

Rapport - IPASS

Jimmy Bierenbroodspot

26 juni 2024

1 Probleembeschrijving

In de toekomst willen wij graag een AI-model hebben, dat de beste CV's op de markt kan maken voor onze klanten. Hiervoor hebben we natuurlijk veel data in de vorm van CV's nodig en gelukkig hebben we deze in grote aantallen. Nu is het probleem dat deze CV's in .pdf-formaat opgeslagen zijn en text hieruit halen is niet makkelijk [5].

We hebben het programma genaamd PyPDF[4] waarmee we in staat zijn om text uit .pdf-bestanden te lezen, het probleem is dat dit ongestructureerd is. Als voorbeeld kunnen we het voorbeeld CV in 1 nemen.

DRIE TREFWOORDEN + PERSOONSOOMSCHRIJVING Rego [WOONPLAATS]	ERVARING <hr/> Functietitel + bij + Bedrijf startmaand startjaar - eindmaand eindjaar [Meest effectieve zin van het hele CV, kort, krachtig, heel erg passend bij de vacature.] Functietitel + bij + Bedrijf startmaand startjaar - eindmaand eindjaar [Optionele uitleg van ervaring en ontwikkeling]
DRIVE [TEKST DRIVE - MIDDEN KOLOM 1 GEPOSITIONEERD]	OPLEIDING <hr/> Titel opleiding + opleidingsinstituut Student startmaand startjaar - eindmaand eindjaar [optionele uitleg studie]
KUNDE [TEKST KUNDE - MIDDEN KOLOM 1 GEPOSITIONEERD]	Titel opleiding + opleidingsinstituut Student startmaand startjaar - eindmaand eindjaar [optionele uitleg studie]
MOGELIJKHEDEN [TEKST MOGELIJKHEDEN - MIDDEN KOLOM 1 GEPOSITIONEERD]	
	<hr/> Contactgegevens kandidaat

Figure 1: Voorbeeld van een Cliq CV

In 2 kunnen we zien wat er gebeurd als we de tekst uit 1 lezen. Dit ziet er uit als een onlogisch rommeltje. We kunnen met PyPDF echter ook de tekstopmaak kopiëren, zoals te zien in 3. Dit ziet er al beter uit maar als we bedenken dat dit allemaal op een regel staat, waarbij de nieuwe regels in dit geval juist worden getoond, maar eigenlijk aangeduid wordt met een \n.

```

Regio [WOONPLAATS]
DRIVE
[TEKST DRIVE - MIDDEN KOLOM
1 GEPOSITIONEERD]
KUNDE
MOGELIJKHEDEN
[TEKST MOGELIJKHEDEN - MIDDEN
KOLOM 1 GEPOSITIONEERD]OPLEIDING
Contactgegevens kandidaatERVARING
Student | startmaand startjaar - eindmaand eindjaarDRIE TREFWOORDEN +
PERSOONSOMSCHRIJVING Functietitel + bij + Bedrijf
[TEKST KUNDE - MIDDEN KOLOM
1 GEPOSITIONEERD]Titel opleiding + opleidinginstituut
[optionele uitleg studie][Meest effectieve zin van het hele CV, kort, krachtig, heel
erg passend bij de vacature.]
Functietitel + bij + Bedrijfstartmaand startjaar - eindmaand eindjaar
[Optionele uitleg van ervaring en ontwikkeling]
Student | startmaand startjaar - eindmaand eindjaarTitel opleiding + opleidinginstituut
[optionele uitleg studie]startmaand startjaar - eindmaand eindjaar

```

Figure 2: Voorbeeld van tekstextractie PyPDF

```

DRIE TREFWOORDEN +
PERSOONSOMSCHRIJVING
Re g io [WOONPL AATS]

DRI VE
[TEKST DRIVE - MIDDEN KOL OM
1 GEPOSITIONEERD]
KUNDE
[TEKST KUNDE - MIDDEN KOL OM
1 GEPOSITIONEERD]
MOGEL IJKHEDEN
[TEKST MOGEL IJKHEDEN - MIDDEN
KOL OM 1 GEPOSITIONEERD]

ERVARING
F u n cti etitel + b ij + Bed rijf
sta rtma a n d sta rtja a r - e in d ma a n d
[Ma e st e Me etie va rin va n h e t h e le
e rg p a s s e n d b ij d e va ca tu re .]
F u n cti etitel + b ij + Bed rijf
sta rtma a n d sta rtja a r - e in d ma a n d
[Op t ie n e le u itle g va n e r va rin g e
OPL EIDING
T itel o p leid in g + o p leid in g in stitu u
Stu d e n t | sta rtma a n d sta rtja a r - e in
[io p t ie n e le u itle g stu d ie ]
T itel o p leid in g + o p leid in g in stitu u
Stu d e n t | sta rtma a n d sta rtja a r - e in
[io p t ie n e le u itle g stu d ie ]
Co n ta ctg e g e ve n s ka n d id a a t

```

Figure 3: Voorbeeld van tekstextractie PyPDF met tekstopmaak

1.1 Waarom willen we dit?

In datawetenschap kennen we het concept GIGO, oftewel garbage in garbage out [3]. Als we aannemen dat hetzelfde zal gelden voor het model dat we gaan gebruiken dan zouden we graag willen dat onze invoerdata gestructureerd is.

Als we aannemen dat de tekst op een regel zou moeten komen te staan om ingelezen te worden door het model, dan zijn beide manieren van text extraheren met PyPDF niet gestructureerd. We willen dus dat de tekst op een logische volgorde komt te staan. Als een CV geen kolommen zou bevatten dan kunnen we met PyPDF de tekstopmaak kopiëren, nu weten we alleen van tevoren niet of een CV kolommen heeft of niet en als dat wel zo is weten we niet hoeveel kolommen er zijn. Om hierachter te komen moeten we naar het CV kijken en dat gaat niet met een geautomatiseerd proces.

2 Eisen

De opdrachtgever heeft in ieder geval de volgende eisen gegeven:

1. Er moet een CV bestand ingelezen worden.
2. Het aantal kolommen moet teruggegeven worden.

Verder zijn er geen functionele of prestatie-eisen.

3 Algoritme

In 3 kunnen we zien dat er een best duidelijke scheiding zit tussen de twee kolommen. We hebben besloten om hier Lloyds algoritme (zie[2]) voor te gebruiken omdat het een makkelijke en snelle oplossing leek te zijn voor dit probleem.

Om erachter te komen hoeveel kolommen het optimale aantal zou zijn gebruiken we de elleboog methode, deze werkt als volgt[1]:

1. Voor elk gewilde aantal clusters wordt het model getraind.
2. Elke keer dat het model getraind is wordt de error opgeslagen.
3. Er wordt een rechte lijn van de eerste observatie naar de laatste observatie getrokken.
4. Bereken de afstand tussen deze lijn en elke score.
5. Neem het aantal clusters van de score die het verst van deze lijn verwijderd is.

4 Resultaat

Na het evalueren van het model op een dataset van 116 bestaande CV's zien we een teleurstellend resultaat. We weten maar van ongeveer 15% van de gevallen het juiste aantal kolommen te vinden. Dit is echt veel te weinig om in een productieomgeving te gebruiken.

5 In het vervolg

We zouden in de toekomst willen kijken of de volgende verbeterpunten het presteren van de clusteringsaanpak zal vergroten:

- het model meerdere keren uitvoeren met verschillende willekeurige clustercentra voor elke keer dat het model getraind wordt.
- Kmeans++ gebruiken.
- De invoerdata standaardiseren.

Ook zouden we dit probleem graag als classificatieprobleem willen proberen aanpakken aangezien dit daadwerkelijk ondersteund wordt door bestaande literatuur.

References

- [1] *Find the "elbow" of a curve*. URL: https://docs.ropensci.org/pathviewr/reference/find_curve_elbow.html (visited on 06/26/2024).
- [2] S. Lloyd. "Least squares quantization in PCM". In: *IEEE transactions on information theory* 28.2 (Mar. 1, 1982), pp. 129–137. DOI: 10.1109/tit.1982.1056489. URL: <https://doi.org/10.1109/tit.1982.1056489>.
- [3] Merriam-Webster. *GIGO*. In: URL: <https://www.merriam-webster.com/dictionary/GIGO> (visited on 05/30/2024).
- [4] *pypdf*. Apr. 7, 2024. URL: <https://pypi.org/project/pypdf/> (visited on 05/30/2024).
- [5] Amit Timalina. *How to extract data from PDF?* May 2, 2024. URL: <https://www.docsumo.com/blog/extract-data-from-pdf> (visited on 05/30/2024).

List of Figures

1	Voorbeeld van een Cliq CV	1
2	Voorbeeld van tekstextractie PyPDF	2
3	Voorbeeld van tekstextractie PyPDF met tekstopmaak	2

List of Tables