

Real-Time Brand Sentiment Analysis of Social Media Text Content

Proposal

Jonathan Agustin

Fernando Calderon

Juliet Lawton



University of San Diego®

Abstract

Businesses increasingly leverage social media platforms to track brand perception using Natural Language Processing (NLP), Classification, and Sentiment Analysis. Machine Learning (ML) models are trained with a large corpus of social media text data using Deep Learning to classify a brand and categorize sentiment towards the brand as positive, negative, or neutral. However, existing systems are either inaccessible or limited in their ability to provide real-time insights into public sentiment about a brand. We introduce an automated system that provides real-time insight on brand sentiment, enabling quick reaction to shifts in public sentiment. The proposed system enhances the ability to maintain a positive brand image and stay ahead of the competition.

Problem

Given a set B of strings representing brand names $B = \{b_1 = \text{"Nike"}, b_2 = \text{"Google"}, b_3 = \text{"Disney"}, \dots, b_i\}$ and a set P of strings representing social media posts

$$\begin{aligned} p_1 &= \text{"Celebrating my birthday"} \\ p_2 &= \text{"I could care less about the new Air Force 1s"} \\ p_3 &= \text{"Happiest place on Earth"} \\ &\vdots \\ p_j &\in P \end{aligned}$$

the goal is to build a $\text{BrandClassifier}(p)$ that predicts whether a social media post p expresses a brand, and returns a brand name b when the post expresses the brand, or No-Brand when the post does not express any brand. The model also returns a confidence score $c \in [0, 1]$ that indicates how confident the model is in its prediction.

$$\text{BrandClassifier}(p) \rightarrow \begin{cases} (b, c) & \text{if } p \text{ expresses a brand } b \in B \\ (\text{No-Brand}, c) & \text{otherwise} \end{cases}$$

We build another model $\text{BrandSentimentAnalyzer}(p, b)$ that predicts the sentiment of a social media post towards a brand, and returns a sentiment label

$$s \in S = \{s^+ = \text{Positive}, s^- = \text{Negative}, s^0 = \text{Neutral}, s^? = \text{Mixed}\}$$

when the post expresses sentiment towards the brand, or No-Sentiment when it does not express any sentiment. The model also returns a confidence score $c \in [0, 1]$.

$$\text{BrandSentimentAnalyzer}(p, b) \rightarrow \begin{cases} (s, c) & \text{if } p \text{ expresses sentiment } s \in S \text{ about brand } b \in B \\ (\text{No-Sentiment}, c) & \text{otherwise} \end{cases}$$

Example. Given the following social media post $p_3 = \text{"I could care less about the new Air Force 1s"}$ and $b_3 = \text{"Nike"}$

$$\begin{aligned} \text{BrandClassifier}(p_3) &\rightarrow (b_3, 0.68) \\ \text{BrandSentimentAnalyzer}(p_3, b_3) &\rightarrow (s^-, 0.82) \end{aligned}$$

Methodology

The work is divided into four stages: assessment, development, deployment, and evaluation. Each stage informs and refines previous stages, creating a feedback loop that continuously improves the overall system.

1. **Assessment.** The stage involves selecting and testing appropriate models and datasets. The specific requirements of the sentiment analysis task, the characteristics of the available data, and the nature of the problem guide this selection process. The chosen models must effectively handle the unstructured and noisy data typical of social media and accurately capture and predict sentiments. The selected datasets must accurately represent the social media environment and provide a diverse view of expressed sentiments.
2. **Development.** This stage preprocesses data in real-time to simulate a live data stream and training the selected models. The unstructured and noisy nature of social media data necessitates the transformation of text into a structured format suitable for modeling. This process involves handling invalid values and discarding non-textual content such as images, videos, or audio. Automation ensures efficiency in this process. The selected models are then adapted to the preprocessed dataset and trained to identify a brand or align with the sentiment expressed in the social media text.
3. **Deployment.** This stage deploys the models to process social media data in near real-time. This stage requires building infrastructure to support the models and integrating the models with social media platforms such as Twitter where data will be streamed.
4. **Evaluation.** This stage measures the models' ability to identify brands and predict sentiment in social media. For a time period, social media posts are collected and manually labeled for brand and sentiment. Standard classification metrics such as Accuracy, Precision, Recall, and F1 Score are used to evaluate the models. The models are then ranked by performance.

Expected Behaviors & Outcomes

We expect our models to demonstrate a nuanced understanding of the English language, accurately identifying brands and sentiments. Despite the inherent challenges associated with sentiment analysis and brand identification in text data, we anticipate that our models will achieve an accuracy rate exceeding 50%. We foresee our models to have real-world applicability, serving as a valuable tool for businesses to accurately analyze sentiment from social media data. The proposed system is expected to provide actionable insights that enable companies to make informed decisions.

Limitations. The model's primary limitation is its unimodal nature. It processes only text data, ignoring non-textual elements like images, videos, and audio clips. Often found in social media posts, these elements can contain significant sentiment-related information. This information can support, contradict, or add nuance to the sentiment expressed in the text. Audio clips, for example, can provide sentiment-related cues through tone, pitch, and other features. The proposed system, however, will identify this multimodal information and discard the input completely. Future work should consider expanding the model to a multimodal framework to provide a more comprehensive sentiment analysis of social media posts.

References

- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805*.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv preprint arXiv:1907.11692*.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language Models are Unsupervised Multitask Learners. *OpenAI Blog*, 1(8).
- Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2019). DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- Hutto, C.J., & Gilbert, E. (2014). VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text. *Eighth International Conference on Weblogs and Social Media (ICWSM-14)*. Ann Arbor, MI, June 2014.