# How Much Noise in Text is too Much: A Study in Automatic Document Classification

Submitted for Blind Review

## Abstract

*Noise is a stark reality in real life data. Especially in the domain of text analytics it has a significant impact as data cleaning forms a very large part (upto $80\%$ time) of the data processing cycle. Noisy unstructured text is common in informal settings such as on-line chat, SMS, email, newsgroups and blogs, automatically transcribed text from speech data, and automatically recognized text from printed or handwritten material. Gigabytes of such data is being generated everyday on the Internet, in contact centers, and on mobile phones. Researchers have looked at various text mining issues such as pre-processing and cleaning noisy text, information extraction, rule learning, and classification for noisy text. This paper focuses on the issues faced by automatic text classifiers in analyzing noisy documents coming from various sources. The goal of this paper is to bring out and study the effect of different kinds of noise on automatic text classification. Does the nature of such text warrant moving beyond traditional text classification techniques? We present detailed experimental results on simulated noise on benchmark datasets viz. Reuters-21578 and 20-newsgroups. We also present interesting results on real life noisy datasets from various contact center domains.*

## 1 Introduction

The importance of text mining applications is growing proportionally with the exponential growth of electronic text. Along with the growth of the Internet many other sources of electronic text have become really popular over the last decade. With the Internet penetrating into the lives of more and more people, email, chat, newsgroups, blogs, discussion fora etc. have become popular and hence generate a huge amount of text data everyday. Other big contributors to the pool of electronic text documents are call centers and CRM organizations in the form of call logs, call transcriptions, problem tickets, complaint emails, electronic text generated by Optical Character Recognition (OCR) process on hand-written or printed documents, conversational data converted automatically to text and mobile text such as Short Message Service(SMS).

Though the nature of these documents is varied, all of them share a common effect - the presence of textual noise. Text produced under such circumstances is typically highly noisy containing spelling errors, abbreviations, non-standard words, false starts, repetitions, missing punctuations, missing letter case information, pause-filling words such as *um* and *uh* and other text and speech disfluencies. More often than not such data requires cleaning and preprocessing before applying any state-of-the-art text analytics technique.

*Noisy Text Analytics* is defined as a process of information extraction whose goal is to automatically extract structured or semistructured information from noisy unstructured text data[1]. However one of the commonly used text mining applications, quite different from extraction of information, is *text classification* or *text categorization*.

The text classification task is one of learning models for a given set of classes and applying these models to new unseen documents for class assignment. Text classification has many important real life applications. For example, categorizing news articles according to topics such as *politics*, *sports*, or *education*; email categorization; building and maintaining web directories like Dmoz[2]; spam filters; automatic call and email routing in contact centers; pornographic material filters and so on. Two types of classifiers are commonly employed viz. statistical and rule based classifiers. In statistical classifiers a *model* is learned on a corpus of already labeled data and once trained the system can be used for automatic assignment of labels to unseen data. Rule based classifiers, on the other hand, are good at finding class boundaries based on presence or absence of words and/or phrases.

In both statistical as well as rule based text classification techniques, the content of the document is the sole determiner of the category to be assigned. However noise in the text distorts the content and hence users can expect the categorization performance to get affected. Classifiers are essentially trained to identify correlations between extracted features (words) and different categories which can be later

---

[1] http://en.wikipedia.org/wiki/Noisy_text_analytics

[2] http://dmoz.org/

utilized to categorize new documents. For example, email containing text like *exciting offer, get a free laptop* might have a stronger correlation with the category *spam* emails than *non-spam* emails. Noise in text distorts this feature space as *excitinng ofer get a tree lap top* will be a new set of features and the categorizer might not be able to relate it to the *spam emails* category. The feature space explodes as the same feature can appear in different forms due to spelling errors, poor recognition and wrong transcription. Noisy text categorization in particular has important practical applications in the form of problem determination in contact centers, call routing, categorization of hand-written customer complaints and automatic SMS routing.

**Our Contribution:** In this paper we will show the effect of different kinds of noise on text classification performance by doing detailed experiments on synthetic as well as real life noisy datasets. Here we are essentially reporting our observations based on experiments and not proposing any new method to combat noise in text for text classification. Our experiments show that text classification algorithms are quite robust even in the presence of a high degree of typographical noise or noise introduced by Automatic Speech Recognition (ASR) systems. We feel this work is a necessary pre-requisite to motivate researchers to look at various issues pertaining to noisy text categorization.

**Organization:** The rest of this paper organized as follows. Section 2 discusses the related work in the noisy text domain and also looks at noisy text classification. Section 3 introduces the various kinds of noisy textual data. Following section describes the settings and systems used in our experiments. In section 5, we describe benchmark and real world datasets used in our experiments. We present detailed results on the datasets, followed by a discussion on the significance of our results. Section 6 sets out our conclusion and explores avenues of future work.

## 2   Previous Work

In this section, we will present some of the relevant work in the following related areas viz. (1) noisy text analytics, (2) text classification and (3) noisy text classification.
**Noisy Text Analytics:** There has been a lot of work on analyzing the conversational data collected in contact centers. These include call type classification for the purpose of categorizing calls [23], call routing [9], obtaining call log summaries [6], agent assisting and monitoring [17], and building of domain models [20]. Wrong spellings can affect automatic classification performance in multiple ways depending on the nature of the classification techniques being used. In the case of statistical techniques, spelling mistakes distort the feature space. A comprehensive survey of techniques pertaining to detecting and correcting spelling errors

can be found in  [12].
**Text Classification:** The two broad types of classification methods used are discriminative and generative methods. Discriminative methods like SVMs [11] or logistic regression (LR) [22] are two-class classifiers that find separators between documents of two classes in some space of representations. Other discriminative models include maximum entropy methods [18] and boosted decision trees in the ADABoost framework [7]. Generative methods are typified by naive Bayes (NB), Latent Dirichlet Allocation [3], and the more recent BayesANIL [19]. Discriminative methods are widely accepted to be more accurate, but generative methods provide intuitive text generation models and have been used in a variety of applications. The industry has also made significant advances in the development and deployment of real-world high-performance text classification systems [15] using combinations of rule-based, hand-tuned, and statistical techniques. However, not all the techniques used in commercial systems are publicly known, and few general principles can be derived from these systems.
**Noisy Text Classification:** Electronically recognized handwritten documents and documents generated from OCR process are typical examples of noisy text. Authors in  [21] have the studied characteristics of noise present in such data and its effects on categorization accuracy. Authors in  [2] proposed a generic system for text categorization based on statistical analysis of representative text corpora. They evaluate their system on the tasks of categorizing abstracts of paper-based German technical reports and business letters concerning complaints. They claim that the tasks achieve recognition scores of approximately 80% and are very robust against recognition or typing errors.
OCR systems produce essentially word substitutions while ASR systems give rise to word substitutions, deletions and insertions. Moreover, ASR systems are constrained by a lexicon and can give as output only words belonging to it, while OCR systems can work without a lexicon (this corresponds to the possibility of transcribing any character string) and can output sequences of symbols not necessarily corresponding to actual words. Such differences are expected to have a strong influence on the performance of systems designed for categorizing ASRed documents in comparison to the systems for OCRed documents. We are not aware of any work dealing with ASR document categorization, relevant issues and reported results though researcher have looked at call-type classification [8].

## 3   Noise in Text

We define noise as *any kind of difference in the surface form of an electronic text from the intended, correct or original text*. We see such noisy text everyday in various forms. Each one has characteristics unique to it and hence requires

special handling. We introduce some such noisy textual data in this section.

- **On line Noisy Documents:** Emails, chat logs, scrapbook entries, newsgroup postings, threads in discussion fora, blogs etc. fall under this category. People are less careful about the lexical accuracy of written content in such informal modes of communication. These are characterized by frequent misspellings, commonly and not so commonly used abbreviations, incomplete sentences, missing punctuations and so on.

- **SMS**: Short Message Services is becoming more and more common everywhere day by day. Language usage over SMS texts significantly differs from the standard form of the language. An urge towards shorter message length facilitating faster typing and the need for semantic clarity, shape the structure of this non-standard form known as the *texting language* [4].

- **Text Generated by ASR Devices:** Automatic Speech Recognition (ASR) is the process of converting a speech signal to a sequence of words. An ASR system takes speech signals such as monologues, discussions between people, telephonic conversations, etc. as input and produces a string of words, typically not demarcated by punctuations, known as a *transcript*. An ASR system consists of an acoustic model, a language model and a decoding algorithm. The acoustic model is trained on speech data and their corresponding manual transcripts. The language model is trained on a large monolingual corpus. ASR converts audio into text by searching the acoustic model and language model space using the decoding algorithm.

- **Text Generated by OCR Devices:** Optical character recognition, or OCR, is a technology that allows digital images of typed or handwritten text to be transferred into an editable text document. It takes a picture of text and translates the text into Unicode or ASCII. For handwritten optical character recognition, the rate of recognition is 80% to 90% with clean handwriting. OCR systems give rise to some typical substitution errors such as *iii* instead of *m*, *5* instead of *s* etc.

- **Call Logs in Contact Centers:** Today's contact centers (also known as call centers) are increasingly contributing to the pool of noisy text by the means of *call logs*. Agents are expected to record summaries immediately after completing interactions with customers and before starting the next one. As the agents work under immense time pressure, the summary logs are very poorly written and sometimes even difficult for humans to interpret. Analysis of such call logs is important to identify problem areas, evaluate agent performance, predict evolving problems etc. They also produce a huge amount of unstructured data in the form of emails, call transcriptions, SMS, chat transcripts etc.

## 4  System Description

### 4.1  Spelling Error Simulation

We developed a program to introduce spelling errors in a text data corpus, *SpellMess*. SpellMess can be customized to introduce *Damerau-type errors*, i.e., insertion, deletion or substitution of a letter or transposition of two letters [5]. It requires two configuration files - (i) *KBMatrix* encoding the keyboard layout in a system understandable format so that the probability of a key getting pressed instead of the intended one can be computed. We assume any of the 8 surrounding letters can be substitute a letter by a wrong keypress, but the two letters on either side in the same row have more chance of getting substituted. (ii) *Weights* containing overall error probability and probability of different types of errors viz. insertion, deletion, transposition, substitution and duplication. For example, one can specify the overall error probability to be 0.25 and individual probabilities of each of the 5 types of errors to be 0.2. In that case, given a text file, 25% of the words (randomly chosen) will be misspelt by any of the 5 equally likely methods.

### 4.2  Automatic Speech Recognition System

We used the automatic speech recognition system developed by IBM Research [1] for generating ASR versions of documents. The acoustic models of the system were built using about 100,000 utterances by 500 speakers which amounted to about 120 hours of speech data. Viterbi alignment was used to generate the labeled vectors for building the initial phone models. A forward-backward algorithm [10] was used to train the HMMs for each arc of a phone. For acoustic front-end processing, 13-dimensional cepstral vectors, each representing a 25 msec duration of speech at every 10 msec were used. First and second-order derivatives are used to capture the dynamics of speech variation and hence a 39-dimensional vector is used to represent speech in the cepstral domain. 9 frames (four previous and four forward frames) of cepstral vectors were concatenated. This forms a 117-dimensional vector on which dimensionality reduction algorithm (LDA) was applied to form a 39-dimensional vector. The Language Model has been trained on a text corpus of 10 million words that represents text from different categories. It consists of a trigram model with an open vocabulary and an unknown word probability of 0.00025.

# 5 Experiments

This section describes our detailed experimental evaluation considering the various aspects of noise. We evaluate the performance of standard text classification algorithms on multiple datasets in different settings. We use *rainbow* from the BOW toolkit [16] for multinomial naive Bayes (NB) classifiers and SVMLight [11] for Support Vector Machine (SVM) classifiers. These two represent the spectrum of generative and discriminative models respectively.

## 5.1 Datasets

We now describe the datasets used in our evaluations. We used real-life datasets from a few contact centers and created some synthetic datasets from benchmark text classification datasets by injecting noise. The objective was to see the variation of classification performance with noise on synthetic as well as validating the propositions on real-life datasets. For each dataset used, we summarize its domain and statistics. We also characterize and/or quantify the noise present or introduced (as described in Section 4)

### 5.1.1 Real-life data

**Contact Center agent summaries:** This dataset is collected from a contact center for a telecommunications company. It contains call-logs for around 25,000 customer calls made to the contact centers; for each call, there are some structured data fields, plus a summary of the call content typed in by a human agent. Each call is also manually classified into two categories; a high-level (chosen from amongst 7 categories), and the other a more precise marker of the complaint (chosen from amongst 100 categories). In this dataset, noise is naturally introduced by the human agents when entering in the call summaries, since these have to be done under great time pressure. Thus the summaries contains many spelling errors and abbreviated forms of words. We denote this dataset **CC-Sum**. An example summary: `(Agent1) /01/06/2005/ - SPK TO (CustName) BILL NOT RECD (PhoneNo) THE COMMUNICATED SLA TO SUBSCRIBER IS 02/06/2005 05:46:00 PM (Place1)COURIER (Place2)/2/6/2005 -D BILL DELVERD & RECIVED BY (Recepient) DATE 02/06/2005......(Agent2)` [3]

**Contact Center customer feedback:** This dataset is collected from multiple contact center business processes for various kinds of companies, such as those offering telecommunications, eCommerce and web services. It contains nearly 10,000 customer feedback records from each of the 3 different business processes; each record has

multiple fields which are entries from a feedback form filled in by a user after concluding an interaction with an agent at the contact center. The key field is the *verbatim*, which is free-form text, and is used by human labelers to classify the customer's complaint under one of a set of around 10 to 40 categories indicating the broad reason for the customer's dissatisfaction. Example categories include *Communication problems* (where the customer is the not happy with say the agent's accent), *Resolution problems* (where the customer complains that her query was not resolved), or *Time Adherence problems* (where the customer complains about the long delays in resolving her issues). In this dataset, there is substantial noise arising out of various spelling and grammatical errors made by customers while filling up the feedback forms. In addition to this, there is also significant *label noise* (i.e., the labels assigned by the human labelers are inconsistent in their definitions), due to substantial vagueness and overlap in the way the semantics for each category are set out. We will discuss about label noise in more detail after presenting the results. We denote this dataset **CCFb**. An example feedback: `Help me identlfy how to verify if my , Want It Now, request was actually posted (could'nt find it in , My mail)`

**Contact Center email:** This dataset is collected from the contact center e-mail process for a financial services company. It contains records of about 30,000 email interactions between customers and the contact center agents. Based on the initial e-mail sent by the customer, each interaction is classified by a human agent into one out of over a hundred different categories, indicating the precise nature of the customer's communication. In this dataset, there is some noise due to typographical and other kinds of errors made whilst typing e-mails. We denote this dataset **CCMail**. An example customer email: `I am moving to (Place1) from (Place2) as i am going to join in FIG commodities division of (BankName) center office.Please send all my statements to the address which i shall confirm u before next week end.  If possible please send a statement dated 24th january by mail to this mail id or to the following address where my parents resides for this jan only.`

### 5.1.2 Benchmark Datasets

**Reuters-21578:** This text classification benchmark dataset is collected from Reuters newswire articles[4]. It contains news articles from different subject categories; articles may belong to multiple categories. The 10 most populated classes of this dataset are typically chosen in

---

[3]We have encrypted some of the confidential details and put inside parenthesis.

[4]Available at `http://www.daviddlewis.com/resources/testcollections/`

literature for supervised learning experiments. We also choose the 90 class subset of this dataset; classes chosen have at least one training and one test document. We denote these sets as R10 and R90. These R10 and R90 subsets of the dataset have emerged as well accepted standards for experiments among researchers.

In this dataset, the base level of noise is virtually zero, since the articles have been revised and proof-read. So, in order to estimate the effect of noise, we artificially introduce varying levels of noise in the data and see how it affects the accuracy of automatic classification. Two kinds of artificial noise are introduced: spelling errors as described in Section 4.1, ranging from 0-100% of the words in the corpus; and noise introduced due to ASR transcription as described in Section 4.2 (these transcriptions were generated only for a subset of 200 documents; 20 from each of the top 10 classes). These generated transcripts are made available for download for the benefit of the noisy text analytics research community[5]. Figure 2 shows an example from the R10 test set a document changing with varying amounts and types of noise.

**20-newsgroups:** This text classification benchmark dataset is collected from on line newsgroup postings; there are about 20,000 documents evenly distributed across the 20 newsgroups[6]. In this dataset, the level of noise is quite low; these postings are typically more carefully written and revised than any of the other real-life datasets mentioned above. Here too, we introduce artificial noise to see how it affects accuracy. We denote this dataset by 20NG.

## 5.2 Results

We report results of our experimental study in this subsection. All results are using the NB and SVM classifiers on specified train-test splits. In a classification problem, the classification system is trained on the training data and effectiveness is measured by accuracy on test data which is the fraction of correctly predicted document–class mappings. We report micro-averaged accuracy in this section, which is sensitive to the skew in class sizes as against macro-average accuracy. Micro-average accuracy is easily computed by dividing the sum of the diagonal elements of a multi-class confusion matrix by the sum of all the elements of the matrix. Our aim here is not to compare algorithms, models, and their effectiveness; rather we want to study the effect of feature noise in detail. Here we do not report other effectiveness measures like precision, recall and $F1$, though we would like to mention that results are similar.

We first report results on the R10 dataset, a benchmark standard with various kinds and amounts of feature noise introduced in the text as described earlier in this paper. In Figure 1 we see the accuracy of the test set containing varying amounts of artificially introduced feature noise as described in Section 4.1. The NB classifierused here was trained on the original training set without any introduced noise. We conducted experiments with the test set corrupted with $0\%$ to $100\%$ noise (in steps of $10\%$); for brevity we report results only at $0, 40, 70, 100\%$ noise.
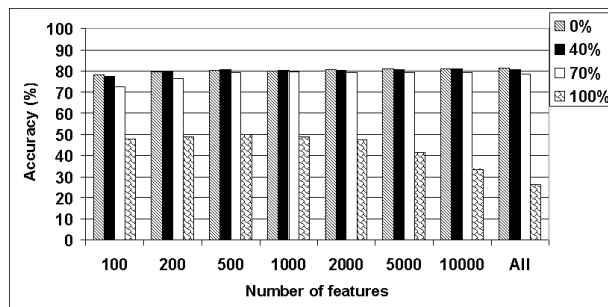


**Figure 1. R10 - trained on clean, tested on noisy**

To our surprise we see that even at $40\%$ noise (emphasizing, on an average 4 out of every 10 words are misspelt), there is little or no drop in accuracy for different numbers of features selected by information gain. The accuracy drops at $70\%$ noise, though only slightly. The accuracy drops significantly at $100\%$ noise – at this level of noise, every word in the test corpus has a spelling error, rendering these words very different from those encountered during training. For this dataset, we also ran SVMs in one-vs-others configuration and achieved very good accuracy numbers. As per traditional use of SVMs we did not perform feature selection and left learning of feature weights to SVM's optimizer. At $0, 40, 70, 100\%$ test noise, the accuracies were $86.2\%$, $85.1\%$, $81.4\%$, and $39.3\%$ respectively – the absolute numbers being higher than NB as per traditional text classification wisdom.

In Figure 3 we repeated the above experiment with the difference that noise (of varying degrees) was also introduced in the training set. The previous experiment is justified in the setting that clean training data for a setting might be available (it is possible to expend resources to build clean domain models), while data to be classified during deployment or testing may be noisy. The current experiment tries to ascertain if there are consistent patterns in the noise that may be learned to help in classification. As we see from the figure, this is not true. Noisy training data leads to worse

```
Original:     Sumitomo Bank Ltd is certain to lose its
status as Japan's most profitable bank as a result
of its merger with the Heiwa Sogo Bank, financial
analysts said.
40% noise:  sumitomo bank ltd is certain to lose its
stxtus as Japan's mozt profitable babk as a ressult of
its merger with the heiwa sogo bank fianncial analysts
said
70% noise:    sumitomo bakn ld is certan to loes is
satus as Jpaan's mpst profitbale bank as a reqult
of its meregr with thye heiwa sogo bakn financial
analystrs sazid
100% noise:   sumitomo bnk ld is ceetain to loes is
sta6us as Japan's mst proifitable bagk as a rexult
of igs mergfer wih thye heiuwa soogo bxnk fnancial
analy5sts sasid
ASR Transcript:   soon is certain lose its status as
chip warns most profitable bank cuts result of its
merger with the high were so woman financial analysis
said
```

**Figure 2. Snippet of a reuters document with varying amounts and types of noise**

off models leading to slightly lower accuracies. This is not unexpected, however, once again, the relationship with the amount of noise in training and test data is interesting. Feature selection proves to be very important in this case. Note how even $40\%$ noise leads to low accuracies at the sub-optimal (small) number of features. At about $5000 - -10000$ features, even $70\%$ noise leaves enough patterns to learn in the training data. One observation comparing these results to the previous set, is that even at $100\%$ noise the accuracy degradation is graceful. We suspect this has to do with the similar nature of noise creeping in during training in this experiment.

For this setting, the four accuracy numbers for SVM were $86.2\%$, $86.4\%$, $84.8\%$, and $83.5\%$. Note that even at $100\%$ training and test noise, SVMs essentially learnt the random pattern in the noise (similar corruptions of short words) for classification.

The second kind of noise that was introduced for this dataset was that caused due to errors made by an Automatic Speech Recognition (ASR) system, as described in Section 4.2. The objective of this experiment was to see effect of ASR (a different kind of noise compared to spelling errors).

A fair comparison can only be done if we create a parallel corpus for which we already have classification accuracy numbers on the clean dataset. The models trained on the training set were then tested both on the original subset of 200 documents, and on the set of their ASR transcripts.
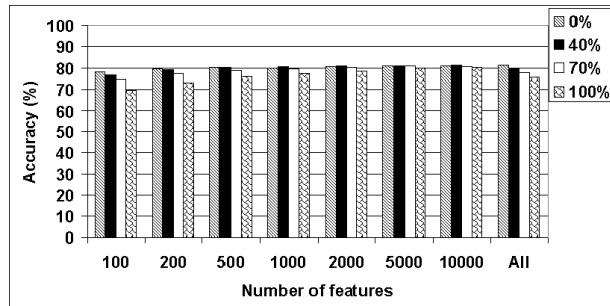


**Figure 3. R10 - trained on noisy, tested on noisy**

The results are shown in Figure 4.

The accuracy of an ASR system is commonly measured as Word Error Rate (WER), which is derived from the Levenshtein distance [13] and works at the word level instead of the character level. WER can be computed as

$$WER = \frac{S + D + I}{N} \qquad (1)$$

where S is the number of substitutions, D is the number of the deletions, I is the number of the insertions, and N is the number of words in the reference.

In this case, even though the word error rate is very high at 66.67%, there is evidently only slight drop in accuracy. This suggests that enough of the key discriminating features between classes get retained in the transcripts, even as a lot of rarer and less relevant words may be corrupted.
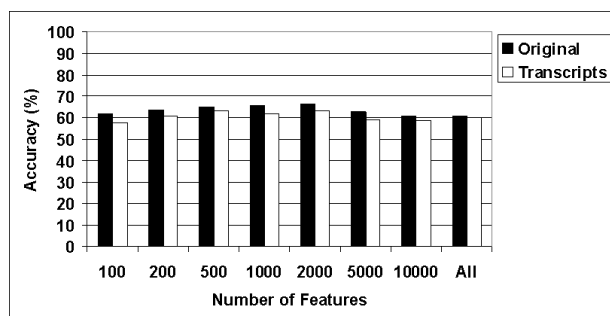


**Figure 4. R10 - trained on clean, tested on transcripts**

These experiments clearly show that text classification does not seem to be very susceptible to feature noise as long as the corpus is large. For small corpora, clearly even a little

noise will disturb the training and test distributions significantly, violating classification's central assumption of similar train and test distributions. These experiments prompted us to investigate the exact relationship between noise, abundance of common features, statistical feature selection, and sparsity of the text classification vector space. *We will return to this investigation after summarizing results for all the other datasets.*
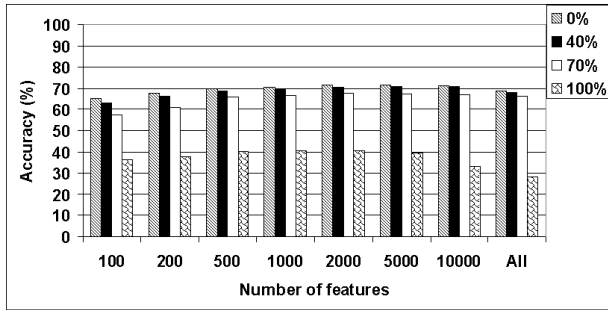


**Figure 5. R90 - trained on clean, tested on noisy**

In Figures 5 and 6, the same two experiments were performed with the Reuters 90-class subset dataset. For this dataset too, the observations were similar. When clean training data was used, there was only a small drop in accuracy at $40\%$ noise; the drop became prominent at $70\%$ and $100\%$ noise as expected. This is also consistent with our discussion above with the R10 dataset and other effects like importance of feature selection when noise is present during training. Again, as is well known, for this dataset too SVM outperformed NB in terms of accuracy. For clean training data, the noisy test accuracies (noise at $0, 40, 70, 100\%$) were $85.6\%, 82.9\%, 77.7\%,$ and $38.3\%$. For noisy training data, the noisy test accuracies were respectively $85.5\%,$ $82.5\%, 79\%,$ and $75.9\%$.

Figures 7 and 8 show the graphs for the same settings for the 20-newsgroups dataset. Once again our observations are similar – the marked difference being the lower absolute accuracy values. The Reuters data is known to be easy to classify given a few terms while the 20NG dataset is a little more noisy. It covers a broader spectrum of topics and has a wider vocabulary because the articles are newsgroup postings, not reviewed for quality.

The main point we would like to stress in this graph is that achievable accuracy levels vary drastically with the domain in question, irrespective of the noise perceived to be present in the domain's documents. It would seem that agent summaries of contact center interactions would be the noisiest to classify since they are written under severe
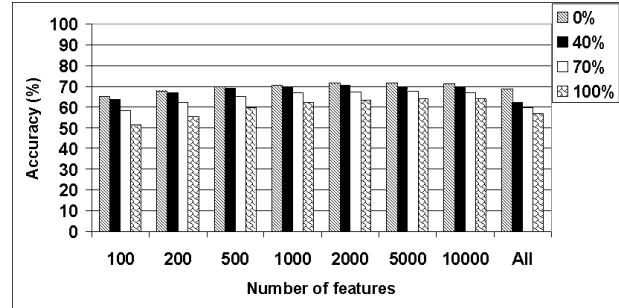


**Figure 6. R90 - trained on noisy, tested on noisy**

time constraints. We achieved text classification accuracy of $85.9\%$ at the first level of the hierarchy of labels (7 categories) and as much as $82.6\%$ accuracy when considering the second level of the hierarchy (100 categories). Accuracy with SVMs for first level touched $88.3\%$.
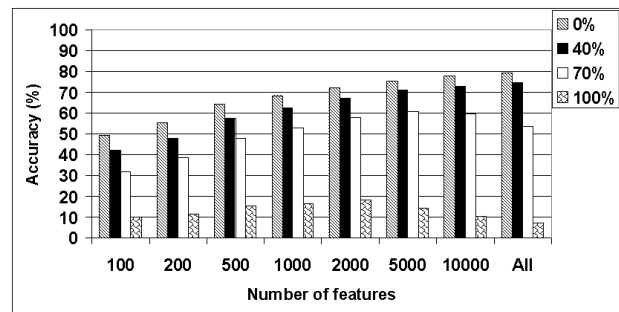


**Figure 7. 20NG - trained on clean, tested on noisy**

We would like to point out an important difference between the classification setting for these datasets against our train-on-noisy and test-on-noisy simulation on the benchmark datasets. In these real-world dataset, noise of at least some kinds tends to be uniform. Customers and agents alike use standard abbreviations and make common spelling mistakes unlike the other situation where spelling errors introduced are random.

The results on these datasets are more instructive, but the best approximation to study such effects in benchmark datasets was to perform experiments in the two settings we described above. We would like to mention that we did not perform hierarchical classification but treated the first and second levels of the hierarchy as flat label-sets. In this do-
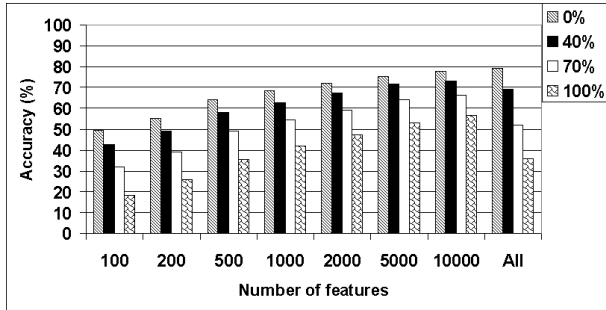
**Figure 8. 20NG - trained on noisy, tested on noisy**

main it was not clear if the hierarchy of labels was constructed for convenience or if it had been factored into designing the label-set. Without losing generality we used the first and second levels of the hierarchy for experiments. We expected the email domain to be the cleanest in terms of quality of language. While this was true, the problem in this domain was the very large number of categories defined. The process of handling email complaints in typical contact centers necessitates on the fly definitions of categories. This left us with over 600 categories. We restricted our attention to only those 50 categories with over a 100 emails. This domain's dataset was not a cleanly defined classification problem. However, we found it instructive to run text classification experiments in this interesting domain from a noise point of view. We achieved 60.1% accuracy with NB for this dataset, and 65.6% with SVMs.
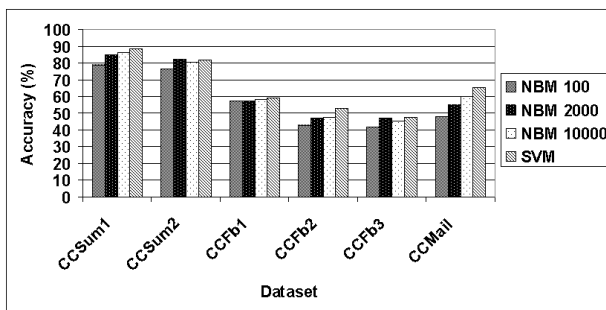


**Figure 9. Real life noisy datasets – accuracy for datasets with NB (100, 2000, 10000 features) and SVMs**

Figure 9 shows the test accuracies of a wide range of real-world noisy text classification datasets. These datasets

have been described and characterized earlier.

The most interesting domain we handled was the contact center customer feedback domain. Feedback to contact centers tends to be short, crisp, and often contains abusive remarks from customers. Many a times the verbatims are very short in length and ambiguous in nature. Also in this domain categories, often known as *call drivers*, may make business sense but seldom have enough data to train models. A harder problem is that the classes defined are often confusing, overlapping, and there is no consistent procedure for labeling comments. This leads to a separability problem to train an automatic classification system driving accuracies down as a whole. For three different client datasets, we got NB accuracies of 58.3%, 47.9%, and 47.6%, and SVM accuracies of 59.1%, 53%, and 47.8%. However the root problem in this domain is not feature noise, which we have been discussing throughout, as much as label noise. With one of the clients, we asked 200 cases to be multi-labeled by two quality analysis domain experts. A week later the same exercise was repeated. A statistical ANOVA Gauge Reproducibility and Repeatability test showed that multi-labeling results were *not reproducible* 53% of the time and the same expert could *not repeat* his own multi-labeling 35% of the time. While multi-labeling has clearly contributed to these very low consistency rates , it points to a larger problem of bad label-set design and the lack of a consistent labeling process. Such an observation is known to some extent to text classification practitioners and about 30% disagreement amongst expert human labelers is accepted [14]. In real-life settings this emerges as a very important kind of noise (label noise) to consider when designing systems. However, we will restrict further discussion on this aspect of noise in this report.

## 5.3 Discussion

In this section we return to inspect the relationship between abundance of terms, sparsity of feature vectors in text classification, statistical (information gain based) feature selection and noise. We noted that corrupting the test set for benchmark datasets like R10 did not lead to large drops in accuracy. This remained true at moderate (40%) and high (70%) levels of noise. We investigated the top 10 most informative features ranked by information gain learned with noisy training data.

In Table 1 we show the top 10 features ranked by information gain with 0, 40, 70, 100% training noise. Note that there is very little difference between the first two sets of features – even 40% training noise finds abundant patterns in the rest of the training data. Even at 70% noise the important words can be still be seen to be occurring though some spelling mistakes (e.g., teh) have now assumed the status of signal-in-the-noise. At 100% noise, as expected,

| Original data | | 40% noise | |
| --- | --- | --- | --- |
| IG | Feature | IG | Feature |
| 0.37063 | lt | 0.22173 | cts |
| 0.27613 | cts | 0.16753 | lt |
| 0.19878 | net | 0.14588 | net |
| 0.16231 | wheat | 0.13415 | wheat |
| 0.14117 | shr | 0.11304 | trade |
| 0.13849 | qtr | 0.10931 | tonnes |
| 0.12909 | trade | 0.10072 | oil |
| 0.12275 | revs | 0.09164 | shr |
| 0.12116 | tonnes | 0.08861 | revs |
| 0.1163 | agriculture | 0.08379 | bank |
| **70% noise** | | **100% noise** | |
| IG | Feature | IG | Feature |
| 0.13281 | cts | 0.10363 | teh |
| 0.09416 | wheat | 0.09123 | cst |
| 0.08846 | trade | 0.0901 | te |
| 0.08594 | tonnes | 0.08862 | cs |
| 0.0852 | teh | 0.07532 | thhe |
| 0.08326 | lt | 0.0622 | nte |
| 0.08104 | te | 0.05835 | ctts |
| 0.07753 | net | 0.05734 | ol |
| 0.07081 | cs | 0.05437 | oli |
| 0.06959 | oil | 0.05046 | tge |

**Table 1. Information gain for most informative features of R10**

| Original data | | 40% noise | |
| --- | --- | --- | --- |
| IG | Feature | IG | Feature |
| 0.09575 | windows | 0.07353 | windows |
| 0.09567 | god | 0.06416 | god |
| 0.08127 | government | 0.06196 | government |
| 0.07828 | dod | 0.05555 | team |
| 0.07013 | team | 0.05331 | people |
| 0.06878 | people | 0.04745 | bike |
| 0.06844 | writes | 0.04745 | game |
| 0.06525 | bike | 0.04671 | jesus |
| 0.06158 | car | 0.04631 | dod |
| 0.06039 | encryption | 0.04562 | encryption |
| **70% noise** | | **100% noise** | |
| IG | Feature | IG | Feature |
| 0.04439 | windows | 0.04072 | gdo |
| 0.04091 | government | 0.03853 | gd |
| 0.03713 | god | 0.02982 | pc |
| 0.03703 | people | 0.02931 | thta |
| 0.03467 | team | 0.02874 | taht |
| 0.03406 | israel | 0.02778 | tat |
| 0.03395 | game | 0.02685 | nto |
| 0.03273 | gdo | 0.0262 | tht |
| 0.03252 | jesus | 0.02487 | cra |
| 0.03197 | bike | 0.02452 | te |

**Table 2. Information gain for most informative features of 20NG**

all words are mangled, and short words (with higher chance of similar corruption due to abundance) emerge as discriminative features. Note the sharp drop in information gain absolute values as noise increases. These numbers are comparable since they are over the same training corpus and document labeling – only feature noise has been introduced in the form of spelling errors. The drop in information is expected because a lot of information is lost as at $40\%$ and $70\%$ noise there is that much probability that each word in the corpus is corrupted. However the abundance of important words repeatedly throws up similar information gain rankings even at high degrees of noise.

Table 2 shows the similar table for the 20NG datasets. Note the consistent drops in the comparable information gain values. The most important feature at higher degrees of noise has even lesser information compared to the $10^{th}$ and even lower ranked features ranked in the clean data.

Coupled with our discussion on label noise in the previous section of real life noisy text classification domains, our observations lead us to believe that feature noise is an important aspect to consider while designing and implementing an operational text classification system. However there are multiple points to consider while designing the systems. An abundance of important features is important in learning robust text classification models. If such an abundance can be confirmed then feature selection needs to be executed

carefully – since the state-of-the-art accuracy achievable on the dataset at hand will be quickly estimated using simple NB models. Consistent with traditional wisdom, SVMs outperform NB, but require more training time and tuning.

The most care needs to be spent in actually tackling label noise, designing a good separable set of classes, and setting up a consistent data labeling process. Feature noise seems to have limited effect in text classification and it can be effectively countered with known feature engineering and feature selection techniques coupled with the choice of a robust classification model.

## 6    Conclusion and Future Work

In this paper, we have studied various aspects of noise as it affects automatic text classification systems. We performed a detailed experimental study introducing spelling and ASR noise in benchmark datasets to study effect on accuracy. The most interesting observation we made for benchmark datasets was that introducing as much as $40\%$ textual noise (spelling mistakes) in documents did not affect text classification accuracy by more than a couple of percentage points. As a contribution to the noisy text analytics community we are making the ASR subset of the Reuters dataset available for further research. We performed experiments on many real-world CRM domains capturing a broad

spectrum of noise (call summaries, customer emails, feedback forms). One of the most striking observations in these real-world datasets was the stark presence of label noise and the pressing need to properly design a separable, non-confusing label-set.

We are intrigued by some of the findings of our experiments. We would like to continue text classification studies with other kinds of noise like time-constrained summaries (of benchmark corpora) of documents. We believe such and other scenarios will be emergent with the growing customer focus of businesses and the ever-growing amount of information present in the real world. We would also like to cover a broader spectrum of real-life noisy datasets depending on their availability for the text classification task.

# References

[1] L. R. Bahl, S. Balakrishnan-Aiyer, J. Bellegarda, M. Franz, P. Gopalakrishnan, D. Nahamoo, M. Novak, M. Padmanabhan, M. Picheny, and S. Roukos. Performance of the IBM large vocabulary continuous speech recognition system on the ARPA wall street journal task. In *Proc. ICASSP '95*, pages 41–44, Detroit, MI, 1995.

[2] T. Bayer, U. Kressel, H. Mogg-Schneider, , and I. Renz. Categorizing paper documents. In *Computer Vision and Image Understanding*, volume 70, pages 299–306, Washington, DC, USA, 1998.

[3] D. Blei, A. Ng, and M. Jordan. Latent dirichlet allocation. In *Proc. of NIPS 14*, 2002.

[4] M. Choudhury, R. Saraf, V. Jain, S. Sarkar, and A. Basu. Investigation and modeling of the structure of texting language. In *Proceedings of the IJCAI-Workshop on Analytics for Noisy Unstructured Text Data*, 2007.

[5] F. J. Damerau. A technique for computer detection and correction of spelling errors. *Commun. ACM*, 7(3):171–176, 1964.

[6] S. Douglas, D. Agarwal, T. Alonso, R. M. Bell, M. Gilbert, D. F. Swayne, and C. Volinsky. Mining customer care dialogs for "daily news". *IEEE Transaction on Speech and Audio Processing*, 13(5):652–660, 2005.

[7] Y. Freund and R. Schapire. A short introduction to boosting. In *Japanese Society for AI 14(5), 771-780. 11*, 1999.

[8] P. Haffner, G. Tur, and J. Wright. Optimizing svms for complex call classification, 2003.

[9] P. Haffner, G. Tur, and J. H. Wright. Optimizing svms for complex call classification. In *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pages 632–635, 2003.

[10] F. Jelinek. *Statistical methods for speech recognition*. MIT Press, Cambridge, MA, USA, 1997.

[11] T. Joachims. Making large-scale support vector machine learning practical. In A. S. B. Schölkopf, C. Burges, editor, *Advances in Kernel Methods: Support Vector Machines*. MIT Press, Cambridge, MA, 1998.

[12] K. Kukich. Technique for automatically correcting words in text. *ACM Comput. Surv.*, 24(4):377–439, 1992.

[13] V. I. Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. Technical Report 8, 1966.

[14] D. Lewis, Y. Yang, T. Rose, and F. Li. Rcv1: A new benchmark collection for text categorization research, 2004.

[15] D. D. Lewis, R. Ghani, D. Mladenic, I. Moulinier, and M. Wasson. In *3rd Workshop on Operational Text Classification (OTC), in conjunction with SIGKDD*, 2003.

[16] A. K. McCallum. Bow: A toolkit for statistical language modeling, text retrieval, classification and clustering. http://www.cs.cmu.edu/ mccallum/bow, 1996.

[17] G. Mishne, D. Carmel, R. Hoory, A. Roytman, and A. Soffer. Automatic analysis of call-center conversations. In *Proceedings of ACM Conference on Information and Knowledge Management (CIKM)*, pages 453–459, 2005.

[18] K. Nigam, J. Lafferty, and A. McCallum. Using maximum entropy for text classification. In *Proc. of IJCAI-99 Workshop on Machine Learning for Information Filtering*, 1999.

[19] G. Ramakrishnan, K. P. Chitrapura, R. Krishnapuram, and P. Bhattacharya. A model for handling approximate, noisy or incomplete labeling in text classification. In *Proc. of ICML*, 2005.

[20] S. Roy and L. V. Subramaniam. Automatic generation of domain models for call centers from noisy transcriptions. In *Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the ACL (COLING/ACL)*, pages 737–744, 2006.

[21] A. Vinciarelli. Noisy text categorization. In *ICPR '04: Proceedings of the Pattern Recognition, 17th International Conference on (ICPR'04) Volume 2*, pages 554–557, Washington, DC, USA, 2004. IEEE Computer Society.

[22] J. Zhang and Y. Yang. Robustness of regularized linear classification methods in text categorization. In *Proc. of SIGIR*, 2003.

[23] G. Zweig, O. Shiohan, G. Saon, B. Ramabhadran, D. Povey, L. Mangu, and B. Kingsbury. Automatic analysis of call-center conversations. In *Proceedings of IEEE International Conference of Acoustics, Speech and Signal Processing (ICASSP)*, pages 589–592, 2006.