

Meepraten over betrouwbare en verantwoorde AI

Wat moet ik weten?

Cynthia C. S. Liem

Multimedia Computing Group, TU Delft

[!\[\]\(666e09182d4cd268646ea700ea60dcdf_img.jpg\) @cynthiacsliem@akademienl.social.nl](mailto:@cynthiacsliem@akademienl.social.nl)

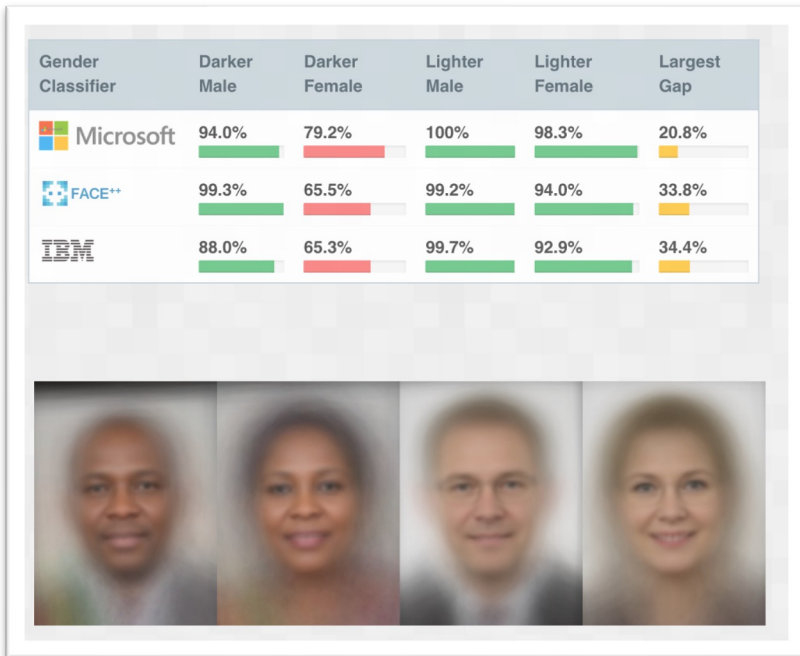
c.c.s.liem@tudelft.nl

Spraakverwarring!

- Imprecies taalgebruik
 - AI, algoritme, systeem, robot bijna synoniem in publiek debat
 - Heel verschillende jargontermen voor informatici
 - ‘Wat is een algoritme dan exact?’ / ‘Wat is het exacte algoritme?’



Schandalen



<http://gendershades.org/overview.html>

Amazon scraps secret AI recruiting tool that showed bias against women

deVerdieping

rouw

WWW.TROUW.NL | 7930E JAARGANG NR 2303 | 1x MAAL VIA WWW.TROUW.NL/DEEPIE | POSTBOS 800 0000 AW AMSTERDAM | REDACTIE DOO 302 3444 | KLAARTEDEEPIE DOE 096 088 | BEZORGER: 096 086 010

TROUW CIRCULO 21 NOVEMBER 2021

Het is de hoogste tijd dat al die onbetrouwbare werknemers een eerlijk loon krijgen

Europaformateur Agnes Jongenius, economie 16

Belastingdienst ging vooral achter lage inkomens aan

Fraudejacht • Om toeslagen te controleren op fouten en fraude gebruikte de Belastingdienst een zelflerende algoritme. Dat selecteerde vooral mensen met lage inkomens voor controle.

Jan Klauwdenhaas
redacteur Inroepen

De Belastingdienst heeft inroepen opgedrukt bij mensen met een laag inkomen. Het was niet de bedoeling. Het was een fout van de Belastingdienst. Het was een fout van de Belastingdienst. Het was een fout van de Belastingdienst.

In het kort

De Belastingdienst gebruikte in de coronatijd op fraude een zelflerende algoritme. In de praktijk werd dat model vooral gebruikt om mensen met een laag inkomen te controleren. Het was geen proces om deze inkomens te controleren.

Er zijn nu vijfde van de controletoeslagen teruggevoerd. Volgens Liem zijn de controletoeslagen teruggevoerd. Volgens Liem zijn de controletoeslagen teruggevoerd. Volgens Liem zijn de controletoeslagen teruggevoerd.

Als de uitkomst zo inzoomt op een specifieke groep dan moeten er alarmbellen afgaan

Cynthia Liem
universiteit hoeddehoken kunstmatige intelligentie

82,3 procent van de groep met de hoogste inkomens had een inkomens van minder dan 20.000 euro

Plannen voor code zwart

Ziekenhuizen in de knel vrijdag 4B

Hoe gaat een agent om met rellen?

Geestelijke zorg bij hulpdiensten vrijdag 7

Criminelen werken op feestdagen

Online oplichters schakelen tactiek bij rotsend 9

Nieuw huis voor De Schreeuw

Oslø heeft een Munchmuseum de verdieping 1213

<https://www.trouw.nl/politiek/hoe-de-belastingdienst-lage-inkomens-profileerde-in-de-jacht-op-fraude~bbb66add/>

Wat nu?

- We spreken momenteel veel over algoritmen, systemen en verantwoordelijkheid.
- Volgens mij gaan deze discussies eigenlijk niet over computersystemen, maar over **hoe wij als mensen met beslisvorming omgaan**.
- Wij (informatici en niet-informatici) moeten heel goed opletten hoe deze discussies gevoerd worden.
- Niet-informatici durven vaak niet mee te doen aan technische/wiskundige discussies. Dat moet anders.

En...

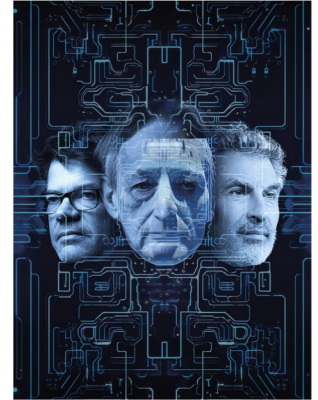
- De AI-discussie neemt momenteel sensationele en bangmakende vormen aan
- Hier lopen veel belangen door elkaar heen
<https://www.nrc.nl/nieuws/2023/06/24/sterft-de-mensheid-uit-door-ai-dat-is-sciencefiction-a4168053>
- Schadelijke beeldvorming

Opnieuw: de kernvragen gaan om **hoe wij als mensen (en instanties) met beslisvorming omgaan.**

The godfathers of AI are at war

Existential threat or humanity's salvation?
The founders of AI don't know what's next.

By Harry Lambert



left to right: Yann LeCun, Geoffrey Hinton and Yoshua Bengio, the three

Mitigating the risk of extinction from AI should be a global priority alongside other societal-scale risks such as pandemics and nuclear war.

Algoritmes 101

Klassiek probleem: sorteren

- Kan op slimme en domme manieren worden gedaan
- Wat is efficiënt bij grote hoeveelheden?
- Wat is toepasbaar op alle problemen waarbij je een volgorde kunt vaststellen?



Bijbehorende waarden

- Efficiëntie
- Generalisatie
- Schaalvergroting
- Exact dezelfde aanpak hanteren

Zie ook Birhane et al. (2022), The Values Encoded in Machine Learning Research,
<https://dl.acm.org/doi/10.1145/3531146.3533083>

Bijbehorende waarden

- Efficiëntie
- Generalisatie
- Schaalvergroting
- Exact dezelfde aanpak hanteren

- Dit kan mensen ‘tot een nummer’ maken

Zie ook Birhane et al. (2022), The Values Encoded in Machine Learning Research,
<https://dl.acm.org/doi/10.1145/3531146.3533083>

Bijbehorende waarden

- Efficiëntie
- Generalisatie
- Schaalvergroting
- Exact dezelfde aanpak hanteren

- Tegelijkertijd kan het ons een spiegel voorhouden

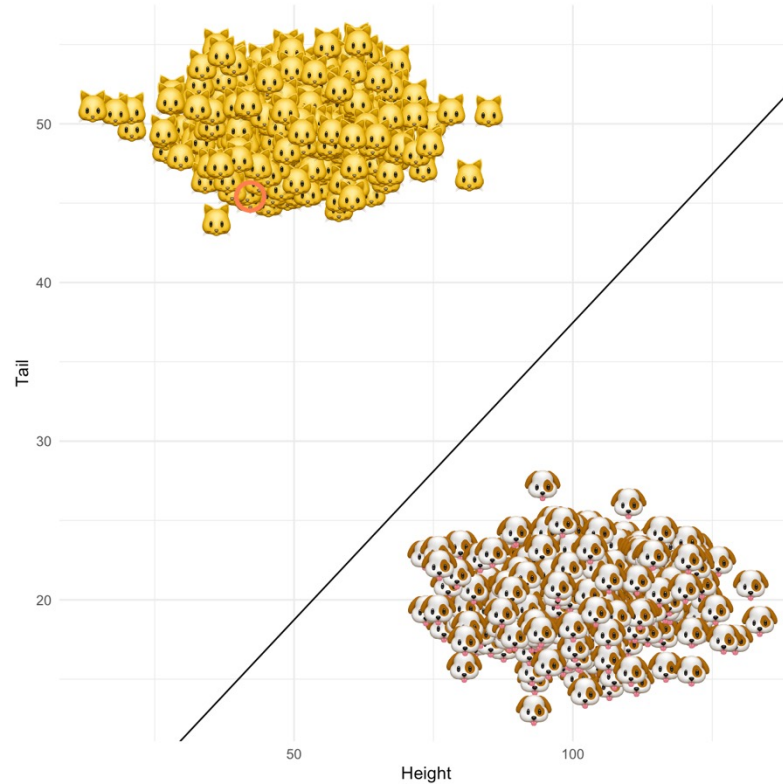
Precies zijn over stappen is moeilijk!



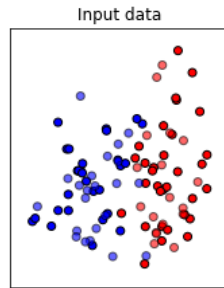
https://www.youtube.com/watch?v=cDA3_5982h8

Machine learning 101

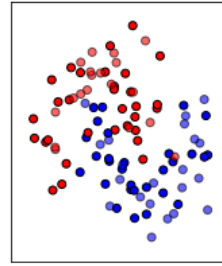
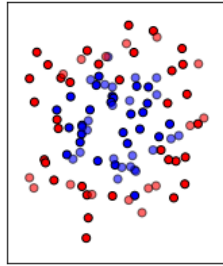
Patronen leren uit data en labels



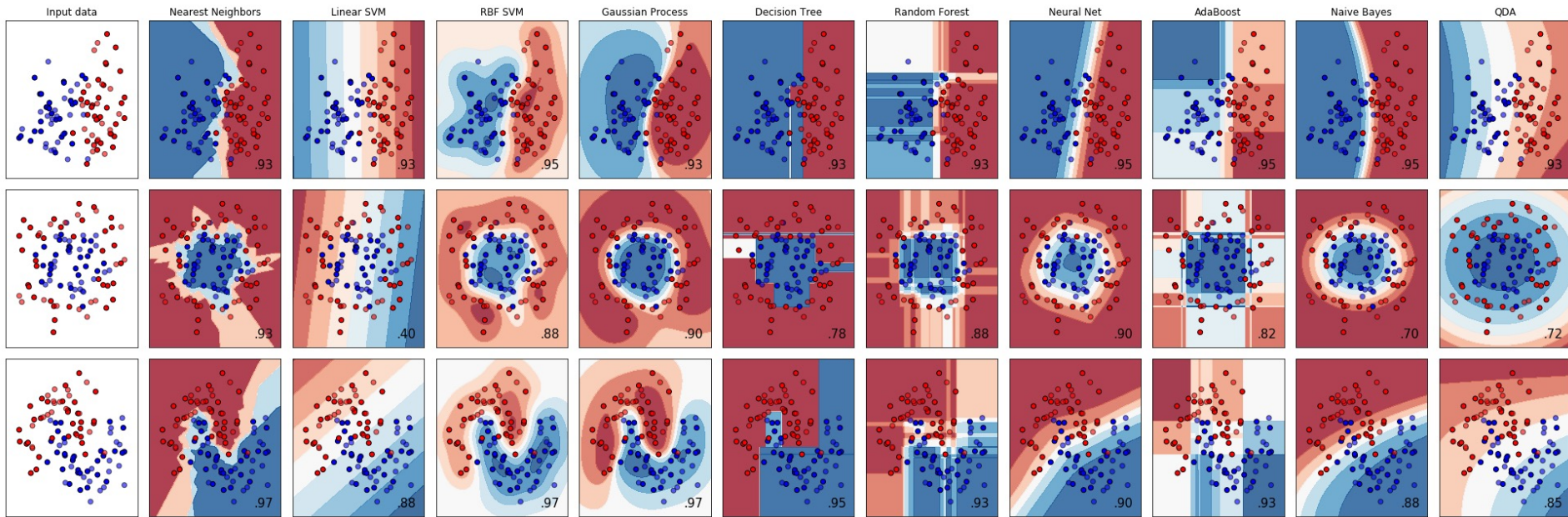
Patronen leren uit data en labels



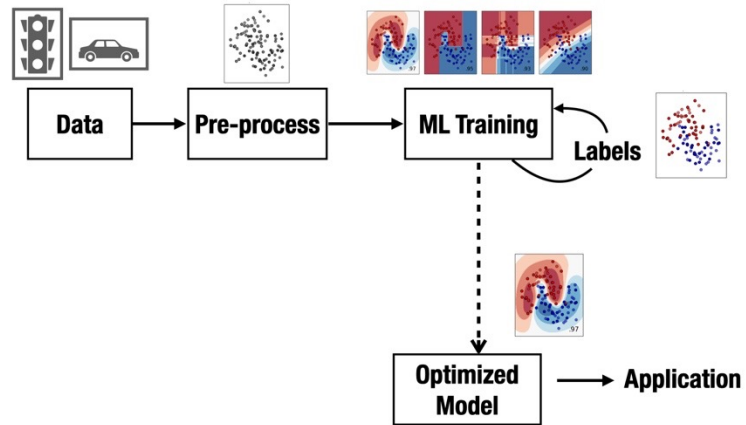
Patronen leren uit data en labels



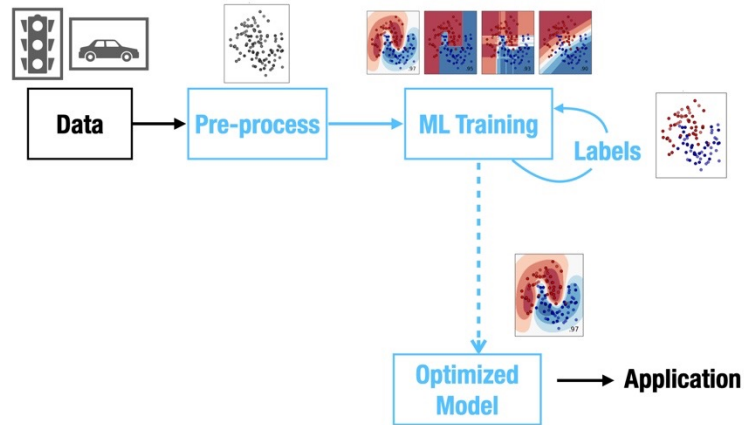
Patronen leren uit data en labels



De opzet

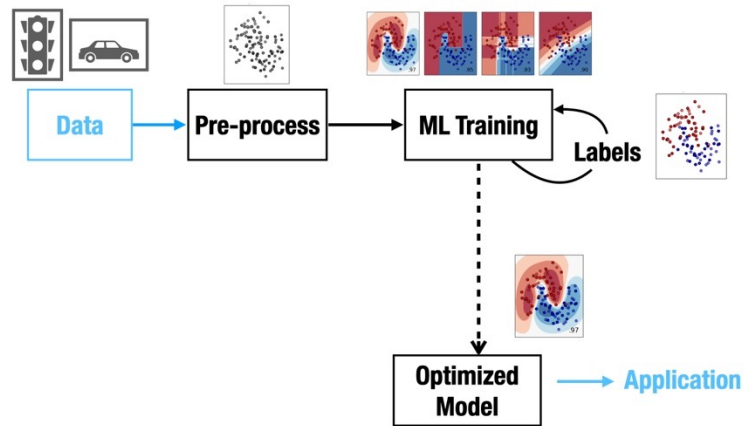


Focus van machine learning-expert



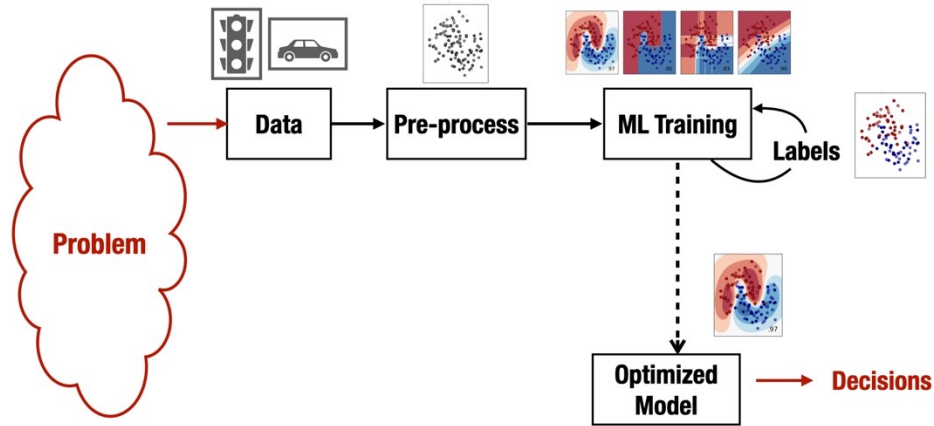
- Wiskundig optimale oplossingen
- We vertrouwen erop dat de data en toepassing geldig zijn

Focus van domeinexpert



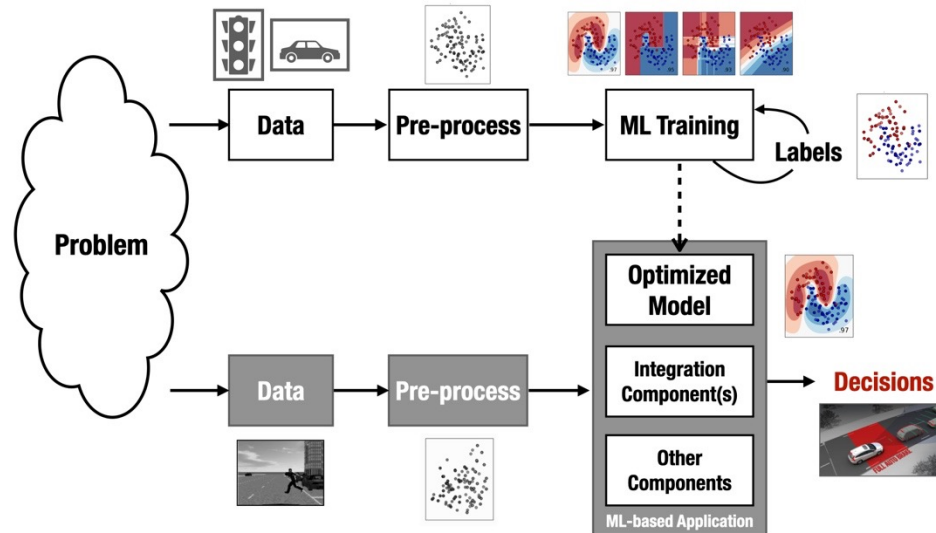
- Zijn de juiste modaliteiten, metingen, kenmerken meegenomen?
- Welk wereldbeeld optimaliseren we?

Waar het eigenlijk om gaat



- Welk probleem willen we werkelijk oplossen?
- Is dit probleem goed ‘vertaald’?
- Trekken we de juiste conclusies?

...of eigenlijk gaat het hierom



- Socio-technische systemen met veel componenten
- Moeilijk een enkele verantwoordelijke aan te wijzen

Patronen als pijnlijke, maar realistische spiegel

Amazon scraps secret AI recruiting tool that showed bias against women

av_nationaliteit	1, _MISSING_, _UNKNOWN_	8
	0	-1

- Patronen worden in context problematisch
- Let hierop: niet meten kan ook ruimte geven voor wegduiken

Hoe 'ziet' een computer de wereld?

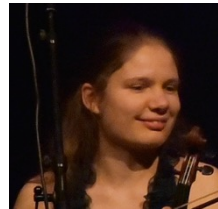
Ruwe multimediate data: heel veel getallen

- 44.1 kHz audio:
44,100 metingen per seconde
- Voor spraak is 8 kHz voldoende. Dat gaat nog steeds om 8000 metingen per seconde. In 45 minuten heb je dan 21.6 miljoen metingen...



← Dit portretje: $224 \times 224 \times 3$ pixels = 150,528 intensiteitswaarden

Wat is het dichtst bij?



Accuraat kunnen voorspellen \neq objectieve natuurwetenschappelijke waarheid

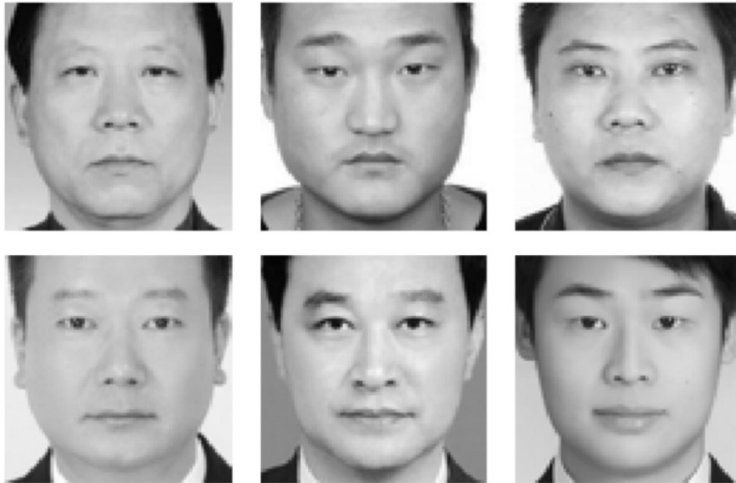


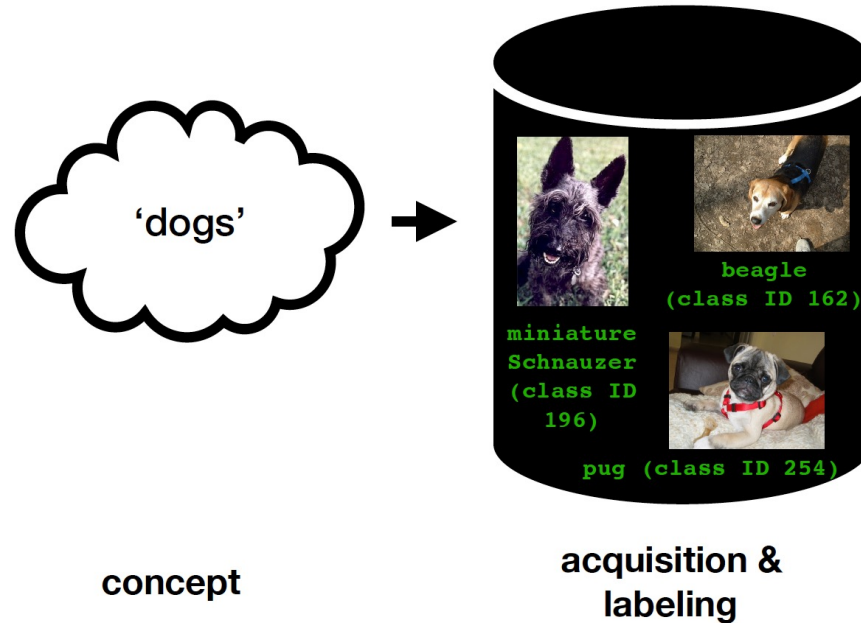
Figure 3. Wu and Zhang's "criminal" images (top) and "non-criminal" images (bottom). In the top images, the people are frowning. In the bottom, they are not. These types of superficial differences can be picked up by a deep learning system.

Deep neural networks are more accurate than humans at detecting sexual orientation from facial images.

Article | [Open Access](#) | Published: 31 March 2023

Using deep learning to predict ideology from facial photographs: expressions, beauty, and extra-facial information

Classificeren en categoriseren

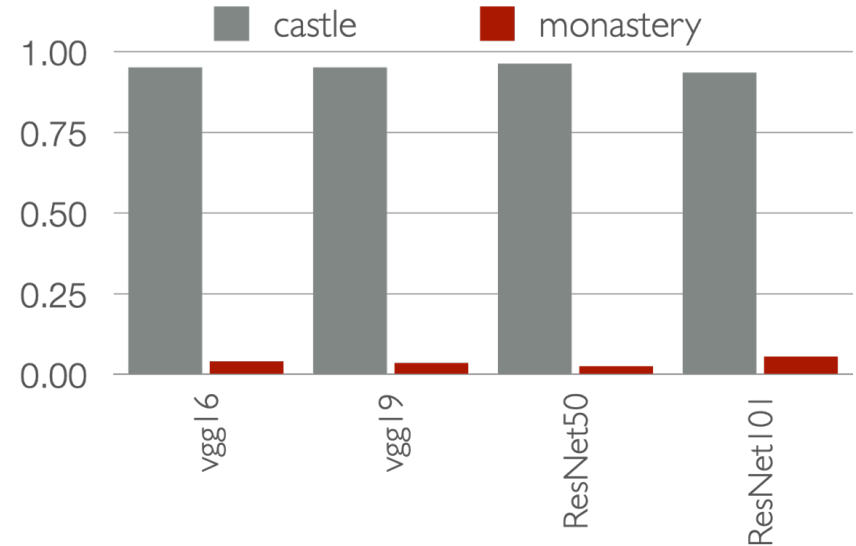


‘Ground truth’: het correcte label

- Is er een kasteel in dit plaatje?

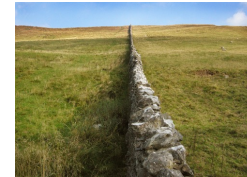
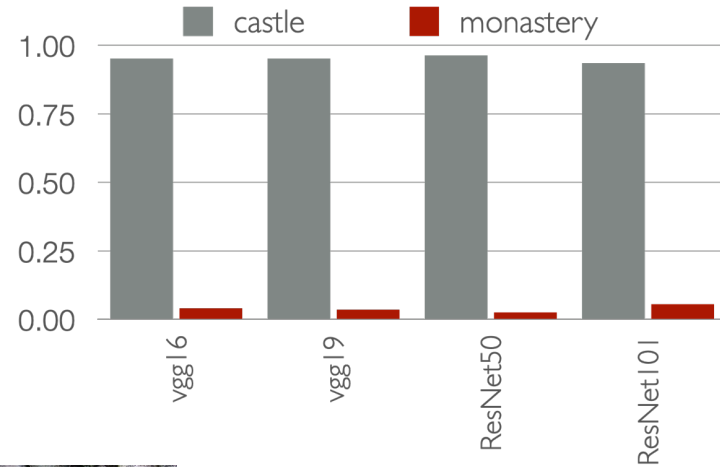


Is er een kasteel in dit plaatje?



Cynthia C. S. Liem and Annibale Panichella, "Oracle Issues in Machine Learning and Where to Find Them," Proc. RAISE 2020.

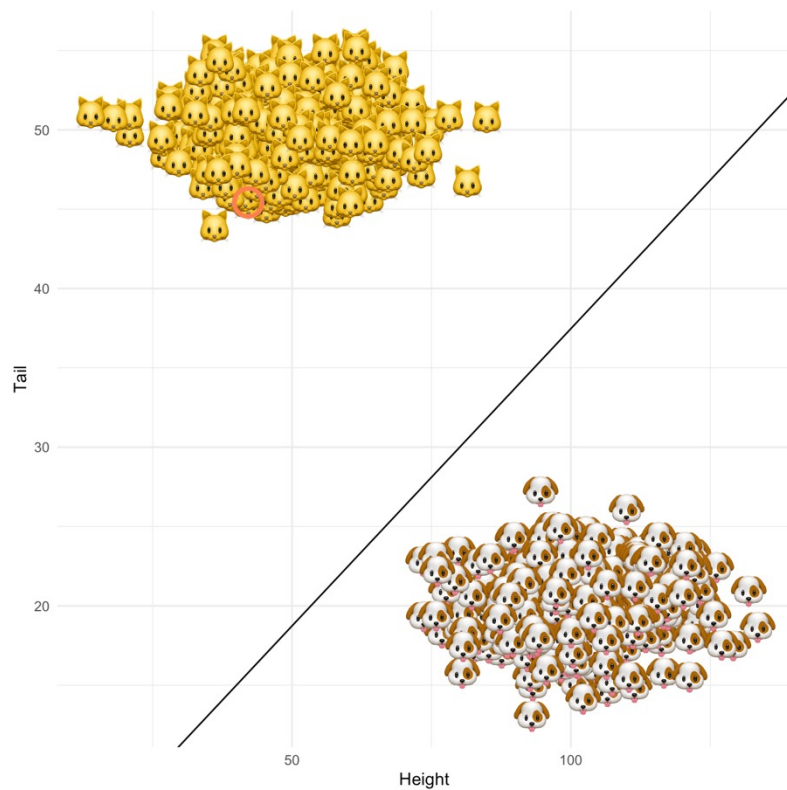
Implicaties van de representatie



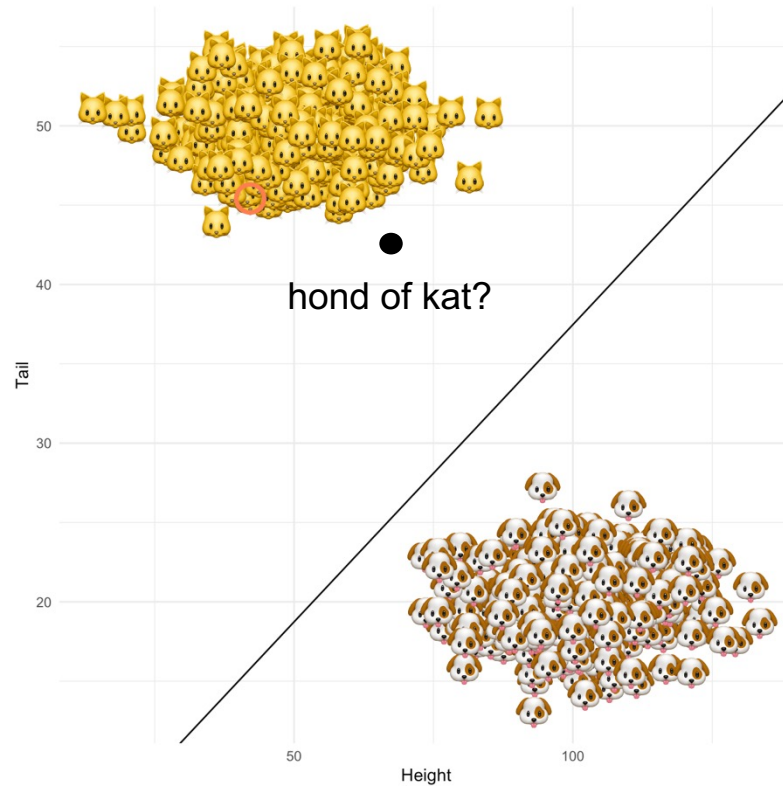
Cynthia C. S. Liem and Annibale Panichella, "Oracle Issues in Machine Learning and Where to Find Them," Proc. RAISE 2020.

Wat doen we met patronen die we vinden?

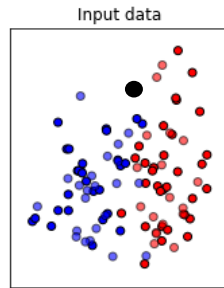
Terug naar de honden en katten



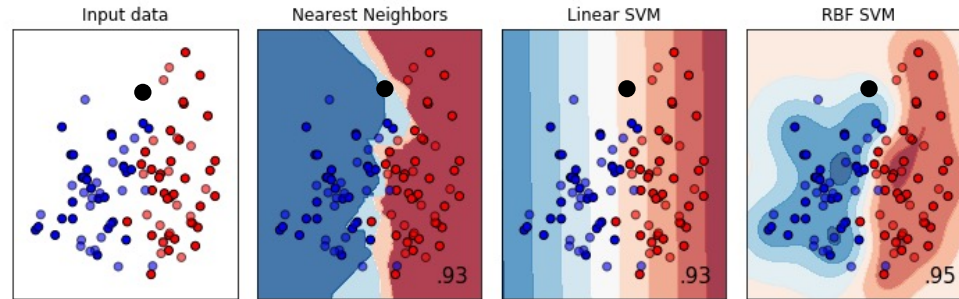
Ongeziene datapunten voorspellen



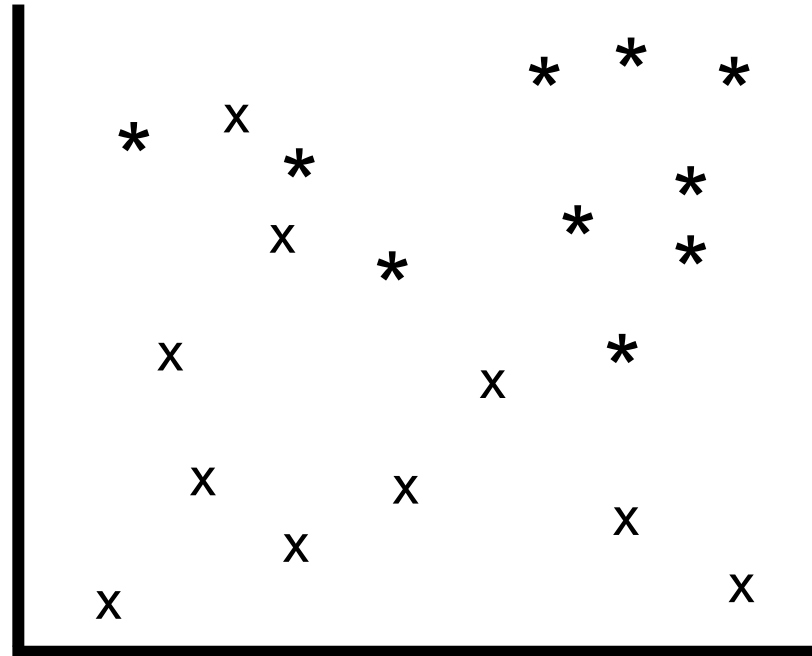
Patronen leren uit data en labels



Patronen leren uit data en labels

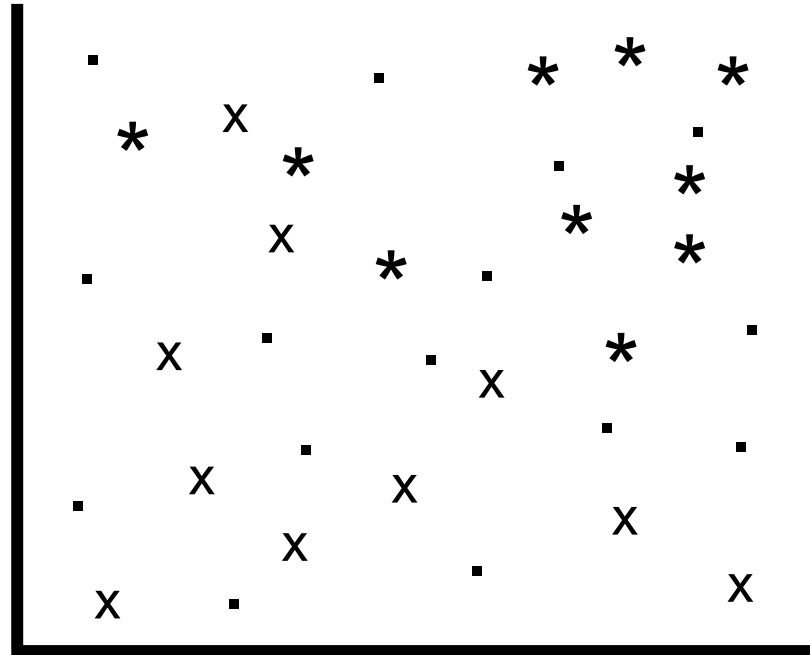


Beslissingen nemen en prioriteren



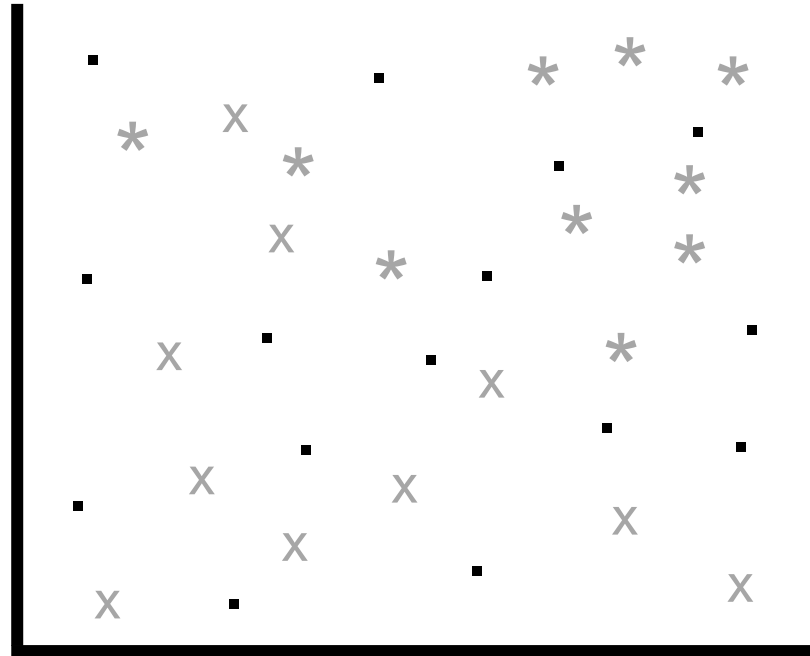
Beslissingen nemen en prioriteren

- Kies 5 datapunten



Beslissingen nemen en prioriteren

- Wat als die datapunten over mensen gaan?



Beslissingen nemen en prioriteren

- Wat als die datapunten over mensen gaan?



https://tudelft.fra1.qualtrics.com/jfe/form/SV_6sqrujNkBAT8N2m

Drie keer hetzelfde?

Drie keer hetzelfde?

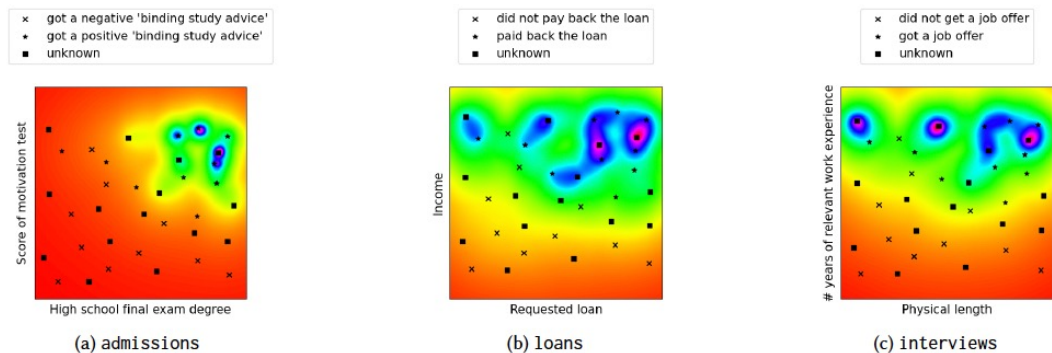


Fig. 6. Aggregated click heatmaps for computer science/data science audiences (python, dataScience and students combined).

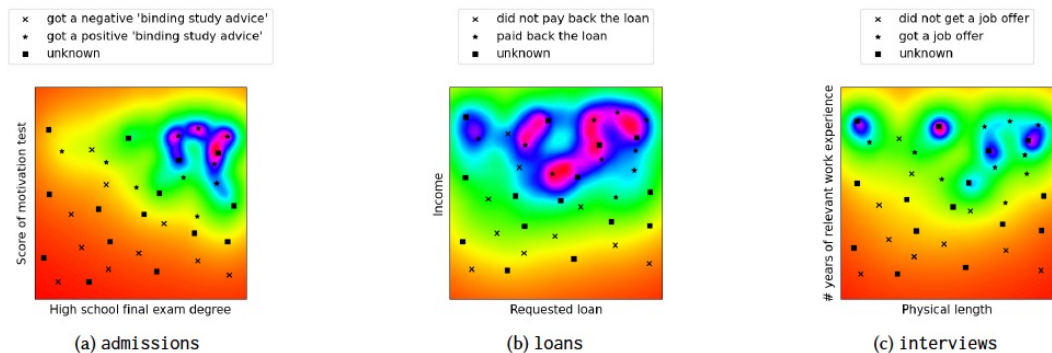


Fig. 7. Aggregated click heatmaps for broader audiences (rotary, policy and librarians combined).

Het goede antwoord?

Perverse prikkels?

“

*"When a measure becomes
a target, it ceases to be a
good measure."*

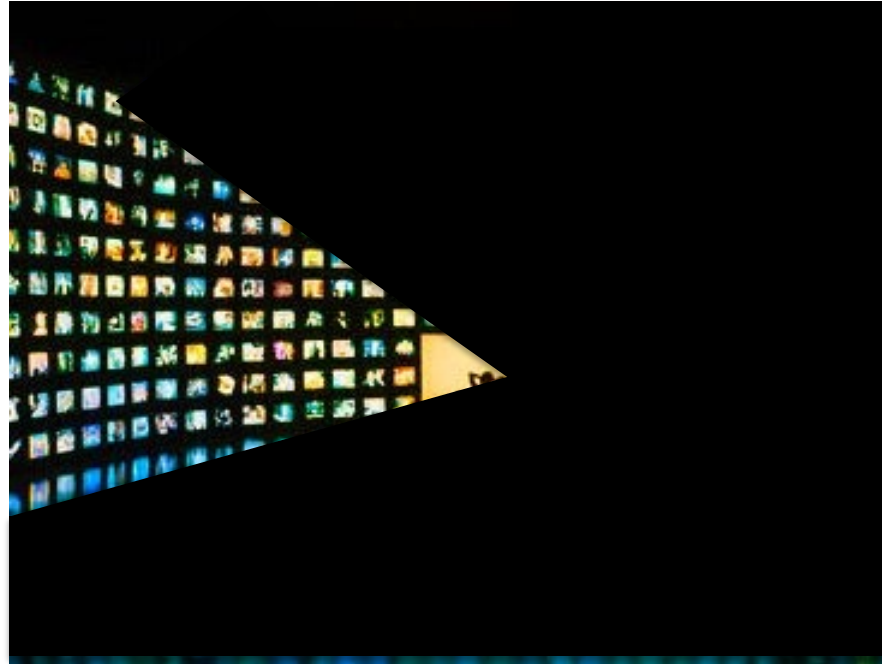
Goodhart's law as
rephrased by Marilyn
Strathern

@dasaptaerwin
CC-0

Een lang bestaand ideaal



Een lang bestaand ideaal



Een lang bestaand ideaal

relevant



Een lang bestaand ideaal



Een lang bestaand ideaal



irrelevant

Een lang bestaand ideaal

voorspellingsfout



irrelevant

Een lang bestaand ideaal

voorspellingsfout



irrelevant

oninteressant

Een lang bestaand ideaal

voorspellingsfout

niet de tijd waard



irrelevant

oninteressant

Een lang bestaand ideaal

voorspellingsfout

niet de tijd waard



irrelevant

oninteressant

niet-bestaand

Een lang bestaand ideaal

voorspellingsfout

niet de tijd waard



irrelevant

oninteressant

niet-bestaand

onbereikbaar

‘Eerlijkheid’

- ‘Eerlijkheid’ / fairness in AI is momenteel een hot topic
- Er zijn echter heel veel mogelijke definities van ‘eerlijkheid’, en wiskundig zijn ze met elkaar in conflict
- Ook hier moet je dus scherp zijn op wat je eigenlijk bedoelt!



Tutorial: 21 fairness definitions and their politics

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

Hoe omgaan met fouten en risico's?

- Tel je hoe vaak je het goed hebt?
- Of geef je aan hoe je zo robuust mogelijk met fouten omgaat?
- Wat weeg je af bij beperkte middelen?

Extreem voorbeeld

- Als we alle criminelen van Nederland willen vangen, kunnen we de hele bevolking arresteren. Dan zijn we 100% succesvol.

Belangen en beloning

- Waar worden verschillende stakeholders voor beloond, en wat zien ze als schade?
- Afreken- en sensatiecultuur in politiek en media

Belangen en beloning

- Genuanceerde afwegingen nodig
 - Goed voorbeeld: afwegingen en moreel beraad in TU Delft rond Kennisveiligheid
<https://www.youtube.com/watch?v=KJTF0blzvKA>
 - (morele) beginselen boven (praktische) gevolgen
 - aandacht voor schadebeperking en moreel leerproces










Een metafoor die mij vaak hielp

Algoritmes zijn als recepten

- Objectief repliceerbaar (als we willen)



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

 <p>Oma's appeltaart Laura's Bakery 4,7 ★★★★★ (96)</p>	 <p>Klassieke appeltaart recept Heel Holland Bakt 3,5 ★★★★★ (2,4K)</p>	 <p>Appeltaart Koopmans 3,5 ★★★★★ (952) 1 uur</p>
 <p>Appeltaart recept Milijuschka 4,1 ★★★★★ (441) 1 uur 15 min</p>	 <p>Appeltaart recept Bakkenderwijs 4,9 ★★★★★ (14)</p>	 <p>Klassieke appeltaart 24Kitchen 2,7 ★★★★★ (9,6K) 2 uur 15 min</p>
 <p>De lekkerste appeltaart – recept Rutger Bakt 4,0 ★★★★★ (1,4K)</p>	 <p>Traditionele Hollandse appeltaart Albert Heijn 3,2 ★★★★★ (1K)</p>	 <p>Appeltaart Leuke Recepten 4,0 ★★★★★ (83) 2 uur 50 min</p>

Algoritmes zijn als recepten

- Eindresultaat is niet 'universeel optimaal' of 'universeel geldend'



Algoritmes zijn als recepten

- Machine learning als een Airfryer?
Generatieve AI als een nieuwe technologie voor een kunstmatige smaakstof?
- Producent moet regels volgen (met toezicht)
- Burger kan niet meepraten over Airfryer-details, en zich toch positioneren
- En ook eigen verantwoordelijkheid nemen qua een 'gezond' dieet

Onze verantwoordelijkheid

- Gan we onze menselijke imperfectie te gretig uit de weg?



<https://www.trouw.nl/tijdgeest/chatgpt-berooft-ons-van-waardevol-denkwerk-wanneer-zijn-we-gestopt-met-ergens-moeite-voor-doen~bda521b4/>

Meepraten over betrouwbare en verantwoorde AI

Wat moet ik weten?

Cynthia C. S. Liem

Multimedia Computing Group, TU Delft

[!\[\]\(4729e517bc6a7cd81c8025b9646574fb_img.jpg\) @cynthiacsliem@akademienl.social.nl](mailto:@cynthiacsliem@akademienl.social.nl)

c.c.s.liem@tudelft.nl

Verdere links

- Nog een aantal presentaties/commentaren over risicomodellering en het toeslagenschandaal:
 - https://tudelft.zoom.us/rec/share/gipC7rMluoijUv82eNEpGKYreTHQyTLM0q3Kv7AX3sIBr-jhuo2mvyeAf-Vk7M6k.HGp0V_X-4-puKPFa?startTime=1648047506000
 - <https://www.vpro.nl/argos/media/luister/argos-radio/onderwerpen/2021/In-het-vizier-van-het-algoritme-.html>
 - <https://www.vpro.nl/argos/lees/onderwerpen/algoritme/2023/fraude-algoritme-op-de-snijtafel-rotterdam-verdenkt-vooral-jonge-moeders.html>
 - https://www.youtube.com/watch?v=t_Mic6Q5A_M