

# Kunstmatige intelligentie

Bert De Saffel

Master in de Industriële Wetenschappen: Informatica Academiejaar 2018–2019

Gecompileerd op 20 februari 2019

# Inhoudsopgave

<b>1</b>	<b>Kunstmatige intelligentie</b>	<b>2</b>
1.1	Kunnen machines denken? . . . . .	2
1.2	Toepassingen van AI en data mining . . . . .	3
1.3	Leren . . . . .	3
1.4	Classificatie . . . . .	4
1.5	Informatie en beslissingsbomen . . . . .	5
1.5.1	Informatie-inhoud . . . . .	5
1.5.2	Beslissingsboom . . . . .	5
1.6	Klasseren zonder leren . . . . .	8
1.6.1	$k$ zwaartepunten . . . . .	8
1.7	Een toepassing: Watson . . . . .	8

# Hoofdstuk 1

## Kunstmatige intelligentie

- Twee doelen van kunstmatige intelligentie:
  - Het laten overnemen, door machines, van taken waarvoor intelligentie vereist is.
  - Studie van natuurlijke intelligentie.
- Twee vormen om kennis in te brengen in een computersysteem:
  - Expliciete kennis.
  - Kennis kan zelf verworven worden.

### 1.1 Kunnen machines denken?

- Twee voorbeelden.
  - ELIZA:
    - ◊ Computerprogramma dat zich voordoet als een psychotherapeut.
    - ◊ Maakt gebruik van simpele vervangingsregels.
    - ◊ Probeert de conversatie zo te sturen zodat de echte persoon het meest moet vertellen.
  - Chinese kamer:
    - ◊ Denkrichting die aantoont dat een entiteit eerst iets moet begrijpen, vooraleer er van intelligentie sprake is.
      1. Iemand die geen Chinees kent wordt in een kamer gebracht.
      2. Door een luik krijgt hij briefjes in het Chinees aangereikt, en de bedoeling is dat hij daar schriftelijk een zinnige antwoord op teruggeeft.
      3. De persoon krijgt handboeken waarin conversieregels staan.
    - ◊ De proefpersoon volgt mechanisch de regels vanuit het handboek, zodat hij wel intelligent gedrag vertoont, maar de berichten niet begrijpt.
- **Denken is elke vorm van complexe informatieverwerking waarvan de onderliggende mechanismen niet volledig gekend zijn.**
- **Turingtest:**
  - Proefpersoon kan contact maken met twee entiteiten: een mens en een machine, maar hij weet niet wie de mens of machine is.
  - De proefpersoon kan eender welke vragen stellen aan beide entiteiten.
  - Als de proefpersoon er niet in slaagt om na zijn vragenronde de entiteit aan te duiden die een machine is, dan is de machine geslaagd voor de Turingtest.

## 1.2 Toepassingen van AI en data mining

- **Classificatie:**
  - Stel een verzameling van  $k$  klassen.
  - Een bepaalde invoer met gelinkt worden aan één van die klassen.
  - Harde classificatie: beperkt aantal duidelijk van elkaar gescheiden klassen. Hier spreekt men ook van patroonherkenning.
  - Zachte classificatie: continue overgang van de klassen.
- Toepassingen:
  - Aanbevelingssystemen.
  - Kwaliteitscontrole.
- Probleemgestuurd: uitgaande van een probleem een oplossing zoeken.
- Datagestuurd: vanuit bestaande informatie problemen zoeken die ermee opgelost kunnen worden. Dit wordt ook data mining genoemd. Vaak moet de data eerst gereorganiseerd worden vooraleer de informatie nuttig wordt.

## 1.3 Leren

- Moderne AI houdt zich bezig met systemen met een zeer groot aantal aanpasbare parameters. Zulke systemen noemt men **massief lerende systemen**.
- **Voorbeelden** van massief lerende systemen:
  - Neurale netwerken: trachten het biologische denksysteem na te bootsen.
  - Hidden Markov Model: wordt gebruikt bij de analyse van allerlei sequenties, waarbij de toestand soms onbekend is.
- Parameters hebben niet noodzakelijk een betekenis, en is daarom ook onmogelijk om ze met de hand in te voeren. Daarom laat men een systeem leren, met behulp van **drie methoden**:
  - Algoritmisch leren: Er wordt gedemonstreerd hoe een bepaalde actie moet uitgevoerd worden. Het systeem kan hierna deze actie inoefenen door het herhalen van deze instructies. Deze vorm komt goed overeen met het programmeren van een computer.
  - Leren met supervisie: Hier wordt er geen gebruik gemaakt van een algoritme maar eerder van voorbeelden. Deze voorbeelden worden een leerverzameling genoemd en bevatten inputgegevens die het systeem moet leren herkennen, met de daarbij horende resultaten. Er wordt een verband opgelegd tussen een bepaalde input en output.
  - Leren zonder supervisie: Dit gebeurt gedeeltelijk algoritmisch aangezien er enige instructies nodig zijn om de machine op gang te krijgen. De machine zal nadien zelf experimenteren wat er gebeurt bij het aanpassen van verschillende parameters. Het leren gebeurt dus niet met voorbeelden, maar uit eigen ervaring. Hier is er dan ook geen verband tussen het resultaat en de verschillende deeltaken, maar er is wel een algemeen idee wat er aangeleerd moet worden.

## 1.4 Classificatie

- Classificatie is het mappen van een bepaalde input op een klasse.
- We spreken van een **item** dat we moeten klasseren.
- Dit item wordt gekarakteriseerd door een aantal **meetwaarden**.
- Twee soorten meetwaarden:
  - Sommige metingen kunnen een groot aantal waarden opleveren die voorgesteld kunnen worden als een getal.
  - Andere metingen hebben maar een beperkt aantal waarden, zoals de indeling van van categorieën.
    - ◊ Deze kunnen omgezet worden zodat ze antwoorden zijn op ja-nee vragen, zoals het gevolg dat ze geconverteerd kunnen worden naar 0 of 1.
    - ✓ Nu zijn alle meetwaarden getallen.
- Aangezien dat alle meetwaarden getallen zijn  $\rightarrow$  standaardvorm: de computer heeft een aantal klassen en moet een getallenrij die de meetwaarden voor een bepaal item bevat toewijzen aan één van de klassen.
- Hoe worden de getallenrijen weergegeven? Een aantal notaties:
  - De  $n$ -dimensionale Euclidische ruimte is de verzameling vectoren met  $n$  reële coördinaten.
  - Zo een vector wordt voorgesteld door een vette letter:  $\mathbf{v} = (v_1, \dots, v_n)$ .
  - Soms vanaf 0 beginnen, zodat we  $n+1$ -dimensionale vectoren hebben:  $\mathbf{v} = (v_0, v_1, \dots, v_n)$ . De waarde  $v_0$  krijgt een speciale betekenis.
  - Reële getallen die niet deel uitmaken van een vector worden weergegeven met hoofdletters: A, B, ...
  - De nulvector:  $\mathbf{0} = (0, \dots, 0)$ .
  - Het inproduct:

$$\mathbf{v} \cdot \mathbf{u} = \sum_i^n v_i u_i$$

Het inproduct is lineair:  $(A\mathbf{v}) \cdot \mathbf{u} = A(\mathbf{v} \cdot \mathbf{u})$  en  $(\mathbf{u} + \mathbf{v}) \cdot \mathbf{x} = \mathbf{u} \cdot \mathbf{x} + \mathbf{v} \cdot \mathbf{x}$ .

Hieruit volgt de norm  $\|\cdot\|$ :

$$\|\mathbf{u}\|^2 = \mathbf{u} \cdot \mathbf{u}$$

- $d(\mathbf{u}, \mathbf{v}) = \|\mathbf{u} - \mathbf{v}\|$  is de lengte van de kortste weg van  $\mathbf{u}$  naar  $\mathbf{v}$ . Hieruit volgt:

$$\|\mathbf{u} + \mathbf{v}\| \leq \|\mathbf{u}\| + \|\mathbf{v}\|$$

Het kwadraat van beide kanten geeft:

$$\mathbf{u} \cdot \mathbf{v} \leq \|\mathbf{u}\| \|\mathbf{v}\|$$

Aangezien dat

$$\cos(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

kan dit omgevormd worden tot de volgende ongelijkheid:

$$-1 \leq \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \leq 1$$

- De afstand en de cosinus geven vaak een goede indruk in hoeverre twee vectoren op elkaar lijken. De cosinus geeft een goede maat voor de afstand tussen twee genormaliseerde vectoren:

$$d\left(\frac{\mathbf{u}}{\|\mathbf{u}\|}, \frac{\mathbf{v}}{\|\mathbf{v}\|}\right)^2 = 2 - 2\cos(\mathbf{u}, \mathbf{v})$$

## 1.5 Informatie en beslissingsbomen

### 1.5.1 Informatie-inhoud

- Een bericht is enkel nuttig indien ontvanger een betekenis kan geven aan het bericht. De belangrijke elementen voor de informatie-inhoud is dus het bericht zelf en de kennis van de ontvanger.
- Met de kennis kan aan elk mogelijk bericht  $B$  een waarschijnlijkheid  $P(B)$  toekennen. De informatie-inhoud wordt dan gedefinieerd door

$$-\log_2(P(B)) \text{ bits}$$

Voor  $P(B) = 1$  is de informatie-inhoud 0 bits, wat logisch is aangezien de ontvanger niets heeft bijgeleerd van dit bericht.

- ! De informatie-inhoud van een bericht is niet altijd een geheel getal.
- ! De informatie-inhoud is nooit negatief.
- Voorbeeld: Stel dat een byte verwacht wordt, maar er is geen idee welke byte. Elke byte is even waarschijnlijk met kans  $1/256$ . De informatie-inhoud van de byte die dan binnenkomt is  $-\log_2(1/256)$  bits = 8 bits.
- Voorbeeld: Stel een alfabet van 4 letters: A, C, G en T. De waarschijnlijkheid dat ze voorkomen wordt weergegeven in tabel 1.1.

A	70,71 %
C	12,50 %
G	8,39 %
T	8,39 %

Tabel 1.1: De waarschijnlijkheden voor de letters A, C, G en T.

Als de ontvanger dit weet dan wordt de informatie-inhoud voor elke letter:

$$\begin{aligned} A &: -\log_2(0,7071) = 0,5 \\ C &: -\log_2(0,1250) = 3,0 \\ G &: -\log_2(0,0839) = 3,575 \\ T &: -\log_2(0,0839) = 3,575 \end{aligned}$$

### 1.5.2 Beslissingsboom

- Elke knoop dat geen blad is bevat een vraag met een beperkt mogelijk aantal antwoorden.
- Elk mogelijk antwoord verwijst naar een kind van de knoop.
- Een item klasseren is een pad vanuit de wortel naar een blad, waarin de klasse staat.
- Hoeveel informatie kan een beslissingsboom geven?
  - Stel  $k$  klassen  $K_1, K_2, \dots, K_k$ .
  - Stel een verzameling  $S$  van items waarbij:
    - ◊  $A(S, i)$  het aantal elementen horend bij  $K_i$  is in de verzameling en,

- ◊  $|S| = \sum_{i=1}^k A(S, i)$  het totaal aantal element is van  $S$ .
- De informatie geleverd door een correcte klassering van alle element is dan:

$$\begin{aligned} E(S) &= \sum_{i=1}^k A(S, i) \left( -\log_2 \left( \frac{A(S, i)}{|S|} \right) \right) \\ &= |S| \log_2(|S|) + \sum_{i=1}^k A(S, i) (-\log_2(A(S, i))) \end{aligned}$$

- Het **Iterative Dichotomiser 3 (ID3)** algoritme is een inhalig algoritme dat een beslissingsboom opstelt vanuit een bepaalde dataset.

- De wortel bevat de vraag die het meeste informatie oplevert.
- Als het  $j$ -de attribuut de leerverzameling  $L$  in de deelverzamelingen  $L_{j,1}, L_{j,2}, \dots, L_{j,n}$  opdeelt, dan is de informatie geleverd door dit attribuut gelijk aan:

$$I(j) = E(L) - \sum_{m=1}^{n_j} E(L_{j,m})$$

- Het attribuut wordt gekozen waarvoor  $I(j)$  maximaal wordt.
- ! Als  $I(j) = 0$ , dan behoren ofwel alle items tot dezelfde klasse, ofwel kan er op basis van het attribuut geen klassering gemaakt worden.
- Na de constructie van de wortel moeten nog  $n_j$  deelbomen geconstrueerd worden.
- Voorbeeld:
  - ◊ Veronderstel volgende informatie:

Beroepscategorie	Jonger dan 20?	Fraudeur?	risico?	frequentie
A	ja	ja	veilig	10
A	ja	nee	riskant	11
A	nee	ja	riskant	18
A	nee	nee	veilig	100
B	ja	ja	veilig	180
B	ja	nee	riskant	8
B	nee	ja	riskant	1
B	nee	nee	veilig	90
C	ja	ja	veilig	50
C	ja	nee	riskant	5
C	nee	ja	riskant	5
C	nee	nee	veilig	50
Totaal veilig:				480
Totaal riskant:				48
Algemeen totaal:				528

- ◊ Voor elk attribuut kan de resulterende verdeling afgeleidt worden (tabel 1.2):
- ◊ Voor elk attribuut kan nu  $I(j)$  berekent worden, hier uitgewerkt voor de beroepsca-

Attribuut	waarde	# veilig	# riskant
(1) Beroepscategorie	A	110	29
	B	270	9
	C	100	10
(2) Jonger dan 20?	ja	240	24
	nee	240	24
(3) Fraudeur?	ja	240	24
	nee	240	24

Tabel 1.2: Resulterende verdeling.

tegorie:

$$\begin{aligned}
 I(1) &= E(L) - \sum_{m=1}^{n_1} E(L_{1,m}) \\
 &= E(L) - (E(L_A) + E(L_B) + E(L_C)) \\
 &\text{met} \\
 E(L) &= 480(-\log_2(480/528)) + 48(-\log_2(48/528)) \\
 &= 232.054 \\
 E(L_A) &= 110(-\log_2(110/139)) + 29(-\log_2(29/139)) \\
 &= 102.702 \\
 E(L_B) &= 270(-\log_2(270/279)) + 9(-\log_2(9/279)) \\
 &= 57.3603 \\
 E(L_C) &= 100(-\log_2(100/110)) + 10(-\log_2(10/110)) \\
 &= 48.3447 \\
 &\text{zodat} \\
 I(1) &= E(L) - (E(L_A) + E(L_B) + E(L_C)) \\
 &= 232.054 - 102.702 - 57.3603 - 48.3447 = 23.647
 \end{aligned}$$

- ◇ Voor  $I(2)$  en  $I(3)$  bedragen beide uitkomsten 0, aangezien er op basis van die attributen geen informatie kan achterhaald worden. Dit is logisch aangezien de verhouding veilige/risicohoudende klanten voor zowel leeftijd als fraudeur 10:1 is.
- ◇ Als wortel van de beslissingsboom wordt als criterium de beroepscategorie genomen.
- ◇ Voor zowel  $L_A$ ,  $L_B$  en  $L_C$  moet apart dezelfde methode uitgevoerd worden. Het kan perfect mogelijk zijn dat op het tweede niveau bij keuze  $A$  eerst wordt gecheckt op fraude, maar bij keuze  $C$  eerst op leeftijd.
- ! Bij  $L_C$  levert geen enkel van de twee attributen informatie op, zodat er random moet gekozen worden:

Attribuut	waarde	# veilig	# riskant
(2) Jonger dan 20?	ja	50	5
	nee	50	5
(3) Fraudeur?	ja	50	5
	nee	50	5

- ◇ Verfijningen:

- \* Invoeren van een drempelwaarde voor de informatiewinst. Als deze te klein is wordt er niet meer opgesplitst. Een blad kan dan meerdere klassen bevatten, elk met een waarschijnlijkheid.



- \* Voor elk mogelijk paar van attributen de informatiewinst berekenen en, als deze te groot is, knopen maken met twee attributen.
- \* Bij effectieve getallen kan ook een drempelwaarde ingevoerd worden. Alle items met een waarde kleiner dan deze drempelwaarde gaan naar links, de andere naar rechts.

## 1.6 Klasseren zonder leren

- Klassen worden niet vooraf gegeven.
- ! Geen leerverzameling aanwezig.
- Een **groep** van punten is wat door een expert als een groep beschouwd wordt.
  1. Twee punten behoren waarschijnlijk tot dezelfde groep als ze zeer dicht bij elkaar liggen. De expert definieert de afstandsfunctie.
  2. Deze eigenschap wordt **transitief** verdergezet. Twee punten  $\mathbf{x}$  en  $\mathbf{z}$  behoren tot dezelfde groep als een rij punten  $\mathbf{y}_1, \dots, \mathbf{y}_n$  bestaat zodanig dat  $\mathbf{x}$  zeer dicht bij  $\mathbf{y}_1$  ligt,  $\mathbf{y}_1$  zeer dicht bij  $\mathbf{y}_2$  ligt, ..., en  $\mathbf{y}_n$  zeer dicht bij  $\mathbf{z}$  ligt.

### 1.6.1 $k$ zwaartepunten

- Op voorhand opgeven dat er  $k$  klassen zijn.
- ! Een zwakte, aangezien men nu moet weten hoeveel klassen er op voorhand zijn. Twee problemen:
  - Het aantal gekozen klassen is te groot zodat samenhangende groepen worden opgesplitst. Dit kan opgelost worden samenhorende groepen op het einde samen te nemen.
  - Het aantal gekozen klassen is te klein zodat verschillende groepen samen worden genomen. Dit wordt opgelost door het algoritme meerdere malen uit te voeren met verschillende initialisaties.
- $k$  **zwaartepunten** poogt een leerverzameling  $L$  op te delen in  $k$  groepen  $G_1, \dots, G_k$ , waarbij  $k$  vooraf opgegeven is. Een klasse wordt voorgesteld door haar zwaartepunt  $m_i$ :

$$m_i = \frac{1}{n_i} \sum_{x \in G_i} \mathbf{x}$$

## 1.7 Een toepassing: Watson