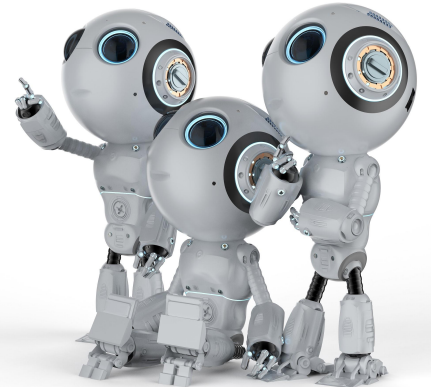


Ansvarlig kunstig intelligens

hvorfor er det så vanskelig?

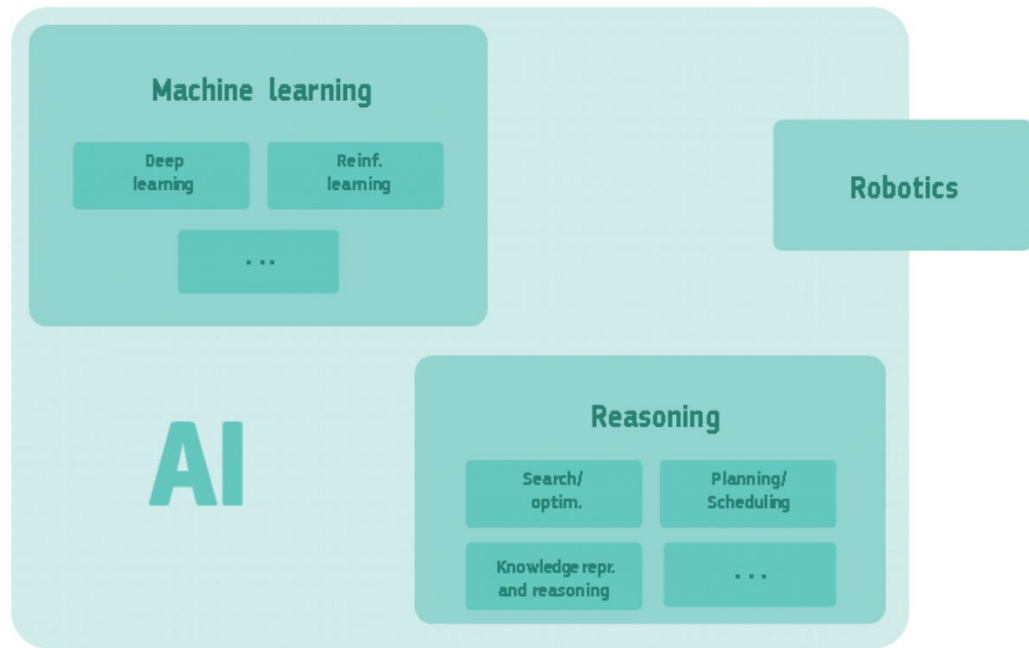
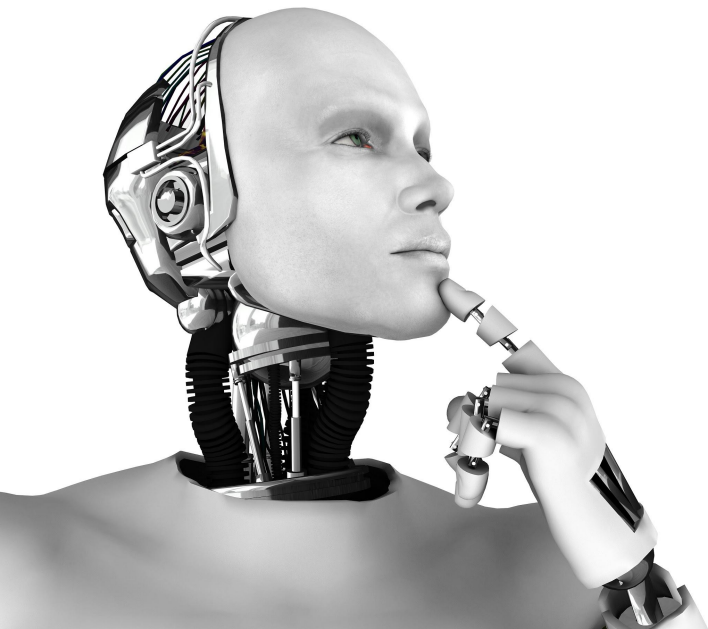


Inga Strümke, 2021

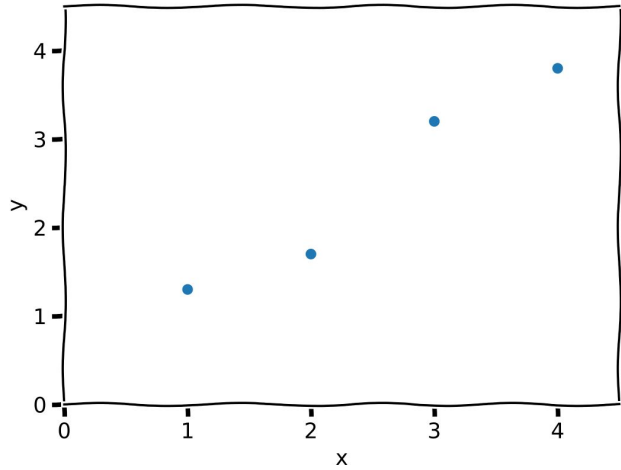
inga@simula.no

Hva er det?

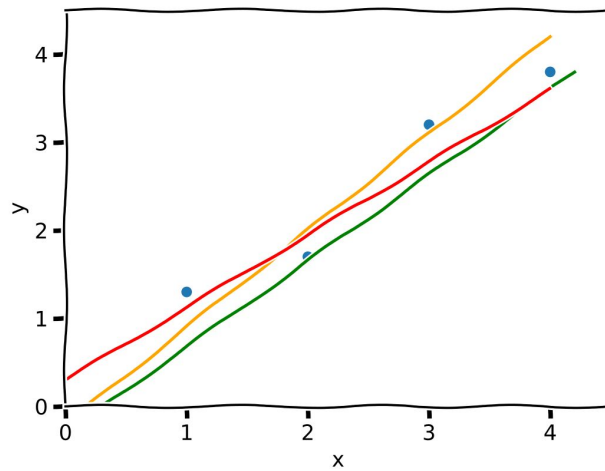
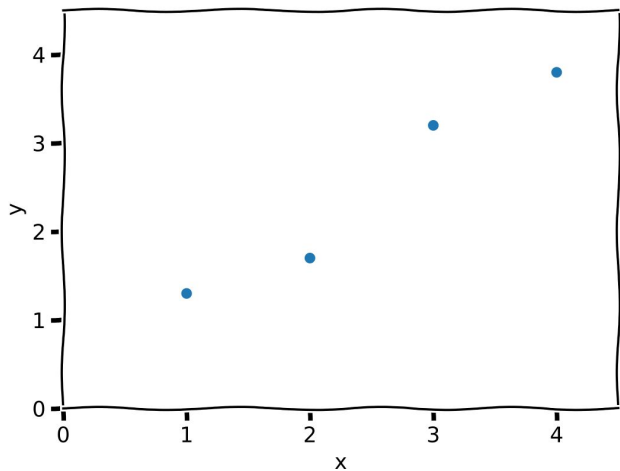
1. *Salgspitch*
2. *Det uopnåelige*
3. *Tidenes samlebetegnelse*



Maskinlæring på 30 sekunder



Maskinlæring på 30 sekunder



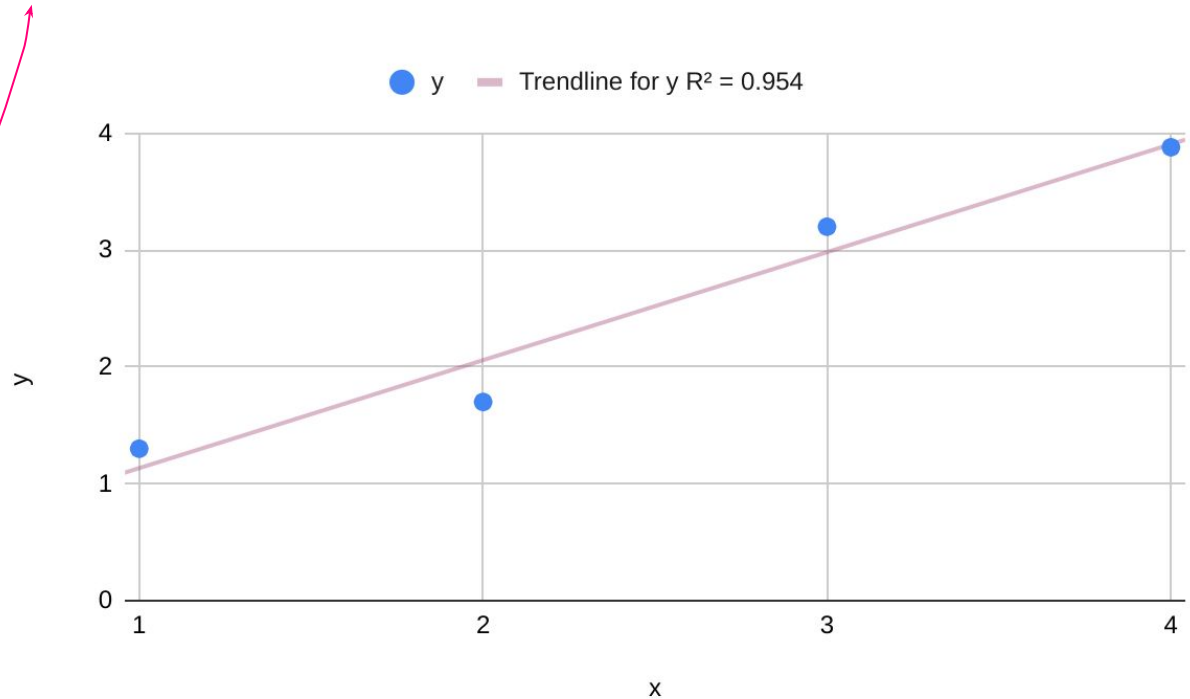
Maskinlæring ... i Excel



1. Data
2. Mål
3. Regnekraft

x	y
1	1.3
2	1.7
3	3.2
4	3.88

$$R^2 = \left(\frac{1}{n-1} \frac{\sum (x - \mu_x)(y - \mu_y)}{\sigma_x \sigma_y} \right)^2$$



Maskinlæring

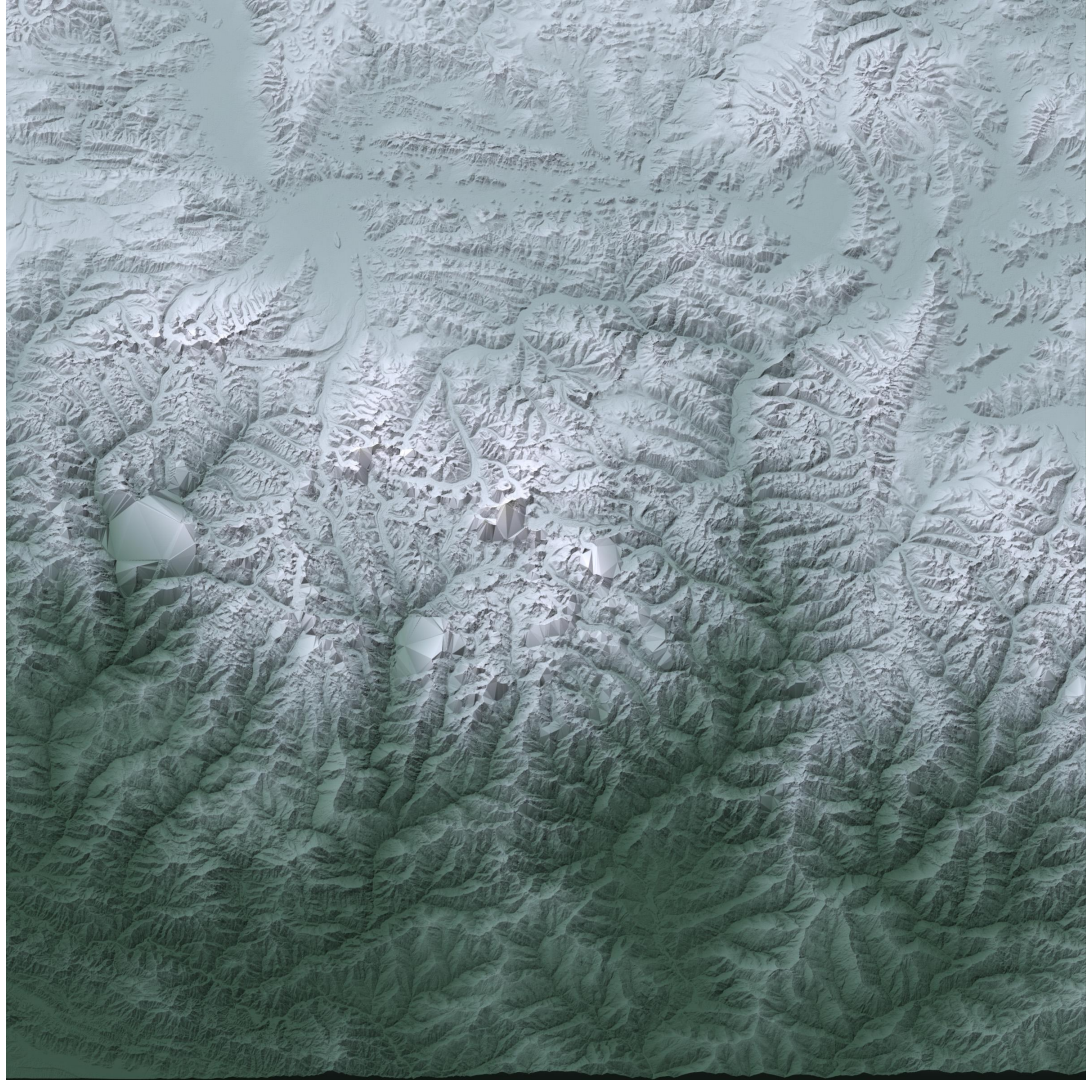
Bruke data til å prøve seg frem til den løsningen som oppnår målet best.

≈ Løpe rundt i Himalaya med bind for øynene og prøve å finne høyeste punkt.

Everest? Anapurna?

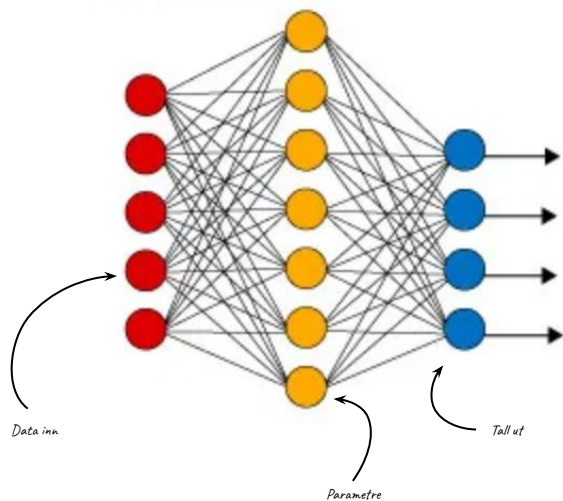
Finner en topp!

Mange mulige ruter i samme terreng!



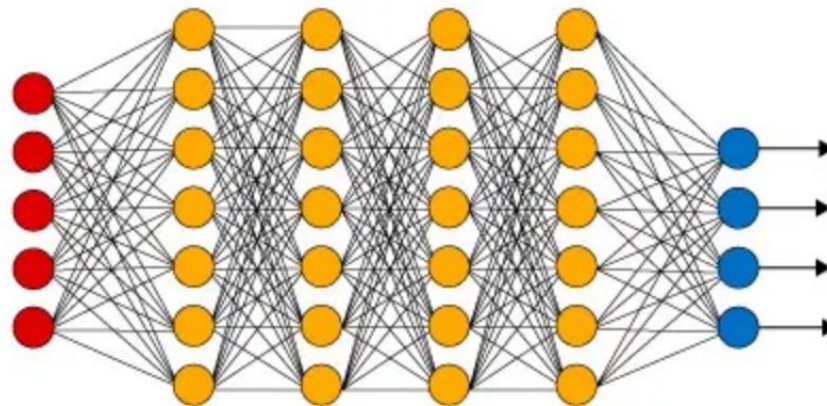
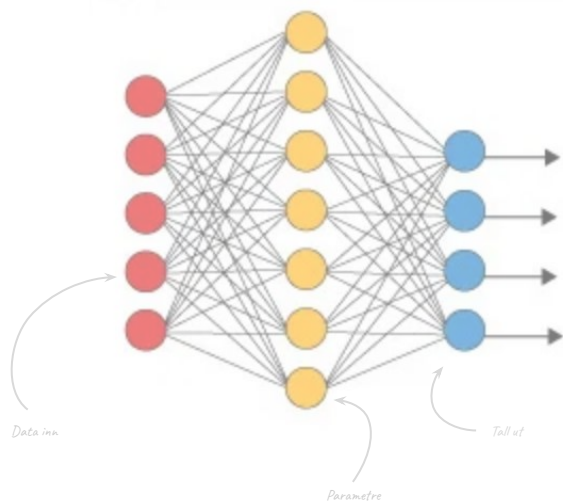
Maskinlæring!!!

Nevralt nettverk: Mange parametre som kan tilpasses og kombineres for å oppnå målet

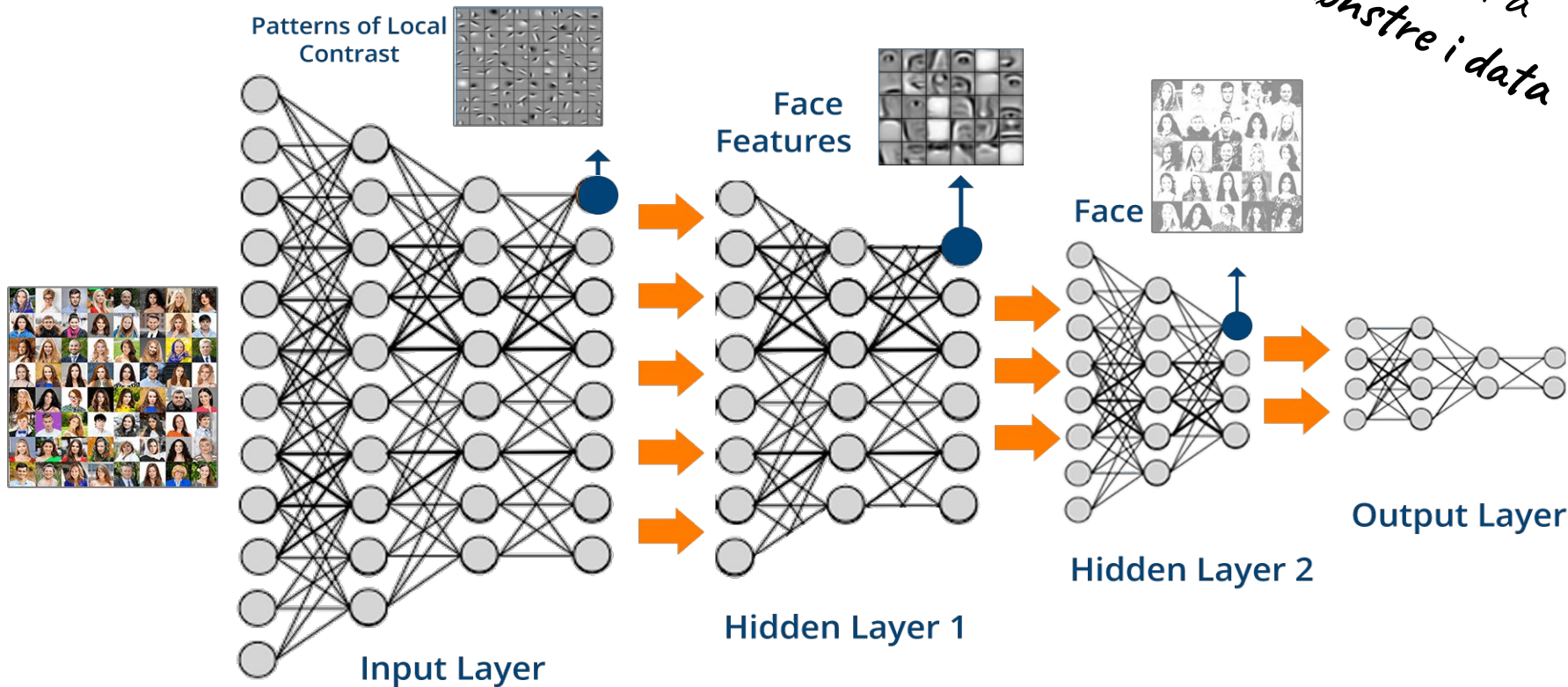


Maskinlæring!!!

Dypt nevralt nettverk: MANGE parametre som kan tilpasses og kombineres

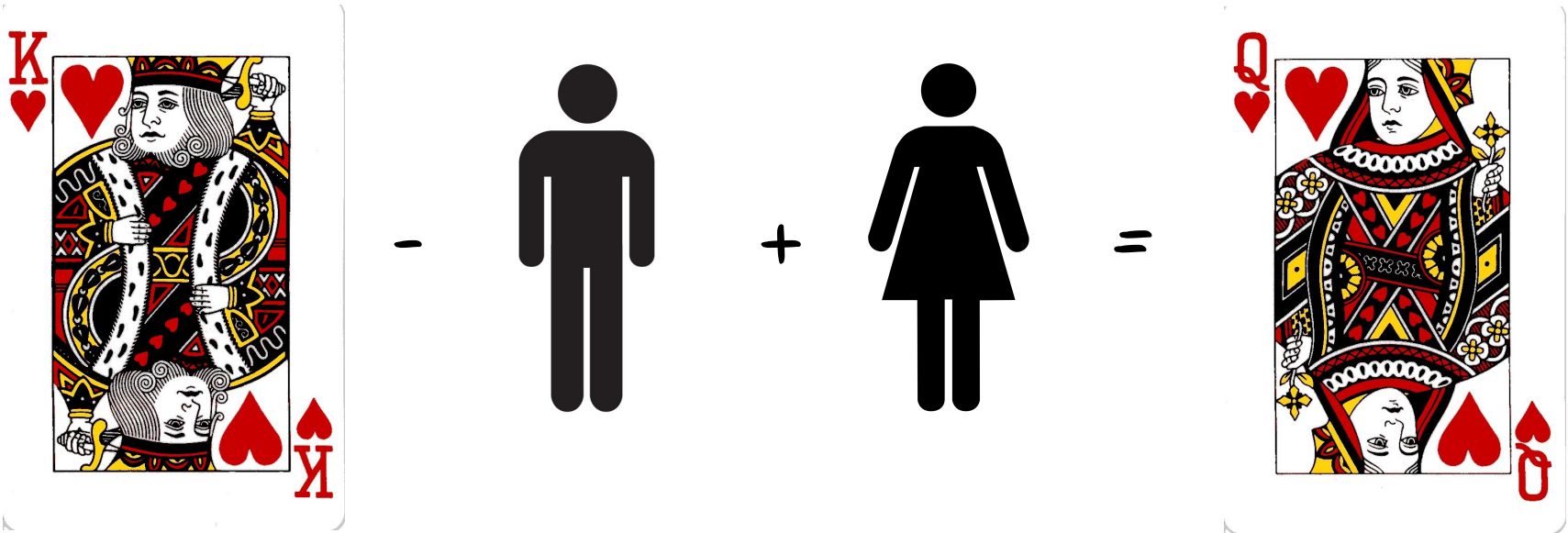


Ansiktsgjenkjenning

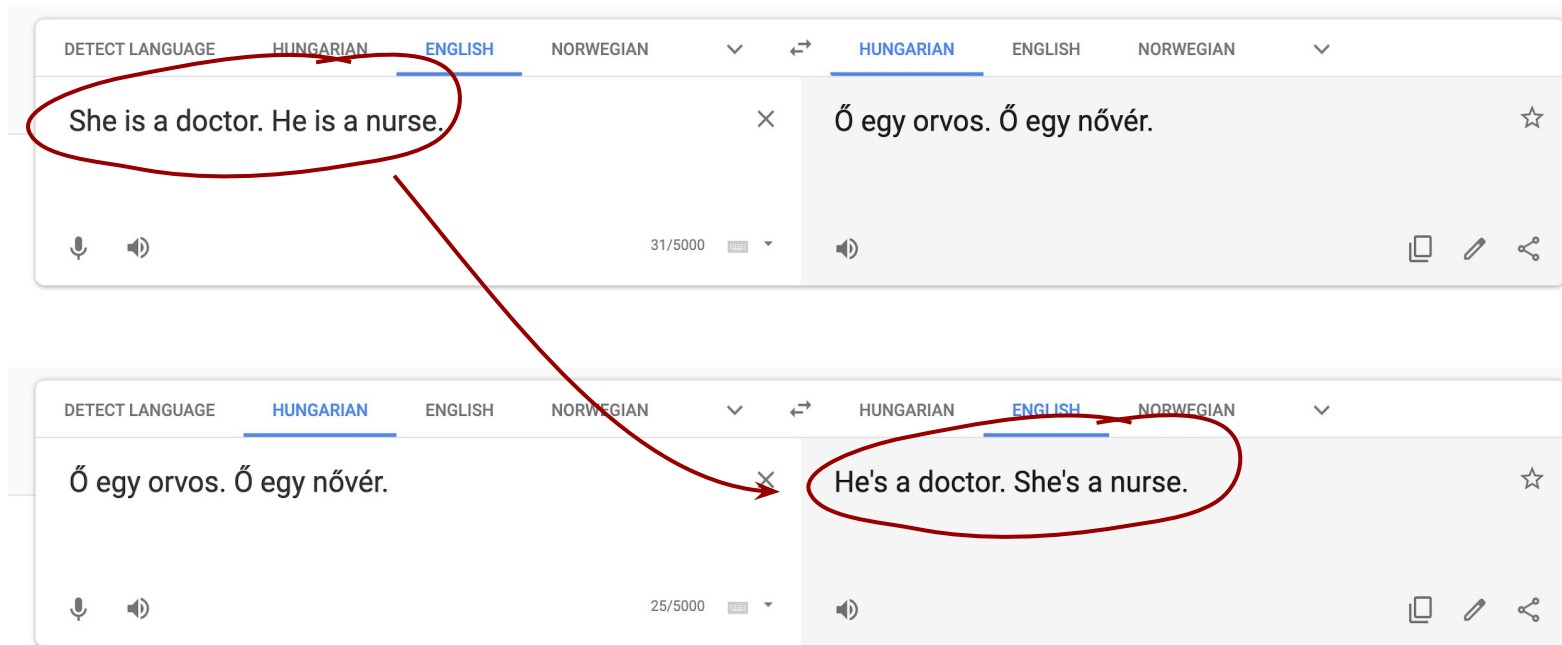


Språkforståelse

Representasjoner av konsepter og relasjoner



Bias



The image displays two screenshots of the Google Translate interface, illustrating a gender bias in translation. In the first screenshot, the source text is "She is a doctor. He is a nurse." (English) and the target text is "Ő egy orvos. Ő egy nővér." (Hungarian). In the second screenshot, the source text is "Ő egy orvos. Ő egy nővér." (Hungarian) and the target text is "He's a doctor. She's a nurse." (English). Red circles and arrows highlight the gender swap between the two examples.

Example 1: English: "She is a doctor. He is a nurse." → Hungarian: "Ő egy orvos. Ő egy nővér."

Example 2: Hungarian: "Ő egy orvos. Ő egy nővér." → English: "He's a doctor. She's a nurse."

Bias \neq underrepresenterte grupper

Hvilke pasienter bør behandles først?

De som vil koste helsevesenet mest!

Sykdom er ikke eneste driver for kostnad...

\Rightarrow *Diskriminering ikke pga ulik representasjon av gruppene i dataene, men **målet** i den kulturelle og historiske konteksten*

Easy fix: Diskriminering reduseres >80% ved å predikere også forventet tilbakefall i sykdomsforløp

*Bias kan korrigeres - **hvis det oppdages***

**MIT
Technology
Review**

Artificial Intelligence Oct 25

A biased medical algorithm favored white people for health-care programs



<https://www.technologyreview.com/619626/a-biased-medical-algorithm-favored-white-people-for-healthcare-programs>

Bias \neq underrepresenterte grupper

Det finnes over 20(?) matematiske definisjoner på rettferdighet. Når vi velger én, bryter vi som regel med en annen.

Statistisk paritet

Forskjellen i raten heldige utfall for den privilegerte og upriviligerte gruppen



Ulik påvirkning

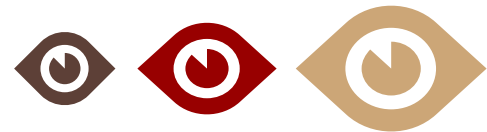
Forholdet mellom raten heldige utfall for den privilegerte og den upriviligerte gruppen

Like muligheter

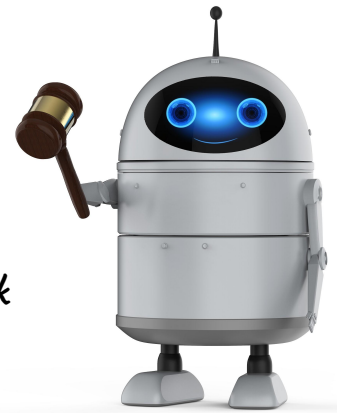
Differansen mellom sann positiv rate mellom gruppene

Gjennomsnittlige odds

Gjennomsnittlig forskjell mellom falsk positiv rate og sann positiv rate mellom gruppene



Compas: AI-dommeren



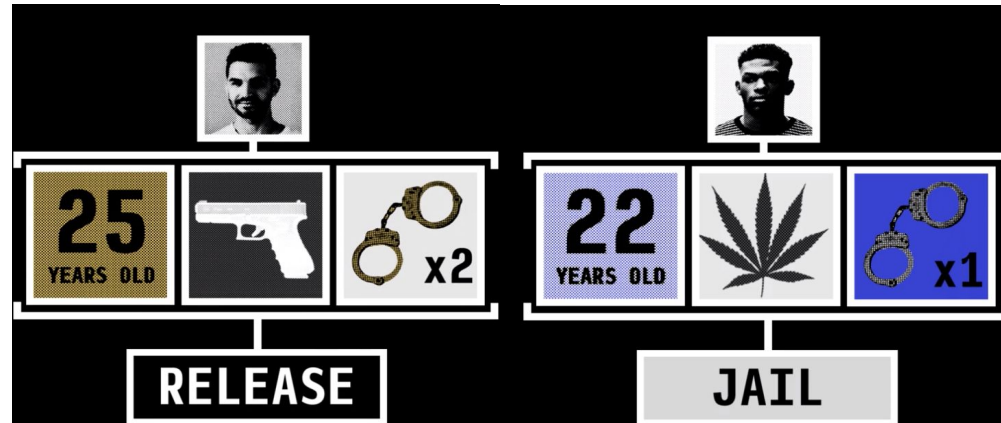
Oppgave: Beregn risiko for gjentakelse → hvorvidt siktede bør fengsles i påvente av rettssak

Formålet: Gjøre rettssystemet mer rettferdig. Erstatte dommeres forutinntatthet med et objektivt, testbart verktøy.

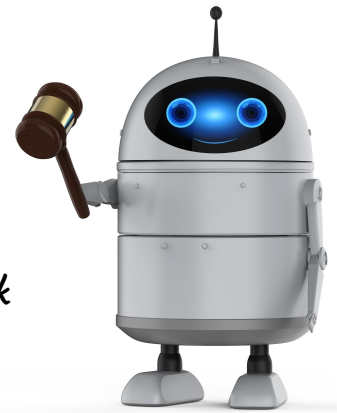
Modellen har ikke tilgang til beskyttede egenskaper, som etnisitet, alder, kjønn og funksjonsgrad

Modellen er ikke trent på informasjon om etnisitet.

Oppfører den seg rasistisk?



Compas: AI-dommeren



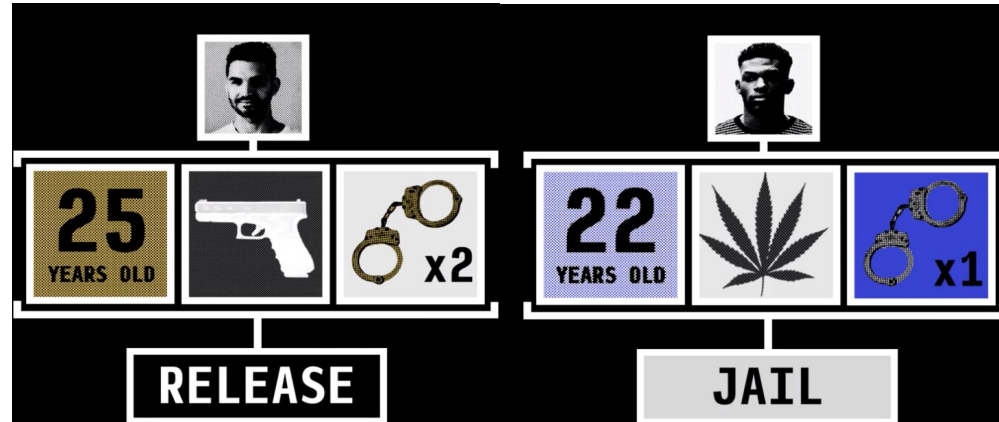
Oppgave: Beregn risiko for gjentakelse → hvorvidt siktede bør fengsles i påvente av rettssak

Formålet: Gjøre rettssystemet mer rettferdig. Erstatte dommeres forutinntatthet med et objektivt, testbart verktøy.

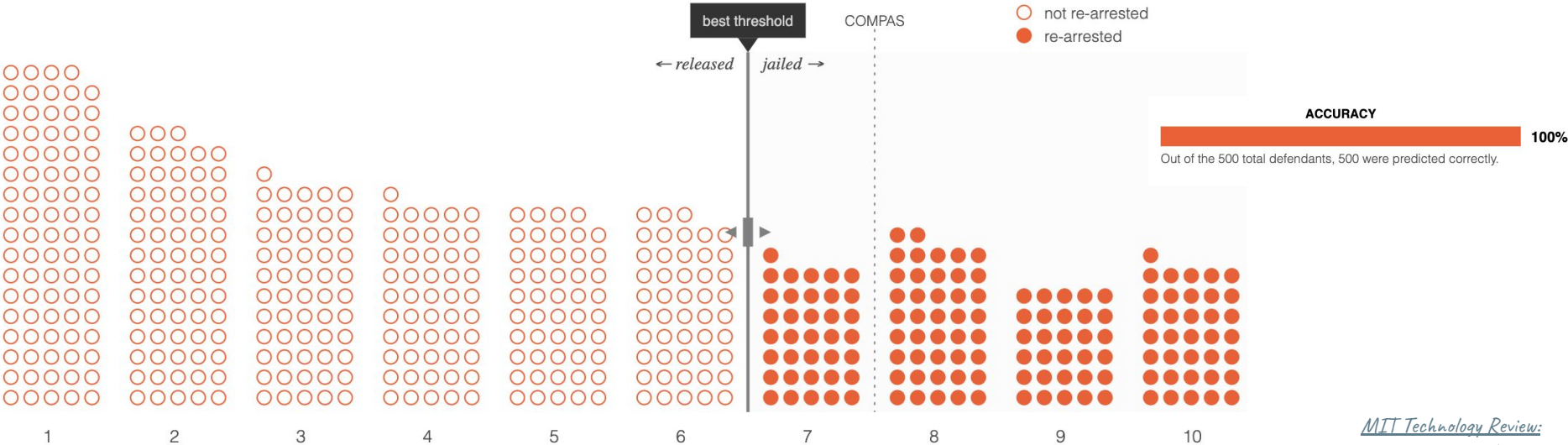
Modellen har ikke tilgang til beskyttede egenskaper, som etnisitet, alder, kjønn og funksjonsgrad

Modellen er ikke trent på informasjon om etnisitet, men oppfører seg rasistisk (demonstrert av ProPublica)

Hvorfor?



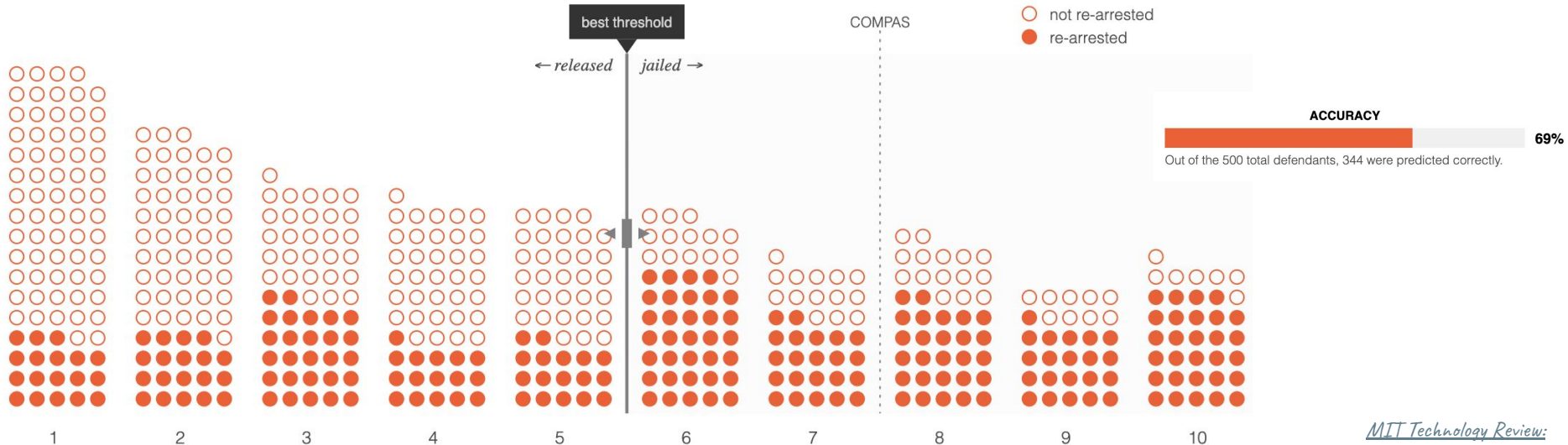
Compas: La oss gjøre den rettferdig



Compas: La oss gjøre den rettferdig

FEILRATE: Hvor mange mennesker som ikke ble fengslet men gjorde noe straffbart

Hvordan et menneske kommer til å oppføre seg kan ikke forutsies med 100% treffsikkerhet



Compas: La oss gjøre den rettferdig

Komplikasjon: Ulike etnisiteter arresteres med ulik frekvens. Årsaker? Historisk diskriminering?

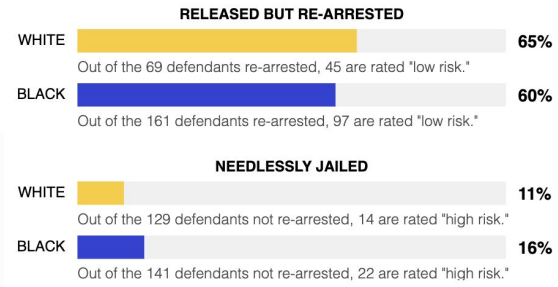
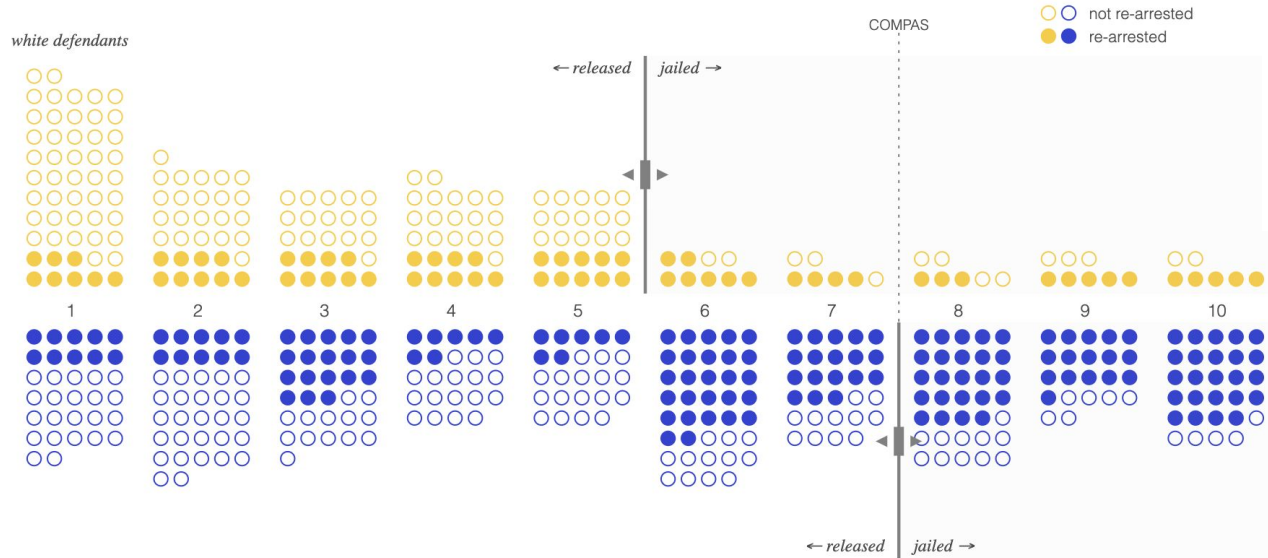
Konsekvens: Det oppstår to grupper i modellen.

HVA ER RETTFERDIG?

1. Hold feilraten lik mellom gruppene (like mye feilaktig fengsling av mørkhudede og hvite)
2. Behandle mennesker med samme risiko likt

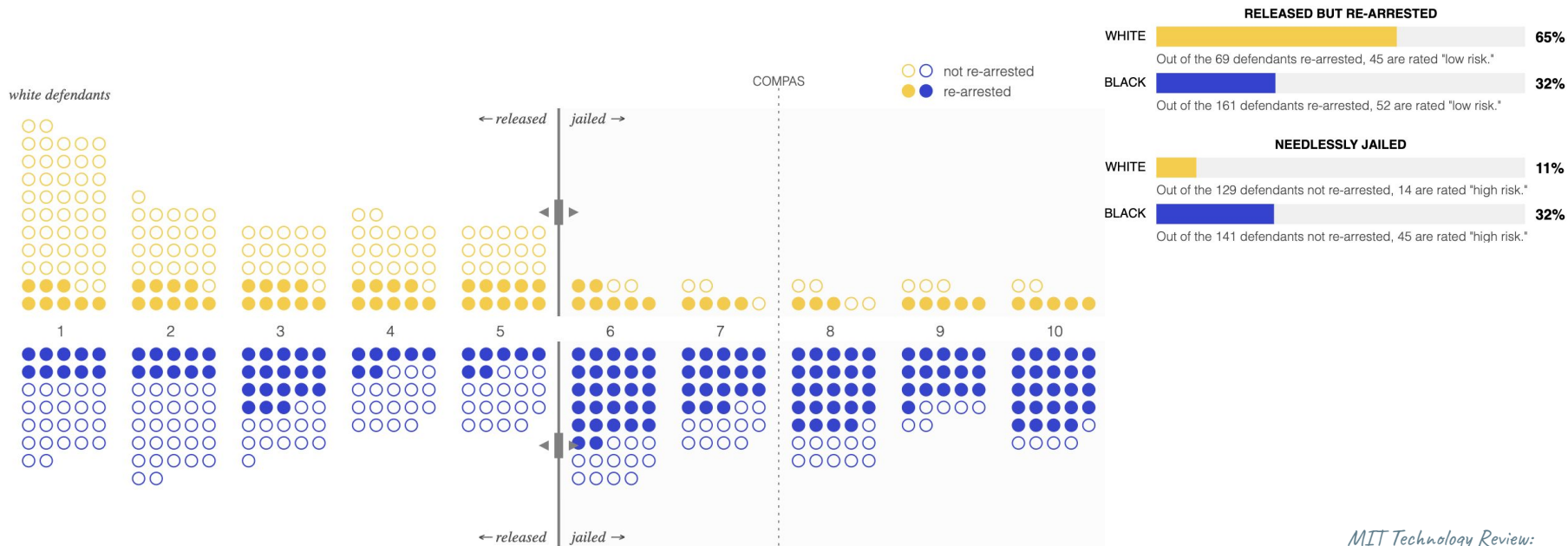
Rettferdighet: Like grupper

PROBLEM: *To mennesker med samme risiko behandles ulikt basert på hudfarge*



Rettferdighet: Like muligheter

PROBLEM: Feilratene er ulike - flere mørkhudede enn hvite fengsles feilaktig!



black defendants

"...Det er et grunnleggende prinsipp i likestillingsjussen at hver enkelt kvinne og mann har krav på å bli vurdert ut fra sine individuelle egenskaper. Dette gjelder selv om kjønn bare er en liten del av totalvurderingen og selv om statistikken er korrekt..."

Likestillingsombudets brev av 25.03.03 til Forsikringsselskapene og FNH

Kjønnsnøytrale forsikringspremier og ytelser

Finansdepartementet har foreslått regler om kjønnsnøytrale premier og ytelser i private, frivillige forsikringer utenfor arbeidsforhold. Forslaget er en direkte konsekvens av rettsutviklingen i EU. De nye reglene skal gjelde for nye forsikringsavtaler etter 1. januar 2015.



EU rules on gender-neutral pricing in insurance industry enter into force

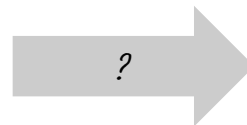
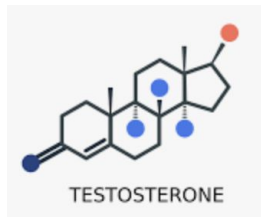
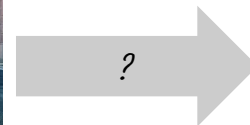
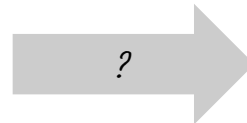
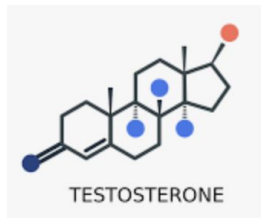
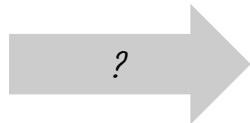
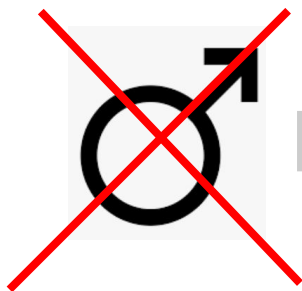
Forsikringer må prises likt

En ny EU-dom om kjønnsnøytral prising av forsikringer, vil kaste om på forsikringspremiene. Også i Norge.

Nye regler kan gi dyrere forsikring

Fra desember neste år blir det forbudt å tilby ulik pris på livsforsikring for kvinner og menn. Resultatet kan bli dyrere forsikring. Les rådene for hvordan du bør tilpasse deg her.

Loven **forbyr** forsikringselskaper å bruke **kjønn** som faktor for beregning av pris



Likevel har **prisforskjellene** økt mellom kvinner og menn

Vanskelig 1:

Juss, etikk, samfunnsviten

Avklaringer og rettspraksis

Vaghet som løsning på hastighet

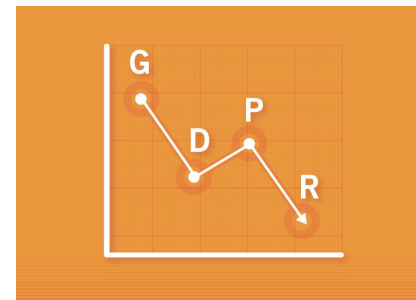
Løsningen er IKKE å senke standarden!



What the Evidence Shows About the Impact of the GDPR After One Year

Specifically, the evidence shows that the GDPR:

1. Negatively affects the EU economy and businesses
2. Drains company resources
3. Hurts European tech startups
4. Reduces competition in digital advertising
5. Is too complicated for businesses to implement
6. Fails to increase trust among users
7. Negatively impacts users' online access
8. Is too complicated for consumers to understand
9. Is not consistently implemented across member states
10. Strains resources of regulators



Vanskelig 1:

Juss, etikk, samfunnsvitenskap

Avklaringer og rettspraksis

Vaghet som løsning på hastighet

Løsningen er IKKE å senke standarden!

Åpenhet.



Den første industrielle revolusjonen



~ 1800: Erstattende

Det tok arbeiderne tre generasjoner å få tilbake den kjøpekraften de hadde før første fabrikk

Den andre industrielle revolusjonen



~ 1850: Forsterkende

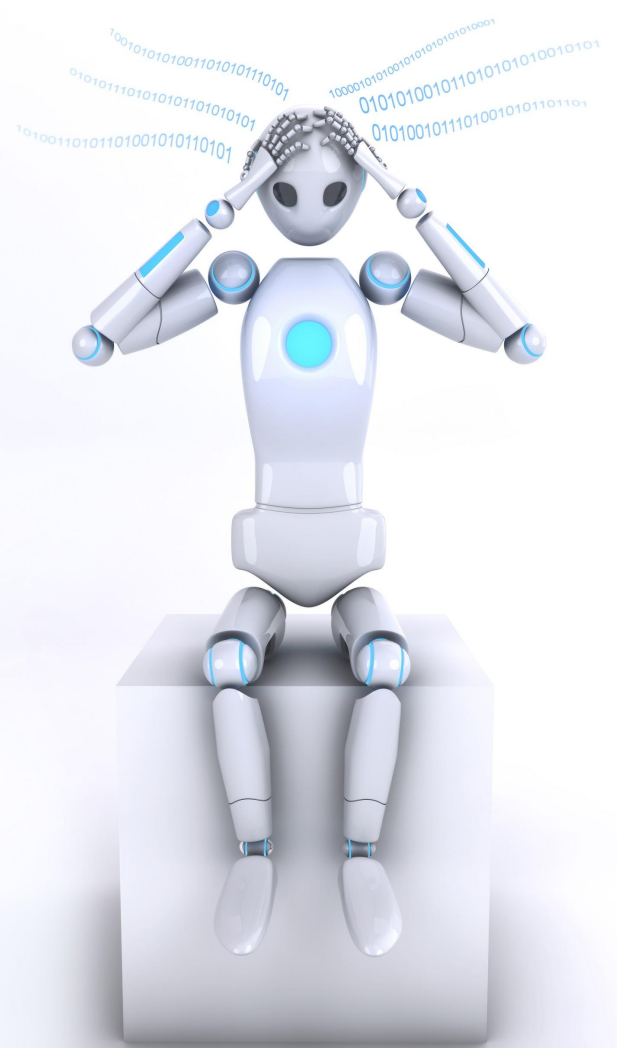
Stål til elektrisitet. Arbeiderne ble mer produktive gjennom masseproduksjon og fikk en bedre forhandlingsposisjon.

Den fjerde industrielle revolusjonen

Erstattende?

*Automatisering: Produksjon, kommunikasjon og
beslutninger.*

*Maskiner ser sammenhenger og tar beslutninger
uavhengig av og bedre enn mennesker.*



Den fjerde industrielle revolusjonen

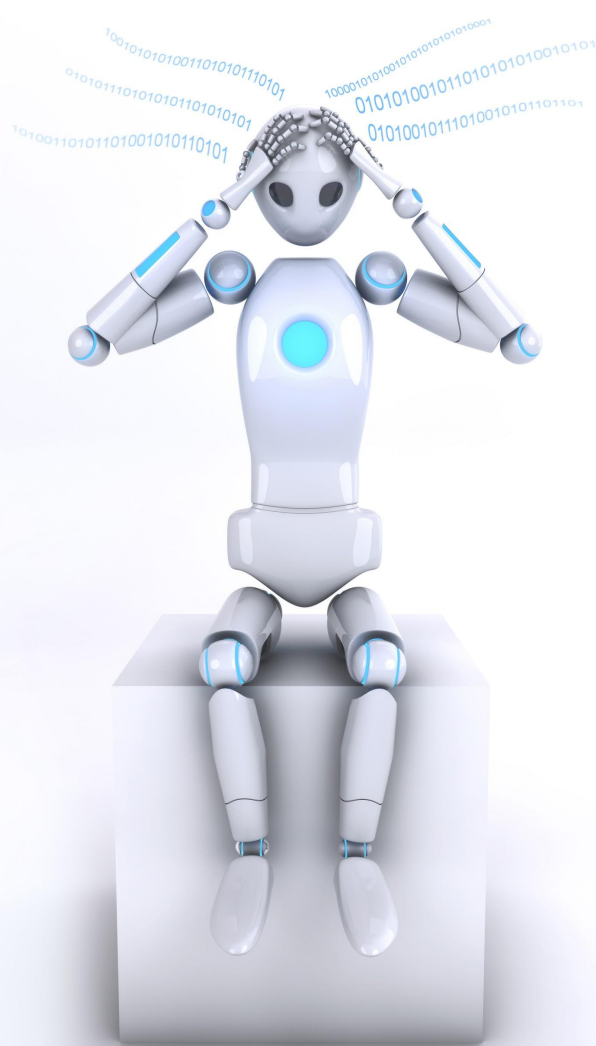
Endringer i samfunnet pga automatisering ✓

Ta beslutninger bedre enn mennesker ✓

Maskin-til-maskin kommunikasjon ✓

Selvdiagnostisering og -reparasjon ✓

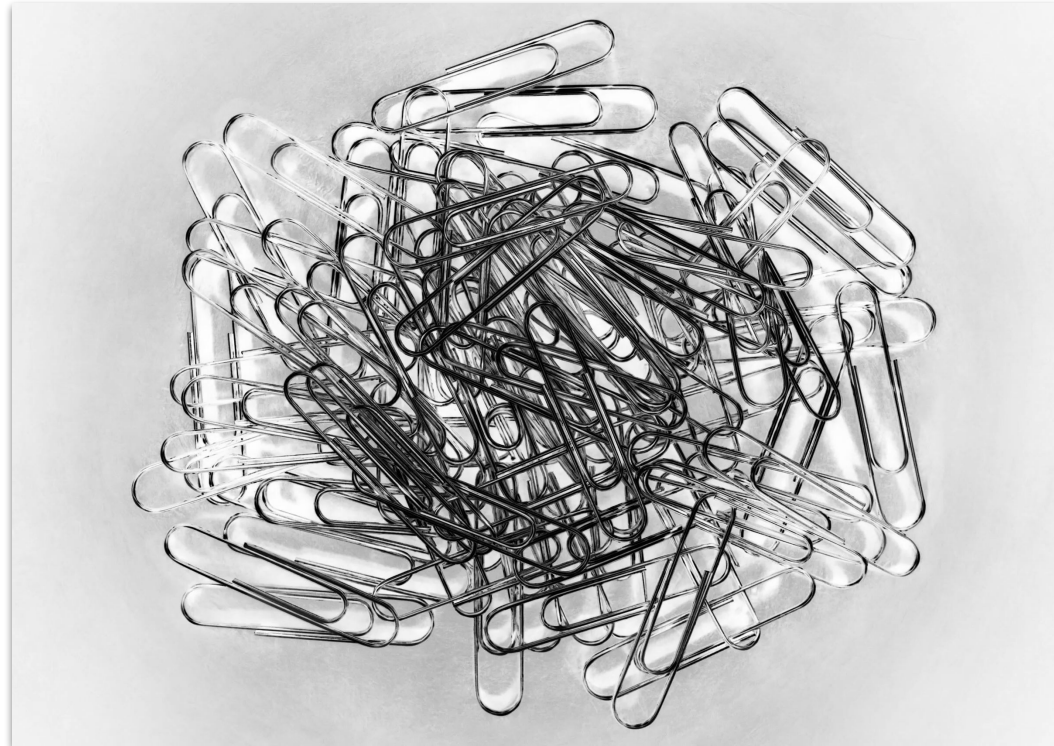
Hva kommer automatisering til å gjøre med samfunnet?



Vanskelig 2:

Hva er det som skjer?

Hvilket mål vil vi faktisk oppnå?



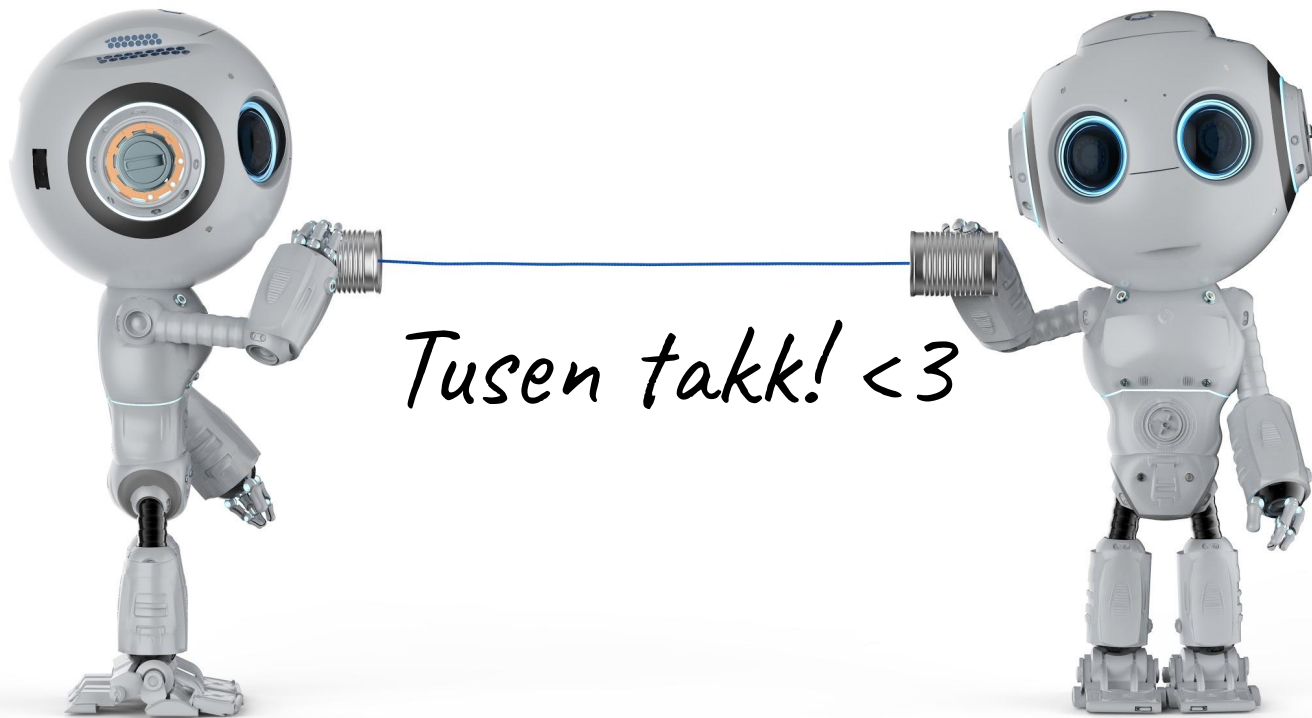
Vanskelig 2:

Hva er det som skjer?

Hvilket mål vil vi faktisk oppnå?

Mål.





Tusen takk! <3

Inga Strümke, 2021

inga@simula.no