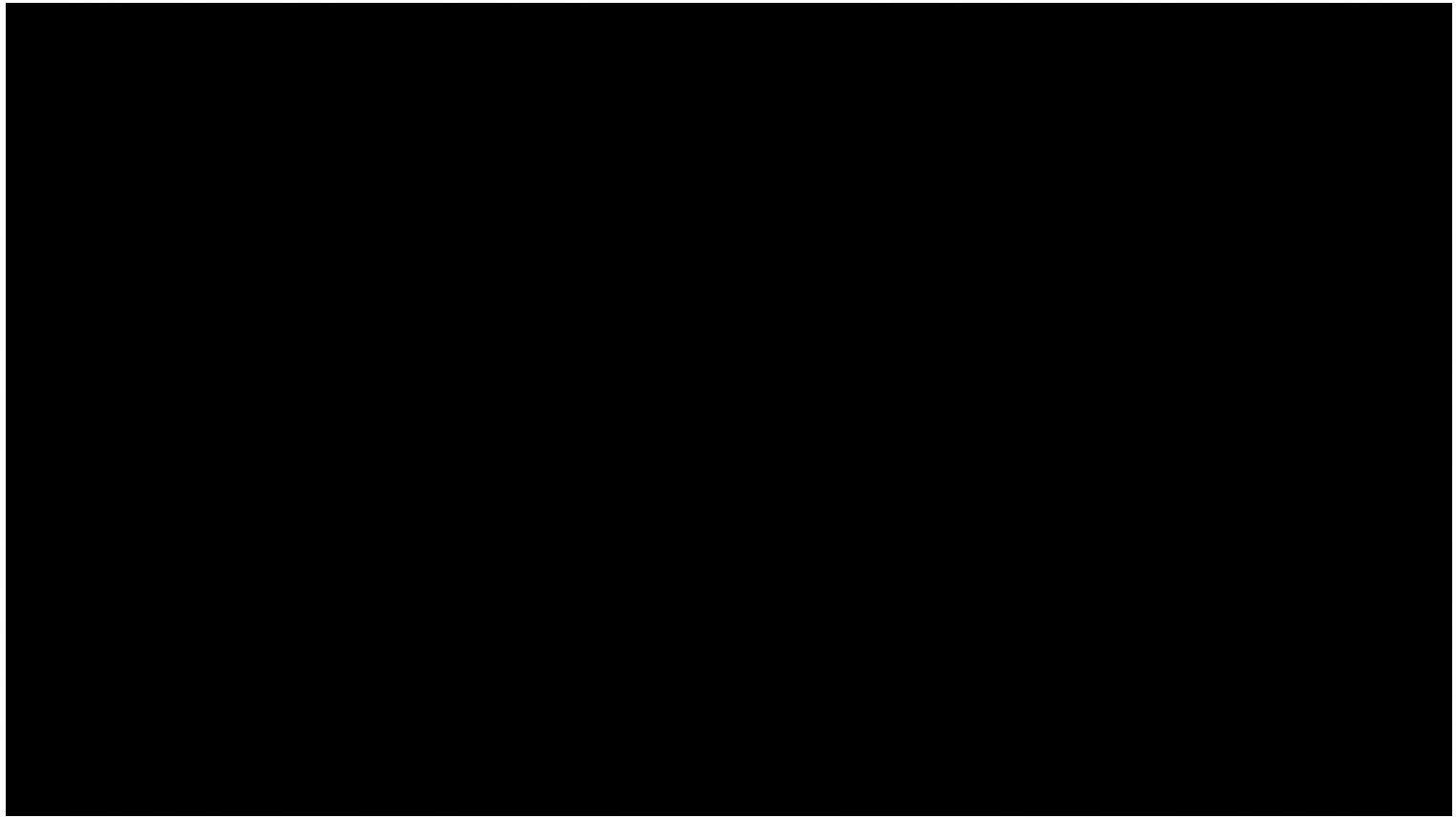


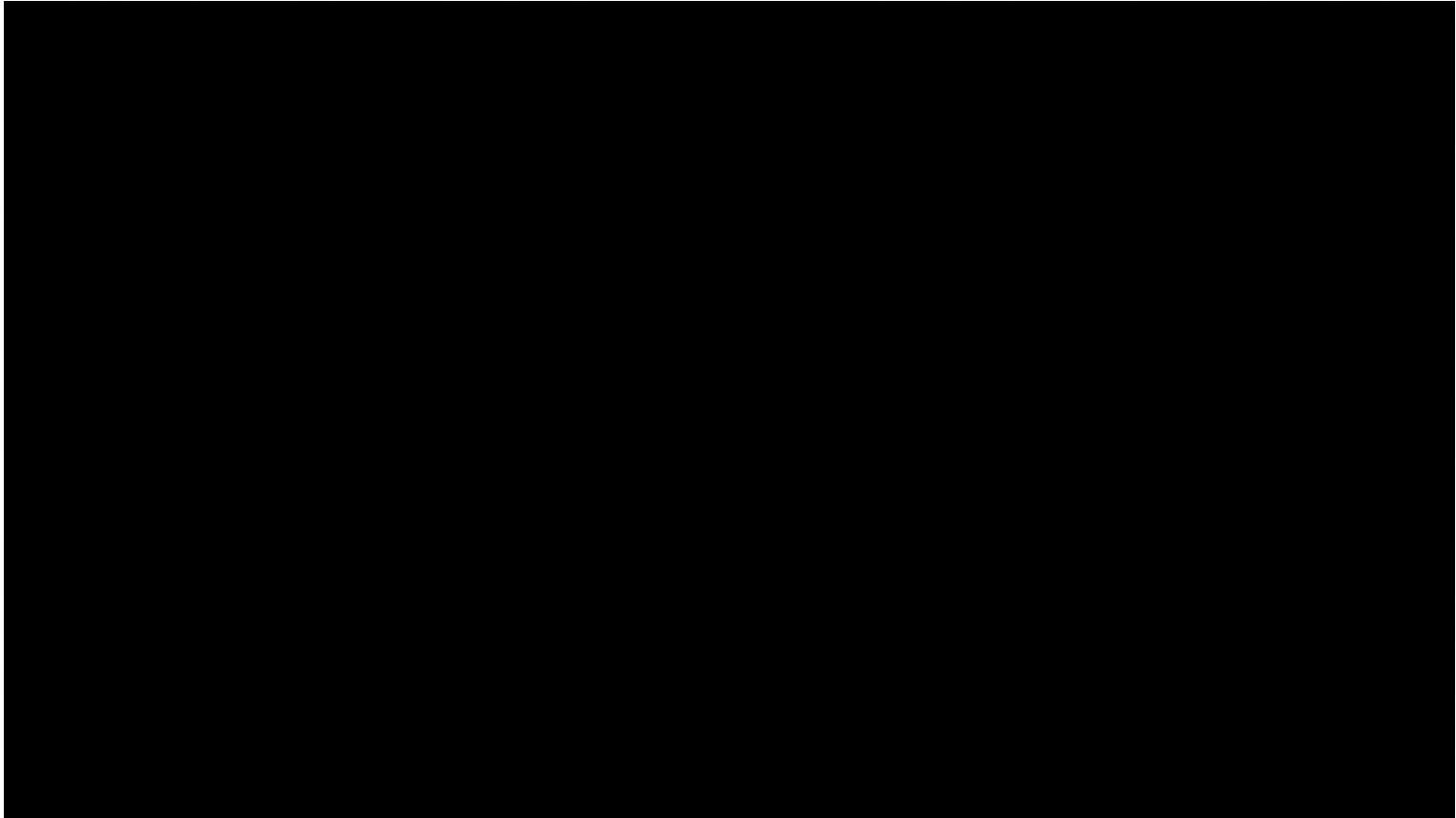
The background features a series of concentric circles in a light gray color, centered on the slide. Overlaid on these circles are stylized circuit board traces in a light blue color. These traces are located in the corners of the slide, with some ending in small circles, resembling electronic components or data paths.

DATA CLEANING

JENS BAETENS



https://www.youtube.com/watch?v=Mp_o_qsyBhA



<https://www.youtube.com/watch?v=-wibyFVm6yg>

KWALITEIT VAN DATASETS IS BELANGRIJK

Slechte data → Slechte $\begin{matrix} \cdot \text{fit}() \\ \cdot \text{predict}() \end{matrix}$
[Garbage in = Garbage out] !

Hoe langer je fouten meesleept in je dataset, hoe kostelijker

- Zelfde als bij software development → *snelere testen is beter*

Fouten in je data kan leiden tot een fout model

PROBLEMEN IN DE DATASETS

Ontbrekende data → *Missing NaN*

Duplicaten

Onmogelijke waarden → *leeftijd is negatief*

Verkeerde dataformaten

→ *columns: ol, mly
ry - m - d
ry - month - d
...*

Onnodige attributen verwijderen

↳ *ID's*
↳ *Foreign Keys*

ONTBREKENDE DATA

Ontbrekende data:

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

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

DUPLICATE DATA

Verwijder exacte duplicaten

Gebruik duplicaten om ontbrekende informatie in te vullen

(Hoe behandel je conflicterende data?

- Zelfde naam maar ander adres bijvoorbeeld

→ rij volledig matchen

→ # kolommen "

→ ID is voldoende

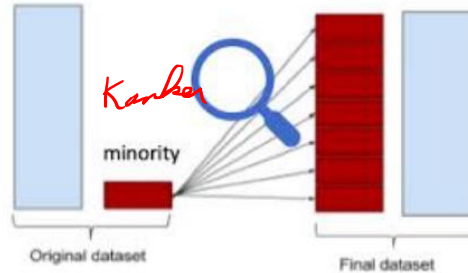
⇒ Welke data
hou je over

DATA BALANCING

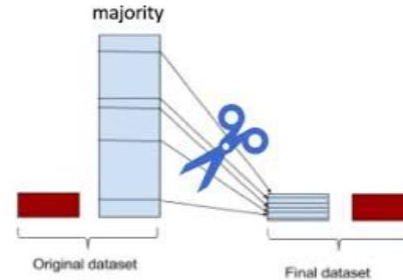
Rijen die waarden bevatten die zelden voorkomen worden best niet verwijderd

Welke kiezen afh van hoeveelheid beschikbare data

geronol



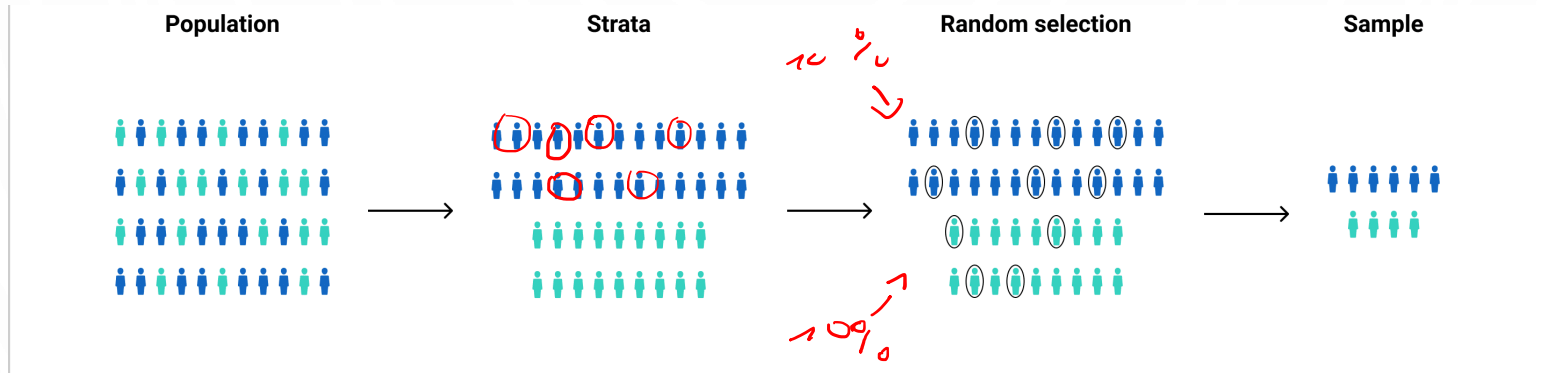
VS



Stratification argument doet dit automatisch in `train_test_split`

Stratificatie ↔ niet bij regressie

DATA BALANCING – STRATIFIED SAMPLING



Stratification argument doet dit automatisch in `train_test_split`

Verdeling blijft behouden

CORRIGEER DATA FORMAT

Typos

Verschillende waarden met dezelfde betekenis:

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

Brussel
Brussels
Bruxelles

Vertalen!

ONE HOT ENCODING VS ORDINAL

classes ~~klassen~~

Kleur
Rood 0
Geel 1
Blauw 2
Blauw 2
Geel 1

3 Klassen

ordinal

Kleur
0
1
2
2
1

→ wat met-afstand?

One Hot Encoding

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

→ tot ±15 klassen

PRIVACY REQUIREMENTS

Zoek naar Personal Identifiable Information (PII)

Deze velden moeten beter afgeschermd zijn dan andere

- Data niet bruikbaar in het geval van hacking

) Big Data

→ data owner control

Data masking of obfuscation

op de data zelf



DATA MASKING – DELETION

Verwijder de gevoelige informatie

Zeer simplistische aanpak

Verwijderde data kan niet meer gebruikt worden in je model

Geeft aan dat er zaken gewijzigd zijn



DATA MASKING - SUBSTITUTION

Wijzig PII door willekeurige zelf gekozen waarden

- Kan afhankelijk zijn van andere velden
- Vb: wijzig naam maar wel afhankelijk van het opgegeven geslacht

Shuffling

- Lijst met mogelijke waarden is de kolom zelf

Naam | - - - -

<i>Jens</i>	<i>A</i>
<i>Leon</i>	<i>B</i>
<i>Jeanine</i>	<i>C</i>



<i>Leon</i>	<i>A</i>
<i>Jens</i>	<i>B</i>
<i>Jeanine</i>	<i>C</i>

DATA MASKING – VARIANTIE TOEVOEGEN

Voeg ruis toe aan de data

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

DATA MASKING – ENCRYPTION

Encrypteer de gevoelige kolommen

Hashing: Indien enkel een identifier moet hebben → *1-way / original way*

Encryption: Indien de waarde nog moet gebruikt kunnen worden

Nadelen:

- Beheer van keys → *gestolen key → alles kan gelezen worden*
- Encryptie en decryptie reken-intensief