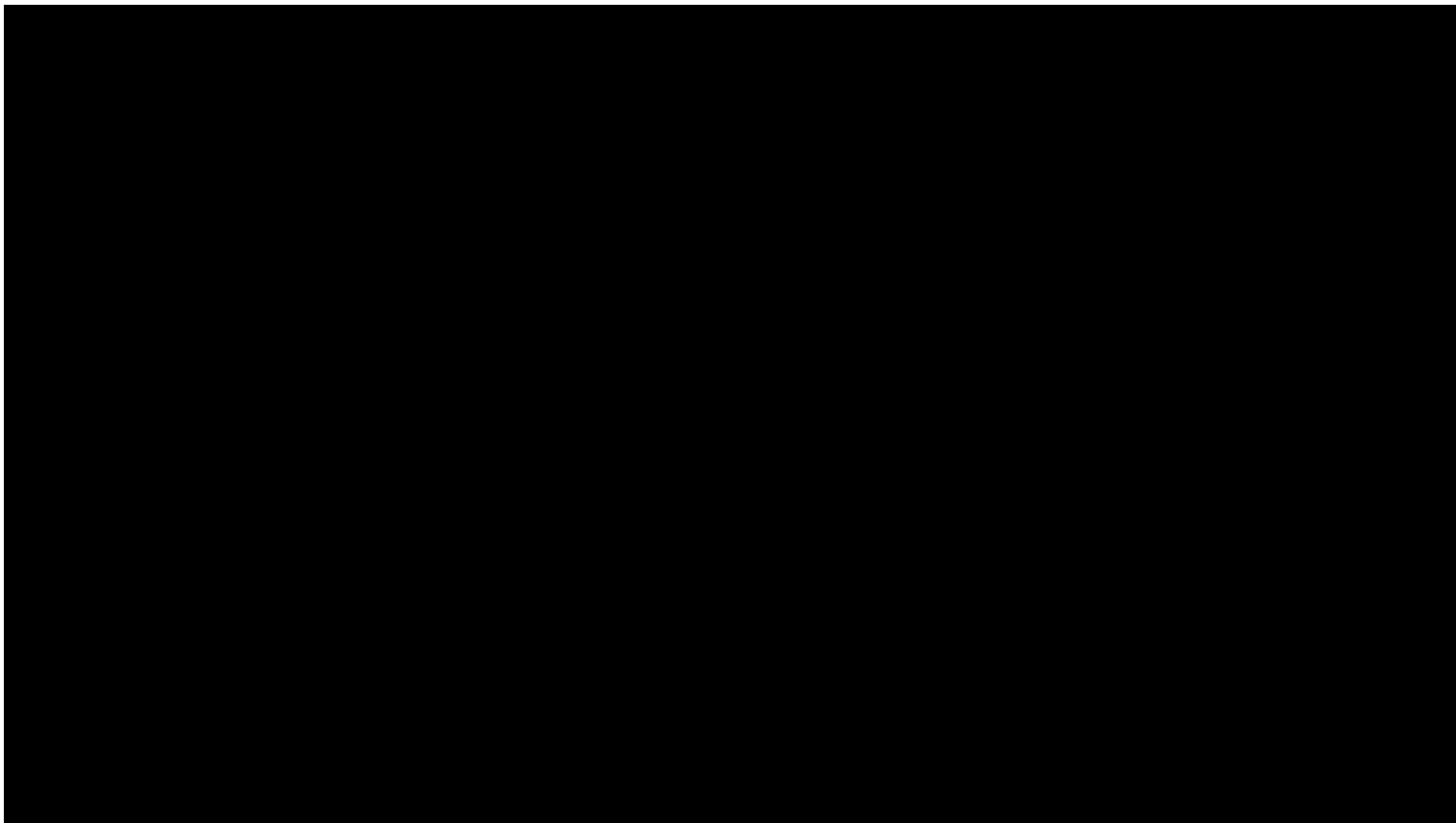


A hand holding a spray bottle, spraying a fine mist against a dark background. The mist is captured in a way that it looks like a soft, glowing cloud. The background is a dark, textured surface, possibly a wall or a large object, which is partially illuminated by the light from the spray.

Data CLEANING

JENS BAETENS





Kwaliteit van datasets is belangrijk

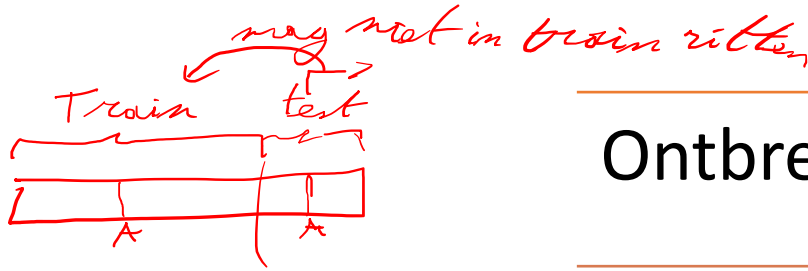
Garbage in = Garbage out

Hoe langer je fouten meesleept in je dataset, hoe kostelijker

- Zelfde als bij software development

Fouten in je data kan leiden tot een fout model





Problemen in de datasets

Ontbrekende data

Duplicaten → Data leakage
↳ Problemen om model te evalueren

Onmogelijke waarden → outliers

Verkeerde dataformaten
↳ lbs/kg ↳ slakken 13030
↳ 23:30
↳ 7330

Onnodige attributen verwijderen



naam age nationaliteit
Jens 32 ?
Jens ? Belg

Ontbrekende data

Ontbrekende data:

- Zoek de data op online
- Bereken de waarde (gemiddelde, gelijkaardige rijen, ...)
- Verwijder duplicaten (rij of kolom)

Niet altijd problematisch: gebrek aan data kan ook een bron van informatie zijn

Duplicate data

	Age
Jens	31
Jens	30

① Verwijder exacte duplicaten

② Gebruik duplicaten om ontbrekende informatie in te vullen

Ⓡ Wat met conflicterende data?

Data balancing

Geen kanker 1%

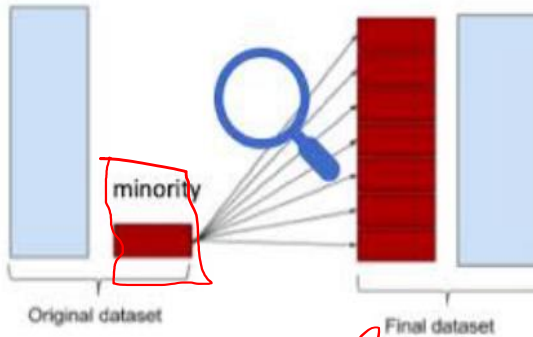
Kanker 99%

model voorspeld
altijd dit \rightarrow 99%

Rijen die waarden bevatten die zelden voorkomen worden best niet
verwijderd *naal = duplicaten*
voordeel = alle data gebruikt

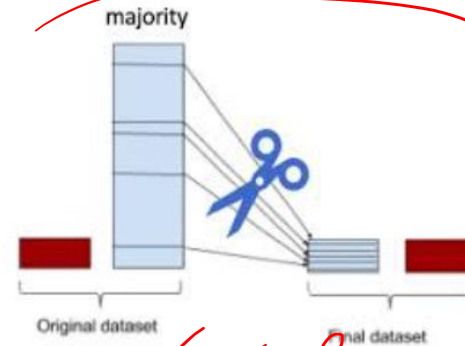
ongeveer

in theorie



oversampling

VS

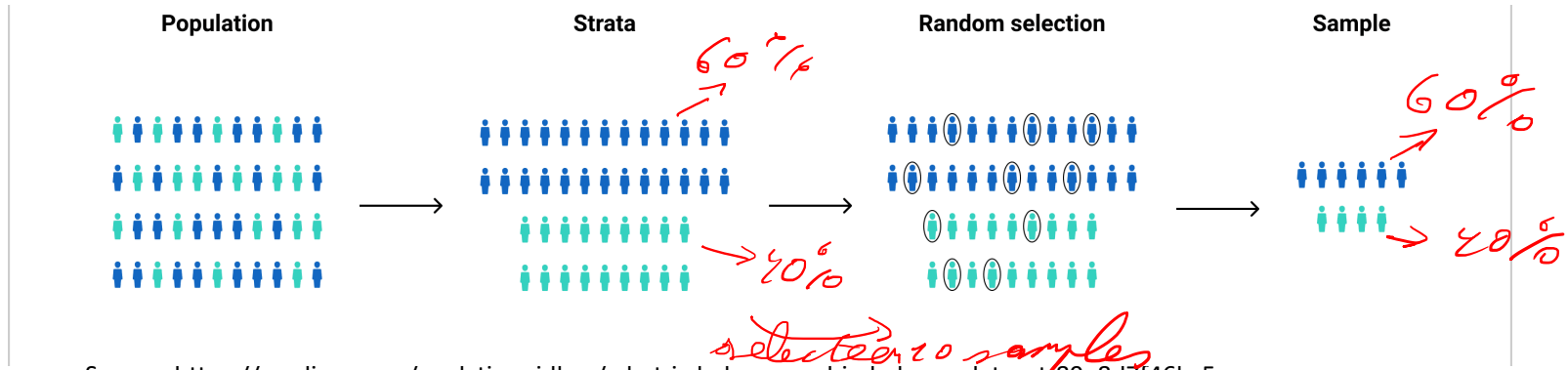


under/sampling

Source: <https://medium.com/analytics-vidhya/what-is-balance-and-imbalance-dataset-89e8d7f46bc5>

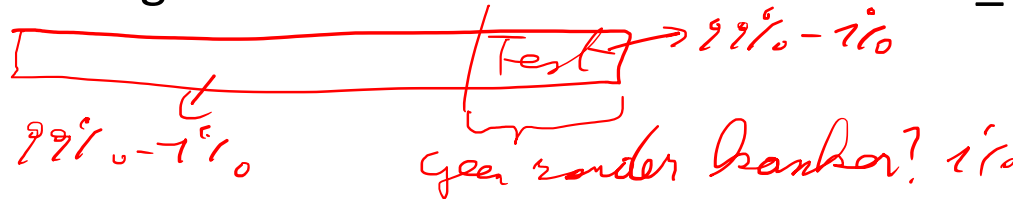
Data balancing – Stratified sampling

↳ kijken naar de verdeling



Source: <https://medium.com/analytics-vidhya/what-is-balance-and-imbalance-dataset-89e8d7f46bc5>

Stratification argument doet dit automatisch in train_test_split





Corrigeer data format

Typos

Verschillende waarden met dezelfde betekenis:

- 0/1 of True/False of ... *17/F → replace*
- Naam van een stad in verschillende talen
- Straat en nummer in 1 veld ipv 2
- Datums: yyyy/mm/dd vs dd/mm/yyyy

10/11/2021
↳ altijd complex problemen
↳ december of oktober?

One hot encoding vs ordinal

categorische data

Rood → 0
Geel → 1
Blauw → 2

Kleur
Rood
Geel
Blauw
Blauw
Geel

voor ML moet dit
numeriek zijn

Ordinal encoding
↳ kan met eenvoudige
replace

Kleur
0
1
2
2
1

one-hot encoding → elke klasse
aparte kolom

Rood	Geel	Blauw
1	0	0
0	1	0
0	0	1
0	0	1
0	1	0

kan problemen
zijn bij
veel klassen
(>50)

Privacy requirements

Zoek naar Personal Identifiable Information (PII)

Deze velden moeten beter afgeschermd zijn dan andere

- Data niet bruikbaar in het geval van hacking

Data masking of obfuscation

verwarring
van of onduidelijkheid
toevoegen

Toegangscontrole

Inhoud onbruikbaar
maken

Data
Masking

Deletion

drop
replace

Verwijder de gevoelige informatie

Zeer simplistische aanpak

Verwijderde data kan niet meer
gebruikt worden in je model

Geeft aan dat er zaken gewijzigd
zijn

*geef je
informatie
aan je klanten*

Data Masking

Wijzig PII door willekeurige waarden

Kan afhankelijk zijn van andere velden

Verwijderde data kan niet meer gebruikt worden in je model

Geeft aan dat er zaken gewijzigd zijn

Substitution

vervangen

→ moeilijker maar nog steeds detecteerbaar

Data Masking – Variantie Toevoegen

Voeg ruis toe aan de data

loon + 10 io

Numerieke waarden

- Tot 10% geeft nog bruikbare data voor prijzen / salarissen

Datums

- Afhankelijk van de toepassing
- tot 120 dagen variantie behoudt de verdeling in de kolom

Indicatie van het origineel

*↳ exacte waarde
niet meer gekend*

loon + 20
↙
~~*loon + 10%*~~
↖
~~*loon + 10*~~

A thick, horizontal, wavy orange line that spans the width of the page, serving as a decorative separator.

Encryption: Indien de waarde nog moet gebruikt kunnen worden

- Beheer van keys
- Encryptie en decryptie reken-intensief

Sleutel veilig \rightarrow data | onleesbaar
| leesbaar

Samenvatting



- De volgende zaken vereisen je aandacht bij data cleaning

- Datatypes van verschillende kolommen \rightarrow astype / replace
 - NaN – waarden \rightarrow fillna / dropna / fillna
 - Outliers \rightarrow volgende week
 - Foutieve dataformaten bij datums/tekst \rightarrow to_datetime / replace
 - Gebruikte categorieën in categorieke kolommen
 - Persoonsgegevens
 - \hookrightarrow deletion
 - substitution
 - variantie
 - encryptie
 - shuffling
- \hookrightarrow unique / nunique / ~~count~~ value_counts
- \hookrightarrow one-hot / ordinal

1 W Willie 3125 pts

6 B Bean 2213 pts

11 M Mike 1713 pts

2 J Jesus 2838 pts

7 R Right 1813 pts

11 m ms 1713 pts

3 l lakaka 2788 pts

7 J Jose 1813 pts

13 J Joske Vermeulen 1675 pts

4 c christian 2475 pts

9 D Dail 1750 pts

14 k kooken12 1488 pts

5 Q Quick 2375 pts

10 L LoroPol 1725 pts

15 H Hey 1263 pts

Click on the projected screen to start the question



5

10

15





5

10

15

1

W

Willie 3125 pts

6

B

Bean 2213 pts

11

M

Mike 1713 pts

2

J

Jesus 2838 pts

7

R

Right 1813 pts

11

m

ms 1713 pts

3

L

lakaka 2788 pts

7

J

Jose 1813 pts

13

J

Joske Vermeulen 1675 pts

4

c

christian 2475 pts

9

D

Dalil 1750 pts

14

k

kooken12 1488 pts

5

Q

Quick 2375 pts

10

L

LoroPol 1725 pts

15

H

Hey 1263 pts

Click on the projected screen to start the question

wooclap



100 %



18 / 23



1 W Willie 3125 pts

6 B Bean 2213 pts

11 M Mike 1713 pts

2 J Jesus 2838 pts

7 R Right 1813 pts

11 m ms 1713 pts

3 l lakaka 2788 pts

7 J Jose 1813 pts

13 J Joske Vermeulen 1675 pts

4 c christian 2475 pts

9 D Daul 1750 pts

14 k kooken12 1488 pts

5 Q Quick 2375 pts

10 L LoroPol 1725 pts

15 H Hey 1263 pts



5

10

15

Click on the projected screen to start the question



1 W Willie 3125 pts

6 B Bean 2213 pts

11 M Mike 1713 pts

2 J Jesus 2838 pts

7 R Right 1813 pts

11 m ms 1713 pts

3 l lakaka 2788 pts

7 J Jose 1813 pts

13 J Joske Vermeulen 1675 pts

4 c christian 2475 pts

9 D Daul 1750 pts

14 k kooken12 1488 pts

5 Q Quick 2375 pts

10 L LoroPol 1725 pts

15 H Hey 1263 pts



5

10

15

Click on the projected screen to start the question



100 %



17



1 W Willie 3125 pts 6 B Bean 2213 pts 11 M Mike 1713 pts

2 J Jesus 2838 pts 7 R Right 1813 pts 11 m ms 1713 pts

3 l lakaka 2788 pts 7 J Jose 1813 pts 13 J Joske Vermeulen 1675 pts

4 c christian 2475 pts 9 D Daul 1750 pts 14 k kooken12 1488 pts

5 Q Quick 2375 pts 10 L LoroPol 1725 pts 15 H Hey 1263 pts

Click on the projected screen to start the question



5

10

15



100 %



15 / 23



1 W Willie 3125 pts 6 B Bean 2213 pts 11 M Mike 1713 pts

2 J Jesus 2838 pts 7 R Right 1813 pts 11 m ms 1713 pts

3 l lakaka 2788 pts 7 J Jose 1813 pts 13 J Joske Vermeulen 1675 pts

4 c christian 2475 pts 9 D Daul 1750 pts 14 k kooken12 1488 pts

5 Q Quick 2375 pts 10 L LoroPol 1725 pts 15 H Hey 1263 pts

Click on the projected screen to start the question



1	W	Willie	3125 pts	6	B	Bean	2213 pts	11	M	Mike	1713 pts
2	J	Jesus	2838 pts	7	R	Right	1813 pts	11	m	ms	1713 pts
3	L	lakaka	2788 pts	7	J	Jose	1813 pts	13	J	Joske Vermeulen	1675 pts
4	C	christian	2475 pts	9	D	Dail	1750 pts	14	k	kooken12	1488 pts
5	Q	Quick	2375 pts	10	L	LoroPol	1725 pts	15	H	Hey	1263 pts



Click on the projected screen to start the question





1

W

Willie

3125 pts

6

B

Bean

2213 pts

11

M

Mike

1713 pts

2

J

Jesus

2838 pts

7

R

Right

1813 pts

11

m

ms

1713 pts

3

L

lakaka

2788 pts

7

J

Jose

1813 pts

13

J

Joske
Vermeulen

1675 pts

4

c

christian

2475 pts

9

D

Dalil

1750 pts

14

k

kooken12

1488 pts

5

Q

Quick

2375 pts

10

L

LoroPol

1725 pts

15

H

Hey

1263 pts

