

The Colour of Cleaning

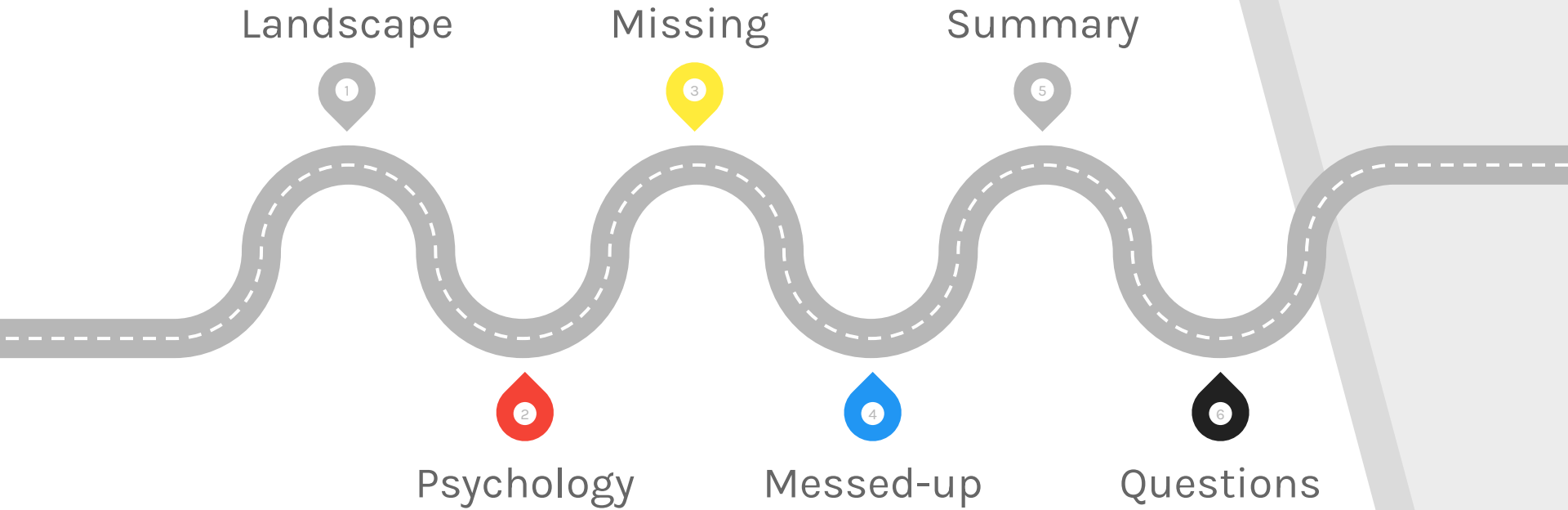
ODSC Europe 2021

Marta Markiewicz

InPost, Wroclaw University of Economics and Business



ROADMAP

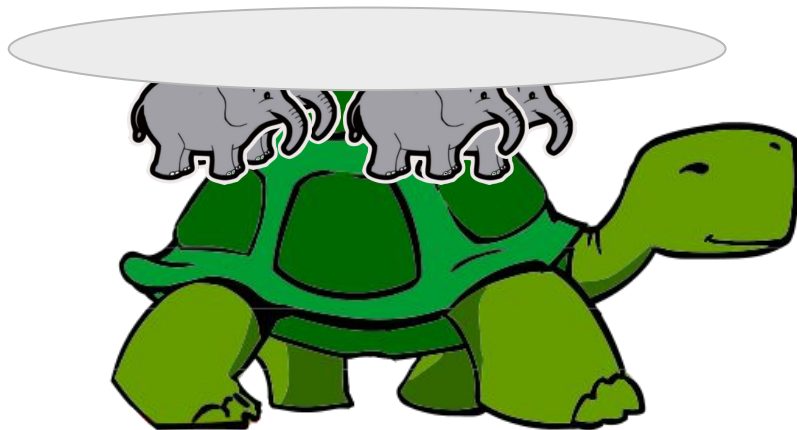




This Thing Called Cleaning

Why care?

Data



Business
Cleaning
Modelling
Deployment

<https://clipartix.com>
<https://imgur.com/gallery/dQcny7b>

Manually introduced data

No proper system for data input

Incorrect data types

Wrong formats

Out-of-domain observations

All possible dates formats

Typos

Single value columns

Lack of uniqueness

Systems integrations

Missing data

Duplicated not bringing value

Improperly gathered data

Improper cross field dependencies

Lack of data input rules

1.

Psychological barriers



DIRTY DATA

US

RoomEscapeArtist.com

Why the dislike?



Pressure

Satisfaction

Never ending

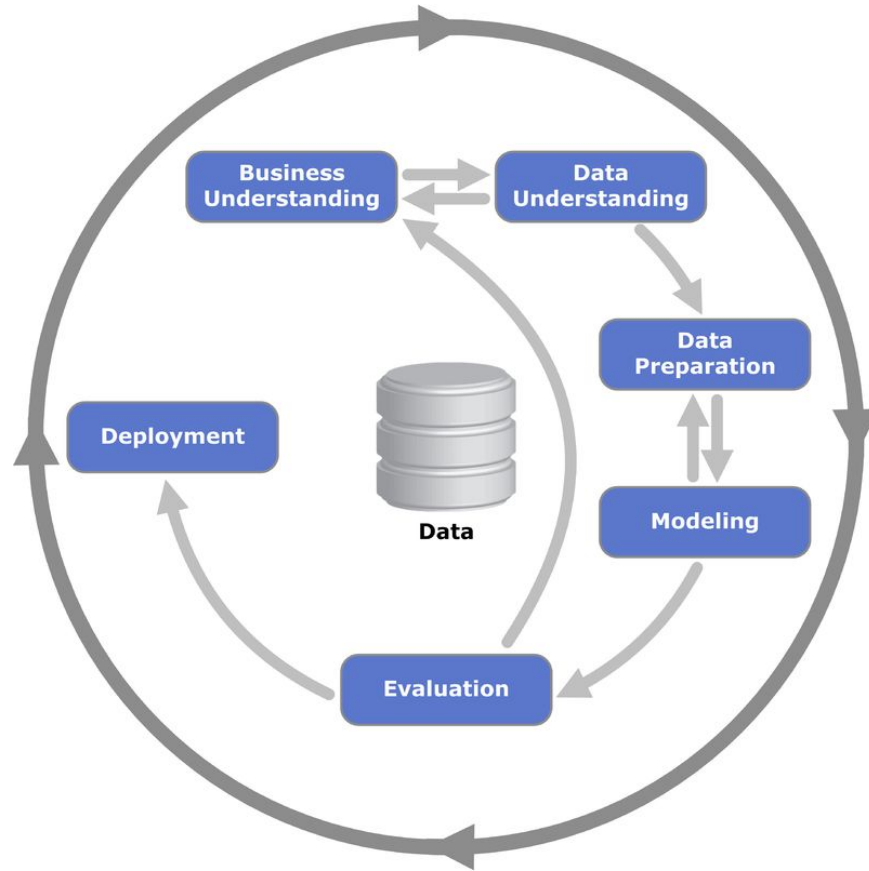
Personal development

Qual & boring

Fame

No fame

Hack your psychology



High iterativeness

Evangelisation

ML for cleaning

2.

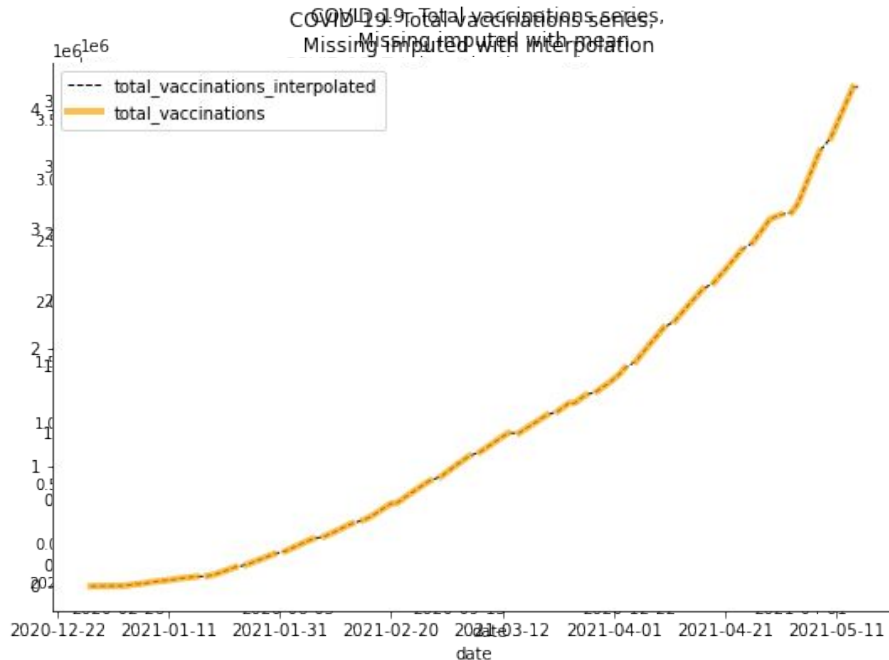
Missing data imputation



Simple approaches

73% empty

1 out of 60



Dropping column

Dropping rows

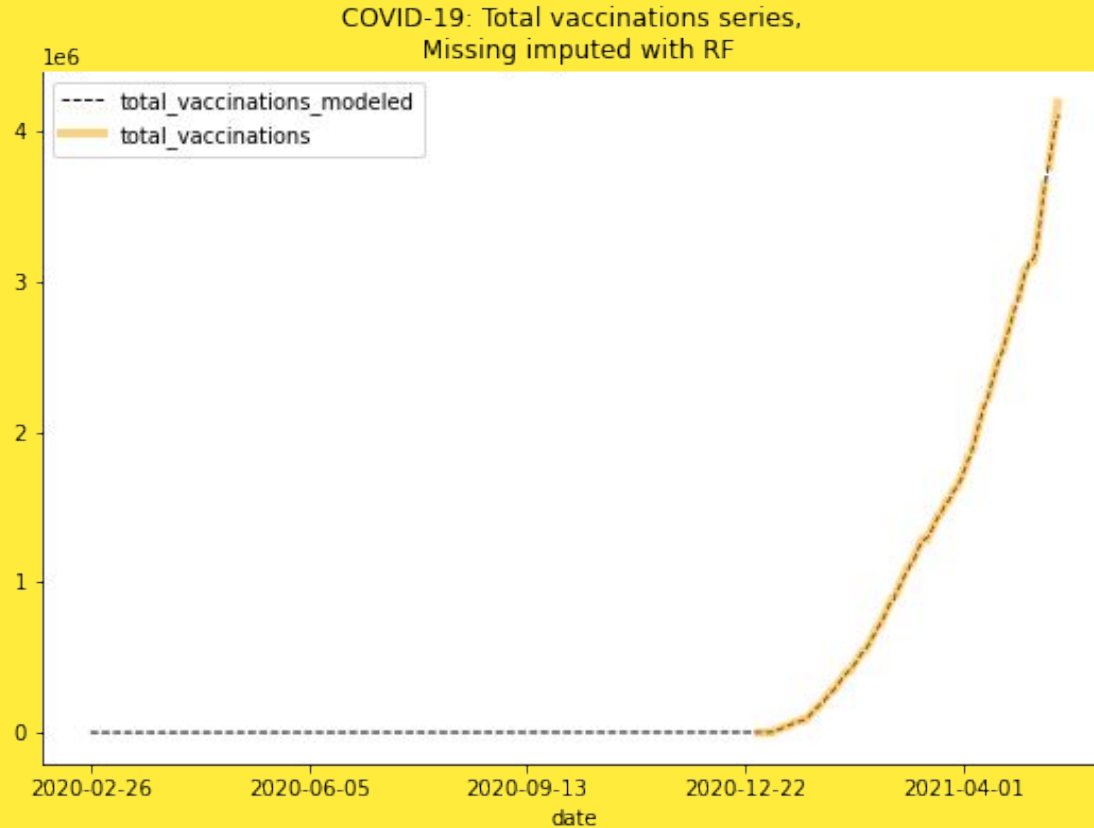
Mean / median / other statistic

Last observation carried forward

Next observation carried backward

Interpolation

Building a model



missForest

Pclass	Gender	Fare		Pclass	Gender	Fare		Pclass	Gender	Fare		Pclass	Gender	Fare		Gender	Fare
1	female	71.2833		1	female	71.2833		1	female	71.2833		1	female	71.2833		0	0
3	female	7.925		3	female	7.925		3	female	7.925		3	female	7.925		0	0
1		263	→	1	male	263	→	1	female	263	→	1	female	263	→	1	0
1	male	27.7208		1	male	27.7208		1	male	27.7208		1	male	27.7208		0	0
3	male	7.8958		3	male	7.8958		3	male	7.8958		3	male	7.8958		0	0
1	female			1	female	32.17		1	female	107.08		1	female	107.08		0	74.91
2		26		2	male	26		2	male	26		2	male	26		0	0

Multiple iterations!

MICE - Multiple Imputation by Chained Equations

Missing data is in red. There is a strong correlation between A and B, so let's try to impute A using B and C.

A	B	C
0.93	1.40	1.53
0.24	0.46	0.76
	0.80	
0.95	1.24	1.46
0.23	0.57	
0.90		1.28
0.15	0.42	
0.47	0.54	0.63
	1.14	
0.89	1.23	1.45

Missing data is filled in randomly. This dilutes the correlations, but allows us to impute using all available data.

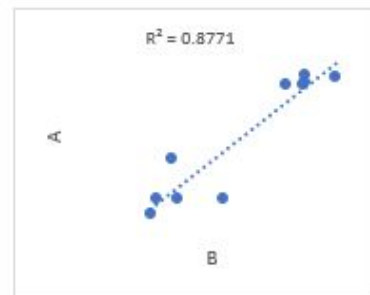
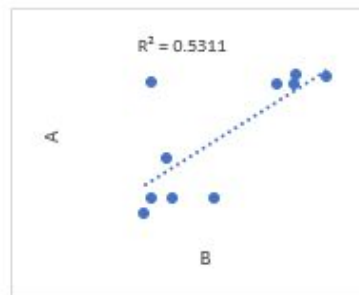
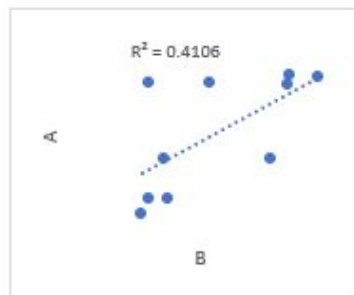
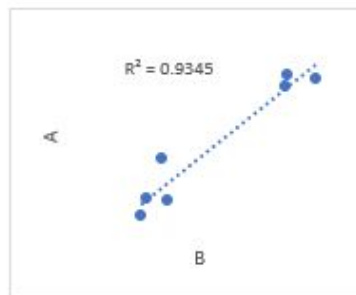
A	B	C
0.93	1.40	1.53
0.24	0.46	0.76
0.90	0.80	1.53
0.95	1.24	1.46
0.23	0.57	1.28
0.90	0.46	1.28
0.15	0.42	1.53
0.47	0.54	0.63
0.47	1.14	1.28
0.89	1.23	1.45

A random forest is used to predict A with B and C. Notice the correlation between A and B improved.

A	B	C
0.93	1.40	1.53
0.24	0.46	0.76
0.24	0.80	1.53
0.95	1.24	1.46
0.23	0.57	1.28
0.90	0.46	1.28
0.15	0.42	1.53
0.47	0.54	0.63
0.89	1.14	1.28
0.89	1.23	1.45

After Imputing B using A and C, we have achieved a correlation between A and B much closer to the original data.

A	B	C
0.93	1.40	1.53
0.24	0.46	0.76
0.24	0.80	1.53
0.95	1.24	1.46
0.23	0.57	1.28
0.90	1.24	1.28
0.15	0.42	1.53
0.47	0.54	0.63
0.89	1.14	1.28
0.89	1.23	1.45



PMM - Predictive Mean Matching

The predicted value of A ($E[A|B,C]$) is shown to the left. We are interested in imputing the bold missing value below

$E[A B,C]$	A	B	C
0.73	0.93	1.40	1.53
0.62	0.24	0.46	0.76
0.60		0.80	1.53
1.39	0.95	1.24	1.46
0.36	0.23	0.57	1.28
1.27	0.90	0.46	1.28
0.15	0.15	0.42	1.53
0.65	0.47	0.54	0.63
1.20		1.14	1.28
1.24	0.89	1.23	1.45

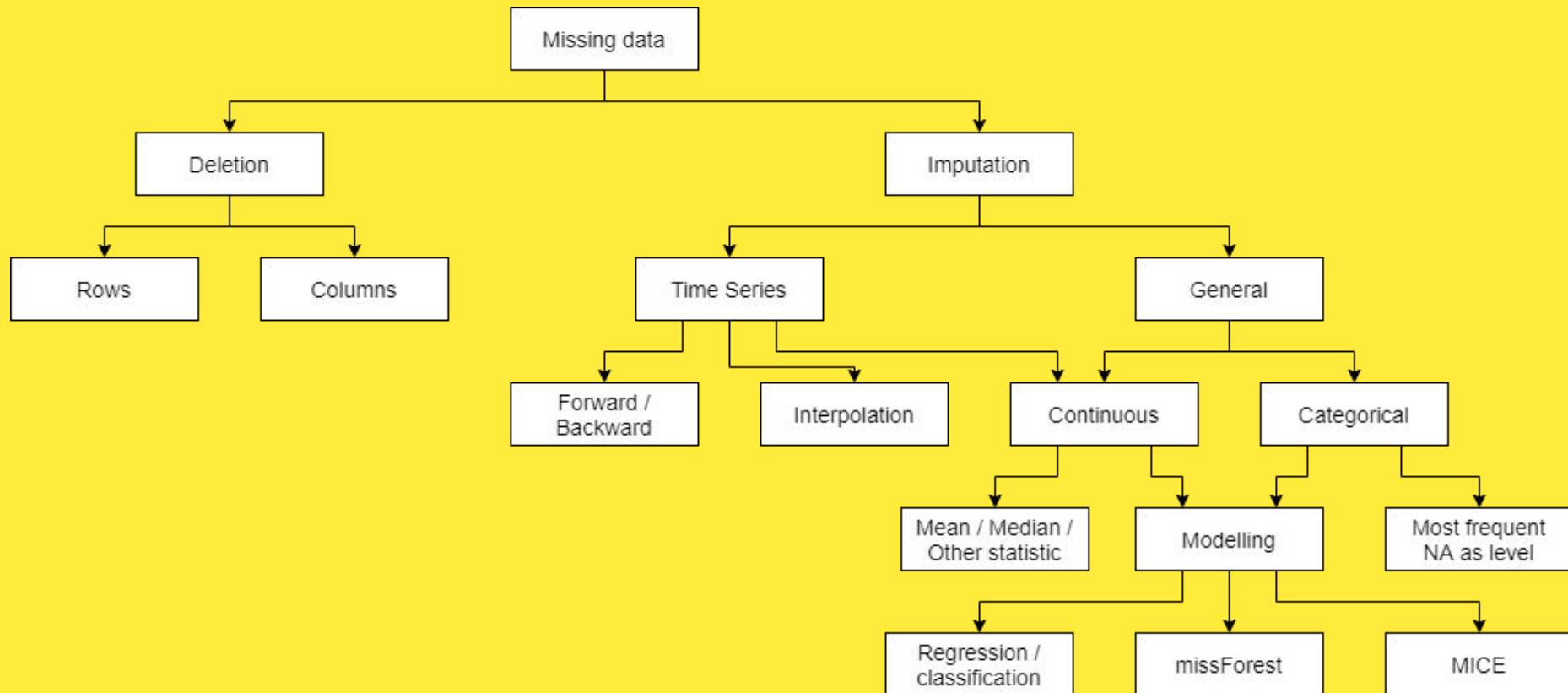
Our predicted value for the first missing sample is 0.60. The closest predicted value is 0.62. We find the closest values for all of our missing samples.

$E[A B,C]$	A	B	C
0.73	0.93	1.40	1.53
0.62	0.24	0.46	0.76
0.60		0.80	1.53
1.39	0.95	1.24	1.46
0.36	0.23	0.57	1.28
1.27	0.90	0.46	1.28
0.15	0.15	0.42	1.53
0.65	0.47	0.54	0.63
1.20		1.14	1.28
1.24	0.89	1.23	1.45

We then impute the value corresponding to the original data.

$E[A B,C]$	A	B	C
0.73	0.93	1.40	1.53
0.62	0.24	0.46	0.76
0.60	0.24	0.80	1.53
1.39	0.95	1.24	1.46
0.36	0.23	0.57	1.28
1.27	0.90	0.46	1.28
0.15	0.15	0.42	1.53
0.65	0.47	0.54	0.63
1.20	0.89	1.14	1.28
1.24	0.89	1.23	1.45

Summary



A thick, solid blue diagonal stripe runs from the top-left towards the bottom-right, separating the white background on the left from the solid blue background on the right.

3.

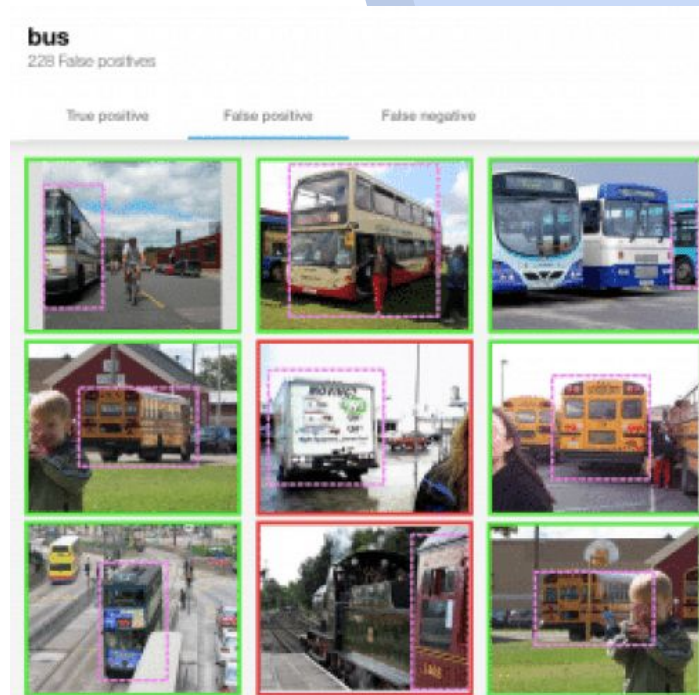
Messed up data

LAST time



Incorrect target labels

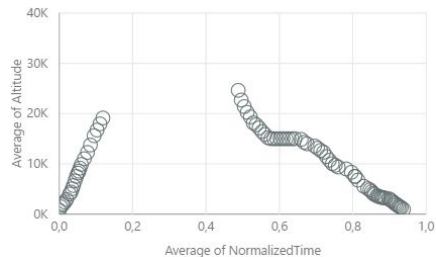
	VOC 2012	COCO 2017	Udacity – Self Driving Car
Training Set Images	17 177	94 439	11 992
Training Set Labels	20	80	9
Training Set Objects	49 834	686 385	78 230
Manual Correction	21.1%	23.6%	25%



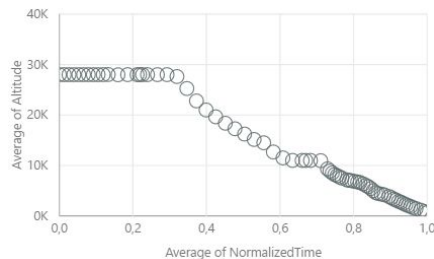
<https://deepomatic.com/en/how-we-improved-computer-vision-metrics-by-more-than-5-percent-only-by-cleaning-labelling-errors>

Breaking quality issues

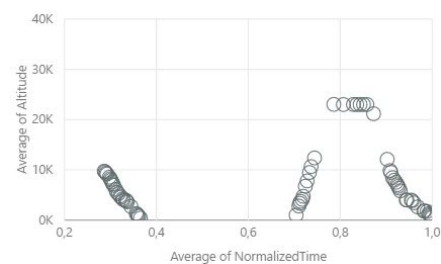
Lack of middle phase



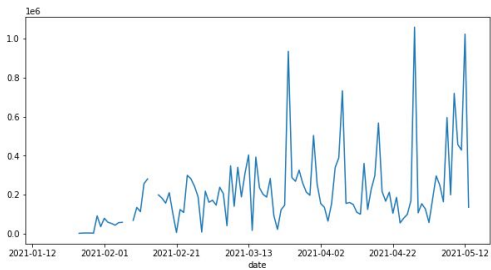
Lacks half of a flight



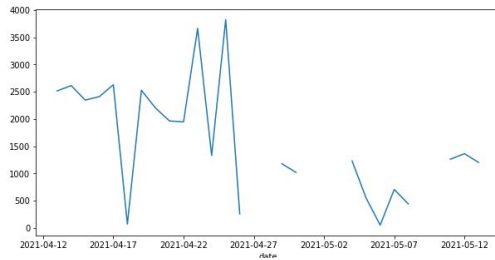
Data gathering error



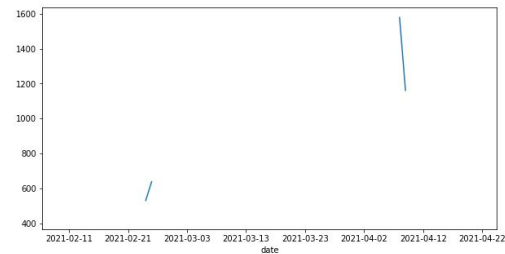
Complete vaccinations series



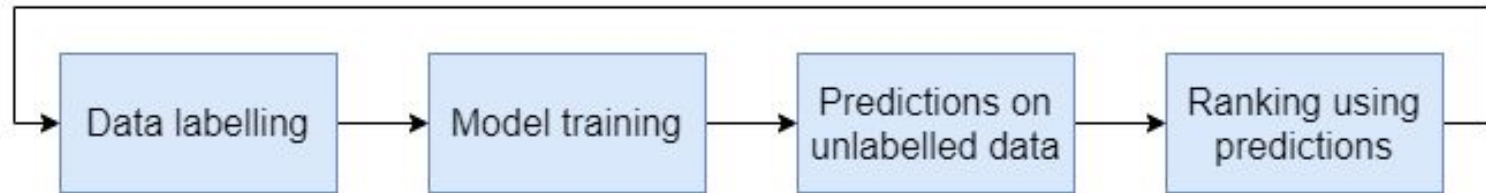
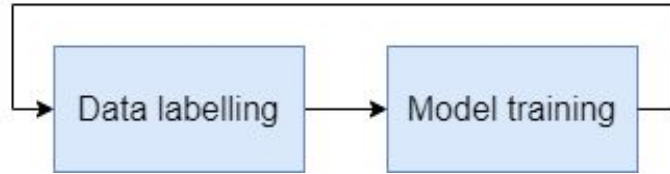
Incomplete vaccinations series



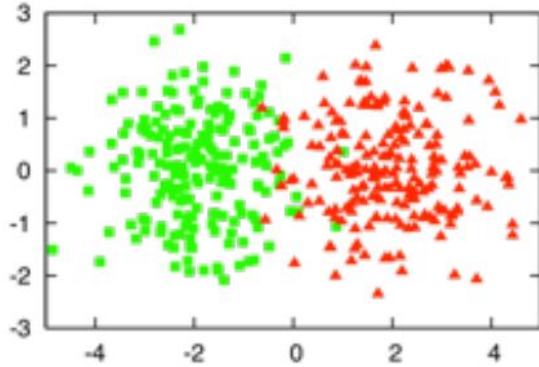
(Almost) absent vaccinations series



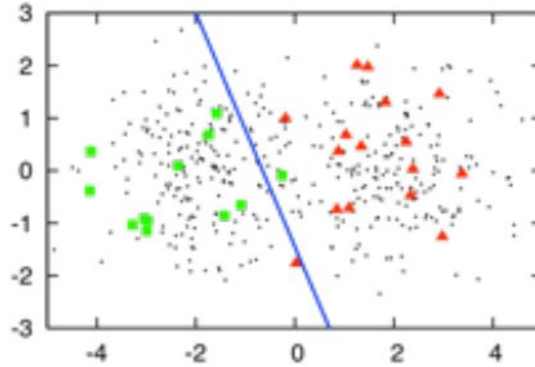
Active learning



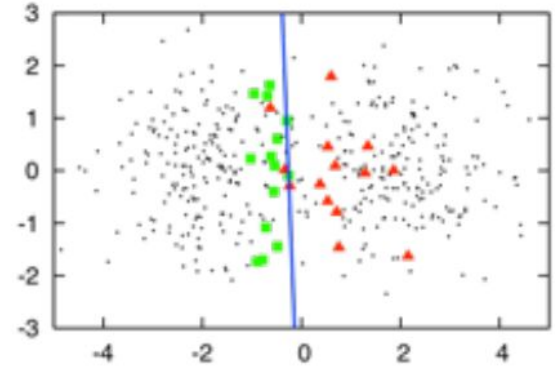
Active learning - intuition



400 instances sampled
from 2 class Gaussians



random sampling
30 labeled instances
(accuracy=0.7)

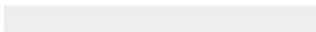


active learning
30 labeled instances
(accuracy=0.9)

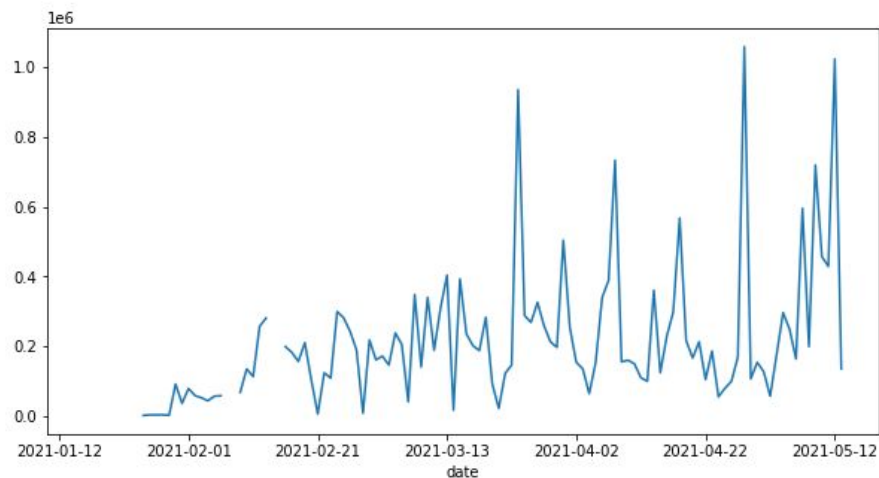
<https://towardsdatascience.com/introduction-to-active-learning-117e0740d7cc>

Active quality learning

Progress:



 Retrain




good

so-so

bad

Other:

Hit enter to submit.

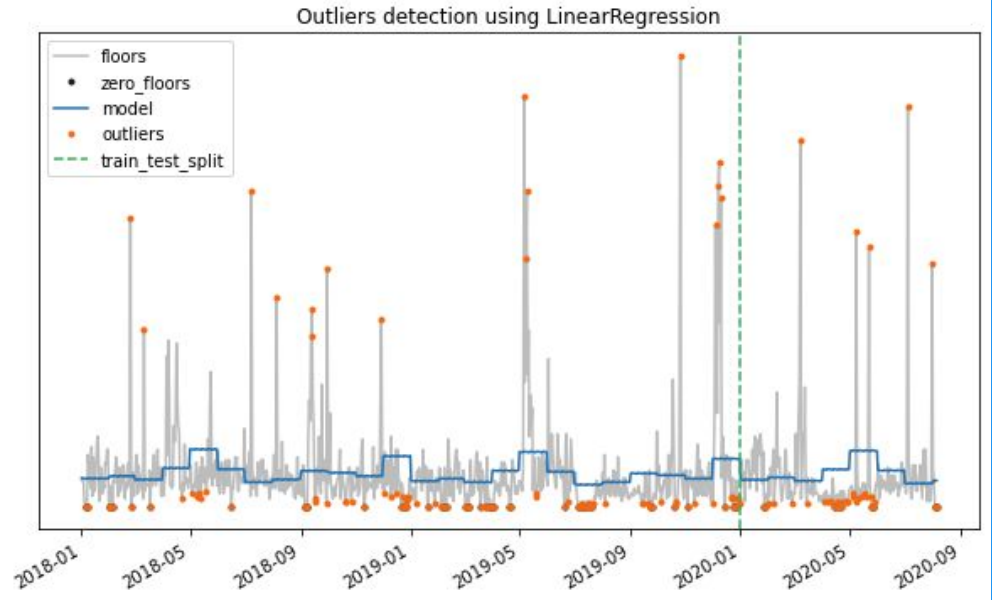
 Sort options

 Skip

 Undo



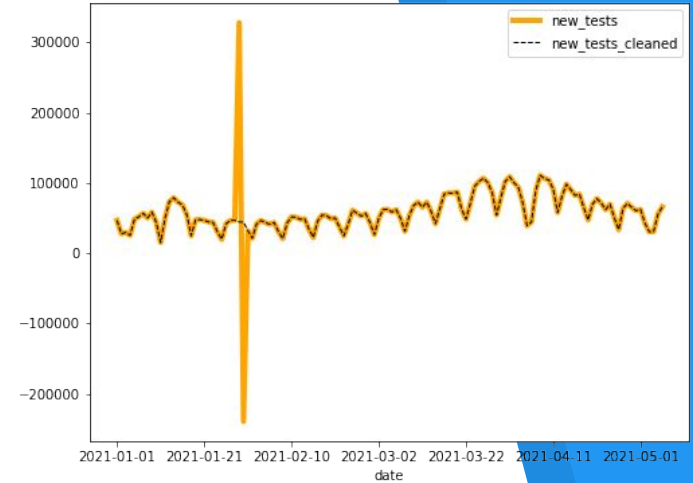
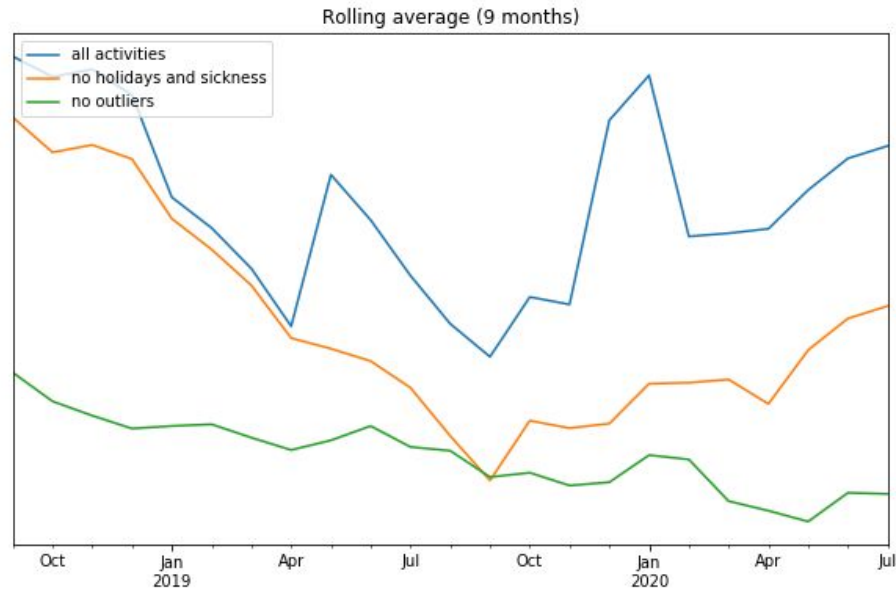
Anomalies - examples

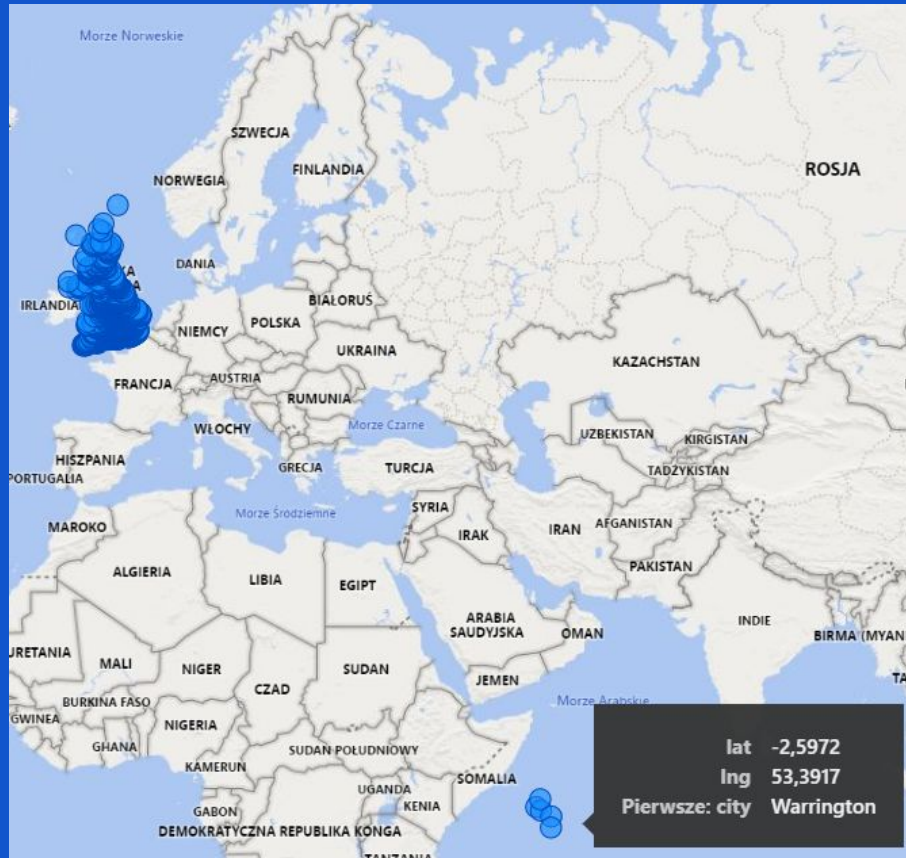


Anomalies

	Intended supervised	A side effect supervised	Unsupervised	Mixed
Manual	A lot	No	No	A little
Immunity to pattern changes	No, unless repeated periodically	Yes	Yes	Partial

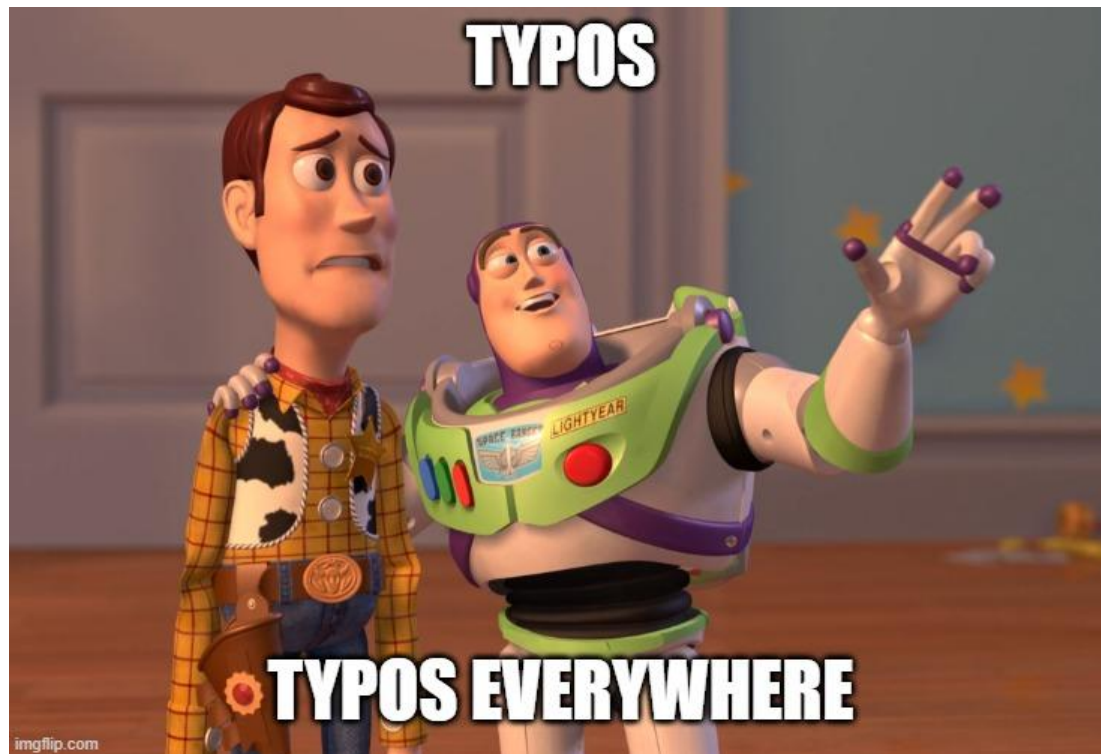
Anomalies - what to do with them?





1. Anomalies detection
2. Swapping
3. Anomalies detection

Typos world



male	375
female	200
Mr.	68
facet	63
monsieur	51
femme	27
kobiet	25
m ale	20
mujer	19
kobieta	18
dziewczyna	16
famale	9

NLP evolution

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez*[†]
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin*[‡]
illia.polosukhin@gmail.com



TF-IDF

love

FastText

RoBERTa

ERNIE 2.0

GPT

GPT-2

GPT-3

TransformerBERT

XLNet

Big Bird

ELMO

Sentence-Bert



Carnegie Mellon University

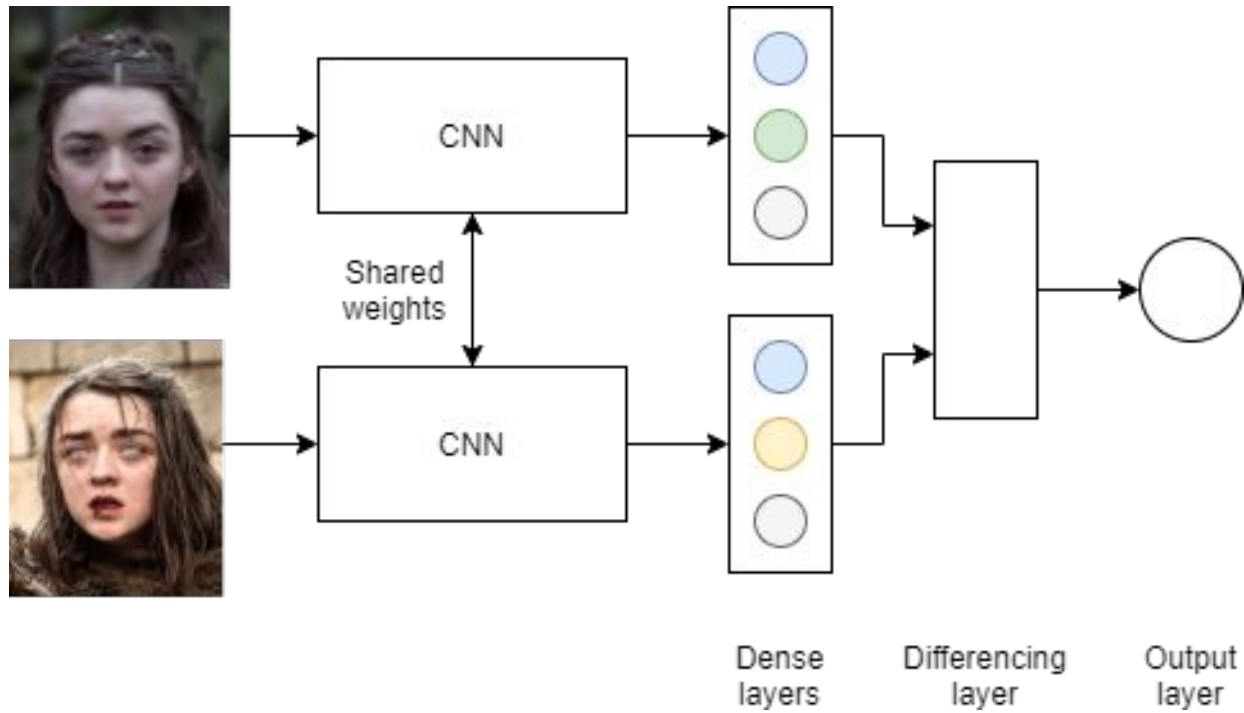


1954

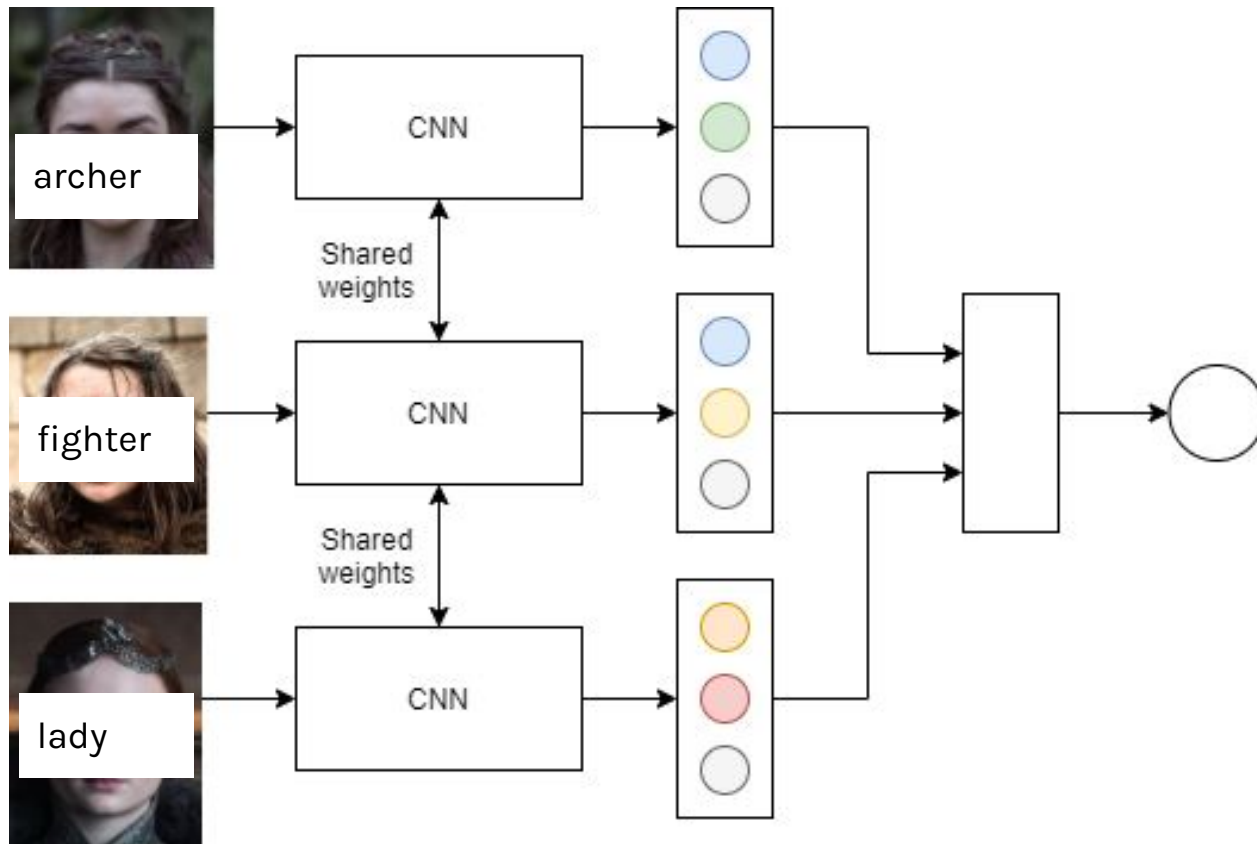
2013

2017

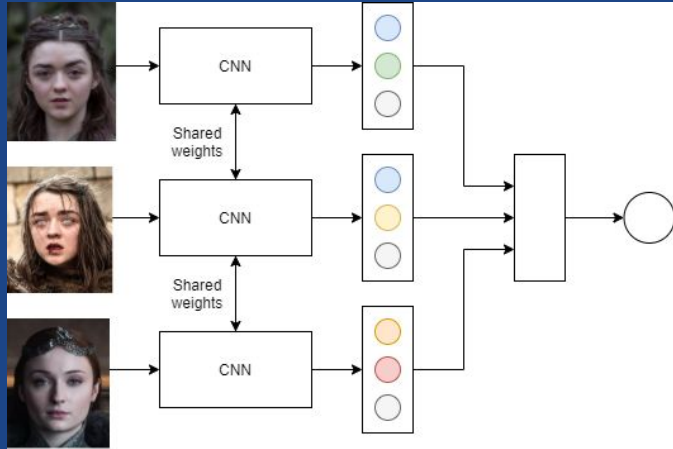
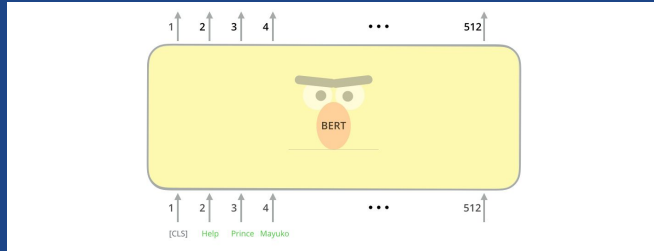
Siamese / Triplets Networks



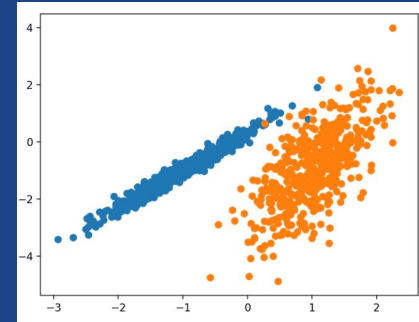
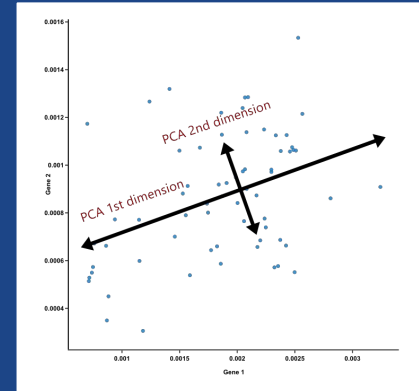
Siamese / Triplets Networks



Embeddings + dimension reduction / clustering



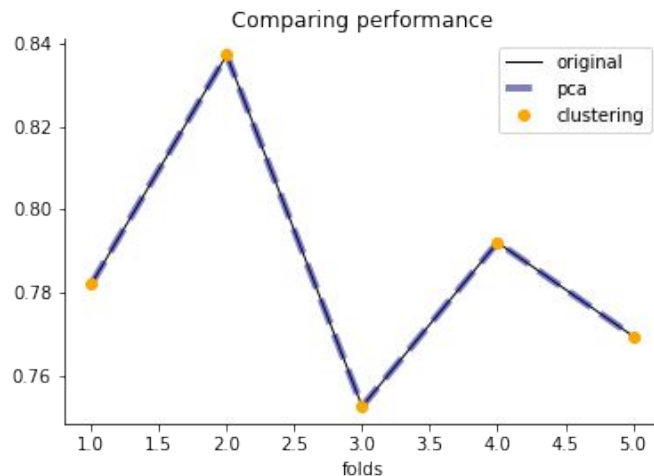
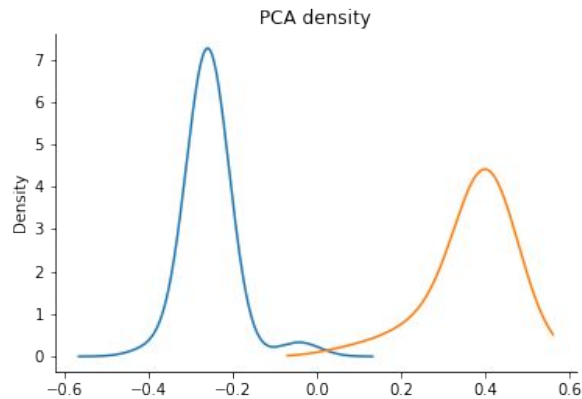
<https://arxiv.org/pdf/1908.10084.pdf>



<https://machinelearningmastery.com/clustering-algorithms-with-python/>

Cleaning categories

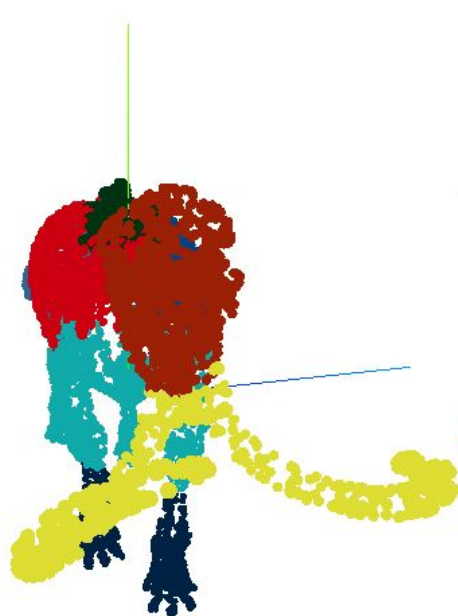
male	462
female	251
mle	24
man	24
monsieur	19
m ale	18
facet	16
mujer	14
Mr.	14
kobiet	11
kobieta	11
femme	10
femalle	10
dziewczyna	7



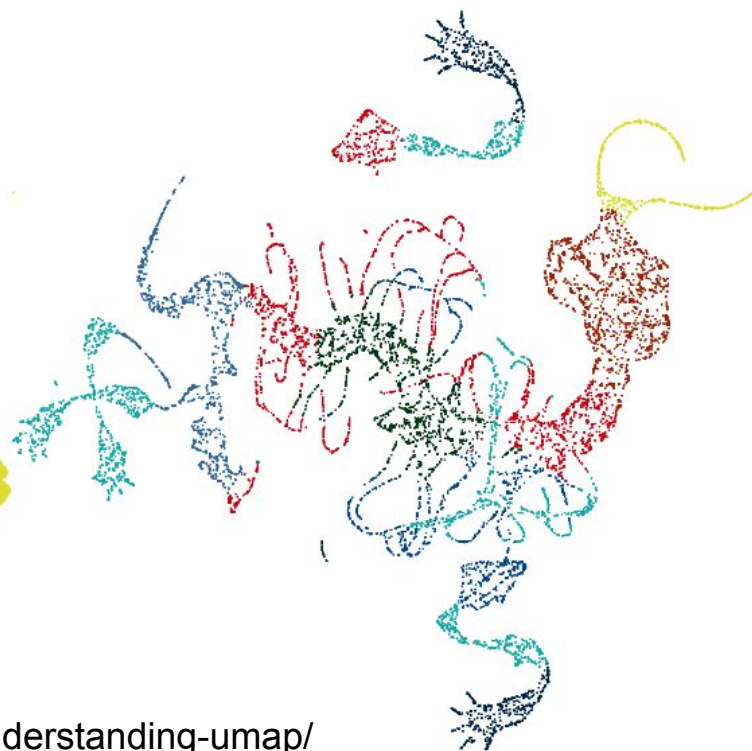
Cleaning categories

male	462
female	251
mle	24
man	24
monsieur	19
m ale	18
facet	16
mujer	14
Mr.	14
kobiet	11
kobieta	11
femme	10
femalle	10
dziewczyna	7

Original 3D Data



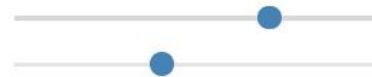
2D UMAP Projection



<https://pair-code.github.io/understanding-umap/>

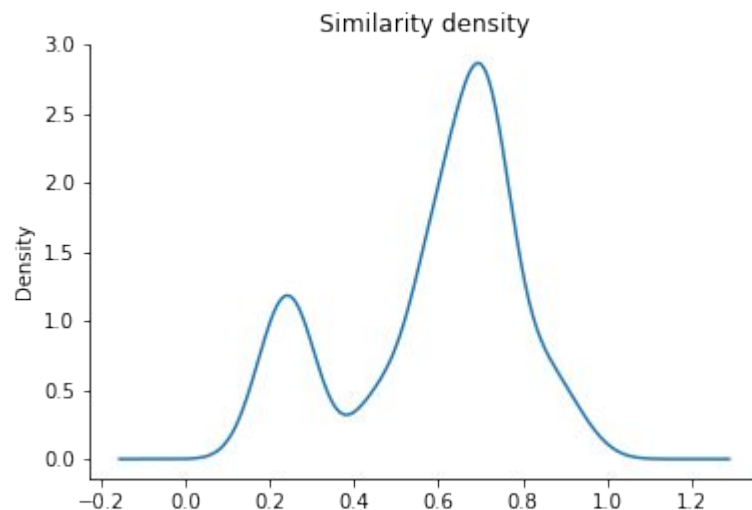
n_neighbors: 50

min_dist: 0.25



Identifying improper categories

category	subcategory
house	bungalow
	bus
	cottage
	log-house
	sem-detached house
	semi-detached house
	short-sleeved
vehicle	calf-length
	car
	caravan
	lorry
	lory
	minivan
	motorcoach
	motorcycle
	truck
	two-storey hous



category	subcategory	similarities
vehicle	two-storey hous	0.204250
house	short-sleeved	0.260327
vehicle	calf-length	0.264653
house	bus	0.450055
vehicle	caravan	0.553825

Key takeaways

- Iterate
- Evangelate
- Leverage ML:
 - Regression / classification
 - Multiple imputation
 - Quality active learning
 - Anomalies detection
 - Embeddings
 - Dimension reduction
 - Clustering
 - ...

Thank you!

Any Questions?

Marta Markiewicz

m.markiewicz.pl@gmail.com

<https://github.com/lady-pandas/cleaning-is-coming>