

UNIVERSITÀ DEGLI STUDI DI
MILANO-BICOCCA

ADVANCED MACHINE LEARNING
FINAL PROJECT

Mercari Price Suggestion Challenge

Authors:

Gabriele Ferrario - 817518 - g.ferrario@campus.unimib.it

Riccardo Pozzi - 807857 - r.pozzi@campus.unimib.it

13 gennaio 2021



Sommario

The ABSTRACT is not a part of the body of the report itself. Rather, the abstract is a brief summary of the report contents that is often separately circulated so potential readers can decide whether to read the report. The abstract should very concisely summarize the whole report: why it was written, what was discovered or developed, and what is claimed to be the significance of the effort. The abstract does not include figures or tables, and only the most significant numerical values or results should be given.

1 Introduction

Il progetto trae origine dalla *Kaggle challenge Mercari Price Suggestion Challenge*[?] aperta a fine Novembre 2017 che come viene reso chiaro dal sottotitolo "*Can you automatically suggest product prices to online sellers?*" pone l'obiettivo di stimare più precisamente possibile il prezzo di determinati prodotti a partire da alcune loro caratteristiche.

Alla base di ciò vi è l'esigenza dell'*ecommerce Mercari*[?] di offrire ai propri venditori un suggerimento sul prezzo di vendita dei prodotti inseriti.

Si tratta quindi di un problema di *regressione* che a partire dalle varie caratteristiche dei prodotti, testuali e non, vuole calcolare il prezzo da suggerire.

Durante lo svolgimento del progetto si valuteranno vari approcci al problema, soprattutto per quanto riguarda i dati di tipo testuale, soffermandosi sulle performance di regressione sia in termini di errore rispetto ai dati di *train* che di costi computazionali.

2 Datasets

Il Dataset consiste in un elenco di 1391082 prodotti descritti tramite le seguenti caratteristiche:

2.0.1 Price

Price rappresenta il prezzo per il quale l'articolo è stato venduto (variabile target). Il prezzo medio è di circa \$26, con un valore minimo pari a \$0 e un valore massimo pari a \$2009; inoltre, presenta una deviazione standard di

circa \$38. Analizzando i percentili ci si accorge che i prezzi sono relativamente bassi in quanto il 75% dei prodotti hanno un prezzo al di sotto di \$29.

2.0.2 Train id

Train_id rappresenta l'identificativo del prodotto nell'elenco.

2.0.3 Name

Name è il nome del prodotto sotto forma di dato non strutturato.

2.0.4 Shipping

Shipping è caratterizzato dal valore 1 se la tassa di spedizione è a carico del venditore, altrimenti 0 se è a carico dell'acquirente. Questo attributo è decentemente ripartito tra i venditori e gli acquirenti, in quanto il 55% dei prodotti prevede un valore di 0. Analizzando il prezzo degli articoli ci si aspetta che quelli i cui costi di spedizione sono a carico del venditore avranno un prezzo più alto. Tuttavia, ci sono una serie di fattori contrastanti. Questo può essere vero all'interno di specifiche categorie di prodotti e condizioni degli articoli, ma non quando si confrontano gli articoli sul totale. Infatti, il prezzo medio pagato dagli utenti che devono pagare le spese di spedizione (circa \$30) è superiore a quello dei prodotti la cui spedizione è pagata dal venditore (circa \$22).

2.0.5 Item condition

Item_condition_id rappresenta lo stato del prodotto fornito dal venditore, questo valore varia da 1 a 5. Il valore più frequente è 1, mentre 4 e 5 sono i più rari. Nei dati non è presente una descrizione dettagliata sul significato di questi valori, analizzando il dataset si può supporre che il valore 1 identifichi la condizione migliore poichè è la più frequente, mentre il valore 5 la condizione peggiore. Tuttavia calcolando i prezzi medi di vendita per ogni condizione non si riesce ad arrivare a una conclusione sicura poichè la condizione 5 è quella con il prezzo medio più alto, mentre la condizione 4 è quella con il prezzo medio più basso e le restanti categorie presentano un prezzo medio molto vicino.

2.0.6 Category Name

Category_name rappresenta la categoria di prodotto a cui appartiene l'articolo. Nel dataset sono presenti 1287 categorie univoche e tra ognuna di esse si vede una categoria principale/generale, seguita da due o più sottocategorie più specifiche (ad esempio: Women/Tops & Blouses/T-Shirts). Inoltre, ci sono 6327 articoli che non hanno una categoria assegnata. Infine, analizzando le dieci categorie più popolari, si nota che l'abbigliamento femminile è molto popolare su Mercari. Infatti, di queste prime dieci categorie 5 sono di abbigliamento femminile; Anche il trucco e l'elettronica sono categorie molto quotate.

2.0.7 Brand Name

Brand_name rappresenta il marchio dell'articolo; nel dataset sono presenti 4809 valori differenti e 632682 valori mancanti.

2.0.8 Item Description

Item_description rappresenta la descrizione del prodotto sotto forma di dato non strutturato. Nel dataset sono presenti 4 istanze senza descrizione e 82494 descrizioni con la stringa "no description yet". Inoltre, non esiste una correlazione tra la lunghezza delle descrizioni e il prezzo, in quanto c'è una correlazione di 0.048. Analizzando le word cloud ottenute tramite i bigrammi delle descrizioni dei prodotti suddivisi in quattro fasce di prezzo: prezzo ≥ 100 (Figura 1), $50 < \text{prezzo} < 100$ (Figura 2), $30 < \text{prezzo} \leq 50$ (Figura 3) e prezzo ≤ 30 (Figura 4); si riescono a notare delle differenze sulle parole più frequenti. Infatti, nella wordcloud di Figura sono molto frequenti bigrammi che danno informazioni sulle buone condizioni dei prodotti, ad esempio: 100 authentic, great condition e good condition. Diminuendo di prezzo questi bigrammi diventano meno frequenti, ma aumentano i bigrammi relativi a descrizioni mancanti.



Figura 1: Word Cloud contenente i bigrammi ottenuti dalle descrizioni dei prodotti con prezzo maggiore o uguale a 100



Figura 2: Word Cloud contenente i bigrammi ottenuti dalle descrizioni dei prodotti con prezzo maggiore di 50 e minore di 100



Figura 3: Word Cloud contenente i bigrammi ottenuti dalle descrizioni dei prodotti con prezzo maggiore di 30 e minore o uguale a 50



Figura 4: Word Cloud contenente i bigrammi ottenuti dalle descrizioni dei prodotti con prezzo minore o uguale a 30

2.0.9 Pulizia Dataset

Dal dataset sono stati eliminati tutti i prodotti con prezzi minori di cinque e maggiori di 2000 poichè sul sito ufficiale di Mercari è specificato che i prezzi possono essere impostati solo nell'intervallo [5,2000] (fonte: https://www.mercari.com/us/help_center/article/69). Durante la fase di analisi si è scoperta la presenza di valori mancanti nei campi: item_description, brand_name e category_name, questi valori sono stati rimpiazzati con il valore NA. Inoltre, è stato effettuato un trattamento dei dati (cito paper sul text preprocessing) non strutturati convertendoli tutti in minuscolo, sono state sostituite le descrizioni mancanti con il valore NA (anche quando era presente la stringa "no description yet"). Sui dati testuali è stata effettuata una fase di lemmatizzazione (perchè usa vocabolario su cui tagliare le parole e cito

articolo). Successivamente, sulle nuove parole ottenute è stata effettuata una fase di pulizia di questi campi non strutturati eliminando le stopwords, la punteggiatura, tutti i caratteri di lunghezza pari a 1 che non sono numeri (è stato deciso di mantenere tutti i numeri poichè molte descrizioni senza di essi perdono di significato) e sono state eliminate le emoji.

I valori del campo `category_name` sono stati codificati in interi tramite la tecnica label encoding.

3 The Methodological Approach

3.1 Rappresentazione del Testo

3.1.1 Bag of Words

Il modello bag-of-words è una rappresentazione semplificata di testi, dove ciascuno di essi è rappresentato come una "borsa" (vettore) delle sue parole. Quindi crea un vettore per ogni documento contenente il conteggio delle occorrenze di ciascuna parola del vocabolario nel documento. Questa tecnica ignora la grammatica e persino l'ordine delle parole, ma mantiene la molteplicità.

3.1.2 Tf-Idf

Tf-Idf (Term frequency-Inverse document frequency) [?] risolve uno dei principali problemi che si hanno con Bag of Words, poiché quest'ultima tecnica non considera relazioni tra i vari documenti. Infatti, Tf-Idf (1) considera l'importanza di una parola per un documento in una collezione di documenti (corpus). Il vettore costruito per ogni documento tramite questa tecnica considera per ciascun termine del documento la frequenza con cui il termine compare in questo documento (Term Frequency, TF (2)) rispetto alla frequenza del termine in tutti i documenti (Inverse Document Frequency, IDF (3)).

$$Tf - Idf_{t,d} = tf_{t,d} \cdot idf_t \quad (1)$$

$$tf_{t,d} = \frac{n_{t,d}}{\text{number of terms in the document } d} \quad (2)$$

$$idf_t = \log \frac{\text{number of documents}}{\text{number of documents with term } t} \quad (3)$$

Dove in (1)(2)(3) t indica il termine e d indica il documento.

3.1.3 Word Embeddings

Le precedenti tecniche sono semplici da realizzare, robuste e funzionali per molti tasks. Tuttavia queste tecniche semplici sono al limite in molti compiti perché trattano le parole come unità atomiche, senza alcuna correlazione poiché vengono rappresentate come indici di un vocabolario. Con lo sviluppo delle tecniche di Machine Learning sono nati modelli in grado di apprendere rappresentazioni vettoriali di alta qualità di parole provenienti da set di dati composti da miliardi di parole e con vocabolari di milioni di parole. Queste nuove tecniche generano rappresentazioni vettoriali delle parole che rispettano la similarità tra le parole considerando più gradi di somiglianza, rappresentando la regolarità semantica e sintattica, permettendo anche operazioni algebriche [?].

3.1.4 GLOVE

GLOVE [?] è un algoritmo di apprendimento non supervisionato che fornisce una rappresentazione vettoriale di ogni parola, considerando le statistiche globali di ricorrenza parola per parola (word-word ?????) da un corpus e le rappresentazioni ottenute mostrano sottostrutture lineari dello spazio vettoriale delle parole. Inoltre, gli ideatori di questa tecnica hanno reso disponibile i word embeddings pre-addestrati. Nel dettaglio in questo lavoro sono stati utilizzati i seguenti embeddings:

- Wikipedia 2014 + Gigaword 5: allenato su 6 miliardi di parole e con un vocabolario di 400 mila parole;
- Common Crawl: allenato su 840 miliardi di parole e con un vocabolario di 2.2 milioni di parole (il più grande disponibile di GLOVE);

3.1.5 Training e valutazione

3.1.6 Dataset split

Al fine di valutare i vari modelli più correttamente possibile il dataset è stato suddiviso in *training-set*, *validation-set* e *test-set*. Quindi il training è stato usato per l'allenamento, il validation per accertarsi che il modello non sia propenso all'overfitting ed il test per la valutazione finale.

3.1.7 Valutazione dell'errore

Per quanto riguarda le performance in termini di errore la funzione di *loss* utilizzata in fase di training è il *Root Mean Squared Logarithmic Error (RM-SLE)*, in quanto è la misura scelta dalla Kaggle challenge per confrontare le performance dei vari partecipanti. Inoltre essa risulta adeguata al problema considerando il vasto intervallo dei valori dei prezzi.

Di seguito la definizione e alcune osservazioni su di RMSLE insieme alle ulteriori metriche utilizzate; ad esempio il *Mean Absolute Error (MAE)*, che fornisce una più immediata comprensione rispetto al RMSLE.

Mean Absolute Error (MAE):

$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (4)$$

Mean Squared Error (MSE):

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (5)$$

Il quadrato fornisce un peso maggiore agli errori dei valori più elevati[?].

Root Mean Squared Error (RMSE):

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (6)$$

La radice dell'MSE, nella stessa scala dei valori.

Mean Squared Logarithmic Error (MSLE):

$$\frac{1}{n} \sum_{i=1}^n (\log(y_i + 1) - \log(\hat{y}_i + 1))^2 = \frac{1}{n} \sum_{i=1}^n \log^2\left(\frac{y_i + 1}{\hat{y}_i + 1}\right) \quad (7)$$

Essendo calcolato a partire da un rapporto riflette l'errore relativo e di conseguenza risulta efficace laddove i valore assoluti presentano variazioni considerevoli.

Root Mean Squared Logarithmic Error (RMSLE):

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (\log(y_i + 1) - \log(\hat{y}_i + 1))^2} \quad (8)$$

La radice dell'MSLE.

3.1.8 Valutazione dei costi computazionali

Il numero dei parametri allenabili dei vari modelli sono stati annotati insieme al tempo richiesto per l'esecuzione di 10 epoche di train sulla piattaforma *Google Colab* abilitando l'accelerazione hardware GPU.

3.2 Categorie

Dire qualcosa sulle categoriche

Inizialmente sono state valutate le performance di regressione dei prezzi a partire dalle sole features categoriche in modo tale da stabilire un punto di partenza per la successiva analisi delle features testuali *name* e *item_description*, più complesse e computazionalmente costose.

Un semplice modello composto da due layer Densi è stato sufficiente a raggiungere le seguenti performance alla decima epoca con l'aiuto di un layer Dropout interposto ad essi per contrastare l'overfitting.

3.3 modello

qualche spiegazione sulla definizione del modello. perchè le lstm, cosa sono, paper, lstm per il testo.

performance gain delle lstm?? perchè non abbiamo provato le gru?

3.4 Rappresentazione del Testo

Descrizione dei vari approcci nella storia (quando sono stati creati)
deep learning e non

3.4.1 clean o raw

abbiamo valutato clean e raw; quali modelli beneficiano del cleanup, quali no? perchè? il deep learning impara meglio se il testo è raw?

- bag of word based models: count vectorizer - tf/idf
- embeddings: keras, glove pretrained. dire perchè non abbiamo provato qualcosa tipo word2vec?
- transformers con fine tuning

4 Results and Evaluation

The Results section is dedicated to presenting the actual results (i.e. measured and calculated quantities), not to discussing their meaning or interpretation. The results should be summarized using appropriate Tables and Figures (graphs or schematics). Every Figure and Table should have a legend that describes concisely what is contained or shown. Figure legends go below the figure, table legends above the table. Throughout the report, but especially in this section, pay attention to reporting numbers with an appropriate number of significant figures.

5 Discussion

The discussion section aims at interpreting the results in light of the project's objectives. The most important goal of this section is to interpret the results so that the reader is informed of the insight or answers that the results provide. This section should also present an evaluation of the particular approach taken by the group. For example: Based on the results, how could the experimental procedure be improved? What additional, future work may be warranted? What recommendations can be drawn?

6 Conclusions

Conclusions should summarize the central points made in the Discussion section, reinforcing for the reader the value and implications of the work.

If the results were not definitive, specific future work that may be needed can be (briefly) described. The conclusions should never contain “surprises”. Therefore, any conclusions should be based on observations and data already discussed. It is considered extremely bad form to introduce new data in the conclusions.