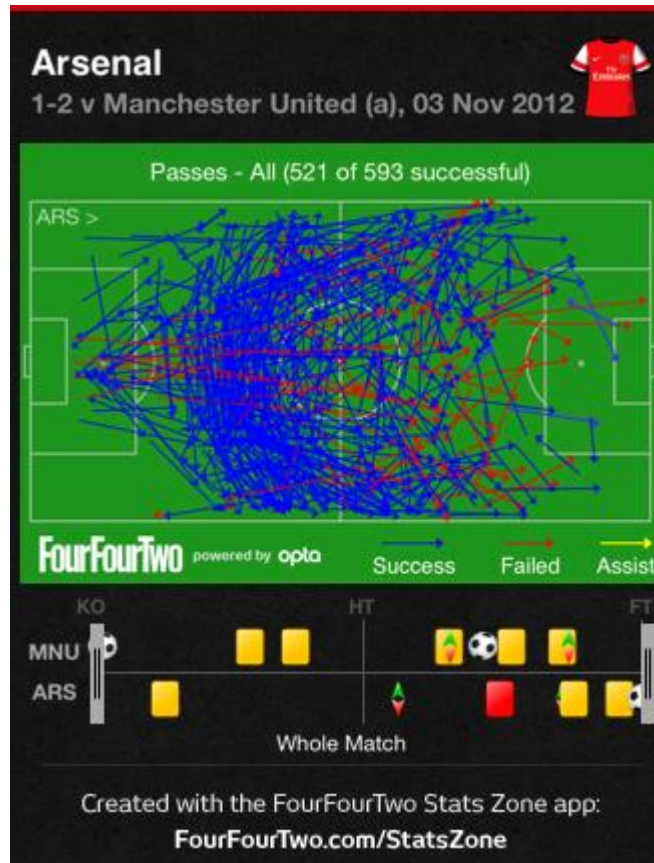
| NAME | CREATOR | DEBUT | METHOD | WINDOW | OFF-BALL INFORMATION |
|------|---------|-------|--------|--------|----------------------|
| Markov Chains | S. Rudd | 2011 | Markov chain | One possession | Defensive states tagged in event data |
| Possession-Based Model | N. Mackay | 2016 | Logistic regression and GAM | One possession | None |
| Expected Threat (xT) | K. Singh | 2019 | Markov-like | Next 5 actions (goal for) | None |
| Valuing Actions by Estimating Probabilities (VAEP) | KU Leuven DTAI | 2019 | Gradient-boosted trees | Next 10 actions (goal for or against) | Possession history proxies |
| Expected Possession Value (EPV) | J. Fernández et al. | 2019 | Multiple models | Next goal (for or against) or end of half | Full tracking data |
| Possession Value (PV) | Stats Perform | 2019 | Gradient-boosted trees | Next 10 seconds (goal for) | Possession history proxies |
| Goals Added (g+) | American Soccer Analysis | 2020 | Gradient-boosted trees | Two possessions | Possession history proxies |
| On-Ball Value (OBV) | StatsBomb | 2021 | Gradient-boosted trees | Two possessions | Broadcast freeze frames (in development) |

FiveThirtyEight

# Beyond Shots

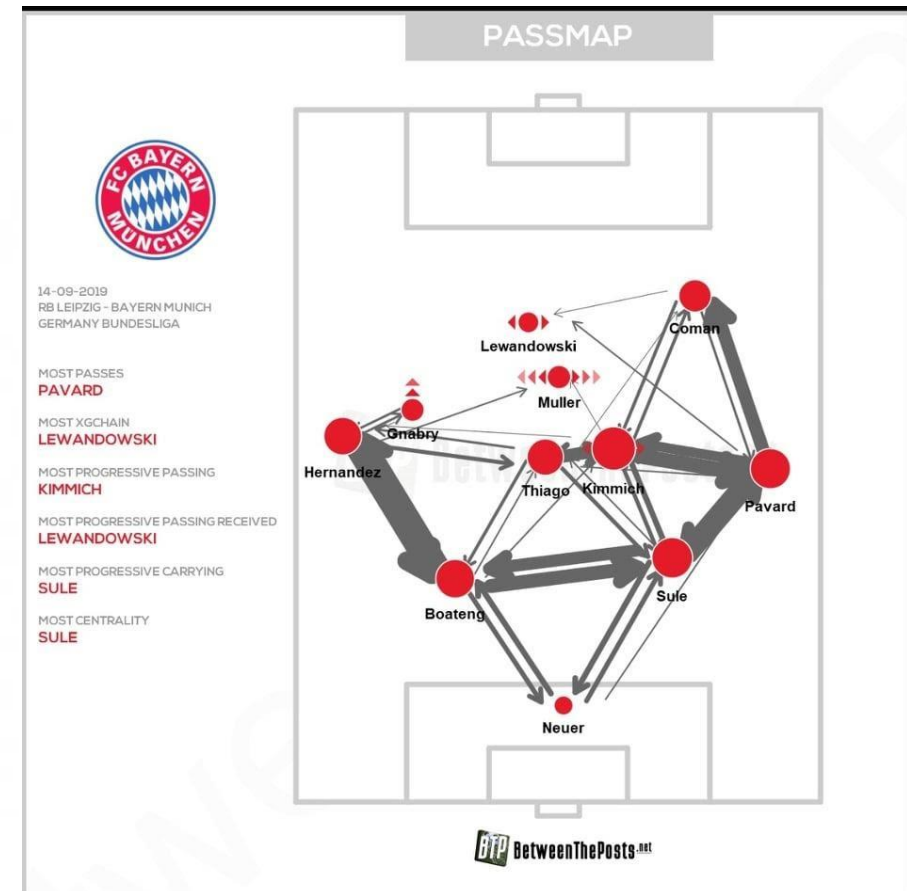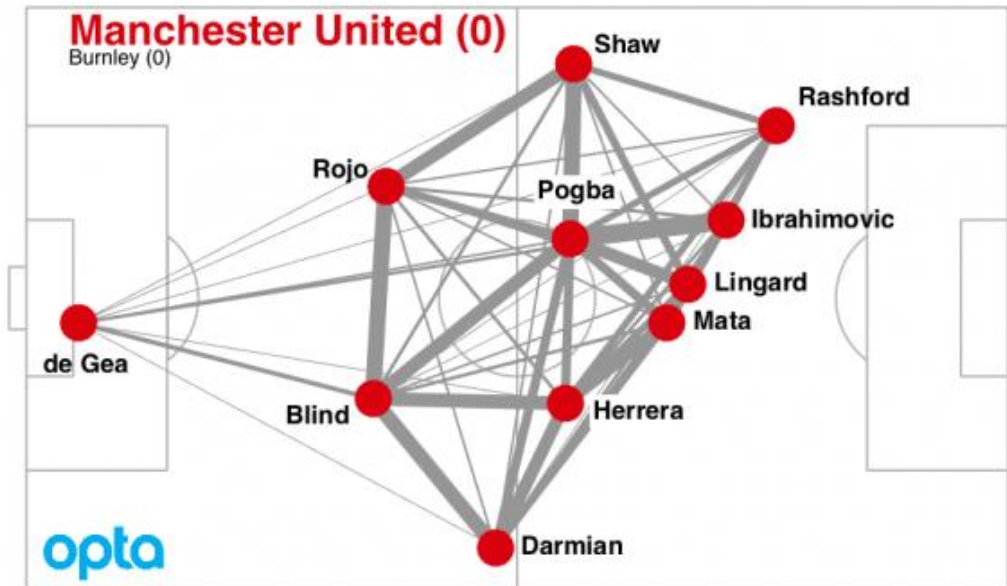# Passing

- "Raw" passing data map

# Passing

- Passing Networks

# Passing

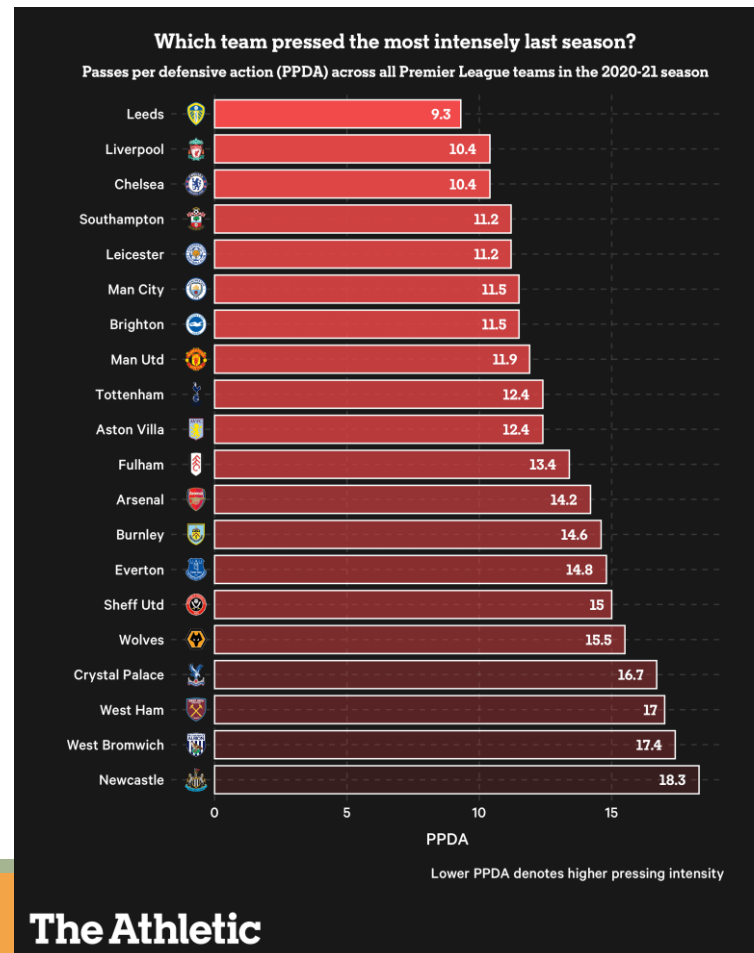- Passing Networks + value of possessions featuring various networks

# Field Tilt and Territorial Dominance

- Field tilt: % of total team + opponent passes in "final third" done by team



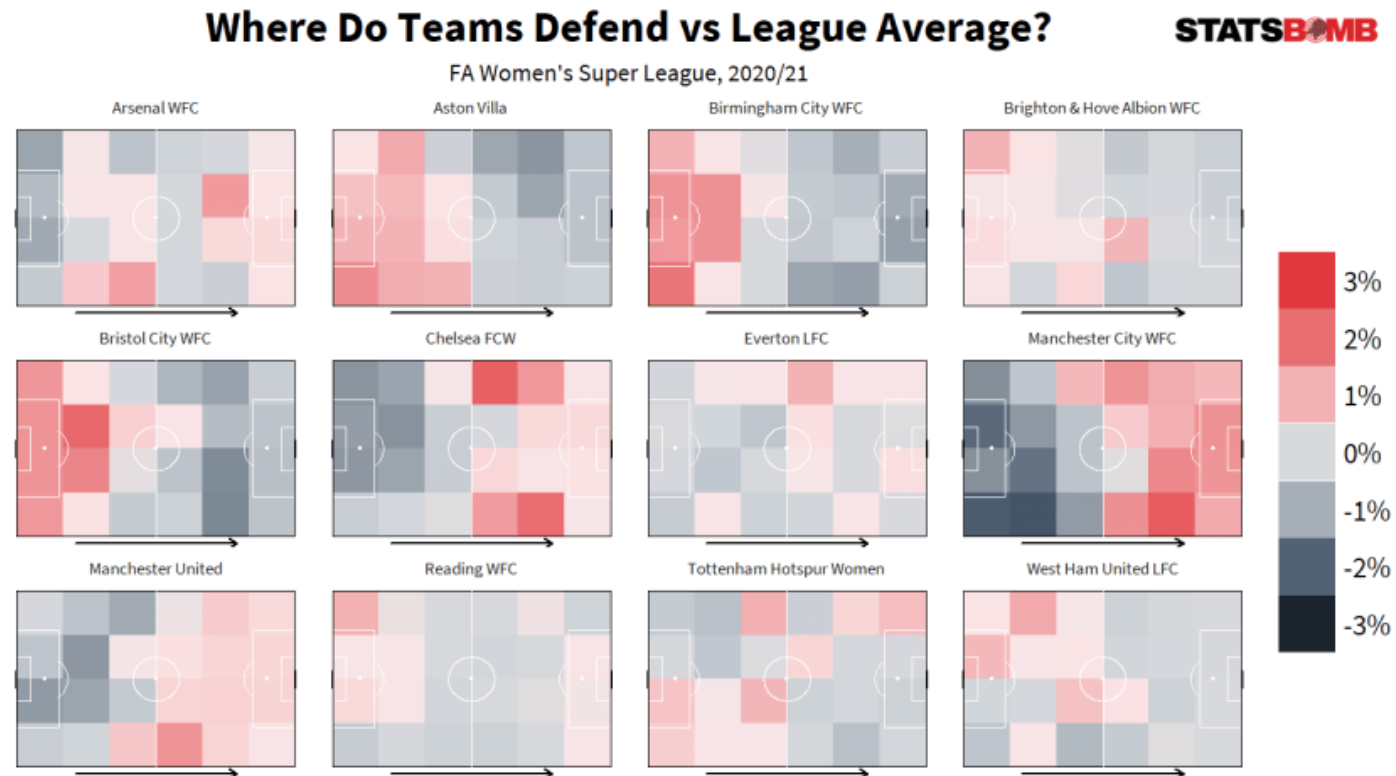**Who had the most territorial dominance last season?**
Field tilt across all Premier League teams in the 2020-21 season

| Team | Field tilt |
|---|---|
| Man City | 70.3% |
| Liverpool | 67.4% |
| Chelsea | 59.9% |
| Man Utd | 59.3% |
| Arsenal | 54.3% |
| Leicester | 53.4% |
| Leeds | 52.5% |
| Brighton | 52.1% |
| Aston Villa | 49.5% |
| Fulham | 49% |
| Sheff Utd | 47.2% |
| Southampton | 46.9% |
| Tottenham | 45.7% |
| Burnley | 45.4% |
| Wolves | 44.7% |
| West Ham | 42.9% |
| Everton | 42.1% |
| Crystal Palace | 39.8% |
| West Bromwich | 38% |
| Newcastle | 35% |

**The Athletic**

# DE-FENSE!

- How aggressive are defenses? **Passes per defensive action (PPDA)**



Which team pressed the most intensely last season?
Passes per defensive action (PPDA) across all Premier League teams in the 2020-21 season

| Team | PPDA |
|------|------|
| Leeds | 9.3 |
| Liverpool | 10.4 |
| Chelsea | 10.4 |
| Southampton | 11.2 |
| Leicester | 11.2 |
| Man City | 11.5 |
| Brighton | 11.5 |
| Man Utd | 11.9 |
| Tottenham | 12.4 |
| Aston Villa | 12.4 |
| Fulham | 13.4 |
| Arsenal | 14.2 |
| Burnley | 14.6 |
| Everton | 14.8 |
| Sheff Utd | 15 |
| Wolves | 15.5 |
| Crystal Palace | 16.7 |
| West Ham | 17 |
| West Bromwich | 17.4 |
| Newcastle | 18.3 |

Lower PPDA denotes higher pressing intensity

**The Athletic**

# DE-FENSE!

- Defending zones, where teams commit defensive actions



**Where Do Teams Defend vs League Average?** STATSBOMB

FA Women's Super League, 2020/21

# Player Evaluation

RADAR CHARTS

# Player Eval: Radars

- Radar Charts from Statsbomb

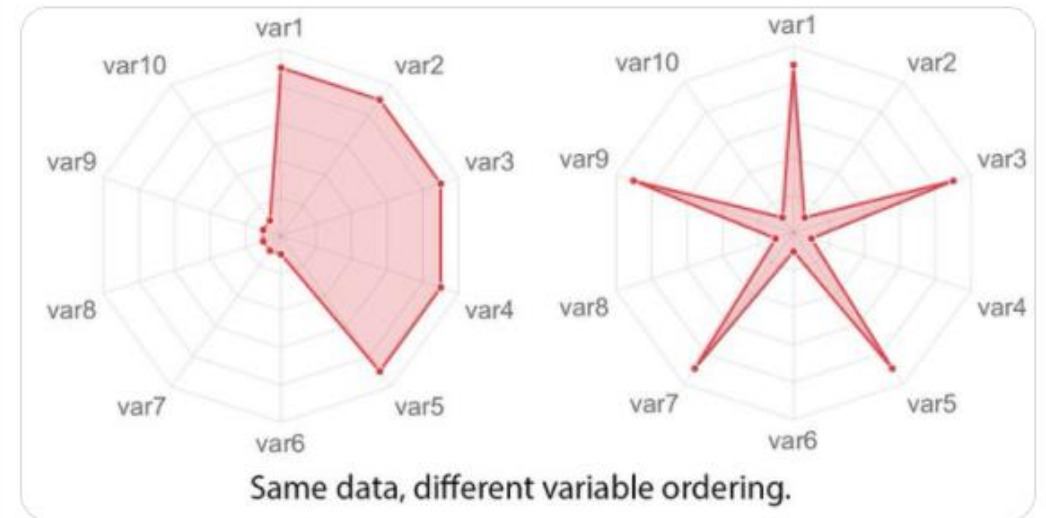- Let's break one of these down…

- Limitations?

# Player Eval: Radars

- Radar Charts from Statsbomb

- Limitations?



Luke Bornn
@LukeBornn

A reminder, blatantly plagiarized from @stat_sam, of why radar plots are misleading. Eye focuses on area, not length.

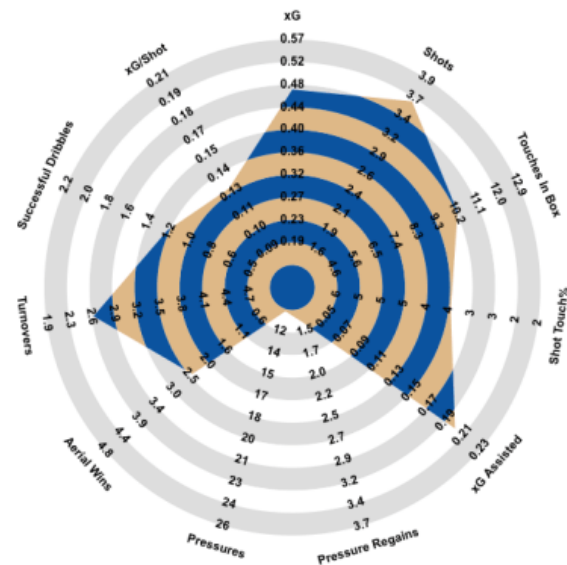Same data, different variable ordering.

10:53 AM · May 17, 2017 · Twitter Web Client

# Player Eval: Radars

- Radar Charts from Statsbomb

# Player Eval: Radars



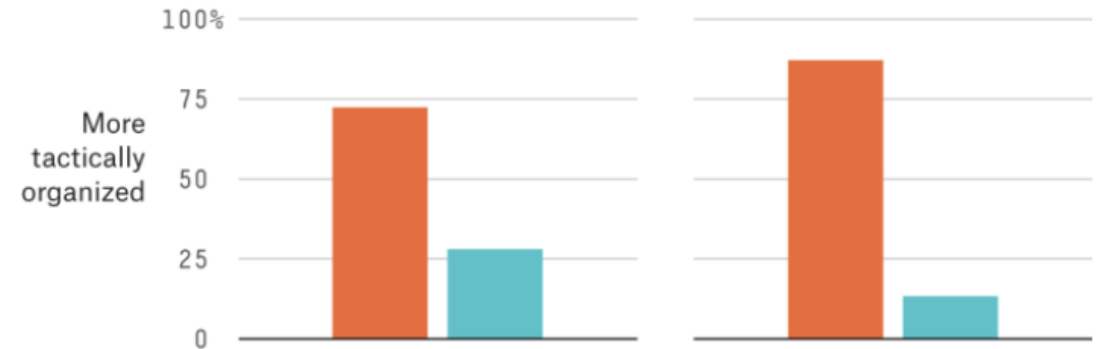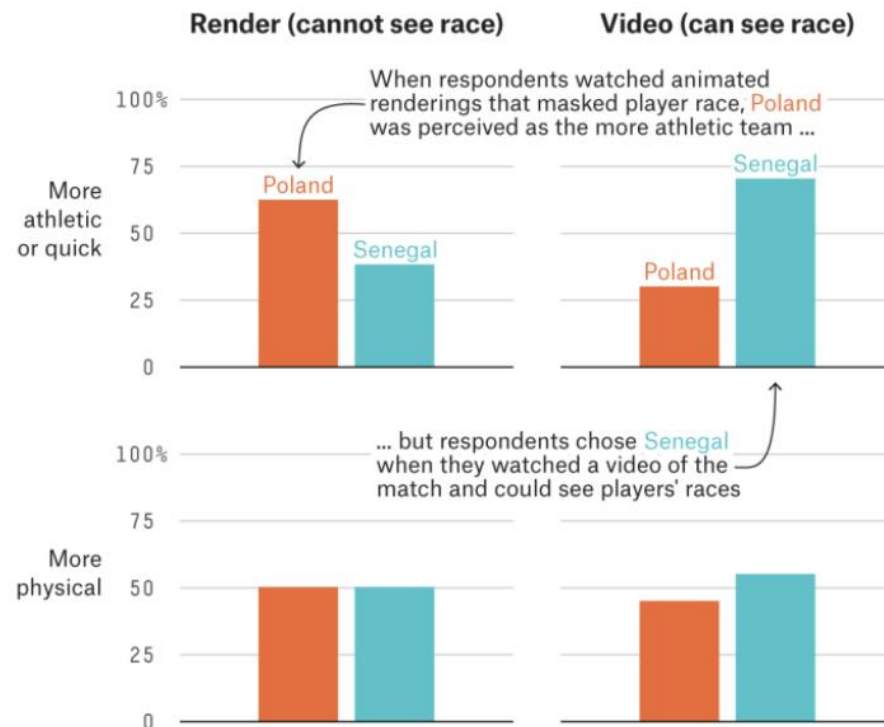Source: StatsBomb
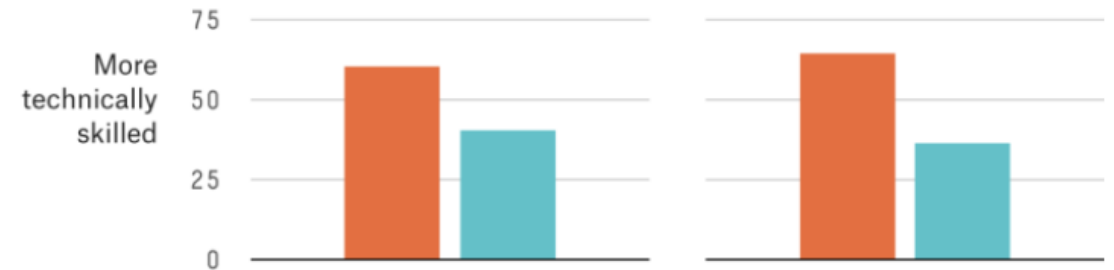
# A Brief Word on Soccer and Racism

# Racial Stereotypes in Soccer (Broadcasting, Scouting)



**Opinions of teams changed when viewers couldn't see race**

Share of respondents reporting whether Senegal or Poland better matched certain playing style characteristics, by whether the respondent watched a broadcast of their game or a two-dimensional render of it

Render (cannot see race) — Video (can see race)

When respondents watched animated renderings that masked player race, Poland was perceived as the more athletic team ...

... but respondents chose Senegal when they watched a video of the match and could see players' races

More athletic or quick / More physical / More tactically organized / More technically skilled

Respondents watched the June 19, 2018, World Cup match between Senegal and Poland.

FiveThirtyEight

SOURCE: GREGORY ET AL.

# More Resources

# More Resources

- Books
  - *The Numbers Game* (Anderson and Sally)
  - *Soccermatics* (Sumpter)

- Companies/Blogs/Video Series
  - Opta
  - Statsbomb (check out their Academy blog posts; also has some public data!)
    - https://statsbomb.com/what-we-do/hub/free-data/
  - Friends of Tracking on YouTube
  - FiveThirtyEight.com, Soccer tag

- Programming
  - @FC_rstats, shaker and worldfootballR package
  - @FCPython, mplsoccer package

- A *million* online resources, analysts, etc. Build trusted network as you would with other sports.

# Thanks!

- Questions? [zbinney@emory.edu](mailto:zbinney@emory.edu), @binney_z on Twitter