



Toepassen van Nederlandse Gevoelsanalyse via sociale media

Yannick Merckx

Vorbereiding op de bachelorproef

Rolnummer: 500294

Promotor: Yann-Michaël De Hauwere
Begeleiders: Maarten Deville
Peter Vrancx



Inhoud

1	Introductie	2
2	Machine Learning	3
2.1	Wat is Machine Learning	3
2.2	Supervised Learning	3
2.3	Regressie Probleem	3
2.4	Lineaire regressie	4
2.5	Kost Functie en Gradule afdaling	4
2.6	Graduele afdaling	5
2.7	Classificatie Probleem	5
2.8	Unsupervised Learning	5
2.9	Data Mining	5
3	Technieken	6
4	Textmining	7
5	Experiment	8
6	Conclusie	9

1 Introductie

Vandaag de dag is sociale media een alledaags gegeven. Het is niet alleen onderdeel van het dagelijkse leven, maar ook een enorme bron aan informatie. Door deze data te analyseren kunnen we nog andere informatie afleiden. Zo kunnen we ook gevoelens afleiden uit de data. Met behulp van computers en meer bepaald

Machine Learning

kunnen we de data verzamelen, verwerken en analyseren zodanig dat men bepaalde informatie kan afleiden uit de data. Hoe het gegeven van verzamelen, verwerken en analyseren juist in elkaar zit, wordt beschreven in deze voorbereiding.

2 Machine Learning

Machine learning is een welgekend begrip in de informatica wereld, maar wat het juist omvat, zijn toepassingen en hoe het helpt om de juiste verbanden te achterhalen uit een enorme datasets wordt uitgelegd in dit hoofdstuk.

2.1 Wat is Machine Learning

Over wat Machine Learning juist is, vindt men nergens een eenduidige definitie. Vele hebben geprobeerd om een eenduidige definitie te definiëren, zoals Arthur Samuel(1959). Hij definieerde Machine Learning als

Field of study that gives computers the ability to learn without being explicitly programmed

. Later heeft Tom Michel(1999) ook een poging ondernomen en stelde een well-posed learning problem als het volgende

A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .

. Als we al deze pogingen proberen te omvatten, kunnen we Machine Learning het best beschrijven als een onderzoeksdomein dat zich bezighoudt met het onderzoeken en de ontwikkeling van zelflerende algorithmes.

Binnen Machine Learning kan men verschillende groepen van lerende algorithmes onderscheiden. Zo heeft men Supervised learning, Unsupervised learning, Reinforcement learning en Recommender systems. In deze voorbereiding gaat men zich enkel opleggen op supervised en unsupervised learning. Deze soorten algorithmen gaan de nodige antwoorden bezorgen om later gevoelsanalyses te kunnen uitvoeren.

2.2 Supervised Learning

Supervised learning is een techniek, waarbij men het algoritme traint met data waarvan men de antwoorden al weet. Algemeen noemt men een dataset waarmee men een algoritme traint een trainingsset. Nadat dit algoritme zijn training heeft ondergaan, kan het zelfstandig keuzes maken aan de hand van de vergaarde kennis.

Bij supervised learning zijn er twee soorten problemen die kunnen optreden: een regressie probleem of een classificatie probleem.

2.3 Regressie Probleem

Het doel dat men wilt bereiken met supervised learning is dat het algoritme na een training antwoorden kan bezorgen. Bij het voorspellen van die antwoorden kunnen we te maken hebben met een regressie probleem. Dit probleem valt het best uit te leggen aan de hand van een voorbeeld.

Neem nu dat men de prijs van een huis wilt voorspellen. Het algoritme traint zich met een trainingsset en bekomt volgend resultaat als men zijn bevindingen zou plotten.

Stel nu dat men aan het algoritme vraagt, wat de prijs is van een huis met 1225 vierkante meter. Deze waarde zat niet in de dataset en moet dus voorspeld worden. Maar welke trend moeten we volgen om onze waarden te voorspellen. We kunnen zowel kiezen voor een rechte of een 2de orde polynoom. Beiden zijn correct, maar geven een verschillend antwoord van elkaar. De situatie, waarbij men een continue waarde moet bepalen en geen echte discrete afbakening bestaat, noemt men een regressie probleem.

Om dit probleem op te lossen, kan men van de techniek

Lineaire regressie

gebruik maken.

2.4 Lineaire regressie

Lineaire regressie is een techniek waarbij het algoritme een hypothese probeert te vormen. De hypothese is een functie die opgesteld is aan de hand van de trainingsset en de gekende en ongekende outputwaarden zo goed mogelijk benaderd.

Als we terug kijken naar het voorbeeld van het huis. Kan men volgende hypothese opstellen.

Gegeven hypothese is een lineaire functie met als parameters θ_0 de nulconditie en θ_1 de richtingscoëfficiënt. Een hypothese met 1 functie noemt men ook wel een 1D lineaire regressie.

Het opstellen van de hypothese introduceert op zijn beurt een

Minimalisatie probleem

. Men moet de hypothese zo goed mogelijk opstellen, zodat de afwijking ten opzichte van de gekende resultaten minimaal is. Als de hypothese minimaal is, kan men er van uit gaan dat de afwijking op ongekende resultaten ook minimaal is.

Het minimalisatie probleem kan opgelost worden met een kost functie en graduele afdaling.

2.5 Kost Functie en Graduele afdaling

Een kost functie is een functie die voor bepaalde waarden van de parameters de gemiddelde afwijking van de hypothese ten opzichten van de resultaten gaat berekenen.

Volgende formule stelt de kost functie voor:

Merk op dat men niet zomaar telkens de som van het verschil tussen het resultaat van de hypothese neemt en de eigelijke waarden. Het kwadraat van het verschil wordt genomen vanwege de negatieve verschillen die ook moeten worden opgenomen als afwijking. Verder vereenvoudigt men het rekenwerk door te delen door twee (De helft van de kleinste waarde, blijft de kleinste waarde).

2.6 Graduele afdaling

2.7 Classificatie Probleem

Een Classificatie probleem is een ander soort van probleem dat zich voordoet bij supervised training. Een classificatie probleem doet zich voor wanneer men data moeten verdelen over verschillende discrete klassen. Ieder element mag maar tot 1 klasse behoren. De classificatie kan gebaseerd zijn op één attribuut, maar ook meerdere.

Hoe men deze classificatie juist aanpakt, wordt later verder uitgelegd.

2.8 Unsupervised Learning

Unsupervised learning is een techniek waarbij het algoritme zelfstandig moet leren hoe het juist moet en deze kennis gebruikt om later patronen en structuren in data te herkennen. De trainingsset bevat niet de juiste antwoorden.

Het herkennen van structuren en patronen is niet voldoende, men moet de data concreet kunnen identificeren. Dit kan men doen door gebruik te maken van cluster algorithmes. Concreet gaat een cluster algoritme de data groeperen of

clusteren"

in groepen.

2.9 Data Mining

3 Technieken

4 Textmining

5 Experiment

6 Conclusie

Men kan besluiten dat Machine learning ons de oplossing biedt. Het bevat alle middelen om de gevoelsanalyse op sociale media toe te passen. Om succesvol een gevoelsanalyse toe te passen moet men volgende stappen ondernemen. Ten eerst moet men de data verzamelen. Aangezien men te maken heeft met een dataset waarvan men geen informatie heeft, bevindt men te maken heeft met een dataset waarvan men geen informatie heeft, bevindt men technieken gebruiken van unsupervised learning. Vervolgens moet met de data preprocessen met het LSA algoritme. Het experiment heeft uitgewezen dat dit een zeer krachtig algoritme is, dat men meer inzicht verwerft in de data. Ten slotte moet men de data verwerken met een cluster algoritme. De sequentie van verzamelen, verwerken en analyseren is een zeer belangrijk gegeven bij onze gevoelsanalyse.

CONCLUSIE KAN NOG SPECIFIEKER, WORDT NOG AANGVULD NAAR-
MATE BP EVALUEERD. VOORAL CLUSTER ALGORITHMEN