

CodeMotion 2017 – Madrid, Spain

Overfitting:

Top 10 Machine Learning Gotchas

@DynamicWebPaige



Paige Bailey

Sr. Cloud Developer Advocate

Work Experience

- Focus at Microsoft is *machine learning* and *artificial intelligence*.
- Prior to joining Microsoft, was a *data scientist* and *geophysical application developer* in the energy industry for 5 years.
- *GIS Technician* (Esri products) for two years.

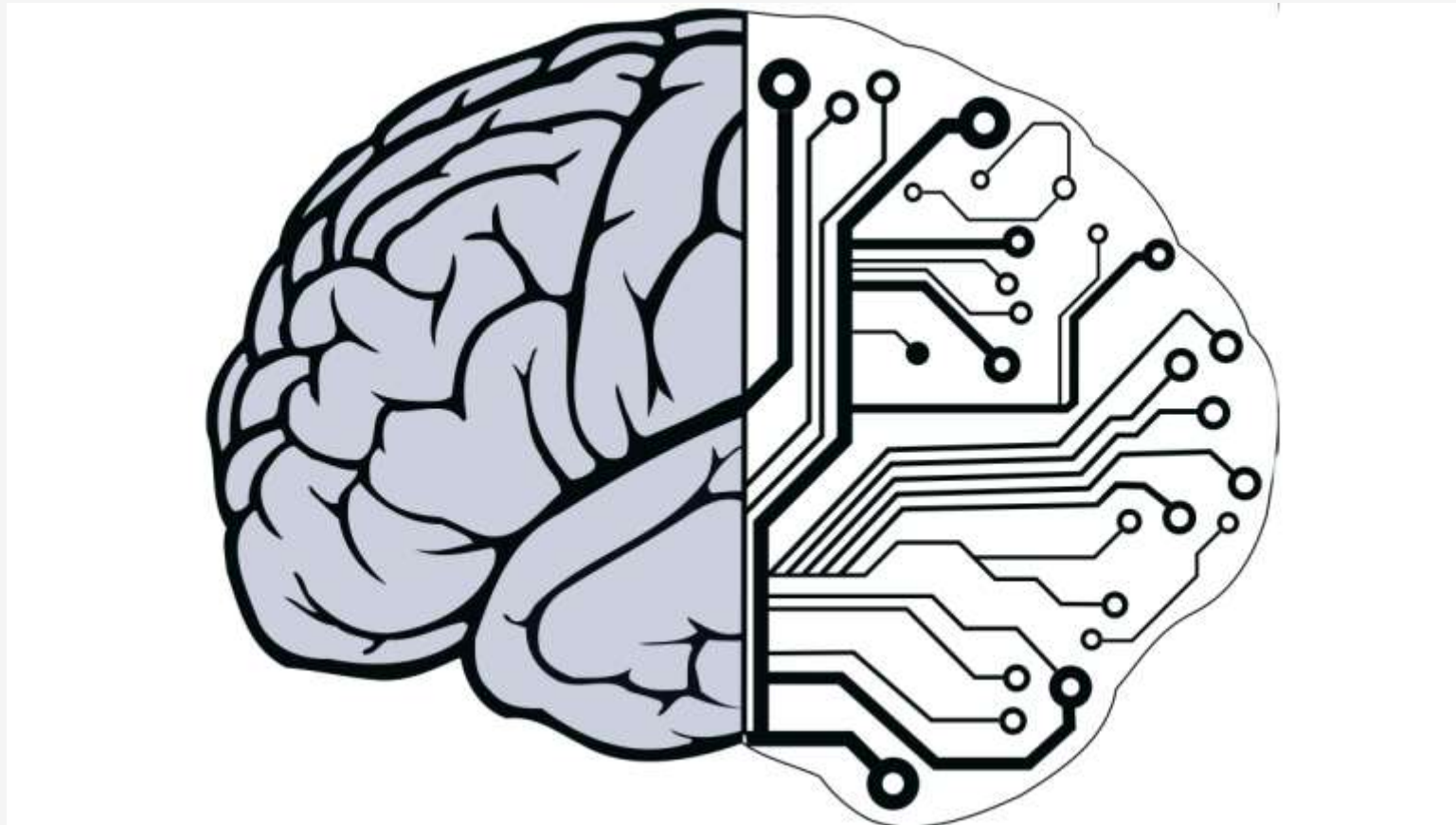
Toolkit

- Python (*10 years*)
- R (*4 years*)
- Spark, Kafka, Hive, HBase (*2 years*)

Location: *Austin, TX*

Twitter: [@DynamicWebPaige](https://twitter.com/DynamicWebPaige)



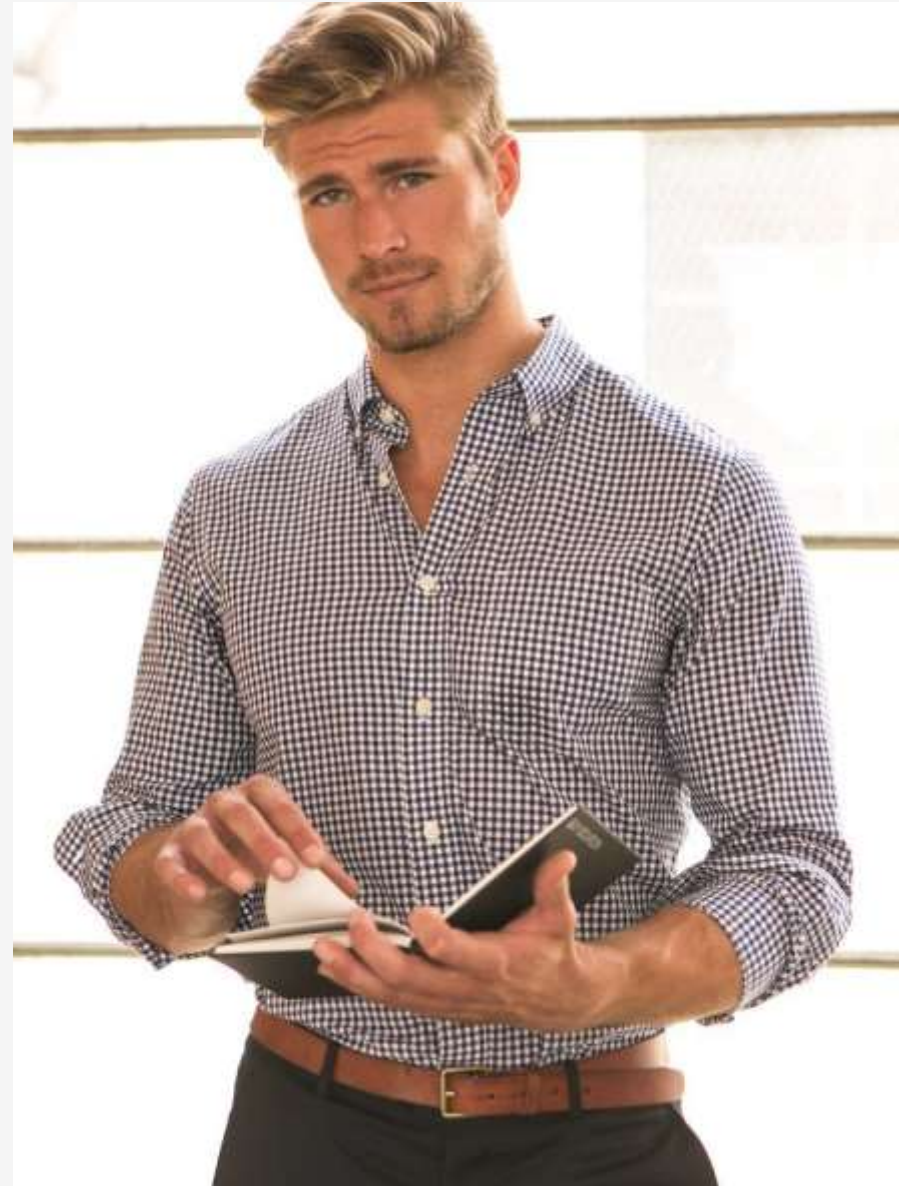


What is Machine Learning?

An example:

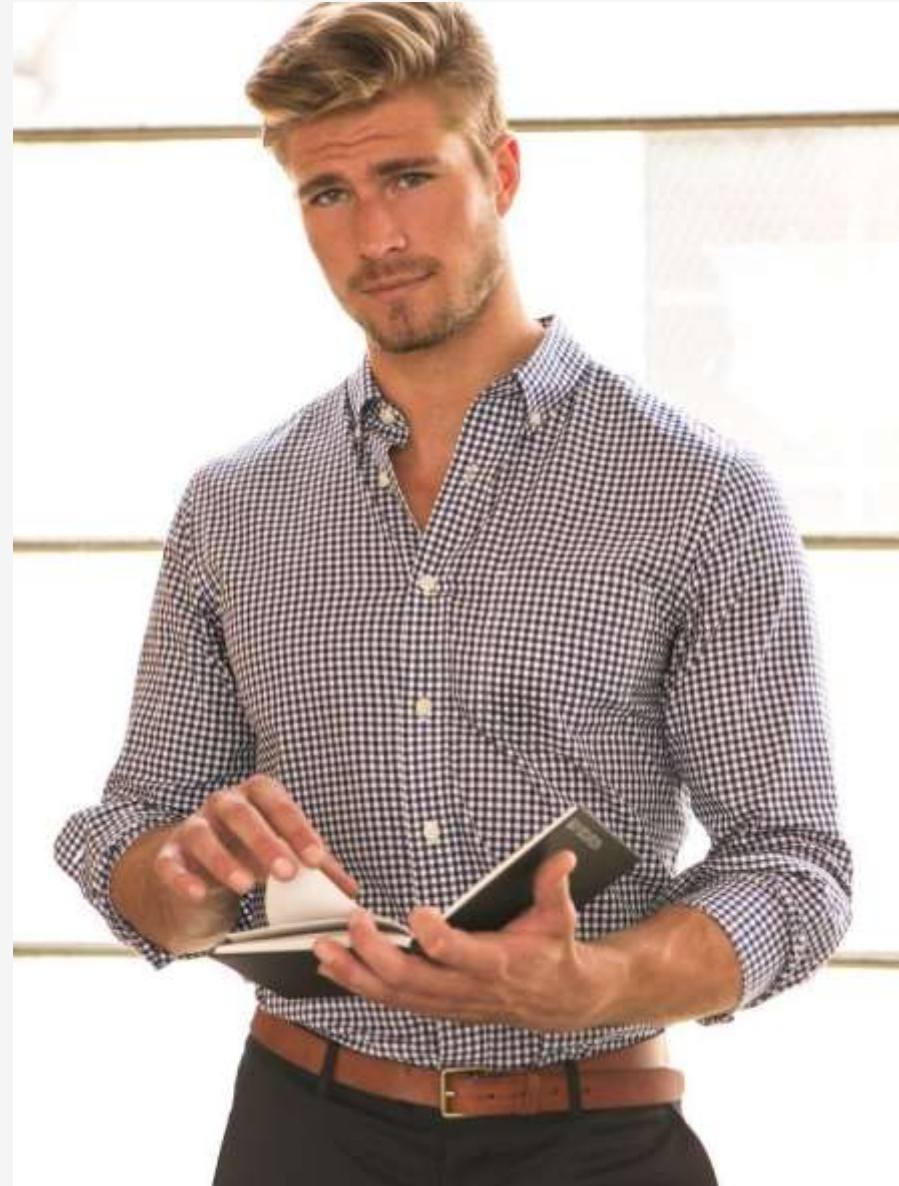
Credit Risk

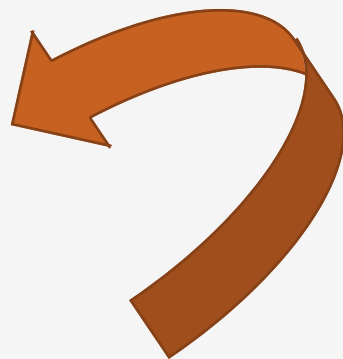
Sven →



We know some things about Sven:

- Age: 32
- Gender: M
- Credit Score: 1
- Homeowner: N
- Account balance: -150
- Employment status: Y
- Education: BS
- Current debt: 2000





Maria



Jeff





Maria

...and we have
similar
information
about each of
them.



Jeff

Name	Age	Gender	Credit Score	Homeowner	Checking Account	Employment Status	Education	Current Debt
Sven	32	M	1	N	-150	Y	BS	2000
Jeff	29	M	4	Y	700	Y	MS	250
Maria	42	F	5	Y	1500	Y	MBA	1500
...
...
...
...
...
...

RECORD



Name	Age	Gender	Credit Score	Homeowner	Checking Account	Employment Status	Education	Current Debt
Sven	32	M	1	N	-150	Y	BS	2000
Jeff	29	M	4	Y	700	Y	MS	250
Maria	42	F	5	Y	1500	Y	MBA	1500
...
...
...
...
...
...

RECORD

FEATURE



Name	Age	Gender	Credit Score	Homeowner	Checking Account	Employment Status	Education	Current Debt
Sven	32	M	1	N	-150	Y	BS	2000
Jeff	29	M	4	Y	700	Y	MS	250
Maria	42	F	5	Y	1500	Y	MBA	1500
...
...
...
...
...
...

RECORD

FEATURE

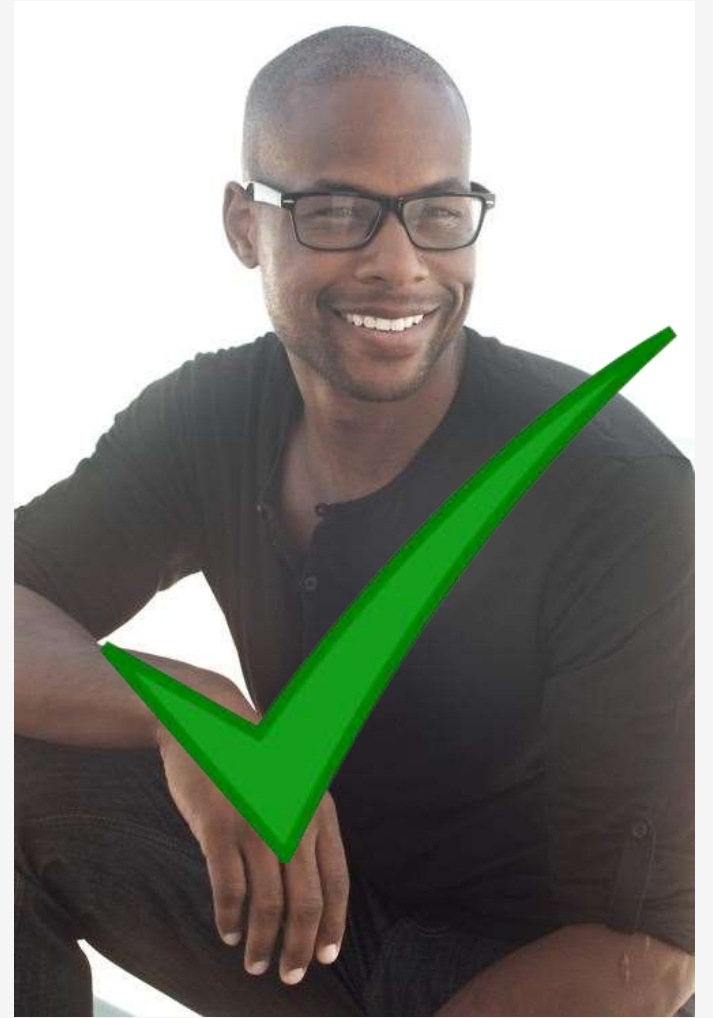
What does this look like in code?



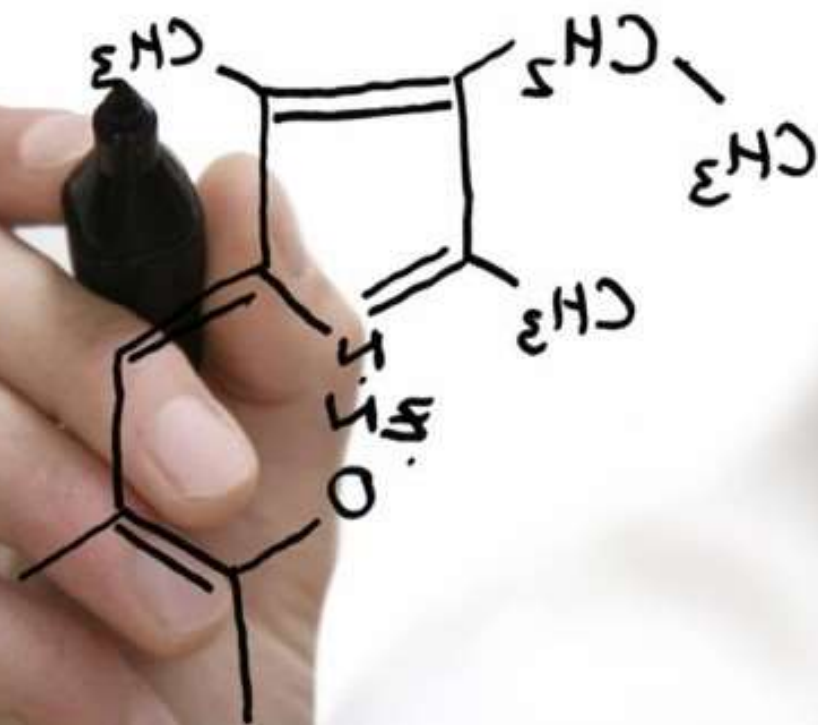
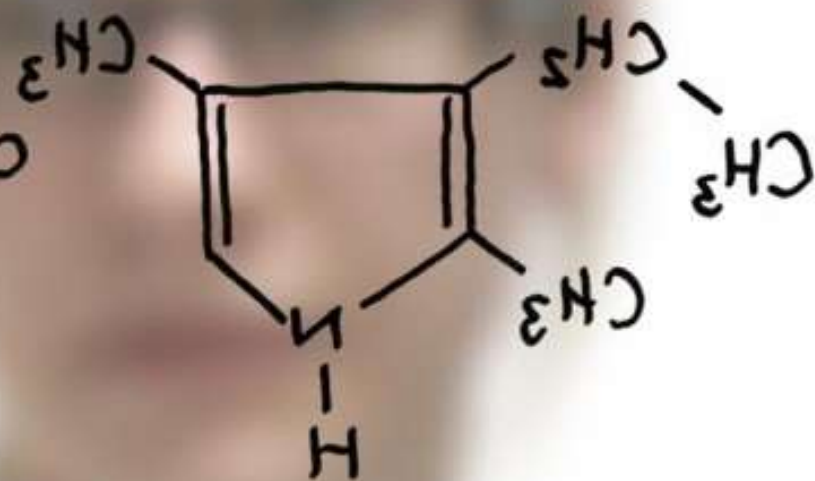
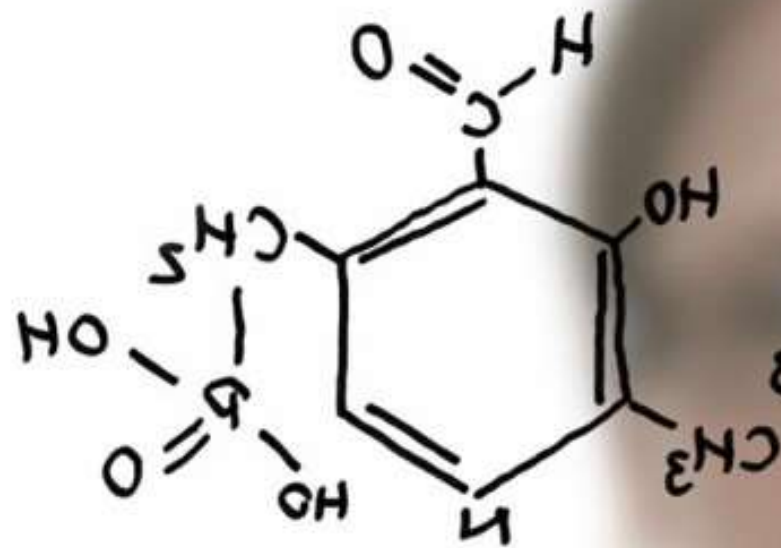
Live Demo

CodeMotion 2017 – Madrid

...so what does that mean?



What is Overfitting?





@DynamicWebPaige

@DynamicWebPaige



An analogy:

Remember those kids in school who were phenomenal at memorization, but could never apply those memorized concepts to a new problem? Every time they were shown a problem they hadn't seen before, their conceptual framework fell apart.

That's **#overfitting**.

2:28 AM - 22 Nov 2017



Live Demo

CodeMotion 2017 – Madrid

Top 10

Machine Learning Gotchas

(and to make it interesting...)

Top 10 Machine Learning Gotchas

(...let's talk about Real Madrid)



The owner of Real Madrid has
heard about your legendary
data science skills...

...and wants you to
predict performance
for the club in 2018.

Not just whether
they'll win

(obviously)



But...

- By how much?
- How will individual players perform?
- How many successful passes?
- How will players grow / deteriorate?
- Will there really be a home field advantage?

etc.

football.db

The screenshot shows the GitHub repository page for **football.db**, an organization that provides open public domain football data. The page layout includes a header with navigation links (Pull requests, Issues, Marketplace, Explore), a repository overview section with pinned repositories (`datafile`, `your-league-starter`, `docs`), and a list of country-specific data repositories. Each repository entry includes a description, a language indicator (Ruby), star count, fork count, and update status. On the right side, there are sections for 'Top languages' (listing Ruby, HTML, and Shell) and 'Most used topics' (listing football and opendata). A 'People' section at the bottom right displays avatars of contributors.

football.db
Open Public Domain Football Data
<http://openfootball.github.io/>

Repositories: 39 | People: 9

Pinned repositories

- datafile**
football.db quick starter datafile templates - worldcup.db, worldcup2014.db, etc.
Ruby ★ 29 🍴 11
- your-league-starter**
football.db league quick starter sample
★ 5
- docs**
football.db documentation incl. notes, articles, tips, guides, etc.
Shell ★ 6 🍴 5

Search repositories... Type: All Language: All

at-austria
Free open public domain football data for Austria (Österreich) incl. Österr. Bundesliga, Erste Liga, Regionalliga (Ost, Mitte, West), ÖFB Cup, etc.
opendata austria football wiener liga bundesliga austria
Ruby ★ 9 🍴 3 Updated 2 days ago

es-espana
Free open public domain football data (football.db) for España (Spain) / Europe - Primera División / La Liga, etc.
★ 15 🍴 13 Updated 2 days ago

it-italy
Free open public domain football data (football.db) for Italy / Europe - Serie A etc.
★ 15 🍴 13 Updated 2 days ago

eng-england
Free open public domain football data for England (and Wales) incl. English Premier League (EPL) etc.
★ 176 🍴 75 📄 CC0-1.0 Updated 2 days ago

de-deutschland
Free open public domain football data for Germany (Deutschland) incl. Deutsche Bundesliga, 2. Bundesliga, 3. Liga, DFB Pokal etc.
★ 15 🍴 13 📄 CC0-1.0 Updated 2 days ago

Top languages
Ruby HTML Shell

Most used topics
football opendata

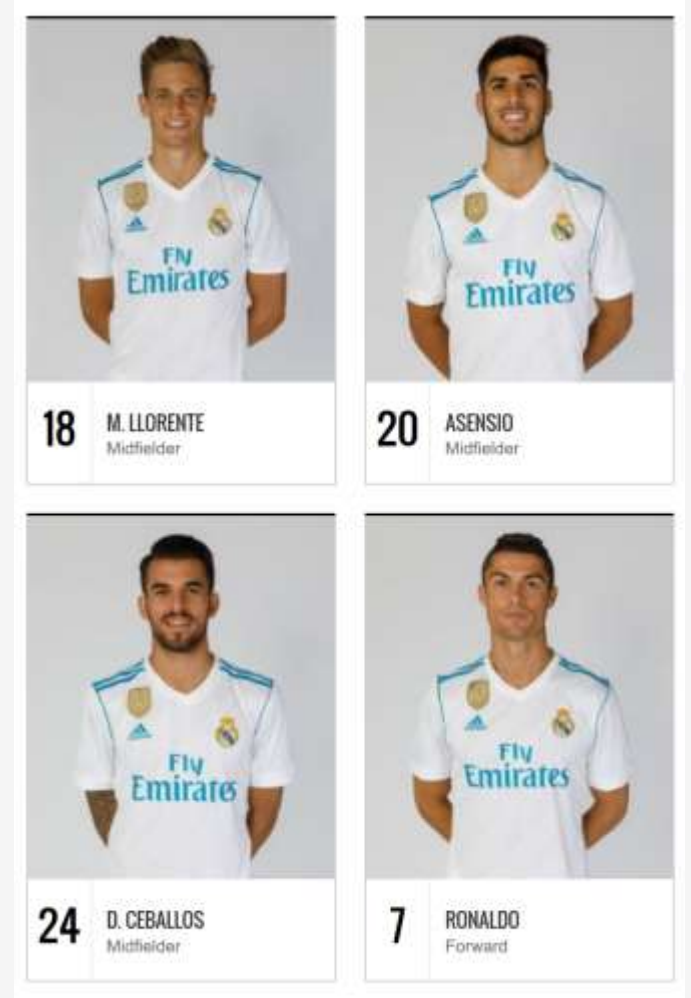
People
9

Example 1:

Too Few Samples

only having data for
a few games...

19	16 January 2011	Almería	1–1	Real Madrid	Almería
20	23 January 2011	Real Madrid	1–0	Mallorca	Madrid
21	30 January 2011	Osasuna	1–0	Real Madrid	Pamplona
22	6 February 2011	Real Madrid	4–1	Real Sociedad	Madrid
23	13 February 2011	Espanyol	0–1	Real Madrid	Cornellà de Llobregat
24	19 February 2011	Real Madrid	2–0	Levante	Madrid
25	26 February 2011	Deportivo La Coruña	0–0	Real Madrid	A Coruña
26	3 March 2011	Real Madrid	7–0	Málaga	Madrid

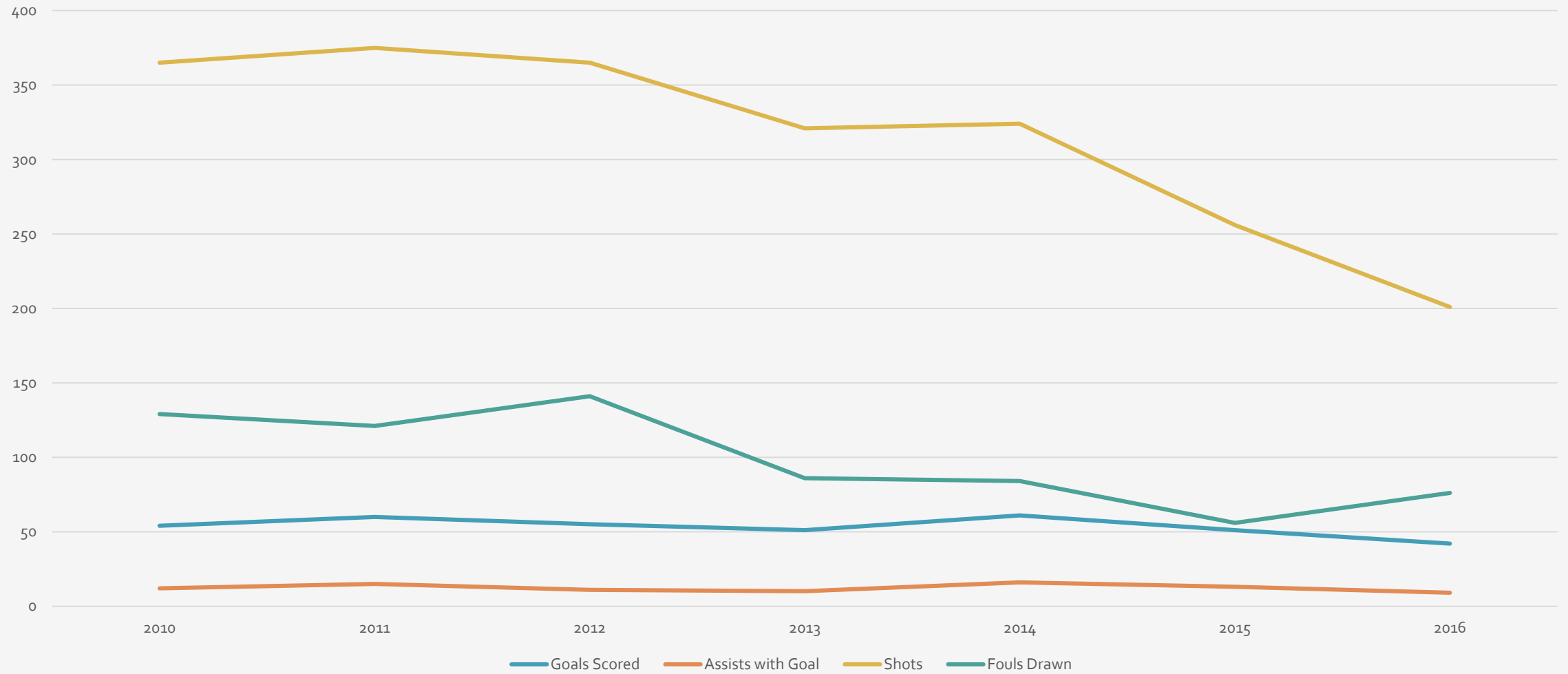


...or a few players.

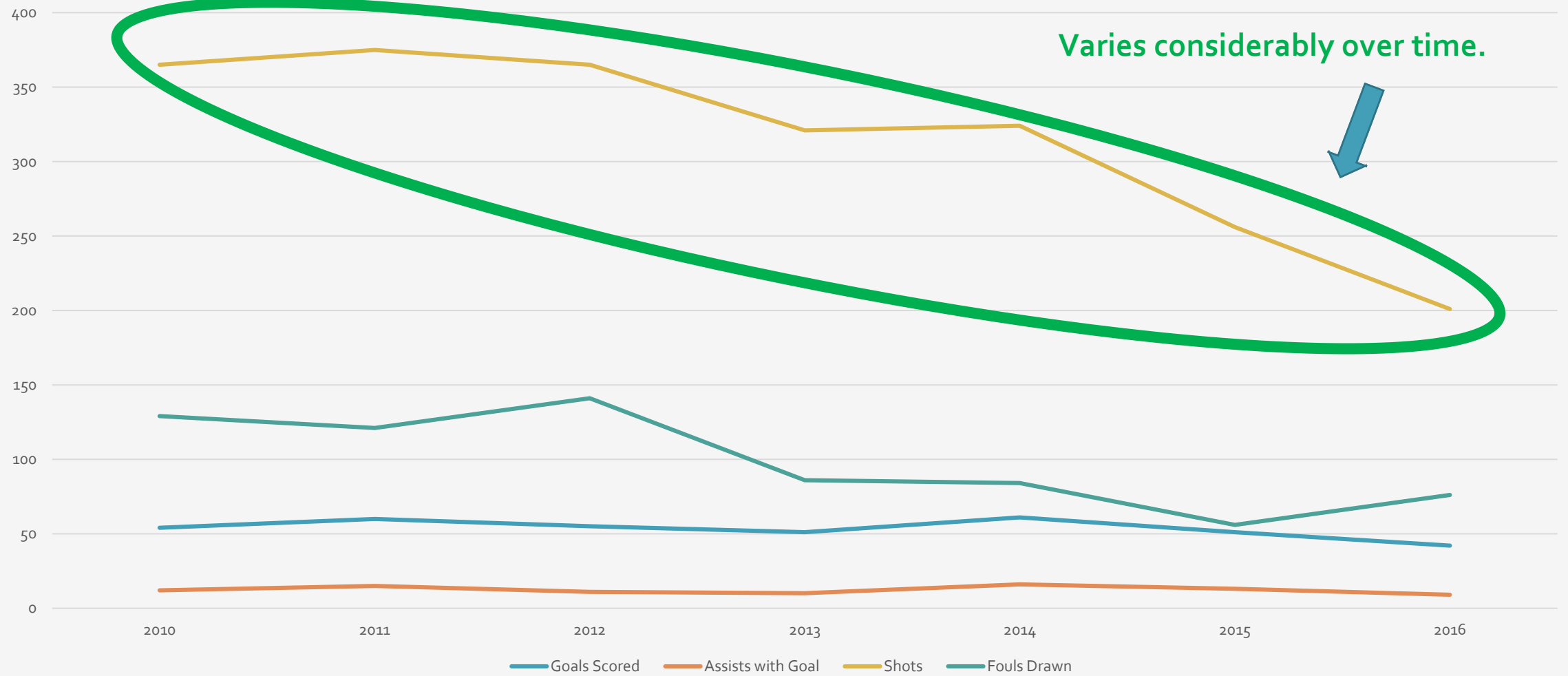
Example 2:

Using Old Data

Peformance Metrics *2010 - 2017*



Peformance Metrics 2010 - 2017



Example 3:

Choice of Measure

What does a “good” season look like?

Season 2012 - 2013

55

Fixtures played

55

Goals Scored

11

Assists w/ Goal

4861

Minutes Played

365

Shots

141

Fouls Drawn

	CHAMPIONS LEAGUE	COPA DEL REY	LIGA ESPAÑOLA 1ª DIVISIÓN	SUPERCOPA DE ESPAÑA
Fixtures played	12	7	34	2
Goals Scored	12	7	34	2
Assists w/ Goal	1	1	9	0
Minutes Played	1141	690	2841	189
Shots	80	42	235	8
Fouls Drawn	30	29	80	2
Matches Started	12	7	30	2
Matches Subbed In	0	0	4	0
Headers Scored	1	1	6	1
Goals w/ Right Foot	9	2	18	1
Goals w/ Left Foot	2	4	10	0
Penalty Kicks Scored	0	1	6	0
Free Kicks Scored	0	0	3	0
Fouls Committed	7	10	27	4
Steals	25	9	49	4
Penalties Received	0	1	1	0
Passes Completed	399	224	884	40

What does a “good” season look like?

Season 2012 · 2013

55
Fixtures played

55
Goals Scored

11
Assists w/ Goal

4861
Minutes Played

365
Shots

141
Fouls Drawn

	CHAMPIONS LEAGUE	COPA DEL REY	LIGA ESPAÑOLA 1ª DIVISIÓN	SUPERCOPA DE ESPAÑA
Fixtures played	12	7	34	2
Goals Scored	12	7	34	2
Assists w/ Goal	1	1	9	0
Minutes Played	1141	690	2841	189
Shots	80	42	235	8
Fouls Drawn	30		80	2
Matches Started	12	7	30	2
Matches Subbed In	0	0	4	0
Headers Scored	1		6	1
Goals w/ Right Foot	9		18	1
Goals w/ Left Foot	2	4	10	0
Penalty Kicks Scored	0		6	0
Free Kicks Scored	0		3	0
Fouls Committed	7	10	27	4
Steals	25	9	49	4
Penalties Received	0	1	1	0
Passes Completed	399	224	884	40

Example 4:

Cherry-picking Data



Example 5:

Reprobleming

Altering the problem so that
your performance improves.

Example 6:

Parameter Tweaking

good results.

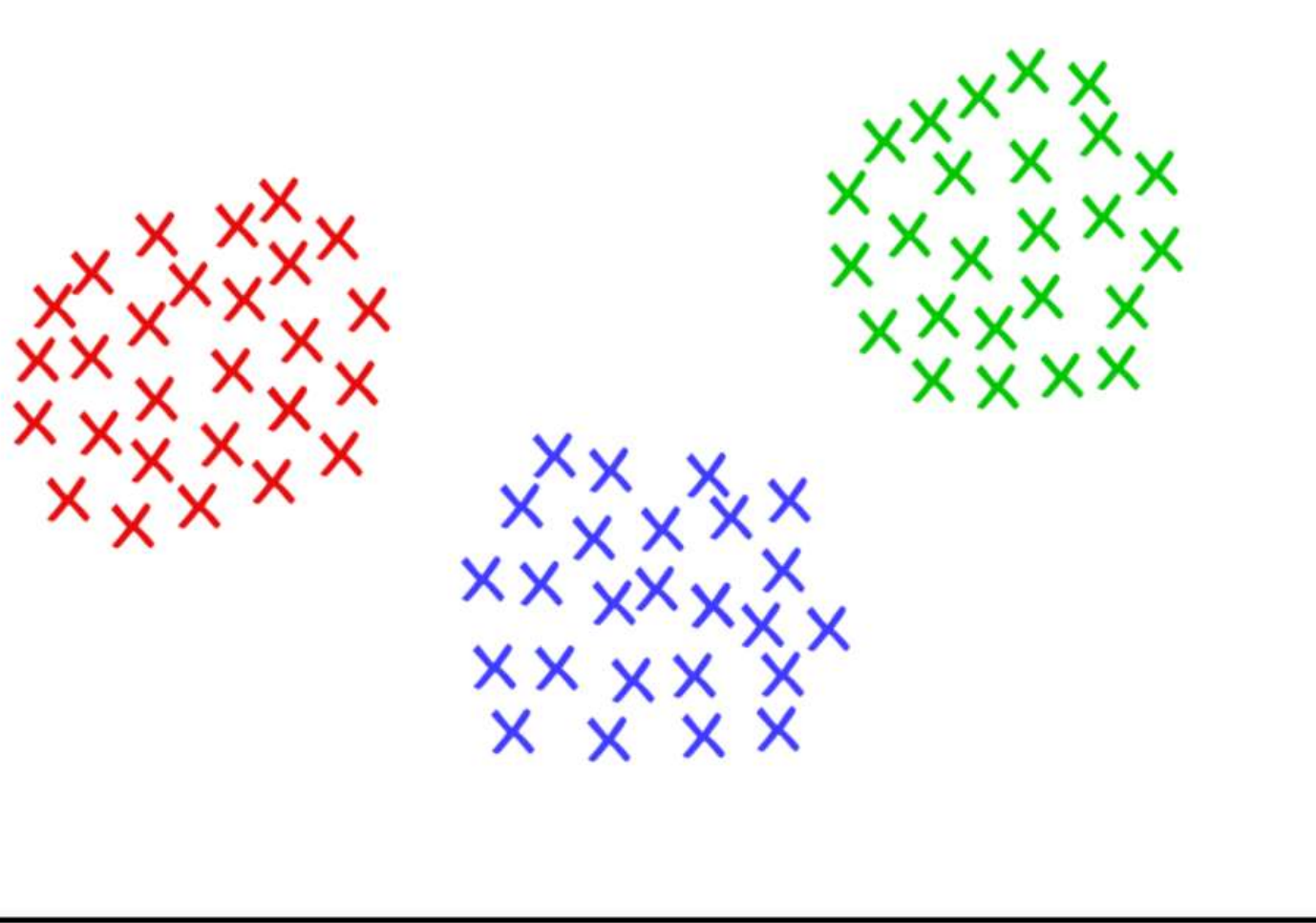
Season 2012 - 2013	
55	55
Fixtures played	Goals Scored
Fixtures played	
Goals Scored	
Assists w/ Goal	
Minutes Played	
Shots	
Fouls Drawn	
Matches Started	
Matches Subbed In	
Headers Scored	
Goals w/ Right Foot	
Goals w/ Left Foot	
Penalty Kicks Scored	
Free Kicks Scored	
Fouls Committed	
Steals	
Penalties Received	
Passes Completed	

good results.

Season 2012 - 2013				
55	55	11	4861	365
Fixtures played	Goals Scored	Assists w/ Goal	Minutes Played	Shots
				Fouls Drawn
	CHAMPIONS LEAGUE	COPA DEL REY	LIGA ESPAÑOLA 1ª DIVISION	SUPERCOPA DE ESPAÑA
Fixtures played	12	7	34	2
Goals Scored	12	7	34	2
Assists w/ Goal	1	1	9	0
Minutes Played	1141	690	2841	189
Shots	80	42	235	8
Matches Started	12	7	30	2
Matches Subbed In	0	0	4	0
Headers Scored	1	1	6	1
Goals w/ Right Foot	9	2	18	1
Goals w/ Left Foot	2	4	10	0
Penalty Kicks Scored	0	1	6	0
Free Kicks Scored	0	0	3	0
Steals	25	9	49	4
Passes Completed	399	224	884	40

Example 7:

Human-Loop Overfitting



Example:

Using a clustering algorithm (on training and test samples) to guide learning algorithm choice.

Example 8:

Collinearity

For example:

Player *age*, *date of birth*, and
number of *years playing*
professionally would probably
give you similar information.

Example 9:

Overfitting by Review




Example:

We **only** see data for the
public games – **not** the
data for **trial** matches.

Example 10:

Deep Learning in General

A common way to find perturbations forcing models to make wrong predictions is to compute **adversarial examples** [SZS13]. They yield perturbations that are very slight and often indistinguishable to humans, yet force machine learning models to produce wrong predictions. For instance, in the illustration below reproduced from [GSS14], the image on the left is correctly classified by a machine learning model as a panda, but adding the noise represented in the middle to that same image results in the image on the right, which is classified as a gibbon by the model.

	$+ .007 \times$		$=$	
x		$\text{sign}(\nabla_x J(\theta, x, y))$		$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“panda”		“nematode”		“gibbon”
57.7% confidence		8.2% confidence		99.3 % confidence

You **will** overfit.

DENIAL

How could my performance on test data be so low? There must be a mistake!

ANGER

What did I do wrong?! This is the [data / algorithm / programming language]'s fault!

BARGAINING

Maybe if I add more data? Or if I do more feature engineering?

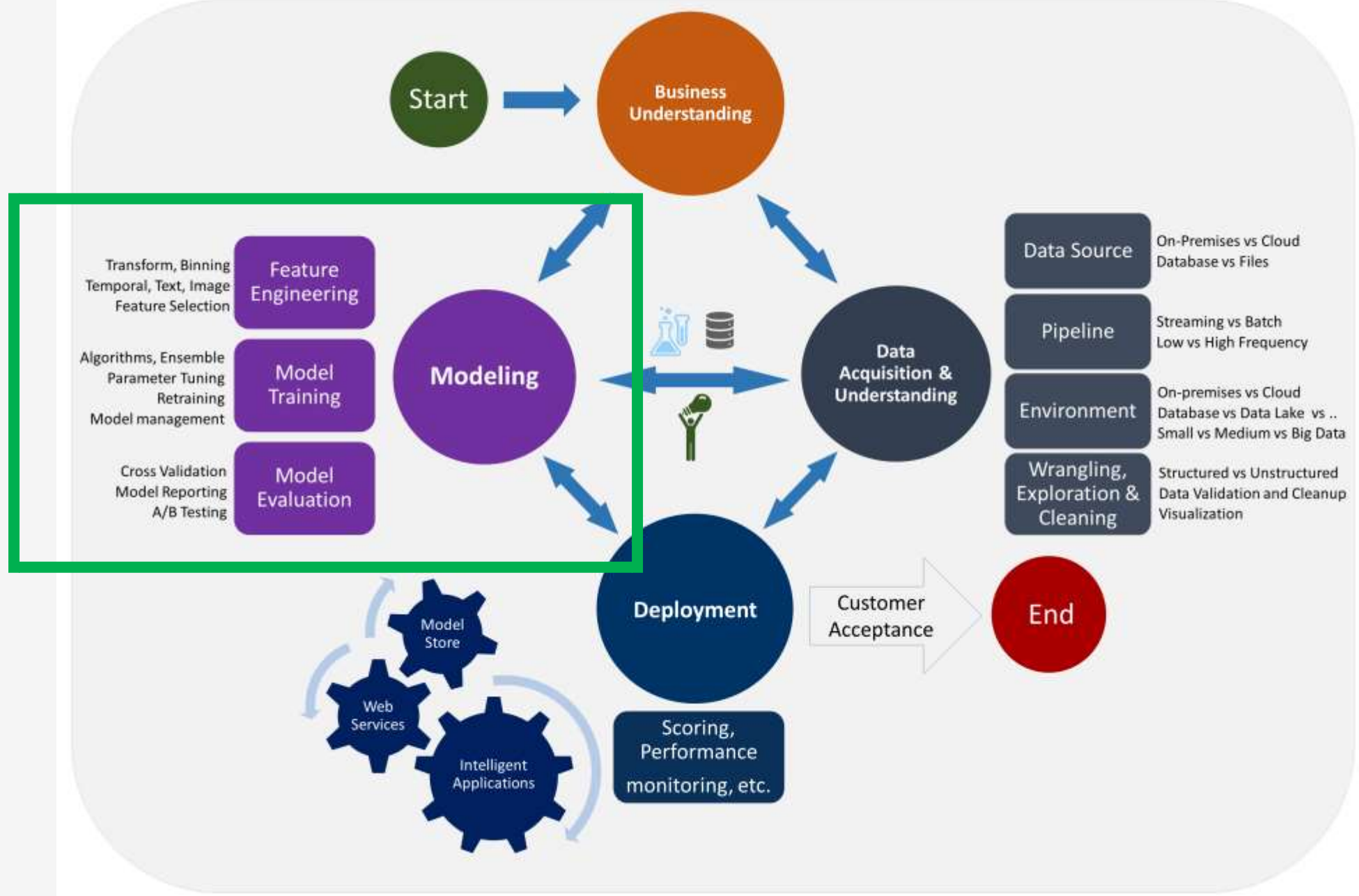
DEPRESSION

Why can't I fix this? Why does machine learning hate me? Maybe I should go be a coffee farmer

ACCEPTANCE

I will overfit. I am one with that reality. I must overfit as little as possible.





RESOURCES (football data):

- [Football.db database](#)
- [Guide to Football and Soccer APIs](#)
- [Sports Open Data Statistics](#)
- [OpenLigaDB](#)

RESOURCES (machine learning):

- [Team Data Science Process](#)
- [Azure Machine Learning Workbench](#)
- [Learn Analytics @ Microsoft](#)
- [Microsoft AI School](#)
- [DataCamp](#)
- [AI Tools for Visual Studio](#)
- [Machine Learning \(Coursera\)](#)
- [Deep Learning \(Coursera\)](#)

Joe Kampschmidt's Code
Blog Projects Code

Guide to Football/Soccer data and APIs

08 Mar 2014
Last updated: Jul

Where can I find football data?
There are three main sources for it or you can find it yourself.
Jump to a specific source

- Open source
 - open
 - joke
 - soccer
 - Other
- Free APIs
 - footb
 - Sport
 - open
 - betli
- Commercial
 - footb
 - Crow
 - SPAR
 - Sport
 - XML
 - Resu
 - opta
 - proce
 - Matz
 - apiFor
- Other Web
 - footb
 - Rec
 - footb
 - euro
 - open

OpenLigaDB
Community

Wer die Idee hat, mach per API abstr
Der hier angebotene welche in einem Liga automatisiert bez Projekten für Tipp
Hier werden nicht bereitgestellt. Spi eigene Liga erstellen

Woher können wir die OpenLigaDB in können eigene Liga erfordert keinerlei
Die Eingabe der S Community - für

Wie funktioniert der Abruf der Daten?
<http://www.openliga.de>
— Spiele des: [http://www.openliga.de](#)

Eine vollständige

football.db
Open Public Domain Football Data
https://github.com/football-db/football-db

Repositories: 29
People: 18

Pinboard repositories

- it-italia**
Free open public domain football data (football.db) for Italy / Europa - Serie A etc.
★ 25
Updated 1 day ago
- your-league-starter**
Football.db league quick starter script
★ 5
- docs**
Football.db documentation and notes, articles, tips, guides, etc.
★ 6
Updated 1 day ago

Top languages

- Ruby
- HTML
- GraphQL

Most used topics

- football
- openliga

People

nächste Spiele:
Fußball - 2. Fußball-Bundesliga 2017/2018

2017/2018

3. Liga

	Punkte
ien	29
	23
	23
hengla...	21
mund	20
fenheim	20
kfurt	19
	18
rkusen	17
	16
	16
15	15
	14
	14



DIG DEEP WITH

AZURE MACHINE LEARNING

Use data analysis to take your business to a whole new level.

Microsoft Azure Machine Learning simplifies data analysis and empowers you to find the answers your business needs.

The question isn't whether you can find the answers.
The question is how.

SELECTING A MODEL

Thank you!