

$$F = G \frac{m_1 m_2}{d^2}$$

$$\phi(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

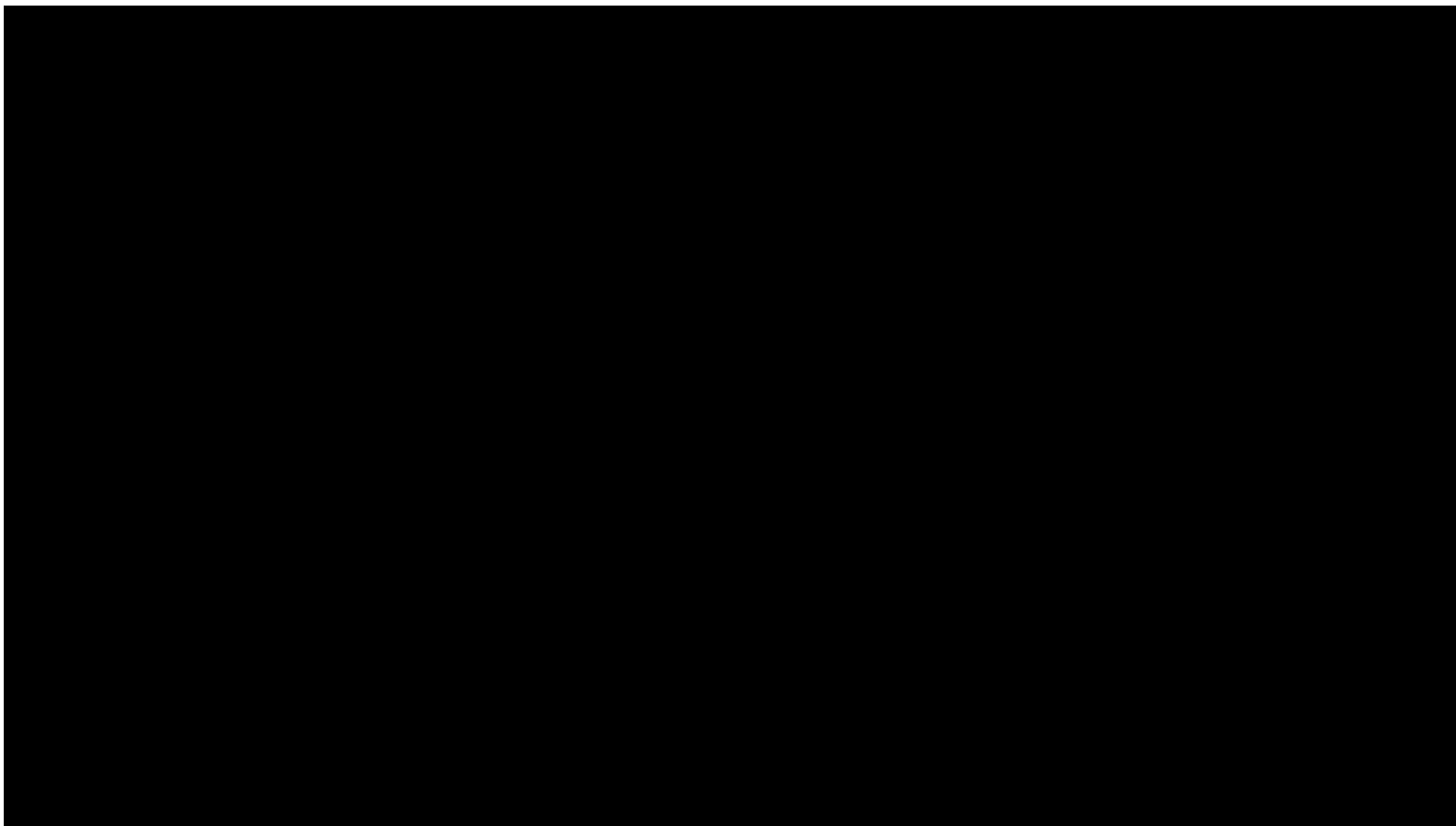
$$i\hbar \frac{\partial}{\partial t} \psi = \hat{H} \psi$$

Data SCIENCE - Lifecycle

JENS BAETENS

$$\frac{\partial^2 u}{\partial t^2} = c^2 \frac{\partial^2 u}{\partial x^2}$$

$$\frac{df}{dt} = \lim_{h \rightarrow 0} \frac{f(t+h) - f(t)}{h}$$



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



Wat is de gestelde vraag of het probleem?

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



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?



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, ...



Fase waarin je de verzamelde data bestudeerd

Zoek naar bestaande patronen en controleer of er een bias aanwezig

Visualiseer en analyseer deze patronen

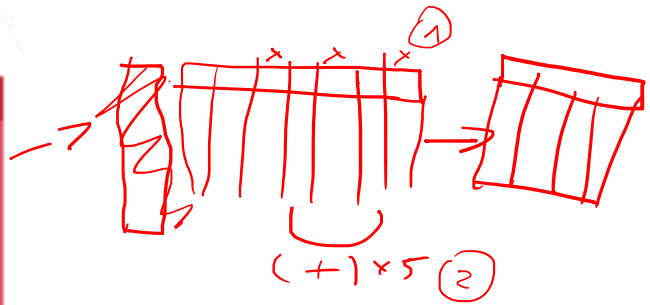
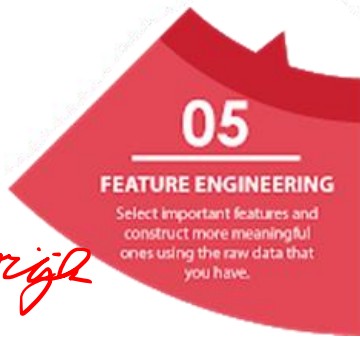
Detecteer outliers

Stel een aantal hypotheses voor

Ook exploratory data analysis genoemd:

https://en.wikipedia.org/wiki/Exploratory_data_analysis

- Kolom/waarvan je denkt dat relevant is



Feature = Een meetbare eigenschap van een geobserveerd datapunt

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

Stappen - Vereist domein kennis om deze te bepalen/berekenen

① Feature Selection \rightarrow erwaring!

- Verwijder onbruikbare features/datapunten

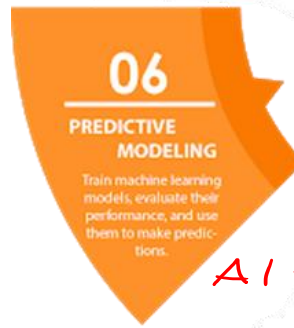
- Curse of dimensionality \rightarrow heel veel features

② Feature Construction

- Nieuwe features op basis van bestaande

- Vaak belangrijk in het geval van beelden

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



AI-step → trial & error

Machine learning model opbouwen

Probeer verschillende varianten en evalueer 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



07

DATA VISUALIZATION

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

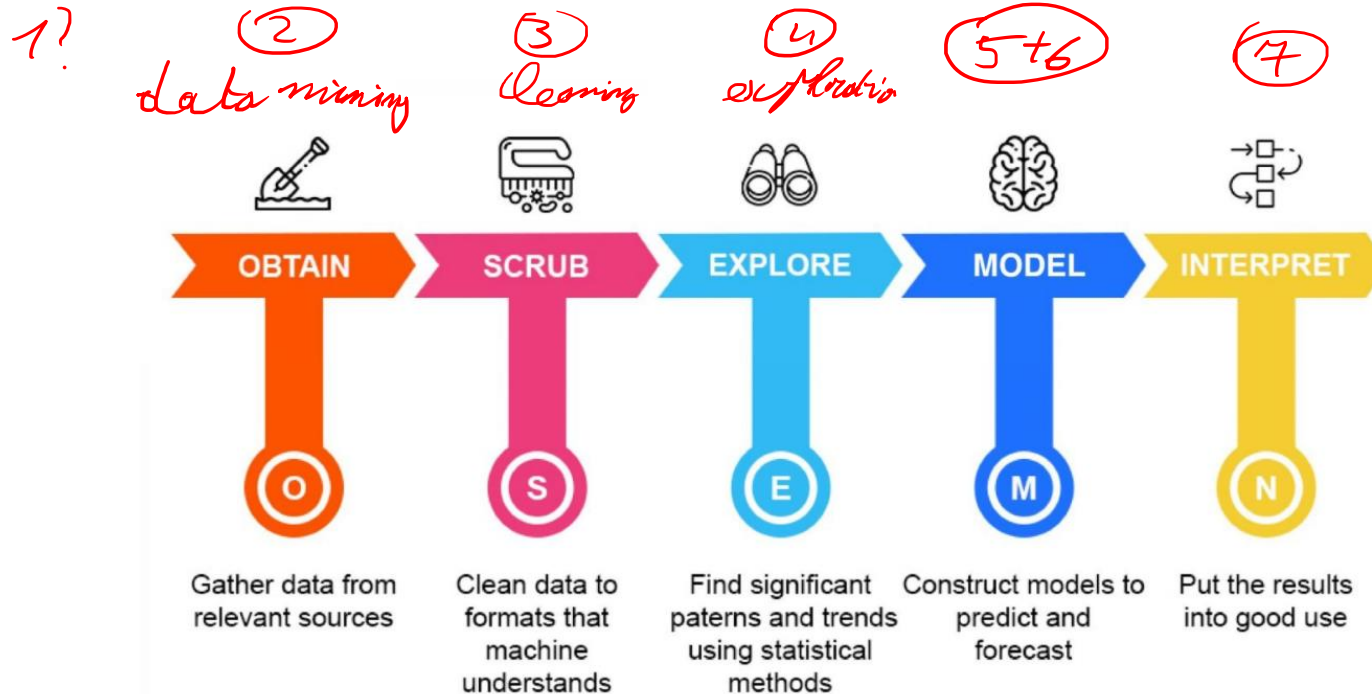
Visualiseer de resultaten van het resulterende model

Ook de behaalde inzichten tijdens het proces zijn belangrijk

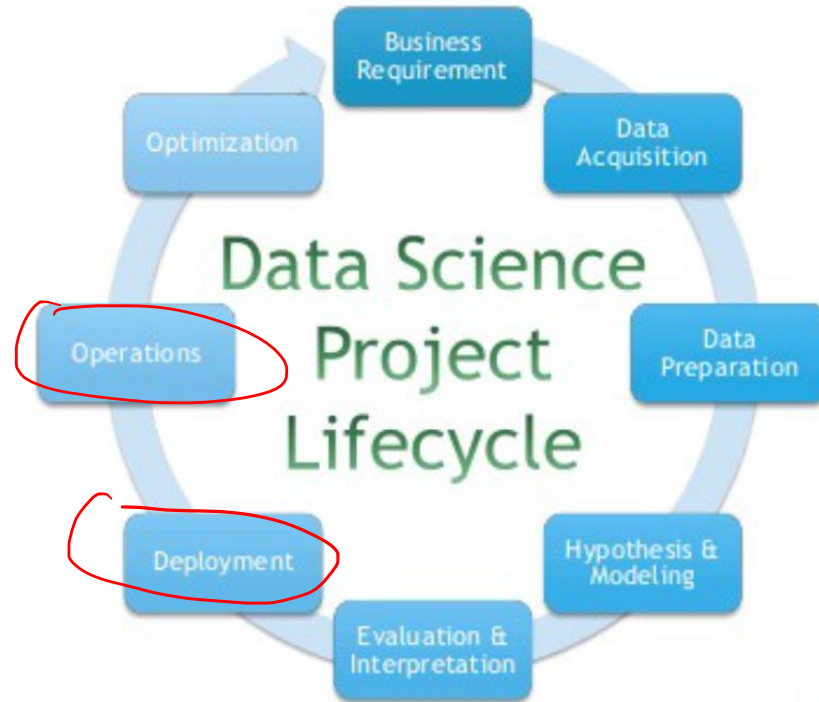
De communicatie moet aangepast zijn aan de verschillende stakeholders



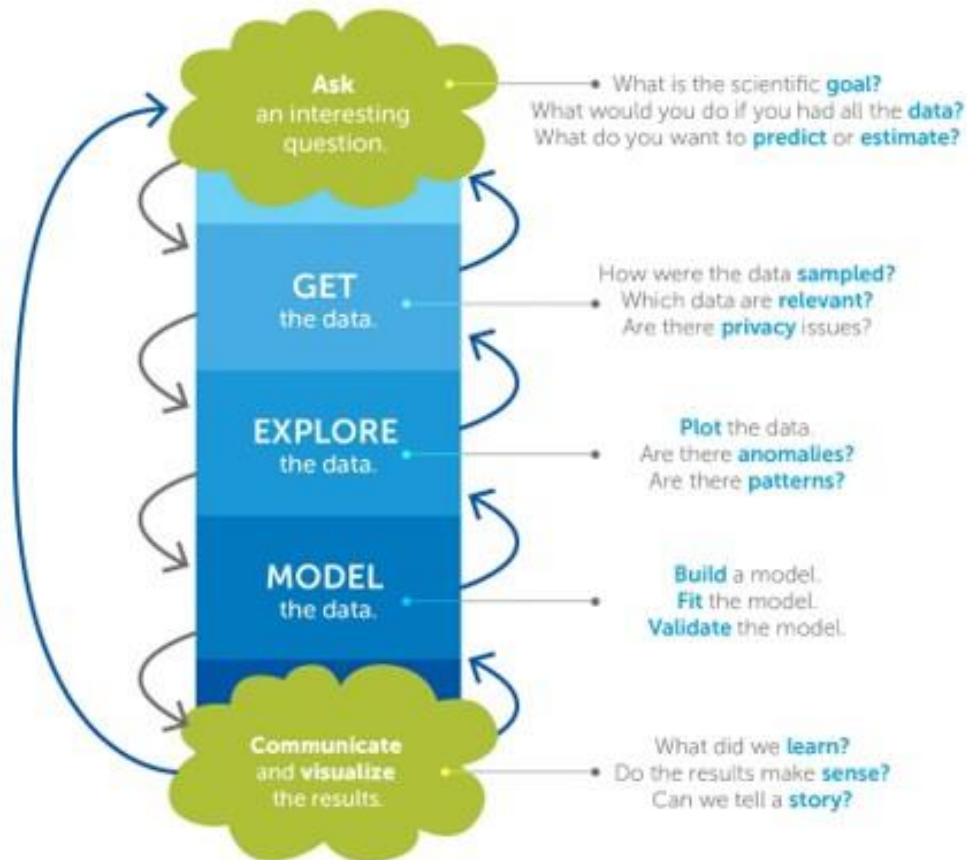
Andere mogelijke lifecycles



Andere mogelijke lifecycles



The Data Science Process



How to participate?



Wat zijn de te volgen stappen in de data science ...



Most frequent combinations:

5

6. Business understanding
2. Data mining
1. Data cleaning
3. Data exploration
- ☒ 5. Feature engineering
4. Predictive modeling
7. Data vizualisation

2

1. Data cleaning
2. Data mining
3. Data exploration
4. Predictive modeling
5. Feature engineering
6. Business understanding
7. Data vizualisation

2

6. Business understanding
2. Data mining
7. Data vizualisation
3. Data exploration
1. Data cleaning
4. Predictive modeling
5. Feature engineering

Wat is het belangrijkste dat er gebeurt in elke stap van de Data Science lifecycle



The most frequent answers are

Business Understanding

← 16 👤 →

Zoeken naar een op te lossen vraag

Data Mining

← 16 👤 →

Verzamelen van data

Data Vizualization

← 15 👤 →

Rapporteer je resultaten

Predictive modelling

← 13 👤 →

Train een model voor de vraag te beantwoorden

Data Cleaning

← 12 👤 →

Oplossen van fouten in de data

Data Exploration

← 9 👤 →

Zoeken naar verbanden in de data

Data Cleaning

← 6 👤 →

Splits de data af die je gaat gebruiken



Resources

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

https://en.wikipedia.org/wiki/Exploratory_data_analysis

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