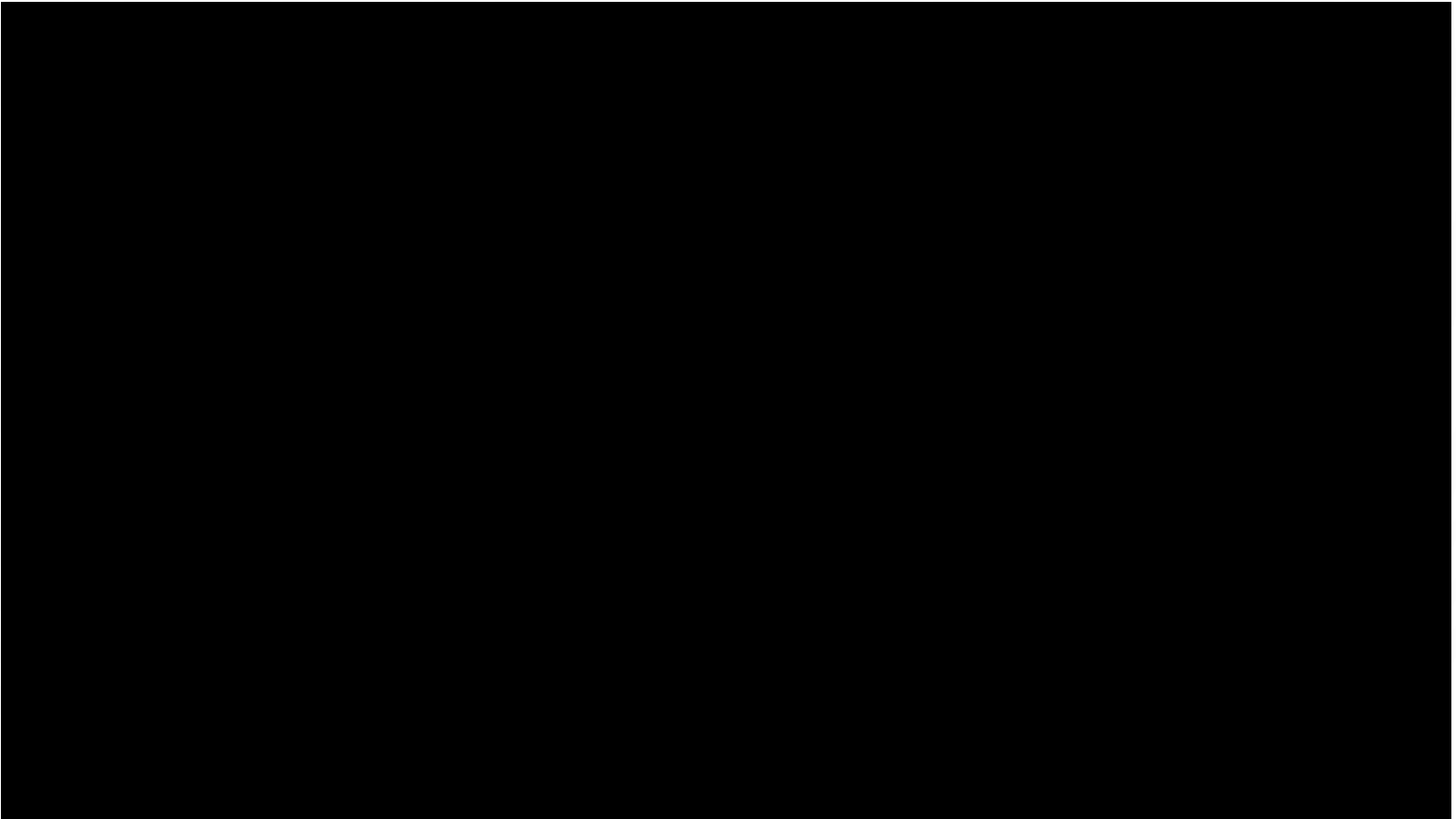


The background features a series of concentric circles in a light gray color, centered on the slide. Overlaid on these are stylized circuit board traces in a light blue color, with small circles at the junctions, located in the corners of the slide.

# **DATA SCIENCE - LIFECYCLE**

JENS BAETENS

# Problem statement



<https://www.youtube.com/watch?v=uO7c2tvrPj0>

01

## BUSINESS UNDERSTANDING

Ask relevant questions and define objectives for the problem that needs to be tackled.

Wat is de gestelde vraag of het probleem? *open-ended*

Formuleer de vragen waarop een antwoord moet gevonden worden

5 soorten vragen:

- Hoeveel? →
- Wat is het? →
- Is het sterk gelijkend op?
- Is het vreemd?
- Welke optie is het beste?

Regressie

Classificatie

Clustering

Anomaly Detection

Recommendation

*getal*  
*keuze voorstellen*



*Welke plane is het*



*Preventive maintenance*

*Netflix*

*gerelateerde artikels*



Verzamel data van verschillende bronnen

Welke data is er nodig?

Hoe geraak ik aan deze data?

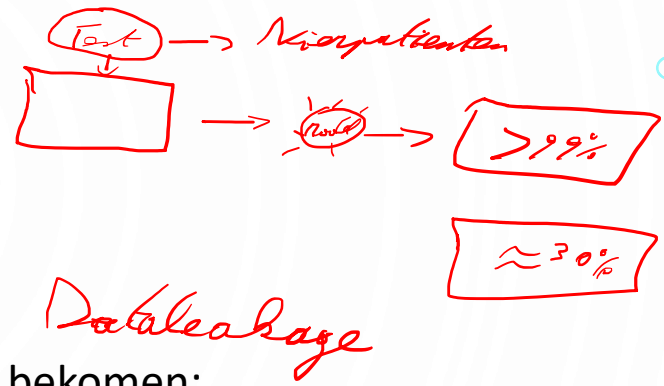
- Lokale databases
- Scraping van webpaginas
- Verzamelen van data van sensoren / apps / satellieten ...

Hoe bewaar ik de verzamelde data? → *lokale of H100 -*  
→ *Nas*  
→ *server*  
→ *cloud* } *Big Data*

03

### DATA CLEANING

Fix the inconsistencies within the data and handle the missing values.



Belangrijke stap voor betrouwbare resultaten te bekomen:

- Garbage In -> Garbage Out

Het doel is om problemen op te lossen in de datasets:

- Ontbrekende data —
- Verkeerd gelabelde data (0/1 vs true/false)
- Verschillende dataformaten (male/m/Male or dates)
- Verbeteren van typos, vertalen van sommige velden, ...

2  
3  
4



11 1% kans  
99% gezond } altijd gezond  
99% juist  
geen andere kans  
getoetst

histogrammen:  
↑ verbanden  
→ correlaties

Fase waarin je de verzamelde data bestudeert

Zoek naar bestaande patronen en controleer of er een bias aanwezig

Visualiseer en analyseer deze patronen

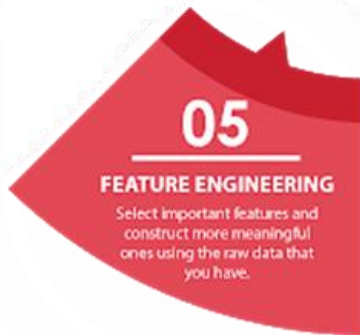
Detecteer outliers } → verwijderen  
→ behouden

Stel een aantal hypothesen voor

Ook exploratory data analysis genoemd:

[https://en.wikipedia.org/wiki/Exploratory\\_data\\_analysis](https://en.wikipedia.org/wiki/Exploratory_data_analysis)

EDA



*kolom/features/variabelen*

*rij 1  
rij 2  
rij 3  
↓  
→ dronatie datapunt*

Feature = Een meetbare eigenschap van een geobserveerd datapunt

Het zoeken naar de beste features van je data om je vraag op te lossen

- Vereist domein kennis om deze te bepalen/berekenen

Feature Selection

- Verwijder onbruikbare features/datapunten
- Curse of dimensionality }

- Feature Construction

*↳ # features*

- Nieuwe features op basis van bestaande

*o.l.*

*→ ontbrek*

*→ opp*

- Vaak belangrijk in het geval van beelden

- vb: Enkel geïnteresseerd of iemand volwassen is en niet de exacte leeftijd.

06

## PREDICTIVE MODELING

Train machine learning  
models, evaluate their  
performance, and use  
them to make predic-  
tions.

### Machine learning model opbouwen

*train  
evalueren*

Probeer verschillende varianten en evaluateer elk model

- Zie cheat sheet voor een aantal mogelijkheden

Beste keuze hang af van:

- Hoeveelheid, type en kwaliteit van de data
- Beschikbare computer-capaciteit
- Gewenste output type → *Kans*  
→ *is het die klasse?*  
:



07

## DATA VISUALIZATION

Communicate the findings with key stakeholders using plots and interactive visualizations.

Visualiseer de resultaten van het resulterende model

Ook de behaalde inzichten tijdens het process zijn belangrijk

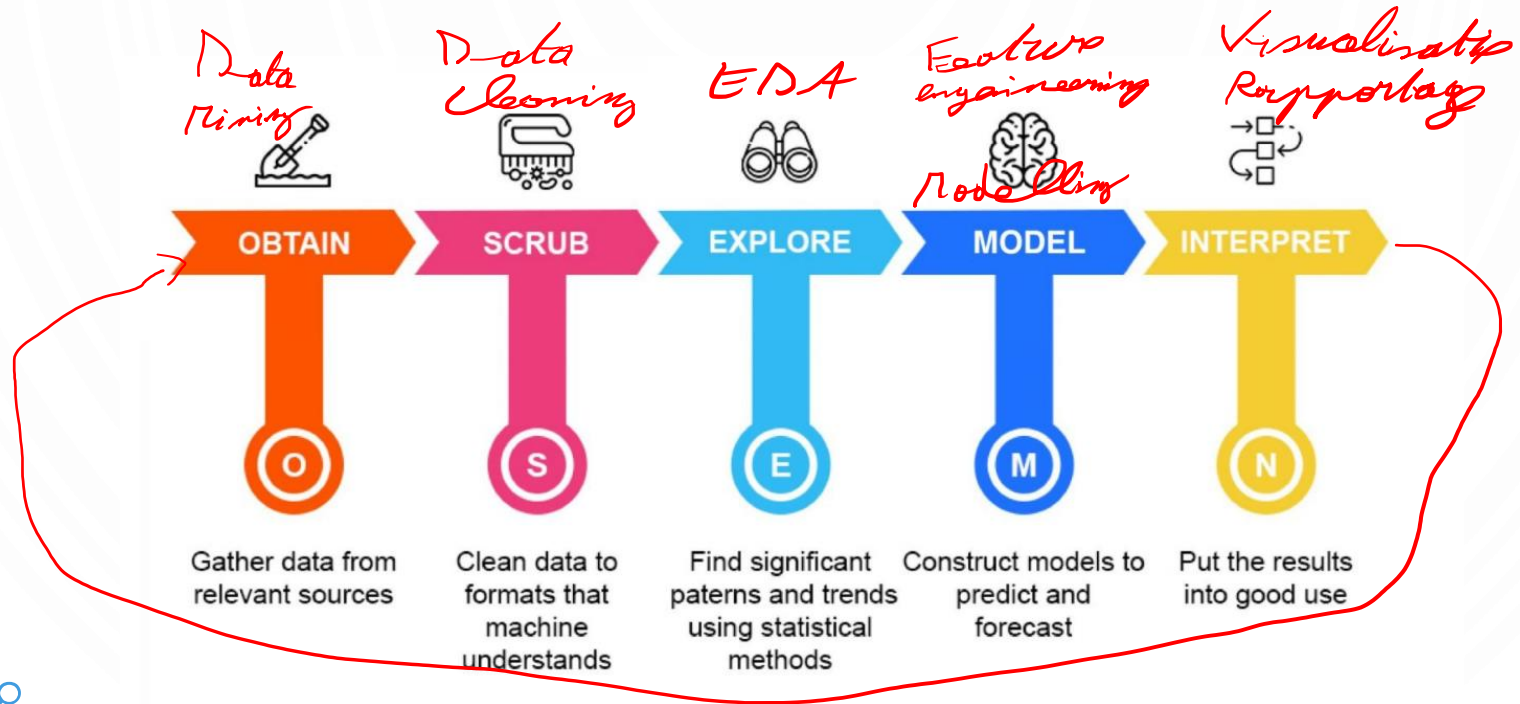
*EDA + modelling*

De communicatie moet aangepast zijn aan de verschillende stakeholders

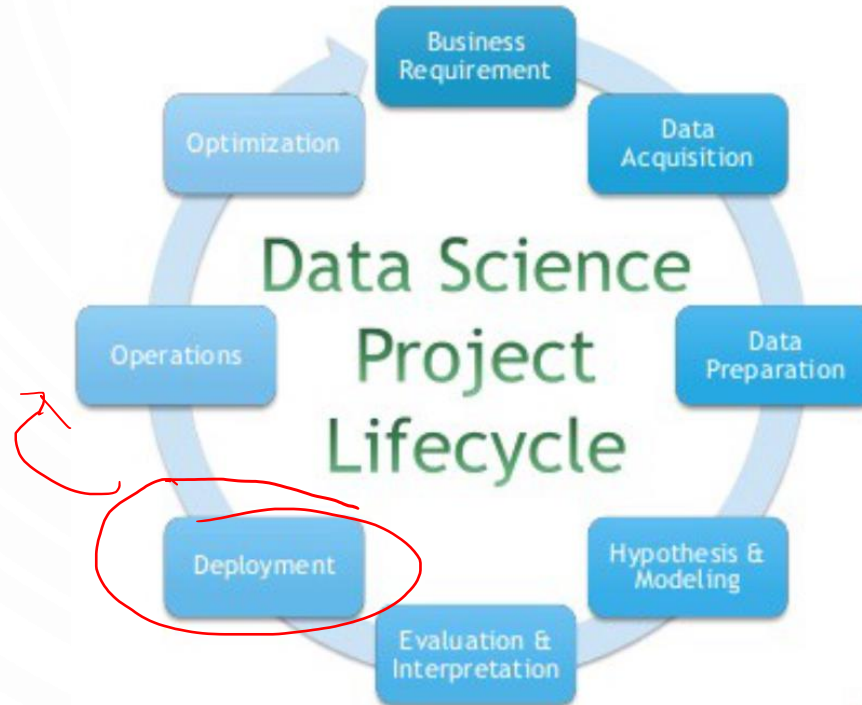
*\$, Cyol      verspaart      Privacy  
                         Rot-de Blaant      - Milieu*



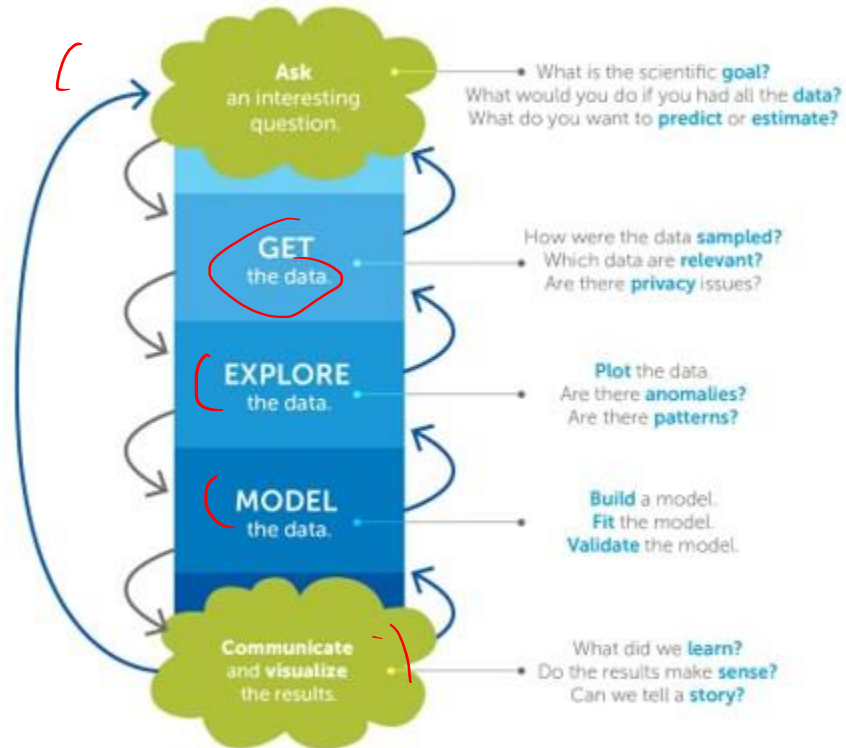
# ANDERE MOGELIJKE LIFECYCLES



# ANDERE MOGELIJKE LIFECYCLES



# The Data Science Process



Derived from the work of Joe Blitzstein and Hanspeter Pfister, originally created for the Harvard data science course <http://cs109.org/>.



# RESOURCES

<http://sudeep.co/data-science/Understanding-the-Data-Science-Lifecycle/>

[https://en.wikipedia.org/wiki/Exploratory\\_data\\_analysis](https://en.wikipedia.org/wiki/Exploratory_data_analysis)

<https://docs.microsoft.com/en-us/azure/machine-learning/algorithm-cheat-sheet>