# Efficient Opinion mining for Multiple Domain.

Mr. Krushna Borude[1], Mr. Rohit Singhal[2]

*[1]M.Tech(CSE Student) at IET, Alwar ( Rajasthan)*
*[2] Associate Professor at IET, Alwar (Rajasthan)*

[1]borude.krushna@gmail.com
[2]mtechrohit@gmail.com

**Abstract:**

**Automatic classification of sentiment is incredibly important for various applications like opinion mining, opinion account, contextual advertising, and marketing research. Typically, sentiment classification has been sculptured as a result of the drawback of work a binary classifier exploitation reviews annotated for positive or negative sentiment. However, sentiment is expressed otherwise in many domains, and increase corpora for every potential domain of interest is pricey. Applying a sentiment classifier trained exploitation labeled knowledge for a particular domain to classify sentiment of user reviews on a special domain typically results in poor performance as a results of words that occur at intervals the train (source) domain might not appear inside the test (target) domain. we have a tendency to propose how to beat this drawback in cross-domain sentiment classification.**

**First, we have a tendency to produce a sentiment sensitive arrangement synonym finder exploitation labeled information for the source domains and unlabeled information for every source and target domains. Sentiment sensitivity is achieved at intervals the synonym finder by incorporating document level sentiment labels at intervals the context vectors used because the basis for measure the arrangement similarity between words. Now, we have a tendency to use the created book of facts to expand feature vectors throughout train and check times in a} very binary classifier. The projected technique significantly outperforms varied baselines and returns results that are comparable previously planned cross-domain sentiment classification ways in which on a benchmark dataset containing Amazon user reviews for varied types of product. we have a tendency to conduct an intensive empirical analysis of the planned technique on single and multi-source domain acceptance, unsupervised and supervised domain acceptance, and diverse similarity measures for creating the sentiment sensitive synonym finder. Moreover, our comparisons against the SentiWord web, a lexical resource for word variations, show that the created sentiment-sensitive synonym finder accurately captures words that specific similar sentiments.**

**Keywords: Cross-Domain Sentiment Classification, Domain Acceptance, Thesauri Creation, sentiwordnet, labeled and unlabeled information**.

## I INTRODUCTION

Graphs and networks actually rank among one amongst the foremost in style knowledge representation models because of their universal relevancy to various application domains. the necessity to research and mine interesting knowledge from graph and network structures has been long recognized, however solely recently the advances in info systems have enabled the analysis of graph structures at immense scales. Analysis of graph and network structures gained new momentum with the arrival of social networks. whereas the analysis of social networks has been a field of intensive analysis, significantly within the domains of social sciences and scientific discipline, economy or chemistry, it's the emergence of big social networking services over the online that spawned the analysis into large-scale structural properties of social networks.. Social networks exhibit a really clear community structure. Such community structure part stems from objective limitations (e.g., internal structure of a corporation will be closely drawn by the ties among a selected social network) or, to some extent, might result from subjective user actions and activities (e.g., bonding with others who share one's interests and hobbies). Social networks area unit extremely effective in bolstering cluster formation of similar people. groups of nodes that share common properties tend to induce connected within the social network. Opinion mining is that the domain of language process and text analytics that aims at the invention and extraction of subjective qualities from matter sources. Opinion mining tasks will be usually classified into 3 varieties. The primary task is named as sentiment analysis and aims at the institution of the polarity of the given supply text (e.g., identifying between negative, neutral and positive opinions). The second task consists in characteristic the degree of judgment and sound judgment of a text (i.e., the identification of factual knowledge as critical opinions).

This task is usually named as opinion extraction. The third task is aims at the invention and/or summarization of specific opinions on elect options of the assessed product. Some authors visit the task as sentiment analysis. All 3 categories of opinion mining tasks will greatly take pleasure in extra knowledge which will be provided from the social network. extra information might include: a node's spatial

relation indexes, a node's cluster membership, terminology utilized among the cluster, average cluster opinion on elect merchandise, group's coherence and cohesion, etc. of these variables enrich opinion mining algorithms and supply extra instructive capabilities to created models.

## II OBJECTIVE

Example: Allow us to think about the reviews shown in Table one for the 2 domains: books and room appliances. Table one show 2 positive and one negative review from every domain. We've got stressed the words that specific the sentiment of the author in an exceedingly review victimization previous face. From Table one, we tend to see that the words wonderful, broad, top quality, attention-grabbing and well researched area unit won't to specific a positive sentiment on books, whereas the word unsuccessful indicates a negative sentiment. On the opposite hand, within the room appliances domain the words excited, top quality, skilled, energy saving, lean, and delicious specific a positive sentiment, whereas the words rust and unsuccessful specific a negative sentiment. though words like top quality would specific a positive sentiment in each domains, and unsuccessful a negative sentiment, it's unlikely that we might encounter words like well researched for room appliances or rust or delicious in reviews on books. Therefore, a model that's trained solely exploitation reviews on books won't have any weights learnt for delicious or rust, that makes it troublesome to accurately classify reviews on kitchen appliances exploitation this model.

One answer to the present feature pair drawback is to use a synonym finder that groups totally different words that specific an equivalent sentiment. as an example, if we all know that each wonderful and delicious are positive sentiment words, then we can use this information to expand a feature vector that contains the word delicious exploitation the word wonderful, thereby reducing the pair between features during a test instance and a trained model. There are two necessary queries that has to be addressed during this approach: the way to automatically construct a synonym finder that's sensitive to the sentiments expressed by words?, and the way to use the synonym finder to expand feature vectors during coaching and classification?. The first question is mentioned in Section 4, wherever we have a tendency to propose a spatial arrangement approach to construct a sentiment sensitive synonym finder exploitation each labeled and unlabeled information from multiple domains. The second question is addressed in Section 5, wherever we have a tendency to propose a ranking score to pick the candidates from the synonym finder to expand a given feature vector.

## III LITERATURE SURVEY

### 3.1 Related work done:

Supervised learning algorithms that need labeled knowledge are successfully used to build sentiment classifiers for a given domain [1]. Still, sentiment is expressed otherwise in many domains, and it's pricey to annotate information for every new domain within which we would prefer to apply a sentiment classifier. As an example, within the electronics domain the words "durable" and light" are wont to specific positive sentiment, whereas "expensive" and "short battery life" usually indicates negative sentiment. On the opposite hand, if we have a tendency to contemplate the books domain the words exciting" and "thriller" expresses positive sentiment, whereas the words "boring" and "lengthy" sometimes specific negative sentiment. A classifier trained on single domain may not execute well on a unique domain as a result of it fails to be told the sentiment of the unseen words. The cross-domain sentiment categorization drawback [7], [8] focuses on the challenge of training a classifier from one or additional domains (source domains) and applying the trained classifier on a special domain (target domain). A cross-domain sentiment arrangement should overcome 2 major challenges.

This is what makes the spatial arrangement synonym finder sentiment sensitive. Unlabeled knowledge is cheaper to gather compared to labeled knowledge and is usually out there in massive quantities. The {use the utilization the employment} of unlabeled knowledge permits us to accurately estimate the distribution of words in supply and target domains. The designed method will learn from an oversized quantity of unlabeled knowledge to leverage a strong cross-domain sentiment classifier. In our projected technique, we have a tendency to use the automatically created synonym finder to expand feature vectors during a binary classifier at train and check times by introducing connected lexical parts from the synonym finder. we have a tendency to use L1 regularized provision regression because the classification algorithmic program.

## IV PROBLEM STATEMENT

### 4.1 Problem Statement:

We outline a domain D as a category of entities within the world or a linguistics conception. as an example, differing kinds of product like books, DVDs, or vehicles are thought of as completely different domains. Given a review written by a user on a product that belongs to a specific domain, the target is to predict the sentiment expressed by the author

within the review regarding the product. We have a tendency to limit ourselves to binary sentiment classification of entire reviews. We have a tendency to denote a supply domain by Dsrc and a target domain by Dtar. The set of labeled instances from the source domain, L(Dsrc), contains pairs (t, c) wherever a review, t, is assigned a sentiment label, c. Here c, $\in$ , and also the sentiment labels +1 and −1 severally denote positive and negative sentiments. Additionally to positive and negative sentiment reviews, there can still be neutral and mixed reviews in sensible applications. If a review discusses each positive and negative aspects of a specific product, then such a review is taken into account as a mixed sentiment review. On the opposite hand, if a review doesn't contain neither positive nor negative sentiment concerning a specific product then it's thought of as neutral. Though this paper solely focuses on positive and negative sentiment reviews, it's not laborious to increase the planned methodology to deal with multi-category sentiment classification issues.

However, the projected methodology is agnostic to the properties of the classifier and might be accustomed expand feature vectors for any binary classifier. As shown later within the experiments, L1 regularization allows us to pick out alittle set of options for the classifier. Our contributions during this work will be summarized as follows.

• We tend to propose a totally automatic methodology to form a synonym finder that's sensitive to the sentiment of words expressed in several domains. we tend to utilize each labeled and unlabeled information on the market for the source domains and unlabeled knowledge from the target domain.

• We tend to propose a method to use the created synonym finder to expand feature vectors at train and check times during a binary classifier.

• We tend to compare the sentiment classification accuracy of our projected methodology against varied baselines and previously projected cross-domain sentiment classification strategies for each single supply and multi-source adaptation settings.

• we tend to study the flexibility of our methodology to accurately predict the polarity of words victimization SentiWordNet, a lexical resource within which every WordNet synset is related to a polarity score.

# V SYSTEM ARCHITECTURE

## 5.1 Dataset

We use the cross-domain sentiment classification dataset1 ready by Blitzer et al. [7] to match the planned technique against previous work on cross-domain sentiment classification. This dataset consists reviews of Amazon product for four completely dissimilar product types: books, DVDs, electronics, and kitchen appliances. Each review is assigned with a rating (0-5 stars), a reviewer name and location, a product name, a review title and date, and also the review text. Reviews with rating &gt; 3 are labeled as positive, whereas those with rating &lt; 3 are labeled as negative. The general structure of this benchmark dataset is shown in Table 6.1. For all domain, there are 100 positive and 100 negative examples, a similar balanced composition because the polarity dataset created by Pang et al. [1]. The dataset additionally contains some unlabeled reviews for the four domains. This benchmark dataset has been utilized in a lot of previous work on cross-domain sentiment classification and by evaluating on that we are able to directly compare the projected technique against existing approaches. Following previous work, we tend to at random choose 800 positive and 800 negative labeled reviews from every domain as coaching instances (total range of coaching instances are 1600×4 = 6400), and also the remainder is used for checking (total range of test instances are 400×4 = 1600). In our experiments, we tend to choose every domain successively because the target domain, with one or a lot of different domains as sources.

Note that {when we tend to once we after we} mix quite one supply domain we limit the overall range of supply domain labeled reviews to 1600, balanced between the domains. as an example, if we tend to mix two supply domains, then we tend to choose 400 positive and 400 negative labeled reviews from every domain giving (400 + 400) × a pair of = 1600. This permits us to perform a good analysis once combining multiple supply domains. We tend to produce a sentiment sensitive synonym finder exploitation labeled knowledge from the supply domain and unlabeled knowledge from supply and target domains as delineated in Section 4.

We then use this synonym finder to expand the labeled feature vectors (train instances) from the supply domains and train an L1 regularized provision regression-based binary classifier (Cassias) 2. L1 regularization is shown to provide a thin model, wherever most incompatible options ar assigned a zero weight [22].

**TABLE 3 The result of using a sentiment sensitive thesaurus for cross-domain sentiment classification**

| Method | Kitchen | DVDs | Books |
|--------|---------|------|-------|
| No. Adapt | 0.7261 | 0.6807 | 0.6272 |
| NSST | 0.7750 | 0.7350 | 0.7146 |
| SST | 0.8518 | 0.7826 | 0.7632 |
| In-Domain | 0.8770 | 0.820 | 0.8040 |

This enables us to pick out helpful features for classification during a systematic manner while not having to preselect options exploitation heuristic approaches. In our beginning experiments, we have a tendency to discovered that the classification accuracy on two development target domains did not vary significantly with totally different L1 regularization parameter values. Therefore, we have a tendency to set the L1 regularization parameter to 1, that is that the default setting in Classias, for all experiments represented during this paper. Next, we have a tendency to use the trained classifier to classify reviews within the target domain. The synonym finder is once more wont to expand feature vectors from the target domain. This procedure is perennial for every domain in Table

The on top of mentioned procedure creates four thesauri (each synonym finder is made by excluding labeled coaching information for a selected target domain).

For example, from the three domains DVDs, electronics and books, we have a tendency to generate 53, 586 lexical elements and 62, 744 sentiment parts to make a synonym finder that's wont to adapt a classifier trained on those 3 domains to the kitchen domain. Similar records of options are generated for the opposite domains similarly. To avoid generating distributed and possibly noisy options, we have a tendency to need that every feature occur in a minimum of two totally different review sentences. we have a trend to use classification accuracy on target domain because the analysis metric. it's the fraction of the correctly classified target domain reviews from the total variety of reviews within the target domain, and is outlined as follows:

## 5.2 Cross-Domain Sentiment Classification

To evaluate the advantage of using a sentiment sensitive synonym finder for cross-domain sentiment classification, we have a tendency to compare the planned technique against three baseline strategies in Table 4. Next, we have a tendency to describe the strategies compared in Table 4.

• No Adapt: This baseline simulates the impact of not playacting any feature enlargement. we have a tendency to simply train a binary classifier exploitation unigrams and bigrams as options from the labeled reviews within the supply domains and apply the trained classifier on a target domain. this will be thought-about as a boundary that doesn't perform domain adaptation.

• NSST (Non-sentiment Sensitive Thesaurus): to judge the advantage of exploitation sentiment options on our planned methodology, we have a tendency to produce a synonym finder only exploitation lexical elements. Lexical elements will be derived from each labeled and unlabelled reviews whereas, sentiment elements are often derived only from labeled reviews. we have a tendency to failed to use rating data within the supply domain labeled knowledge during this baseline. A synonym finder is made utilization those options and afterwards used for feature growth. A binary classifier is trained exploitation the enlarged options.

• planned (SST: sentiment sensitive thesaurus): this can be the planned methodology represented during this paper. we have a tendency to use the sentiment sensitive synonym finder created exploitation the procedure represented in Section 4 and use the synonym finder for feature growth in a binary classifier.

• In-Domain: during this methodology, we have a tendency to train a binary classifier exploitation the labeled knowledge from the target domain. This methodology provides an bound for the cross-domain sentiment analysis. This higher baseline demonstrates the classification accuracy we are able to hope to get if we have a tendency to had labeled knowledge for the target domain. Note that this can be not a cross-domain classification setting. Table four shows the classification accuracy of the abovementioned strategies for every of the four domains within the benchmark dataset because the target domain. Moreover, for every domain we've got shown in boldface the simplest cross-domain sentiment classification results. Note that the In-Domain baseline isn't a cross-domain sentiment classification setting and acts as a bound. From the leads to Table 4, we have a tendency to see that the planned (sentiment sensitive thesaurus) returns the simplest cross-domain sentiment classification accuracy for all four domains. The analysis of variance (ANOVA) and Tukey's honestly important variations (HSD) tests on the classification accuracies for the four domains show that our planned methodology is statistically considerably higher than each the no synonym finder and non-sentiment sensitive synonym finder baselines, at confidence level zero.05. This shows that exploitation the sentiment sensitive

synonym finder for feature enlargement is helpful for cross-domain sentiment classification.

### 5.3 Effect of Relatedness Measures

The choice of the connection live is an important call in a very thesauri-based approach. totally different completely different} connection measures can list different lexical elements as neighbors for a selected lexical element.

Therefore, the set of growth candidates are going to be directly influenced by the connection measure accustomed produce the synonym finder. to check the result of the connection measure on the performance of the projected methodology, we have a tendency to construct four sentiment sensitive thesauri exploitation four totally different connection measures. we have a tendency to then conduct feature growth and coaching within the same manner as represented in Section 5 with all four connection measures. we have a tendency to use the 3 domains at a time because the sources and therefore the remaining domain because the target during this experiment. The classification accuracies obtained victimization the various connection measures are shown in Table 5. Next, we have a tendency to describe the four connection measures compared in Table 5.

**TABLE 4.Comparison of different relatedness measures**

| Method | Kitchen | DVDs | electronics | Books | Overall |
|--------|---------|------|-------------|-------|---------|
| Cosine | 0.8342 | 0.7826 | 0.8363 | 0.7657 | 0.8047 |
| Lin[19] | 0.8367 | 0.7826 | 0.8438 | 0.7632 | 0.8066 |
| Proposed | 0.8518 | 0.7826 | 0.8386 | 0.7632 | 0.8091 |
| Reversed | 0.8342 | 0.7852 | 0.8463 | 0.7632 | 0.8072 |

• **Cosine Similarity**: This is the cosine of the angle between the two vectors that represent two lexical elements $u$ and $v$. Using the notation introduced in Section 4, it can be computed as follows:

$$r(v,u) = \frac{\sum \omega \epsilon\, r(v) f(u,\omega)}{\| u \| \| v \|}$$

$$\| v \| = \sqrt{\sum_{\omega\epsilon\, r(\upsilon)} (f(v,\omega))^2}$$

$$\| u \| = \sqrt{\sum_{\omega\epsilon\, r(u)>0} (f(u,\omega))^2}$$

Here, $\Gamma(v) = \{x/f(v, x) > 0\}$, is the set of features $x$ that have positive pmi values in the feature vector for the element $v$. Cosine similarity is widely used as a measure of relatedness in numerous tasks in natural language processing [23].

• **Lin's Similarity Measure**: We use the similarity measure proposed by Lin [19] for clustering similar words. This measure has shown to outperform numerous other similarity measures for word clustering tasks. It is computed as follows:

$$r(v,u) = \frac{\sum_{\omega\epsilon\, r(v)\cap r(u)}(f(v,\omega) + f(u,\omega))}{\sum_{\omega\epsilon\, r(v)} f(v,\omega) + \sum_{\omega\epsilon\, r(u)} f(u,\omega)}$$

• **Proposed**: This is the relatedness measure proposed in this paper and is defined by Equation 2. Unlike the **Cosine Similarity** and **Lin's Similarity Measure**, this relatedness measure is asymmetric.

• **Reversed**: As a baseline that demonstrates the asymmetric nature of the relatedness measure proposed in Equation 2, we exchange the two arguments $u$ and $v$ in Equation 2 to construct a baseline relatedness measure. Specifically, the reversed baseline is computed as follows:

$$r(v,u) = \frac{\sum \omega\epsilon\{x|f(\mathbf{v},x) > 0\}f(u,\omega)}{\sum \omega\epsilon\{x|f(\mathbf{u},x) > 0\}f(u,\omega)}$$
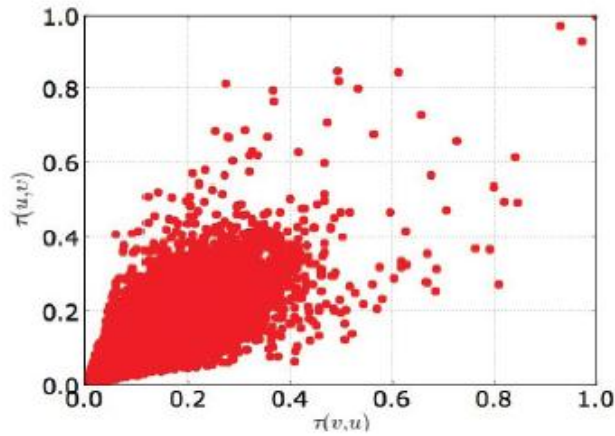
Fig. 2. Correlation between relatedness scores.

From Table 5 we tend to see that the planned connection measure reports the best overall classification accuracy followed by the Reversed baseline, Lin's Similarity measure, and therefore the cosine Similarity in this order. However, it should be prominent that the variations in performance among those connection measures don't seem to be statistically important. This result implies that a wide-range of connection measures is accustomed produce a sentiment sensitive wordbook to be used with the feature enlargement methodology planned within the paper. any investigations into the unfitness of the planned methodology to the connection measures discovered three important reasons that we'll discuss next. First, remind that the planned feature enlargement methodology (Section 5) doesn't use absolutely the worth of connection scores, still only uses the relative rank among the enlargement candidates. Therefore, two connection measures that turn out totally different absolute scores will acquire similar performance if the relative rankings among enlargement candidates are similar. Second, as a posterior step to feature enlargement we tend to train a binary classifier with L1 regularization exploitation source domain labelled knowledge.
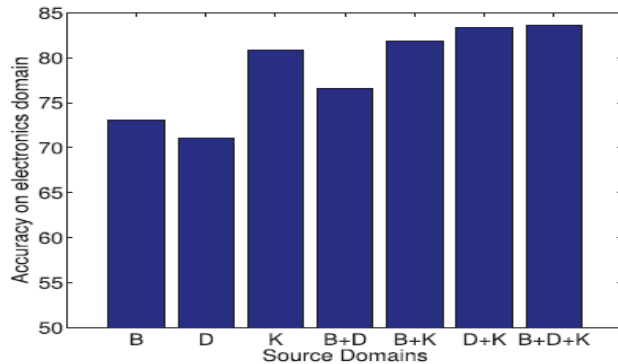
Therefore, if we tend to introduce any incorrect enlargement candidates that don't properly mirror sentiment, those enlargement candidates are appointed zero weights. Consequently, invalid enlargement candidates are cropped out from the ultimate model learnt by the binary classifier. However, it should be emphasised that though this posterior classifier coaching step will take away incorrect expansions, it cannot introduce the proper expansions. Therefore, it's important to the performance of the planned methodology that a connection measure identifies correct enlargement candidates throughout the feature enlargement step. to review the degree of spatiality within the connection measure planned in Equation 2, and its outcome on the performance of the planned cross-domain sentiment classification methodology, we tend to conduct the subsequent experiment. For word pairs (u, v) inside the sentiment sensitive wordbook, we have a tendency to plot the relation scores $\tau$ (u, v) next to $\tau$ (v, u) as shown in Figure two. There ar one, 000, 000 such word pairs (data points) in Figure 2. From Figure 2, we have a tendency to see that $\tau$ (u, v) is very related to $\tau$ (v, u). In reality the Pearson parametric statistic for Figure 2 is as high as 0.8839 with a decent confidence interval of [0.8835, 0.8844]. This experimental result indicates that, though by definition Equation two is uneven, its level of spatiality is incredibly little in observe. each the planned methodology and its Reversed baseline (Equation 8) coverage similar accuracy values in Table five any supports this finding. we have a tendency to take into account this perceived low level of spatiality to be a third reason that explains the similar performance among symmetric and asymmetric connection measures compared in Table 5.
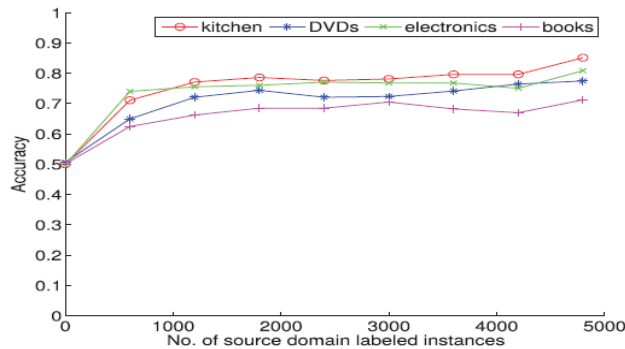
## 5.4 Effect of using Multiple Sources

In real-world cross-domain sentiment classification settings usually we've more than one supply domains at our disposal. Choosing the right provide domains to adapt to a given target domain could be a difficult drawback [24]. To check the result of exploitation multiple supply domain within the planned technique, we tend to choose the electronics domain because the target and train a sentiment classifier exploitation all possible mixtures of the three source domains books (B), kitchen appliances (K), and DVDs (D). Note that we tend to fix the entire number of labeled coaching instances after we combine multiple domains as sources to avoid any performance gains just because of the magnified variety of labeled instances as already explained in Section 6.1. Specifically, once employing a single source domains we tend to take 800 positive and 800 negative labeled reviews, once exploitation two source domains we tend to take 400 positive and 400 negative labeled reviews from every source domain, and once exploitation all 3 source domains we tend to take 266 positive and 266 negative labeled reviews. Moreover, we tend to use all obtainable unlabeled reviews from every supply domain and also the target domain. Figure three shows the result of mixing multiple supply domains to make a sentiment classifier for the physical science domain. we tend to see that the room domain is that the single Fig. 4. result of supply domain labeled information. best input domain once adapting to the electronics target domain. This behavior is

explained by the actual fact that generally kitchen appliances and electronic things have similar aspects. however a additional interesting observation is that the accuracy that we tend to get {when we tend to|once we|after we} use two source domains is usually bigger than the accuracy if we use those domains severally. the best accuracy is achieved after we use all three source domains. though not shown here for house limitations, we tend to ascertained similar trends with different domains within the benchmarkdataset.



Graph5.1: Effect of using multiple source domains.



Graph5.2.Effect of source domain labeled data.

# VI ALGORITHMS

## 6.1 Lemmatisation

- Lemmatisation is a process of identifying the lemma of a word. Algorithms for performing this operation typically use dictionaries, where they look up the primary form of the word.
- Lemmatisation may find several different lemmas for a given word, if the word is the inflected form of a lot of various lemmas.
- The use of lemmatisation reduces the number of terms present in the corpus and allows matching of words in documents, even if words tend to occur in different grammatical forms.
- The use of lemmatisation may result in deterioration of the classification accuracy, due to the possible occurrence of words in different forms derived from one lemma, depending on the document affiliation to one of the classes.
- Pseudo Code for Lemmatisation :

1) IF suffix("") THEN change(""-->"") EXCEPT
2) IF suffix("ote") THEN change("ote"-->"ite")
3) ELSE IF suffix("ten") THEN change("ten"-->"e")
4) ELSE IF suffix("s") THEN change("s"-->"")
5) ELSE IF suffix("g") THEN change(""-->"") EXCEPT
6) IF suffix("hing") THEN change("ing"-->"e")
7) ELSE IF suffix("d") THEN change("d"-->"")
8) ELSE IF suffix("e") THEN change(""-->"")
9) ELSE IF suffix("r") THEN change(""-->"")
10) ELSE IF suffix("f") THEN change(""-->"")
11) ELSE IF suffix("t") THEN change(""-->"")

## 6.2.Stopwords

- A *stop-list* is a set of words that should be removed at early stage of text processing. In the majority cases, these are conjunctions and other words which do not contribute additional information to the content of the sentence.
- Often, stop-list words are present in the sentence solely due to the requirements of language's grammar. In many cases the use of stop-lists improves accuracy and performance of text document processing
- Pseudo Code for Stop words:

1   For (i=0;i<=length;i++)
2   {
3   If(word == removerword[i])
4   {
5   Removeword();
6   }
7   End for }

## 6.3 Stemming:

- Stemming is a process similar to lemmatization. It aims to extract the core of the word, referred to as the stem, from the inflectional word forms. Stemming typically involves removal and replacement of prefixes and

suffixes. The result of stemming does not need to be and often is not a proper lemma

- Pseudo Code  Stemming:

- **Step 1a:**
  Remove matching instrumental case if preceded by double Consonants and remainder is a suitable word:
  al el
  Then remove one of the twice consonants (digraphs included)

- **Step 1b**
  Remove the following noun cases if the remainder is a valid word: ra re nak nek ban ben ba be tól t®l ról r®l
  ból b®l hoz hez nál nél ként ig val vel
  Then if the last letter of the new word is á change it to a
  Else if the last letter of the new word is é change it to e

- **Step 2**
  Remove longest matching personal suffix for owned nouns if remainder is valid word: oké öké 'aké eké

- **Step 3a**
  Search for the longest among the following singular owner suffixes and do the action indicated:
  em om am m ünk unk nk juk jük uk ük od ed ad öd d ja je a e o Remove if the remainder is a valid word ánk ájuk ám ád á Replace with a if remainder is a valid word énk éjük ém éd é Replace with e if remainder is a valid word

- **Step 3b**
  Search for the longest among the following plural owner suffixes and perform the action indicated:
  im jaid jeid aid eid id jaim jeim aim eim jai jei ai ei i jaitok jeitek aitok eitek itek jeik jaik aik eik ik jaink jeink eink aink ink Remove if the remainder is a valid word áim áid ái áink áitok áik Replace with a if remainder is a legal word éim éid éi éink éitek éik Replace with e if remainder is a legal word

- **Step 4**
  Search for the longest among the following plural suffixes and perform the action indicated: ök ok ek ak Remove if the remainder is a valid word ák Replace with a if the remainder is a valid word ék Replace with e if remainder is a valid word

**6.4 social opinion algorithms:**

The method proposed in this paper for determining term's semantic orientation is a variant of the method used in [1].

The drawback of the original method is that it assigns maximum or minimum value to all terms if they occur in only one class, regardless of the number of occurrences. Therefore, we have proposed an alternative way of calculating the semantic orientation of a term. Our method is based on the ratio of phrase occurrence frequency in documents assigned to positive and negative classes. According to our approach the scoring function for assigning positive and negative scores to terms to terms becomes

$$\text{Score}(t) = \begin{cases} p_t - 1, iff\ \ pi \geq 1 \\ \left(\frac{1}{pt} - 1\right) iff\ pi < 1 \end{cases}$$

Where

$$p_t = \begin{cases} p(t|C_p) + \varepsilon \\ p(t|C_N) + \varepsilon \end{cases}$$

- The score value of a term determined as above increases or decreases with changing frequency of term occurrences in positive or negative class, even if the term occurs in only one class.
- The score value of a term resolute as above increases or decreases with changing frequency of term occurrences in positive or negative class, even if the term occurs in only one class. Similarly to the score method, the disadvantage of the proportional method is the noise resulting from an insufficient number of term's instances in the training set. However, when proportional method is used, the influence of the noise is limited in comparison to the score method.
- This limitation results from the use of the scaling value .The score value assigned to a term which occurs only once in the training set is limited by the ratio of cardinalities of classes, whereas the semantic orientation of terms characteristic to positive or negative documents is often orders of magnitude greater. To further reduce the impact of the noise on the effectiveness of the algorithm, we plan to add filtering by removing from the dictionary terms that occur in fewer than $\beta$ documents

$$\beta = \left\lfloor \frac{|C_*|}{C_\#} \right\rfloor + 2$$

the minority class in the training set. Setting the threshold $\beta$ of term occurrences in the training set allows to eliminate terms that are not characteristic for any of the document classes, i.e. these terms for which conditional probabilities of term occurrences are similar for both classes, but which occurred too rarely in the training set, to have their evaluation been determined to be equal or close to zero.

Experiments:

- Test sets

The main objective of experiments was to test the accuracy of the classification algorithm proposed in Section IV. We used collections of opinions harvested from the e-commerce site Merlin, and two social networks Znany lekarz and Ceneo. The first dataset is the collection of movie reviews from the Merlin website. The reviewers were grading movies using the scale from 1 to 5, where the reviews with grades 1 or 2 are considered negative, and the reviews with grades 4 and 5 are considered positive. We have discarded neutral reviews with grade equal to 3. The dataset consists of 1055 negative reviews and 9068 positive reviews. The second dataset contains opinions on consumer products aggregated by the website Ceneo. Among the reviews graded from 0 to 5, we have chosen 793 negative reviews with grades 0 or 1 and 16 674 positive reviews graded 4 or 5,. Again, we have discarded all neutral reviews. The third dataset comes from the website Znany lekarz which gathers opinions about physicians. We have assumed that opinions associated with grades 1 and 2 on a scale 1-6 are negative, and opinions with grades 5 and 6 indicate a positive feedback. The dataset contains 2380 -ve opinions and 11 764 positive opinions. In addition we have performed tests using an aggregated dataset created by merging the three datasets. The aggregated dataset contains 4228 negative opinions and 37 506 positive opinions.

Pseudo Socio-opinion Algo:

1. For(i=0;i<total comments ; i++)
2. {
3. If(points > 1)
4. {
5. Scorer ++
6. }
7. Else
8. {
9. Scorer –
10. }
11. Counter = 1/scorer
12. }end for
13. Total counter = counter;
14. Update Scorer

6.5 Cross domain analysis algorithm:

- Input:
  - $-D_{src}$ (source domain)
  - $-D_{tar}$ (target domain)
- Output:
  $$D_{src} \quad \cap \quad D_{tar}$$
  $$L(D_{src}) \quad \cap \quad L(D_{tar})$$
- Process:
  1) Get input source domain
  2) Genrate $L(D_{src})$ with pair (t,c) review t with label c
  3) Genrate $( (\cup(D_{src})) \cup (D_{tar}))$
  4) For $t \in D_{src}$
     If(+ve sentiment)*P
     If(-ve sentiment)*N
  5) Compute a
     For all i: 1to n
     A=C(i,w)
     a=A/N

  6) calculate b

     For all j: 1to m
     B=C(j,u)
     b=B/N
  7) Compute

     $$F = \log \frac{c(u,w)/n}{a*b}$$

  8) Find out target domain

     for 'i' having 0 to w length

     if $(f_N > 0)$ a + = w(x)

     else a=a

     for j having 0 to w length

     if (f(u)>0)

     b+=w(x)

  9) if (a/b > 0.5)

     {

$$Dtar = w$$

}

Else

Continue

10) update sentiment to $D_{tar}$



# VII RESULT's

**Feedback Analysis Report**

College Name :College Name 1 2
Year          : SE
Department  : comp
Subject       : OSA

**Feedback Details**

| Orientation | Feedback Details |
|---|---|
| Positive | hi this is very good subject |
| Positive | osa is very good subject and also dsps is very good |

Total Feedback            :2
Positive Orientation Count  : 6
Negative Orientation Count : -2
Result                     : Positive Orientation

Signature of Authority

# VIII CONCLUSION

We planned a cross-domain sentiment classifier exploitation an automatically extracted sentiment sensitive synonym finder. to overcome the feature mis-match drawback in cross-domain sentiment classification, we tend to use labelled data from multiple source domains and unlabeled data from source and target domains to calculate the connection of features and construct a sentiment sensitive synonym finder. we tend to then use the created synonym finder to expand feature vectors throughout train and check times for a binary classifier. A relevant set of the options is selected exploitation L1 regularization. The planned technique considerably outperforms many baselines and reports results that ar comparable previously planned cross-domain sentiment classification strategies on a benchmark dataset. Moreover, our comparisons against the SentiWordNet show that the created sentiment-sensitive synonym finder accurately teams words that specific similar sentiments. In future, we decide to generalize the planned technique to resolve different kinds of domain adaptation tasks. we implement it and compare the results with expected results.

## IX REFERENCES:

[1] B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs up? Sentiment classification using machine learning techniques," in EMNLP 2002, 2002IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING

[2] P. D. Turney, "Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews," in ACL 2002, 2002

[3] B. Pang and L. Lee, "Opinion mining and sentiment analysis," Foundations and Trends in Information Retrieval, vol. 2, no. 1-2, 2008.

[4] Y. Lu, C. Zhai, and N. Sundaresan, "Rated aspect summarization of short comments," in WWW 2009, 2009,

[5] T.-K. Fan and C.-H. Chang, "Sentiment-oriented contextual advertising," Knowledge and Information Systems, vol. 23, no. 3, 2010.

[6] M. Hu and B. Liu, "Mining and summarizing customer reviews," in KDD 2004, 2004,.

[7] J. Blitzer, M. Dredze, and F. Pereira, "Biographies, bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification," in ACL 2007, 2007

[8] S. J. Pan, X. Ni, J.-T. Sun, Q. Yang, and Z. Chen, "Cross-domain sentiment classification via spectral feature alignment," in WWW 2010, 2010.

[9] H. Fang, "A re-examination of query expansion using lexical resources," in ACL 2008, 2008,

[10] G. Salton and C. Buckley, Introduction to Modern Information Retreival. McGraw-Hill Book Company, 1983.

[11] D. Shen, J. Wu, B. Cao, J.-T. Sun, Q. Yang, Z. Chen, and Y. Li, "Exploiting term relationship to boost text classification," in CIKM'09, 2009

[12] T. Briscoe, J. Carroll, and R. Watson, "The second release of the rasp system," in COLING/ACL 2006 Interactive Presentation Sessions, 2006.

[13] T. Joachims, "Text categorization with support vector machines: Learning with many relevant features," in ECML 1998, 1998, pp. 137–142.

[14] V. Hatzivassiloglou and K. R. McKeown, "Predicting the semantic orientation of adjectives," in ACL 1997, 1997, pp. 174–181.

[15] J. M. Wiebe, "Learning subjective adjective from corpora," in AAAI 2000, 2000, pp. 735–740.

[16] Z. Harris, "Distributional structure," Word, vol. 10, pp..

[17] P. Turney, "Similarity of semantic relations," Computational Linguistics, vol. 32, no. 3, 2006.

[18] P. Pantel and D. Ravichandran, "Automatically labeling semantic classes," in NAACL-HLT'04, 2004,

[19] D. Lin, "Automatic retrieval and clustering of similar words," in ACL 1998, 1998,

[20] J. Weeds and D. Weir, "Co-occurrence retrieval: A flexible framework for lexical distributional similarity," Computational Linguistics, vol. 31, no. 4, pp. 439–475, 2006.[21] S. Sarawagi and A. Kirpal, "Efficient set joins on similarity predicates," in SIGMOD '04, 2004, pp. 743–754.

[22] A. Y. Ng, "Feature selection, l1 vs. l2 regularization, and rotational invariance," in ICML 2004, 2004.

[23] C. D. Manning and H. Sch ¨ utze, Foundations of Statistical Natural Language Processing. Cambridge, Massachusetts: The MIT Press, 2002.

[24] M. T. Rosenstein, Z. Marx, and L. P. Kaelbling, "To transfer or not to transfer," in NIPS 2005 Workshop on Transfer Learning, 2005.

[25] J. Blitzer, R. McDonald, and F. Pereira, "Domain adaptation with structural correspondence learning," in EMNLP 2006,

[26] A. Esuli and F. Sebastiani, "Sentiwordnet: A publicly available lexical resource for opinion mining," in LREC 2006,

[27] G. Qiu, B. Liu, J. Bu, and C. Chen, "Expanding domain sentiment lexicon through double propagation," in IJCAI 2009,

[28] N. Kaji and M. Kitsuregawa, "Building lexicon for sentiment analysis from

massive collection of html documents," in EMNLP 2007, 2007.