

EMD Python

Łukasz Eckert

February 2021

Contents

1	Wprowadzenie	2
2	Analiza zbioru danych	2
3	Przebieg eksperymentów	2
4	Modele	3
4.1	Baseline	3
4.2	Target Encoding	4
4.3	Modele oparte na tekście	5
4.3.1	TF-IDF	5
4.3.2	FastText	6
4.3.3	LSTM	7
5	Wnioski	8

1 Wprowadzenie

Celem projektu jest próba zaproponowanie klasyfikatora starającego się przewidzieć wystawioną ocenę przez użytkowników. Oceniane są aplikacje.

2 Analiza zbioru danych

Zbiór treningowy posiada 555791 rekordów. Każdy rekord opisany jest przy pomocy 9 atrybutów. Są nimi:

- reviewerID - string określający użytkowników
- asin - identyfikator aplikacji
- reviewerName - nazwa użytkownika
- helpful - czy opinia została oceniona jako pomocna
- reviewText - tekst recenzji
- summary - streszczenie recenzji
- unixReviewTime - moment wystawienia oceny w formacie unixowym
- reviewTime - timestamp recenzji
- score - liczba od 1 do 5 określająca ocenę tej aplikacji przez użytkownika

Niektóre z tych cech w naturalny sposób nadają się do dodania do modelu. Jednak niektóre z nich nie mogą być dodane w sposób naturalny do zbioru uczącego. Dodanie reviewerID oraz asin spowodowałoby mocne przeuczenie.

3 Przebieg eksperymentów

Wraz z tym sprawozdaniem dostarczony jest kod wszystkich modeli wraz z wytrenowanymi modelami. W celu ewaluacji na nowym zbiorze należy uruchomić wybrany moduł podając ścieżkę do pliku csv. Tak jak to zostało zauważone wcześniej, zbiór danych jest niezbalansowany. Istnieje ocena o wartości 5 stanowiąca niemal 50% wszystkich danych. Dodatkowo ilość dostarczonych rekordów pozwoliła by na ich zbalansowanie. Jak strategię balansującą zdecydowano się wykorzystać undersampling. Liczność przykładów z każdej z klas została zredukowana w sposób losowy do klasy o najmniejszej liczności. Jednak eksperymenty na tak przygotowanym zbiorze okazały się generować gorsze modele. Co okazało się być sporym zaskoczeniem. Modele nauczone na pełnym zbiorze dawały lepsze wyniki nie tylko na pełnym zbiorze walidacyjnym, ale także na tym zrównoważonym. Dlatego też zdecydowano się uczyć modele na oryginalnych danych.

4 Modele

Ta sekcja ma za zdanie opis wszystkich nauczonych i przetestowanych modeli. Cechy tekstowe zostały wstępnie oczyszczone. Wszystkie wielkie litery zostały przetransformowane to swoich "małych" odpowiedników. Wszystkie znaki które nie należą do zbioru a-z zostały usunięte. Niektóre modele wymagały specjalnego preprocesingu i znajdują się w sekcji poświęconej temu modelowi. Jako miare do porównania modeli zdecydowano się wykorzystać:

- Accuracy - dokładność klasyfikacji. Ze względu na zbalansowanie danych ta miara jest tylko używana do porównania się z Baslinem.
- L1 - wartość bezwzględna pomiędzy najbardziej prawdopodobną predykcją a oczekiwaną wartością.
- E[L1] - oczekiwana wartość miary L1.

miara $E[L1]$ jest zdefiniowana następująco:

$$E[L1] = E(proba, target) = abs(\sum_{i=1}^5 proba[i] * i - target)$$

Jest to naturalne przekształcenie wyjścia klasyfikatora na model regresji. Idea tej miary jest następująca. Jeżeli prawdziwa ocena jest ocena 3. Model przyznaje ocene 4 z 51% prawdopodobieństwem a ocene 3 z 49% prawdopodobieństwem. Dokładność wynosi w takim przypadku 0, L1 1, a E[L1] 0.49. Pomyłka pomiędzy 3 a 4 jest mniejsza pomyłka niż pomiędzy 5 a 1.

4.1 Baseline

Pierwszym zaproponowanym i przetestowanym modelem jest Baseline. Jest to model który za każdym razem zwraca najczęstsza klasę (5). Taki model osiąga accuracy równe 0.50. Jest to częstość występowania najczęstszej klasy. L1 oraz E[L1] wynoszą 1.04. Całość zostało zaimplementowana przy pomocy DummyClassifier z pakietu sklearn. Obraz 1 pokazuje macierz pomyłek.



Figure 1: Macierz pomyłek dla modelu Baseline

4.2 Target Encoding

Ten model reprezentuje pierwszą i jedyną próbę stworzenia modelu, który nie używa informacji znajdujących się w tekście. W danych znajdują się dwie kolumny, które są identyfikatorami: reviewerID oraz asin. Pierwsza określa osobę, która wystawiła ocenę. Druga określa aplikację. Dodanie ich bezpośrednio do danych uczących byłoby błędem. Mogłoby spowodować przeuczenie się modelu. Dlatego też zdecydowano się zastosować inną strategię. ID każdego użytkownika zostało zastąpione średnią ocen danej osoby. Podobnie zostało to wykonane dla aplikacji. Na tak stworzonych danych została nauczona regresja logistyczna. Ostatecznie udało się uzyskać następujące wyniki:

- Accuracy: 0.49
- L1 1.11
- E[L1] 1.11

Obraz 2 prezentuje macierz pomyłek dla tego modelu. Występuje tutaj ciekawe zjawisko. Model odkrył, dwie skrajne klasy 5 oraz 1. Dodatkowo nigdy nie wybiera innych klas.



Figure 2: Macierz pomyłek dla modelu Target Encoding

4.3 Modele oparte na tekście

W tej sekcji opisane są modele oparte na tekście zawartym w danych. Istnieją dwie kolumny posiadające informacje o typie tekstowym: reviewText oraz summary. Wykorzystanie obydwóch niestety nie jest możliwe przez ograniczenia obliczeniowe. Tekst został wstępnie oczyszczony, tak jak to zostało opisane we wcześniejszej sekcji.

4.3.1 TF-IDF

Model opisywany w tej sekcji został zaimplementowany przy pomocy sklearn. Model składa się z następujących komponentów:

- Stemizacja - tokeny z tekstu zostały poddane skracaniu przy pomocy SnowballStemmer z pakietu nltk
- TfidfVectorizer - tworzy wielowymiarową macierz cech. Jest to jedno z klasycznych podejść na ekstrakcję cech z tekstu
- TruncatedSVD - wymiarowość TFIDF może sięgać ilości unikatowych słów w korpusie. To jest zdecydowanie zbyt dużo, aby w sensownym czasie nauczyć jakiś model. Ten krok ma na celu redukcję wymiarowości. Zachowywane są 500 pierwszych wymiarów. Opisują one 0.83% wariancji w danych
- RandomForestClassifier - klasyfikator. Inne rodzaje klasyfikatorów zostały sprawdzone (SQCClassifier, KMeansClassifier), jednak radziły sobie one zdecydowanie gorzej.

Na tak przygotowanym modelu został zastosowany GridSearch w celu znalezienia optymalnych parametrów. Liczba wymiarów TruncatedSVD została dobrana ręcznie. Optymalizowanymi parametrami była liczba drzew w lesie (100, 500, 1000]) oraz maksymalna głębokość drzew (4,8). Jako strategię podziału na kubeczki do cross walidacji został wykorzystany TimeSeriesSplit. Tak przygotowane model był uczony raz na kolumnie summary raz na reviewText. Ostatecznie okazało się, że model oparty o reviewText jest lepszy. I to dla niego raportowane są metryki.

Obraz 3 pokazuje macierz pomyłek

- Accuracy: 0.54
- L1 0.9
- E[L1] 0.92

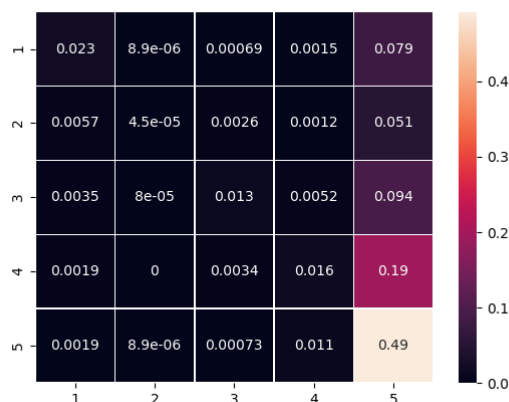


Figure 3: Macierz pomyłek dla modelu opartego na tf-idf

4.3.2 FastText

Pierwszy testowany model oparty o uczenie głębokie. Jest to model bazujący na architekturze FastText. Słowa były zamieniane na wektory przy pomocy już nauczonych wektor glove. Każdy token w tekście zostaje zamieniony na wektor o długości 100. Możliwość wektorów słów powodowała przeuczenie się do zbioru treningowego (0.96% acc dla zbioru treningowego przy 0.3% dla zbioru walidacyjnego). Dlatego też wartości wektorów ostatecznie nie były uczone.

Wektory dla słów danej recenzji są uśredniane a następnie przepuszczane przez sieć neuronową o 2 warstwach ukrytych. Na wyjściu znajduje się Log-Softmax. Taka funkcja aktywacyjna jest bardziej stabilna numerycznie od standardowej funkcji Softmax w czasie obliczania błędu. Jako optyimizator został

użyty Adam. Całość modelu została zaimplementowana przy pomocy frameworka pytorch. Model trenowany był tylko przez 10 epok przez ograniczenia czasowe.

Minusem tego rozwiązania jest utrata informacji o kolejności słów występujących w recenzji. Następny model będzie się starał poprawić tą wadę.

- Accuracy: 0.56
- L1 0.74
- E[L1] 0.79

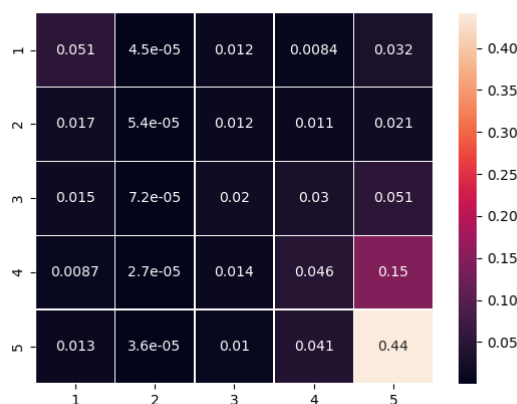


Figure 4: Macierz pomyłek dla modelu opartego na modelu FastText

4.3.3 LSTM

Model oparty na architekturze FastText ma jedną wadę. Tracona jest informacja o kolejności słów w recenzji. Zaproponowany tutaj model zastępuje obliczanie średniej z wektorów słów przez rekurencyjną sieć neuronową. Dokładnie przez LSTM. Wszystkie pozostałe założenia pozostają bez zmian. Dodatkowo została dodana warstwa Dropout aby ograniczyć efekt przeuczenia.

Model na 10, ostatniej epoce nie wykazywał jeszcze efektów przeuczenia. Możliwe, że wydłużenie czasu uczenia mogłoby poprawić wyniki.

- Accuracy: 0.64
- L1 0.49
- E[L1] 0.55

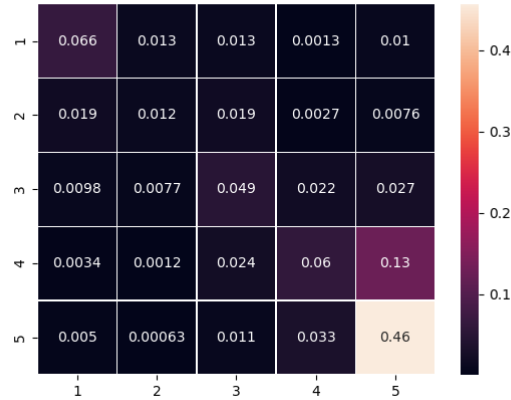


Figure 5: Macierz pomyłek dla modelu opartego na LSTM

5 Wnioski

Tabela 1 pokazuje zebrane wyniki dla każdego modelu. Problem okazał się dużo trudniejszy niż to mogło się wydawać. Jednym z problemów występującym w danych jest dość spore niezbalansowanie dystrybucji ocen. Jednak najprostsza z taktyk możliwa do zastosowania w takim wypadku, undersampling, nie poprawiała wyników a wręcz je pogorszyła.

Kolejnym problemem jest to w jaki sposób ludzie wystawiają oceny. Różnica pomiędzy dwiema sąsiednimi decyzjami jest zazwyczaj dość mała. Szczególnie to widać na macierzy najlepszego modelu LSTM. Wyniki predykcji są położone w pobliżu przekątnej.

Ostatecznie najlepszy model nie ma dobrej dokładności predykcji. Jednak pocieszające są wartości dla L1 oraz E[L1]. Analizując jest możemy dojść do wniosku, że zazwyczaj predykcja modelu jest oddalona maksymalnie o 1 od właściwej oceny.

Praktycznie wszystkie cechy dostępnym w zbiorze danych zostały w jakiś sposób wykorzystane.

- reviewerID oraz asin - model Target Encoding
- znaczniki czasowe - podział na zbiór treningowy oraz walidacyjny.
- Cechy tekstowe - w 3 modelach opartych o tekst

Jedna rzecz która mogła zostać rozwinięta to jest preprocessing danych. Jednak ostatecznym celem było stworzenie modelu opartego o uczenie głębokie i już wytrenowane embedding. Mocny preprocessing nie jest w takim momencie zalecany. Podejście zastosowane w tym projekcie wydaje się być zazwyczaj wystarczające.

	Baseline	Target Encoding	Tf-IDF	FastText	LSTM
Accuracy	0.50	0.49	0.54	0.56	0.64
L1	1.04	1.11	0.9	0.74	0.49
E[L1]	1.04	1.11	0.92	0.79	0.55

Table 1: Zbiorcze wyniki