

Classificazione di testo con Perceptron e Naive Bayes

Carlo Baronti

Università degli Studi di Firenze

Email: carlo.baronti@stud.unifi.it

Giulio Baronti

Università degli Studi di Firenze

Email: giulio.baronti@stud.unifi.it

Sommario—Il lavoro consiste nell'implementazione e analisi delle performance di classificatori di testo tramite algoritmi Bernoulli Naive Bayes, Multinomial Naive Bayes e Perceptron. Come dataset di riferimento è stato usato il ModAdpte split di Reuters-21578 di cui sono state scelte le 10 categorie con più documenti.

1. Introduzione

Inizialmente si sono divisi i documenti di training e i documenti di test, dopo di che si è estrapolato un insieme di parole detto vocabolario da cui sono state eliminate le stopwords; per fare ciò ci si è avvalsi dello stemmer di Porter. Tali funzioni preliminari sono svolte all'interno del modulo `manager` e vengono eseguite in circa un minuto su Google Colab. In seguito, si continuerà a fare riferimento ai tempi di esecuzione su Google Colab. La struttura del programma comprende, oltre al `main.py`, i moduli `bernoulliNB.py`, `multinomialNB.py` e `perceptron.py` relativi ai tre algoritmi omonimi [2], [3]. I risultati presentati in questo lavoro si riferiscono solo alla categoria `interest`, mentre i risultati relativi alle altre categorie sono accessibili presso [4]

2. ModAdpte split Reuters-21578 dataset

Delle 90 categorie presenti nel dataset Reuters-21578 sono state selezionate le 10 più popolose in termini di documenti. Si riportano il numero di documenti di training e test per ognuna delle categorie selezionate:

- acq: train 1650, test 719
- corn: train 181, test 56
- crude: train 389, test 189
- earn: train 2877, test 1087
- grain: train 433, test 149
- interest: train 347, test 131
- money-fx: train 538, test 179
- ship: train 197, test 89
- trade: train 368, test 117
- wheat: train 212, test 71

3. Definizioni

Nel testo si farà uso delle seguenti definizioni:

- d : documento
- $|D|$: numero di documenti
- C : categoria
- V : vocabolario
- V_c : vocabolario della categoria
- w : parola
- N_{cw} : numero di documenti nella categoria c che contengono la parola w
- N_c : numero documenti nella categoria c

4. Bernoulli Naive Bayes

L'algoritmo ha richiesto la suddivisione del vocabolario precedentemente estratto in 10 dizionari contenenti le parole appartenenti a ciascuna categoria; analogamente si sono divisi i documenti per categoria di appartenenza. Si è quindi proceduto alla definizione di una stima di probabilità circa l'appartenenza a priori di un documento a una categoria:

$$P(d \in C) = \frac{|\{d : d \in C\}|}{|D|} \quad (1)$$

Si procede contando il numero di documenti per categoria in cui occorre una certa parola e si calcola la probabilità condizionata che data la presenza di una parola w nel documento d tale documento appartenga alla categoria C (sempre in un ottica frequentista):

$$P(d \in C | w \in V_c) = \frac{N_{cw} + 1}{N_c + 2} \quad (2)$$

Il processo di training consiste quindi nel calcolare, considerando l'intera collezione di documenti di training, tali probabilità rispettivamente per ogni categoria e per ogni coppia categoria-parola. Il processo di test consiste nel calcolare la probabilità condizionata che il documento stia in una o più categorie considerando quali parole compaiono o meno nel documento. Poichè gli eventi di occorrenza delle parole nei documenti sono considerati **indipendenti** gli uni dagli altri, la probabilità condizionata si semplifica in una produttoria delle probabilità dell'occorrenza delle singole parole.

4.1. Risultati e prestazioni

I risultati in termini di predizione ottenuti non risultano particolarmente efficaci, come si può evincere dal grafico.

$$t_{train} + t_{test} \approx 12 \text{ min}$$

I processi di training e test richiedono complessivamente circa 12 minuti.

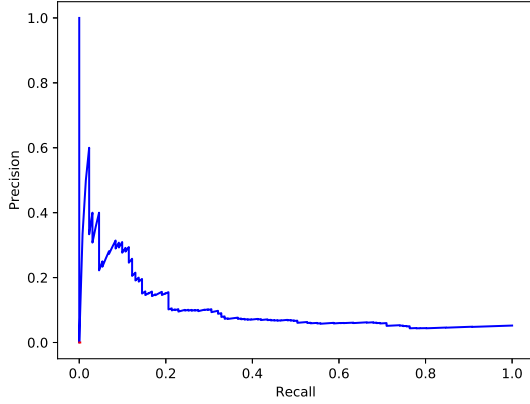


Figura 1. Precision-recall curve categoria interest

Categoria	F1 Measure
acq	0.13670886
corn	0.02985075
crude	0.28440367
earn	0.69054054
grain	0.36363636
interest	0.19314642
money-fx	0.18581907
ship	0.1512605
trade	0.0
wheat	0.0

5. Multinomial Naive Bayes

Multinomial Naive Bayes differisce dall'algoritmo precedentemente per il semplice fatto che non si limita a controllare se una parola è presente o meno nel documento, ma conta anche quante volte essa occorre nei documenti. Quindi, per quanto la stima di probabilità definita della sezione precedente rimanga valida, la probabilità condizionata risulta:

$$P(d \in C | \dots) = \frac{n_{w_i, c} + 1}{\sum_{w \in V} n_{w, c} + 1} \quad (3)$$

dove $n_{w_i, c}$ rappresenta il numero di occorrenze della parola w_i nella categoria c e $\sum_{w \in V} n_{w, c}$ tutti i tokens (ovvero le parole, con ripetizioni) nella categoria.

5.1. Risultati e prestazioni

Il conteggio delle occorrenze delle parole nei documenti, di contro alla variante di Bernoulli, comporta un aumento deciso dei tempi di esecuzione.

$$t_{train} \approx 8 \text{ min}, \quad t_{test} \approx 10 \text{ min}$$

Tuttavia all'aumento del tempo di esecuzione non sembra essere collegato un aumento della precisione di classificazione: osservando le curve di recall-precision non si osservano miglioramenti sostanziali e in alcuni casi un lieve peggioramento.

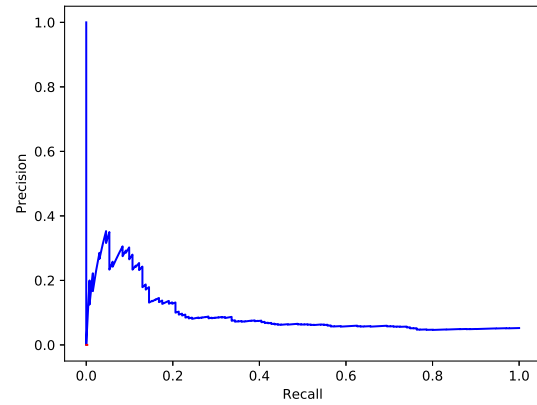


Figura 2. Multinomial NB precision-recall curve categoria interest

Categoria	F1 Measure
acq	0.68218527
corn	0.5
crude	0.0
earn	0.61764706
grain	0.94614265
interest	0.277777778
money-fx	0.0
ship	0.18987342
trade	0.0
wheat	0.0

6. Perceptron

Perceptron rappresenta una prima rappresentazione del concetto di neurone. Utilizza un approccio totalmente diverso dai precedenti algoritmi in quanto non sfrutta il teorema di Bayes bensì il concetto di regressione lineare. Occorre prima però dare un'adeguata rappresentazione ai documenti.

6.1. Rappresentazione documenti

Dopo aver proceduto all'estrazione del vocabolario, alla suddivisione dei documenti di training e test, e dall'ulteriore divisioni di tali documenti per categorie, per ogni documento viene estratta una sua rappresentazione TF-IDF

(Term Frequency-Inverse Document Frequency). Tale rappresentazione consiste nel pesare la rilevanza di un dato documento nei confronti di una certa parola prendendo in considerazione due fattori:

- 1) la frequenza di una parola nel documento (TF)
- 2) l'importanza generale di una parola nell'insieme dei documenti considerati (IDF)

Moltiplicando la IDF relativa a una parola w con la corrispondente TF nel documento d , si ottiene un vettore di lunghezza pari al numero di parole presenti nel vocabolario, in cui ogni elemento del vettore è il prodotto $tf[w] * idf[w]$.

6.1.1. Frequency. Nello specifico la frequenza della i -esima parola nel j -esimo documento d_j risulta

$$tf_{i,j} = \frac{n_{i,j}}{|d_j|} \quad (4)$$

dove $n_{i,j}$ rappresenta il numero di occorrenze della parola i -esima nel j -esimo documento.

6.1.2. Document Frequency. La IDF, come detto, misura l'importanza di una parola nell'insieme dei documenti considerati:

$$idf_i = \log \frac{|D|}{|\{d : w_i \in d\}|} \quad (5)$$

dove $|D|$ rappresenta il numero di documenti di training e il denominatore rappresenta il numero di documenti che contengono la i -esima parola w_i del vocabolario V .

6.1.3. TF-IDF. Di conseguenza la rappresentazione TF-IDF per la parola i -esima w_i e il documento j -esimo risulta

$$tfidf_{i,j} = tf_{i,j} \cdot idf_i \quad (6)$$

6.2. Algoritmo

L'algoritmo suddivide lo spazio dei documenti in due parti, positiva e negativa, tramite un iperpiano. Tale iperpiano è formato dai vettori dei pesi (*weights*) e dei *bias*. L'aggiornamento viene effettuato in base al risultato di

$$y_i(w \odot x) \quad (7)$$

dove

- 1) y_i ha valore in $\{-1, 1\}$ a seconda che il documento i -esimo appartenga o meno alla categoria esaminata
- 2) w è il vettore dei pesi corrente
- 3) x è la rappresentazione TF-IDF del documento considerato

Pesi e bias vengono quindi aggiornati se la classificazione commette un errore, cioè quando l'equazione (7) ha valore minore di zero. Per i dettagli delle regole di aggiornamento si fa riferimento a [1].

6.3. Risultati e tempi di esecuzione

Nell'osservare i tempi di esecuzione del Perceptron c'è da tenere conto che

$$t_{\text{perceptron}} = t_{\text{init}} + t_{\text{tfidf}} + t_{\text{train}} + t_{\text{test}}$$

dove

- t_{init} : tempo di inizializzazione (come Naive Bayes)
- t_{tfidf} : tempo per il calcolo della TF-IDF
- t_{train} : tempo per training
- t_{test} : tempo per il test

Il training richiede un tempo variabile a seconda del numero massimo di iterazioni impostogli, tuttavia effettuando 10 iterazioni, poichè $t_{\text{init}} \approx 1 \text{ min}$, $t_{\text{tfidf}} \approx 4 \text{ min}$, $t_{\text{train}} \approx 6 \text{ h}$ e $t_{\text{test}} \approx 1 \text{ min}$ il tempo complessivo risulta

$$t_{\text{perceptron}} \approx 6 \text{ h}$$

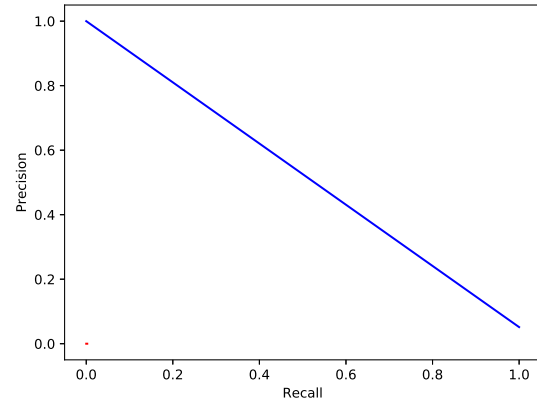


Figura 3. Perceptron precision-recall curve categoria interest

7. Conclusioni

Si è notato che gli algoritmi di Bernoulli e quello Multinomial hanno dei tempi d'esecuzione accettabili, ma un pessimo rapporto precision recall (Fig.1),(Fig.2). Nel caso del Perceptron si è riscontrato (Fig.3) un comportamento opposto vale a dire, dei lunghi tempi d'esecuzione a favore di buone precision recall.

Riferimenti bibliografici

- [1] N.Cristianini and J.Shawe-Taylor, An Introduction to Support Vector Machines and Other Kernel-based Learning Methods, Cambridge University Press, 2002
- [2] <https://nlp.stanford.edu/IR-book/html/htmledition/naive-bayes-text-classification-1.html>
- [3] <https://nlp.stanford.edu/IR-book/html/htmledition/the-bernoulli-model-1.html>
- [4] <https://github.com/HJuls2/perceptron-naive-bayes-text-classification>