

Targetowanie reklam z machine learning

Czyli słowa kluczowe i trochę Adtechu

Spis treści

1. Analiza danych
2. Parsowanie danych
3. Zdanie jako wektor
4. TFIDF
5. LSA
6. Clustering + k-means
7. Wynik + co dalej

Uwaga - dużo nerdbelkotu!

Ten nudny slajd ze mną

1. Jestem programistą Pythona
2. Należę do Hackerspace Silesia
3. W pracy zajmuję się Adtech-em
4. Zajmuję się zliczaniem wejść i zakupów na stronę oraz raportowaniem



Geneza problemu

- Klient zajmuje się reklamą wielu sklepów
- Klient bardzo chciał machine learning
- Jedynie jakie mamy dane to z logi tracker-ów
- Dane muszą zostać przerobione do postaci linia - słowa kluczowe

Analiza danych

- Jakie to są dane?
- Jak one wyglądają?
- (dokąd zmierzamy?)
- Czy wszystkie słowa są nam potrzebne?

Parsowanie danych

1. Mamy tylko logi z NGINX
2. Interesuje nas tylko adres (bez host-a)
3. Pozbywamy się stron typu admin / login
4. Pozbywamy się liczb (nic nie mówiące ID)
5. Pozbywamy się niepotrzebnych słów oraz znaków specjalnych
6. Słowa kluczowe grupujemy po IP/cookieID
7. Dokument = słowa kluczowe per IP/cookieID

Uczenie Maszynowe

- Dane musimy przerobić tak, by stały się “przeliczone”
- Jeżeli uda nam się zamienić na macierz liczb - możemy stworzyć model dla uczenia maszynowego

Zdanie jako wektor

- Dokument = 1 linijka zgrupowanych słów (per IP/cookie)
- Dokument zawiera “Term-y”
- Term-em może być znak, słowo, para słów, link etc.
- Dokument możemy przedstawić jako wektor częstotliwości słów
- Wiele dokumentów = wiele wektorów = macierz $^{\wedge}_{_}^{\wedge}$

1. Ala ma kota
2. Kot ma Alę
3. Basia ma kota

Wiersz	Ala	ma	kota	kot	Alę	Basia
1	1	1	1	0	0	0
2	0	1	0	1	1	0
3	0	1	1	0	0	1

1. Ala ma kota
2. Basia ma kota

Wiersz	Ala	Ala ma	ma kota	kota	Basia	Basia ma
1	1	1	1	1	0	0
2	0	0	1	1	1	1

TFIDF

- TF = **T**erm **F**requency
- (częstotliwość term-u w dokumencie)
- IDF = **I**nverse **D**ocument **F**requency
- (im rzadziej się term pojawia w dokumentach tym bardziej ma wyższy ranking)

$$\text{TFIDF}(\text{term}, \text{doc}) = \text{TF}(\text{term}, \text{doc}) \times \text{IDF}(\text{term})$$

$$\text{TF} = \text{count}(\text{term}) / \text{count}(\text{terms in doc})$$

$$\text{IDF} = \log_{10} (\text{count}(\text{docs}) / \text{count}(\text{docs with term}))$$

term	TF	IFD	TFIDF
Ala	0.33	$\log(3/1) \approx 1.1$	0.363
ma	0.33	$\log(3/3) = 0$	0
kota	0.33	$\log(3/2) \approx 0.4$	0.132

1. Ala ma kota
2. Kot ma Ale
3. Basia ma kota

term	TF	IFD	TFIDF
Kot	0.33	$\log(3/1) \approx 1.1$	0.363
ma	0.33	$\log(3/3) = 0$	0
Ale	0.33	$\log(3/1) \approx 1.1$	0.363

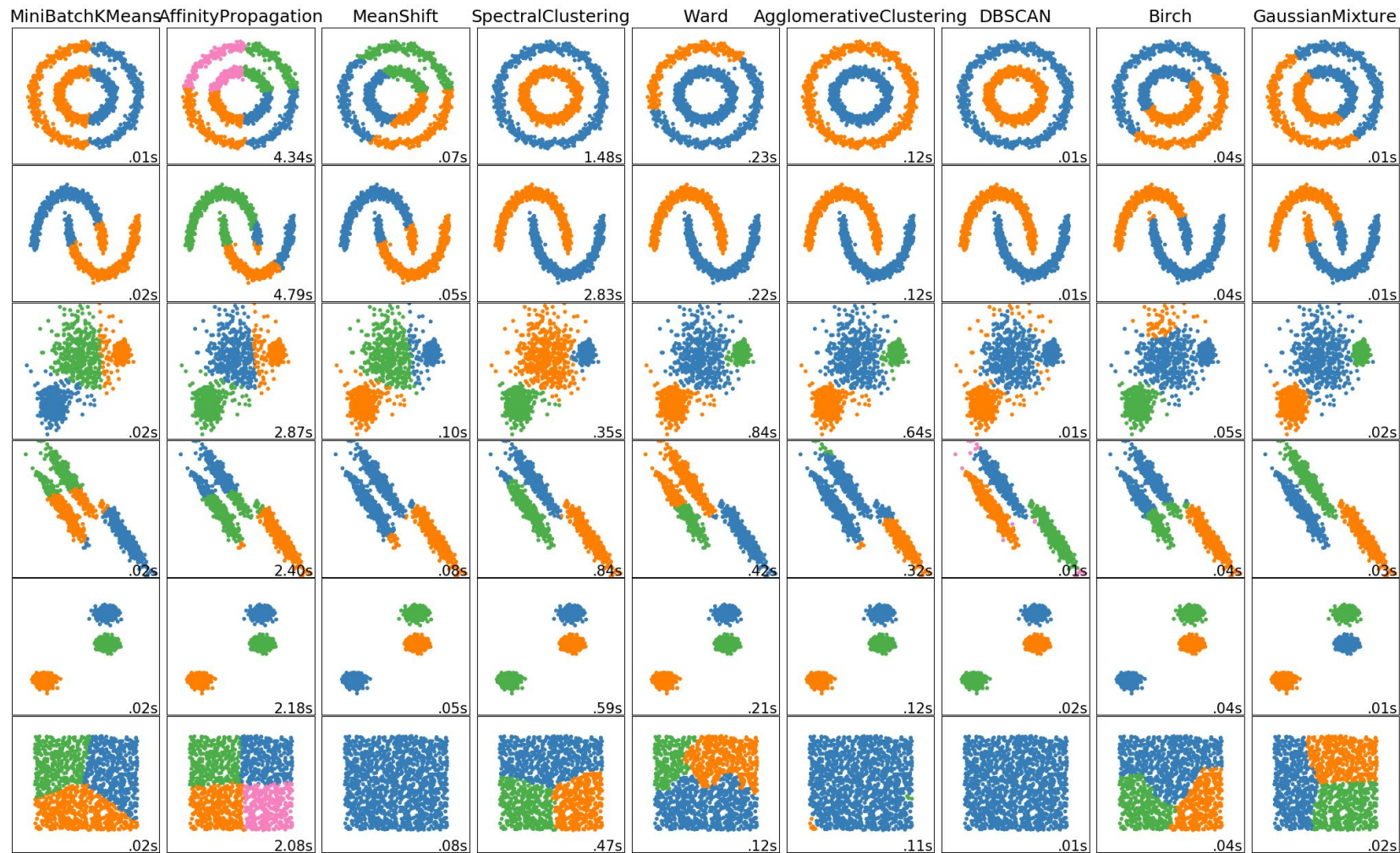
term	TF	IFD	TFIDF
Basia	0.33	$\log(3/1) \approx 1.1$	0.363
ma	0.33	$\log(3/3) = 0$	0
kota	0.33	$\log(3/2) \approx 0.4$	0.132

LSA

- LSA - Latent semantic analysis
- Grupuje wektory na podstawie wagi (np. TFIDF) term-ów
- Grupuje term-y na podstawie podobnych wag
- Na podstawie grup jesteśmy w stanie określić relację termów
- Animacja bo ja tego nie wytłumaczę :)

Clustering

- Nie wiemy nic co możemy uzyskać z danych
- Musimy użyć z uczenia nienadzorowanego
- Clustering; Klasteryzacja - na podstawie wartości wektorów jesteśmy w stanie sklasyfikować grupy term-ów

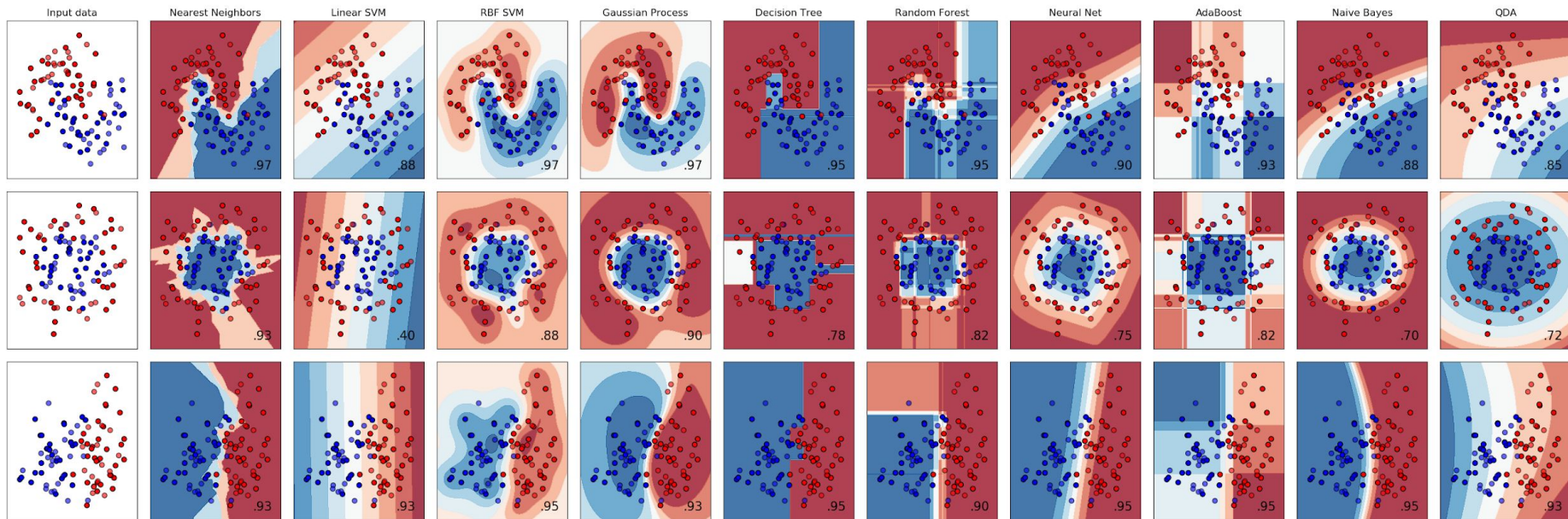


Clustering

- Każdy jeden wiersz w macierzy możemy potraktować jako wektor w przestrzeni N wymiarowej
- Algorytm k-means: szukanie N najbliższych wektorów tworząc grupę

Klasyfikacja

- Na podstawie stworzonych grup możemy zbudować klasyfikator
- Zbudowany klasyfikator będzie w stanie ocenić ‘grupę’ na podstawie wejścia (w naszym przypadku, zawartości stron po których się porusza)



Wynik + co dalej

- Dzięki klasteryzacji słów jesteśmy w stanie stworzyć odpowiednie kampanie reklamowe czy zmienić kategorię produktów na stronie
- Na podstawie cookieID oraz słów kluczowych których używa gość na stronie mogą wyświetlić potencjalne reklamy lub zaproponować odpowiedni produkt
- **WIĘCEJ DANYCH** - z wielu dni + inne dane o klientach

Źródło

1. https://en.wikipedia.org/wiki/Latent_semantic_analysis
2. http://scikit-learn.org/stable/auto_examples/text/document_clustering.html
3. http://scikit-learn.org/stable/auto_examples/classification/plot_classifier_comparison.html
4. http://scikit-learn.org/stable/auto_examples/cluster/plot_cluster_comparison.html
5. <https://pandas.pydata.org/>
6. <https://github.com/firemark/word-analyser>