

PRESENTACIÓ CAS KAGGLE

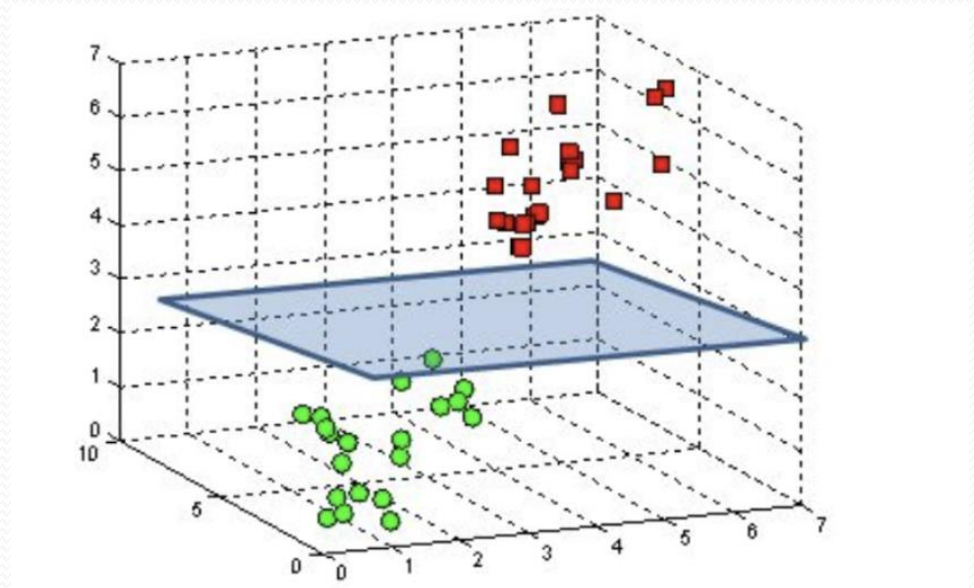


Autor: Òscar Pocurull Rodríguez. NIU: 1496348

Dataset: Classification of websites

INDEX:

- 1. INTRODUCCIÓ
- 2. PREPROCESSING
- 3. SELECCIÓ DE MODELS
- 4. RESULTATS
- 5. CONCLUSIÓ



1. INTRODUCCIÓ

- En aquesta practica analitzarem un dataset que classifica les webs segons la categoria a la qual pertanyen.

Unnamed: 0		website_url	cleaned_website_text	Category
0	0	https://www.booking.com/index.html?aid=1743217	official site good hotel accommodation big sav...	Travel
1	1	https://travelsites.com/expedia/	expedia hotel book sites like use vacation wor...	Travel
2	2	https://travelsites.com/tripadvisor/	tripadvisor hotel book sites like previously d...	Travel
3	3	https://www.momondo.in/?ispredir=true	cheap flights search compare flights momondo f...	Travel
4	4	https://www.ebookers.com/?AFFCID=EBOOKERS-UK.n...	bot create free account create free account si...	Travel
...
1403	1403	http://www.oldwomen.org/	old nude women porn mature granny sex horny ol...	Adult
1404	1404	http://www.webcamslave.com	bdsb cams bdsm chat bondage cams free bdsm vid...	Adult
1405	1405	http://www.buyeuroporn.com/	porno dvd online european porn dvd cheap adult...	Adult
1406	1406	http://www.analdreamhouse.com/30/03/agecheck/i...	anal dream house anal dream house anal dream h...	Adult
1407	1407	http://www.world-sex-news.com/	world sex news daily sex news adult news eroti...	Adult

1. INTRODUCCIÓ

- L'objectiu es aprendre quina es la categoria ('Category') a la qual pertanyen les diferents pàgines web del nostre data set.



2. PREPROCESSING

- - Eliminem atribut 'Unnamed: 0'
- + Afegim atribut 'Category_ID'
- Eliminem les webs repetides amb `drop_duplicates('cleaned_website_text')` on passem de 1408 a 1384 webs.
- No te valors nulls



2. PREPROCESSING

- Tf-idf: Convertim el text en dades numèriques.
- Paràmetres `TfidfVectorizer`-> `sublinear_tf=True`, `min_df: 6`, `ngram_range(1,2)` (unigrams i bigrams), `stop_words='english'`.
- l'utilitzarem per trobar les paraules (unigrams i bigrams) més correlacionades de cada categoria amb la funció `chi2()` de `sklearn.feature_selection`.

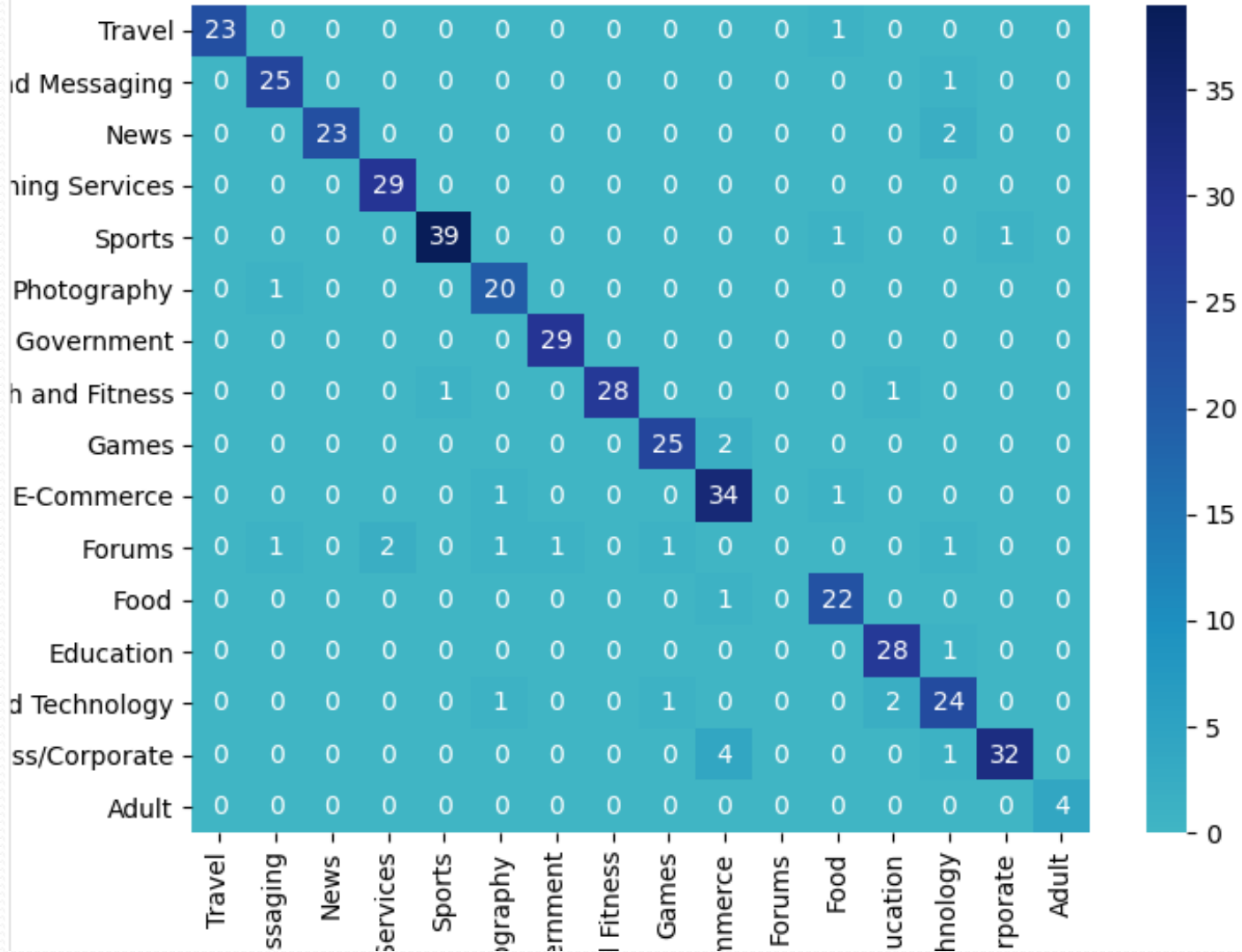
3. SELECCIÓ DE MODELS

- Models: 70% train, 30% test i cross-validation k-fold, k=5.

Model	Hiperparametres	Mètrica	Temps
Random Forest	100 Trees, max_dept: 5	68%	5227 ms
Random Forest	100 Trees, max_dept: inf	83%	17310 ms
LinearSVC	penalty: l2, multi_class: ovr, max_iter: 1000	91%	2304 ms
MultinomialNB	alpha: 1.0	85%	812 ms
GaussianNB	var_smoothing: 1e-09	73%	6173 ms
Logistic Regression	penalty: l2, C: 1, max_iter: 100	89%	17900 ms

4. RESULTATS

- LinearSVC: model amb millor accuracy.
- LinearSVC calibrat = 94% accuracy.
- LinearSVC usant Pipeline (no cv) = 89% acc. ($89 < 91$).
- SVC amb kernel: rgb = 87%. Time = 257 seg.
- SVC amb kernel: linear = 90% acc. Time = 242 seg.
- SVC amb kernel: linear, C:100 = 83% acc.
- SVC amb kernel: poly = 43% acc. Time = 252 seg.

[illegible]

4. RESULTATS

- Applicant només unigrames al tf-idf:

Model	Hiperparametres	Mètrica	Temps
Random Forest	100 Trees, max_dept: inf	84%	13776 ms
LinearSVC	penalty: l2, multi_class: ovr, max_iter: 1000	91%	1511 ms
MultinomialNB	alpha: 1.0	85%	420 ms
GaussianNB	var_smoothing: 1e-09	64%	3365 ms
Logistic Regression	penalty: l2, C: 1, max_iter: 100	89%	13563 ms

- Applicant només bigrames al tf-idf:

Model	Hiperparametres	Mètrica	Temps
Random Forest	100 Trees, max_dept: inf	65%	17830 ms
LinearSVC	penalty: l2, multi_class: ovr, max_iter: 1000	76%	555 ms
MultinomialNB	alpha: 1.0	71%	195 ms
GaussianNB	var_smoothing: 1e-09	74%	1700 ms
Logistic Regression	penalty: l2, C: 1, max_iter: 100	73%	4926 ms

5. CONCLUSIONS

- El millor model es el LinearSVC: max. Accuracy (91%) i ràpid al executar-lo (< 2 seg).
- SVC i NuSVC (implementats amb **libsvm**) tenen pitjor rendiment que el LinearSVC (amb **liblinear**) per a un gran numero de mostres (1384).
- Seria interessant optimitzar més les dades de text per trobar el balanç de millor temps i rendiment possible.