

Contents

Abstract	3
1 Introduction	5
1.1 Sentiment analysis and its applications	5
1.2 Definitions	7
2 Classification Techniques	8
2.1 k Nearest Neighbors	8
2.2 Naive Bayes	10
2.3 Logistic regression	15
3 Feature Extraction	17
3.1 Motivation	17
3.2 Important Words Selection	19
3.3 Context Mining	20
4 Effectiveness of trained models	21
Conclusion	22
References	23

Abstract

In this paper I would like to investigate various methods of text analysis to improve quality of predictions based on text messages. Text analysis, which has other names, in particular, text mining and text analytics, is continuously developing. Labor-intensive manual text mining approaches first surfaced in the mid-1980s, but technological advances have enabled the field to advance during the past decade. Text mining usually involves the process of structuring the input text (usually parsing, along with the addition of some derived linguistic features and the removal of others, and subsequent insertion into a database), deriving patterns within the structured data, and finally evaluation and interpretation of the output. Typical text mining tasks include text categorization, text clustering, concept/entity extraction, production of granular taxonomies, sentiment analysis, document summarization, and a plethora of others. The term text analytics describes a set of linguistic, statistical, and machine learning techniques that model and structure the information content of textual sources for business intelligence, exploratory data analysis, research, or investigation.

Sentiment analysis, which refers to the use of natural language processing, text analysis and computational linguistics to identify and extract subjective information in source materials, is the most crucial problem for businesses related to social media and market advertisement. An average client may produce about 100 messages per day. What is more, sometimes sentiment data is generated by web crawlers which aggregate data from various sources. Given the amount of data those enterprises poses, it is almost infeasible to process all the messages manually. Therefore, automatic text analytics is indispensable. Even those methods with low precision can generate profits.

A basic task in sentiment analysis is classifying the polarity of a given text at the document, sentence, or feature/aspect level - whether the expressed opinion

in a document, a sentence or an entity feature/aspect is positive, negative, or neutral. Advanced, "beyond polarity" sentiment classification looks, for instance, at emotional states such as "angry," "sad," and "happy." In this work only the basic task is considered.

Existing approaches to sentiment analysis can be grouped into three main categories: knowledge-based techniques, statistical methods, and hybrid approaches. Knowledge-based techniques classify text by affect categories based on the presence of unambiguous affect words such as happy, sad, afraid, and bored. Some knowledge bases not only list obvious affect words, but also assign arbitrary words a probable "affinity" to particular emotions. Statistical methods leverage on elements from machine learning such as latent semantic analysis, support vector machines, "bag of words" and Semantic Orientation - Pointwise Mutual Information.

In this paper I consider using text analysis and statistical methods to perform sentiment analysis of the data referred to specific domain. I leverage different common approaches and suggest improvements to them for solving the particular problem I work on. Moreover, I analyze some of the common problems like feature vector selection. Finally, I provide a performance comparison of the detailed algorithms in practice.

Chapter 1

Introduction

1.1 Sentiment analysis and its applications

A marketing campaign is efforts of a company or a third-party marketing company to increase awareness for a particular product or service, or to increase consumer awareness of a business or organization. It has a limited duration. A marketing campaign consists of specific activities designed to promote a product, service or business. A marketing campaign is a coordinated series of steps that can include promotion of a product through different mediums (television, radio, print, online) using a variety of different types of advertisements. The campaign doesn't have to rely solely on advertising, and can also include demonstrations, word of mouth and other interactive techniques.

After a particular campaign is finished, the business is interested in aggregating people's opinion about the products which were being promoted. The vast majority of the reviews come from social media and shopping websites. Basically a review is a short text which is either positive or negative. As was mentioned before, because of large quantities of customers and review sources, manual review analysis is not feasible. The solution is to use an automated review analysis. A classifier can be built on the data that have been already analysed by humans. The larger the corpus the better.

In general, the problem of sentiment analysis depends heavily on the application's domain and can benefit from additional metadata available along with the text data. For example, in case of reviews related to political domain, the analyst might expect some amount of sarcasm and indirect implications. More-

over, a language which is used also plays important role; similar expressions in different languages have different shades of meaning. Therefore, it often is beneficial to target a specific domain and a specific language in text sentiment analysis problem.

In this paper I am going to target customer reviews of baby products in English available on amazon website, the biggest retailer in the world. The data is available at https://s3.amazonaws.com/static.dato.com/files/coursera/course-3/amazon_baby.csv.zip. This data contains 183532 labeled products reviews. The polarity grade is 1-5. 1 - means the most negative review and 5 - the most positive.

1.2 Definitions

Let us denote a set of all possible words W and a set of all possible permutations G . And let R be a set of all the reviews available. So $r \in R$ - product review. $R = \{(w, g) : w \in W, g \in G\}$. P - review polarities set. Each review has corresponding polarity $p = \{1, 2, 3, 4, 5\}, p \in P$. $|R| = |P| = m$.

X feature matrix for machine learning. X is an $m \times n$ matrix, where m is a number of training samples and n is a feature vector cardinality.

Y - label vector for machine learning. In this case it is a mapped polarity of a review. $Y = \{-1, 1\}$, where -1 correspond to negative opinion and 1 - to the positive. $|Y| = m$.

$X_{train} \subseteq X$ - training data.

$X_{test} \subseteq X$ - test data.

$Y_{train} \subseteq Y$ - training labels.

$Y_{test} \subseteq Y$ - test labels.

$extr : R \rightarrow X$ - feature extraction function. Given a review r , this functions maps it to the feature vector x .

$pol : P \rightarrow Y$ - polarity map function. Maps polarity to binary vector.

$f : X \rightarrow Y$ - classification function. Given a feature vector x , this function predicts its label y .

Now text analysis process can be defined. The sets R and P are available. Firstly, an extraction function $extr$ must be defined. After feature extraction process we should be able to form a feature matrix X . Also the representation of the polarities set P is changed by applying a map function pol . Hence, X and Y are defined.

At this moment, I shall split the data into training and cross-validation set. Or in other words training and test data. Now f must be found by a classification algorithm on the training set.

The ultimate goal is to obtain a decision function $F : R \rightarrow \{extr, f\}$ that would minimize $\frac{1}{m_{test}} \sum_{i=1}^{m_{test}} (f(x_i) - y_i)^2$ for training data.

TODO precession and recall

Chapter 2

Classification Techniques

2.1 k Nearest Neighbors

k Nearest Neighbors (KNN) is part of supervised learning that has been used in many applications in the field of data mining, statistical pattern recognition, image processing and many others. The idea is to predict class for a given feature vector using so called majority voting. Suppose we want to predict label for a feature vector $x_i \in X_{test}$. Then we look on labels of k training vectors which are "the closest" to x_i and use the most frequent label as a prediction. To define "closeness" some distance metric is used. There are several options for distance functions.

A commonly used distance metric for continuous and discrete variables is Euclidean distance:

$$\sum_{i=1}^n (a_i - b_i)^2$$

, where a, b - feature vectors. $a_i - b_i$ - vector difference.

For discrete variables, such as for text search, another metric can be used, such as the overlap metric (or Hamming distance).

The best choice of k depends upon the data; generally, larger values of k reduce the effect of noise on the classification, but make boundaries between classes less distinct. A good k can be selected by various heuristic techniques. In our case though I am going to use a fixed value.

This algorithm does not require training step. Having sets X_{train} and Y_{train} and distance function $dist$ is enough for making predictions. On the other hand, the

classification of test data takes time, because for every review without label, the whole training corpus must be processed to calculate distance. For one prediction expected time complexity is $O(m_{train} * n)$. Performing indexing on the training set can decrease the running closer to $O(n)$. However, indexing is a computationally heavy operation and in case data is updated often, this approach becomes infeasible. Moreover, precision of this classification algorithm is usually low compared to others. Another drawback of the basic "majority voting" classification occurs when the class distribution is skewed. That is, examples of a more frequent class tend to dominate the prediction of the new example, because they tend to be common among the k nearest neighbors due to their large number.

In practice KNN algorithm performs worse than other algorithms. Therefore KNN model is mostly used as a baseline to compare it with other models.

2.2 Naive Bayes

Naive Bayes algorithm is a classification algorithm based on Bayes' theorems, and can be used for both exploratory and predictive modeling. The word naive in the name Naive Bayes derives from the fact that the algorithm uses Bayesian techniques but does not take into account dependencies that may exist. This algorithm is less computationally intense than most other classification algorithms, and therefore is useful for quickly generating mining models to discover relationships between input columns and predictable columns

Naive Bayes relies on very simple representation of the document - Bag of words.

Input:

- A document d
- A fixed set of classes $C = \{c_1, c_2, \dots, c_j\}$
- A training set of m hand-labeled documents $(d_1, c_1), (d_2, c_2), \dots, (d_m, c_m)$

Output:

- A learned classifier $y : d \rightarrow c$.

According to the Bayes theorem, for a document d and class c :

$$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

.

This probability can be used to predict a class of a document.

$$c_{best} = \operatorname{argmax}_{c \in C} P(c|d) = \operatorname{argmax}_{c \in C} \frac{P(d|c)P(c)}{P(d)} = \operatorname{argmax}_{c \in C} P(d|c)P(c)$$

Taking into account the fact that d is a bag (set) of words:

$$c_{best} = \operatorname{argmax}_{c \in C} P(d|c)P(c) = \operatorname{argmax}_{c \in C} P(w_1, w_2, \dots, w_n|c)P(c)$$

Now let us find the way to compute two probabilities in the equation above. Let us start with $P(c)$. This probability can be estimated by computing how often this class occurs. It can be computed by counting relative frequencies in a corpus.

Computing $P(w_1, w_2, \dots, w_n|c)$ is more subtle. There is $|W|^n * |C|$ parameters to be computed, which is impossible in practice. In Naive Bayes classifier two simplifying assumptions are made:

- **Bag of Words assumption:** Assume position does not matter
- **Conditional Independence:** Assume the feature probabilities $P(x_i|c)$ are independent given a class c .

Despite these simplification, in practice a problem can be solved with a high degree of accuracy. The result of the simplifying assumptions is the following equation:

$$P(w_1, w_2, \dots, w_n|c) = P(w_1|c) \times P(w_2|c) \times \dots \times P(w_n|c)$$

Therefore:

$$c_{best} = \operatorname{argmax}_{c \in C} P(w_1, w_2, \dots, w_n|c)P(c) = \operatorname{argmax}_{c \in C} P(c) \prod_{w \in W} P(w|c)$$

Now we have to find the way to compute these probabilities. The most natural way to implement this is simply to use frequencies in the data.

$$P(w_i|c) = \frac{\text{doccount}(C = c)}{N_{doc}}$$

$$P(c) = \frac{\text{count}(w_i, c)}{\sum_{w \in W} \text{count}(w, c)}$$

So I am going to compute the fraction of times word w_i appears among all words in documents of class c . To do that it is possible to concatenate all available documents of a particular class c into one mega-document. In case of this paper there is two mega-documents D_{pos} and D_{neg} , which correspond to positive and negative reviews.

Unfortunately, there is a problem with the approach described above. What should happen if we have not seen training example with the word **fantastic** and classified in the positive?

$$P(\text{"fantastic"} | \text{positive}) = \frac{\text{count}(\text{"fantastic"}, \text{positive})}{\sum_{w \in W} \text{count}(w, \text{positive})} = 0$$

As can be seen, even though the word **fantastic** is usually associated with positive sentiment, the probability of a document which contains this word belongs to D_{pos} is zero. Here is why:

$$c_{best} = \operatorname{argmax}_{c \in C} P(c) P(\text{"fantastic"} | c) \prod_{w \in W / \{\text{"fantastic"}\}} P(w | c)$$

For positive class this product is equal to zero.

The solution to this problem is adding Laplace (add-1) smoothing for Naive Bayes. Simply one is added to each of the counts:

$$P(w_i | c) = \frac{\text{count}(w_i, c) + 1}{(\sum_{w \in W} (\text{count}(w, c) + 1))} = \frac{\text{count}(w_i, c) + 1}{\sum_{w \in W} \text{count}(w, c) + |W|}$$

Now the algorithm can be defined:

- From training corpus, extract Vocabulary - the list of words.
- Calculate $P(c)$ terms

For each $c_j \in C$ do

$docs_j \leftarrow$ all documents with class c_j

$$P(w_i | c_j) = \frac{|docs_j|}{|total \# documents|}$$

- Calculate $P(w_k | c_j)$ terms

$Text_j \leftarrow$ single document containing all $docs_j$

For each $w_k \in W$

$n_k \leftarrow$ number of occurrences of w_k in $Text_j$

n <- total number of words in the class j .

$$P(w_k|c_j) = \frac{n_k + \alpha}{n + \alpha \times |W|}$$

Note: In the last equation α is used instead of one to generalize the model.

It turns out Naive Bayes has a very close relationship to language modeling. Each class in Naive Bayes classifier is a unigram language model. [TODO link on unigram language model.](#)

Probaility of word: $P(word|c)$

Probability of sentence: $P(s|c) = \prod P(word|c)$

Let us see how it works. Imagine that we have a class positive: And we have table of likelihoods of each word.

Likelihood	Word
0.05	I
0.2	love
0.002	this
0.05	beautiful
0.03	dress

Then the probability of sentence "I love this fun film" will be $P(s|pos) = 0.00000003$.

A model negative can look like:

Likelihood	Word
0.05	I
0.002	love
0.002	this
0.0001	beautiful
0.03	dress

Then the probability of sentence "I love this fun film" will be $P(s|neg) = 0.00000000000006$.

As we can see $P(s|pos) > P(s|neg)$. Therefore, this sentence is classified as positive, which is correct. Hence, Naive Bayes can be viewed as a unigram language model conditioned on the class.

Summary:

- Very fast algorithm, low storage requirements
- Robust to irrelevant features: Irrelevant features cancel each other without affecting results.
- Very good at domains with many equally important features
- Optimal if independence assumption is true. Of course it is rare.

2.3 Logistic regression

Linear regression can easily be used for classification in domains with numeric attributes. Indeed, we can use any regression technique, whether linear or nonlinear, for classification. The trick is to perform a regression for each class, setting the output equal to 1 for training instances that belong to the class and 0 for those that do not. The result is a linear expression for the class. Then, given a test example of unknown class, calculate the value of each linear expression and choose the one that is largest. This scheme is sometimes called multiresponse linear regression.

One way of looking at multiresponse linear regression is to imagine that it approximates a numeric membership function for each class. The membership function is 1 for instances that belong to that class and 0 for other instances. Given a new instance, we calculate its membership for each class and select the biggest.

Multiresponse linear regression often yields good results in practice. However, it has two drawbacks. First, the membership values it produces are not proper probabilities because they can fall outside the range 0 to 1. Second, least-squares regression assumes that the errors are not only statistically independent but are also normally distributed with the same standard deviation, an assumption that is blatantly violated when the method is applied to classification problems because the observations only ever take on the values 0 and 1.

A related statistical technique called logistic regression does not suffer from these problems. Instead of approximating the 0 and 1 values directly, thereby risking illegitimate probability values when the target is overshoot, logistic regression builds a linear model based on a transformed target variable

In logistic regression the probability of classes is estimated in the following way:

$$P(1|w_1, w_2, \dots, w_n) = \frac{1}{1 + e^{(-a_0 - a_1w_1 - \dots - a_nw_n)}}$$

- where a_i is weight.

$$P(0|w_1, w_2, \dots, w_n) = 1 - P(1|w_1, w_2, \dots, w_n)$$

The goal function, which must be minimized, is:

$$\sum_i^m (1 - x_i) \log(1 - P[1|w_1, w_2, \dots, w_n]) + x_i \log(P[1|w_1, w_2, \dots, w_n])$$

In comparison with Naive Bayes algorithm Logistic regression is better at Domains with many non-equally important features. Moreover, Logistic regression is as well robust to irrelevant features when regularization is properly set

There are also drawbacks, for instance, Logistic regression is really fast algorithm, still Naive Bayes is faster.

Chapter 3

Feature Extraction

3.1 Motivation

Consider the simple case of text classification based on the presence or absence of just one word $w \in W$. Suppose we know that the word w only occurs in reviews with negative opinion and never is present in reviews with positive opinion. This gives us confidence that any message containing w is negative. This approach can be generalized to the probability of a message feature vector occurring in the message.

The entities, I need to classify, are text messages that are given in the form of strings. Raw strings are not convenient objects to handle by classifiers. Most machine learning algorithms can only classify numerical objects or otherwise require a distance metric or other measure of similarity between the objects. Moreover, string by itself doesn't represent a polarity of the message in any way. Therefore, set of features should be extracted.

Before proceeding with machine learning we have to convert all messages to numerical vectors called *feature vectors*, and then classify these vectors. The simplest example of a feature vector is the vector of the numbers of occurrences of all the words in the dictionary in a message. The problem with this approach is that not all the words affect a sentiment. In particular there are so-called STOP words which are completely useless in this respect. This causes significant performance degradation. Additionally, this approach does not analyze the context of the word, which causes prediction precision to be lower.

Extraction of features usually means that some information from the original

message is lost. On the other hand, the way feature vector is chosen is crucial for the performance of the filter. There is a trade-off between prediction quality and training speed. Usually, the difference between performance of a really detailed model and a reasonably simple model is not significant. Yet the training speed of a simple model is considerably swifter. Therefore, in most practical applications the most basic vector of word frequencies or its modification is used.

There are plethora of machine learning algorithms which are suitable for review classification task. Choosing the best algorithm may improve a quality of prediction. In practice, however, it is much more important what features are chosen for classification than what classification algorithm is used. If the features are chosen badly, then any machine learning algorithm will fail to classify data correctly no matter how good it is. A reasonable choice of features can make classification both precise and quick.

3.2 Important Words Selection

As already was mentioned in the previous section, selecting all the words as features is not very practical. Yet in performance comparison I will use this approach as well because it can be a baseline for other extraction functions quality. In case of using all the words let us define a word dictionary $D \subseteq W$. Let us suppose $|D| = k$. Then for matrix X which is $m \times n$ dimensional, now $n = k$. Now $x_{ij} \in X$ represents number of occurrences of the word $d_j \in D$ in the review $r_i \in R$.

Now let us consider extracting a smaller dictionary. The problem is to decide which words to include not to lose important features. As it turns out adjectives used most often define an opinion. There is almost no use in adverbs, because of a domain. It is possible that in restaurant reviews there are lots of them, yet for baby products it is not. Also I am going to use link verbs. Common link verbs are: be, appear, look, seem, become, get. For instance, in the sentence "My baby seemed happy". Another reasonable idea is including some verbs such as love, like, enjoy, recommend, etc.

3.3 Context Mining

In the previous section I described in detail the process of words selection. However, before beginning classification there is still a work to be done. Consider the following example: "My baby likes this toy", "My baby doesn't like this toy". Using our previous approach there is no difference between those two sentences in classification matrix. But obviously there is a significant difference - polarity is exactly the opposite. One solution would be to include "not" in the dictionary and hope that classifier will give it significant negative weight. The issue with this approach though is that in more complex reviews, a negation can bear positive sense. Consider an example "I really like this cradle. It looks nice and my baby does not cry as often as before". Therefore, another solution should be found. I suggest using antonyms.

Chapter 4

Effectiveness of trained models

Words, tables, graphs, pictures, code.

Conclusion

Words.

References

- [1] Manning, Chris and Hinrich Schütze. *Foundations of Statistical Natural Language Processing*. Cambridge [England]: MIT Press, 1999.
- [2] Wake, Lisa. *NLP : principles in practice*. St. Albans, Herts [England]: Ecademy Press, 2010.
- [3] Murphy, Kevin. *Machine Learning A Probabilistic Perspective*. Cambridge, Massachusetts [USA]: Massachusetts Institue of Technology, 2012.
- [4] Leskovec, Jure, Anand Rajaraman, and Jeff Ullman. *Mining of Massive Datasets SECOND EDITION*. Cambridge [England]: Cambridge University Press, 2014.
- [5] Witten, Ian, and Eibe Frank. *Data Mining: Practical Machine Learning Tools and Techniques THIRD EDITION*. Burlington [USA]: Morgan Kaufmann Publications, 2010.
- [6] McKinney, Wes. *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython 1st Edition*. Sebastopol [USA]: 2012.
- [7] <https://msdn.microsoft.com/en-us/library/ms174806.aspx> TODO