

# コーパスを用いた中国語ネット語の判定システム

寶 梓瑜 (東京農工大学 工学府 情報工学専攻)  
古宮 嘉那子 (東京農工大学 工学研究院 先端情報科学部門)  
小谷 善行 (東京農工大学 工学研究院 先端情報科学部門)

## A Detection System of Chinese Netspeak Using Text Corpus

Ziyu Dou (Graduate School of Engineering, Tokyo University of Agriculture and Technology)  
Kanako Komiya (Institution of Engineering, Tokyo University of Agriculture and Technology)  
Yoshiyuki Kotani (Institution of Engineering, Tokyo University of Agriculture and Technology)

### 1. 研究背景

現在、中国では、インターネット利用者が爆発的に増え、それに伴って大量な中国語ネット語（以下、ネット語と書く）が現れた。ネット語の独特な言葉やその使い方は、インターネットだけではなく、徐々に人々の生活にも浸透してきている。しかし、中国の人口はおよそ13.5億であるのに対して、中国のインターネット利用者は5.13億と言われており、中国人の半分以上はインターネットを利用していない。そのような人々にとって、ネット語は理解しにくく、意味が分からなかったり、または意味の誤解から、トラブルになったりすることがある。こうした事態を避けるため、我々はコンピューターで自動的にネット語か書き言葉かどうかを区別するシステムを作成した。本システムは任意の中国語の一つ以上の文の入力に対して、ネット語であるかどうかの判断結果を出力する。

### 2. 中国語ネット語特徴の検出システムの構成

中国語ネット語特徴検出システムの構成を、図1に示す。

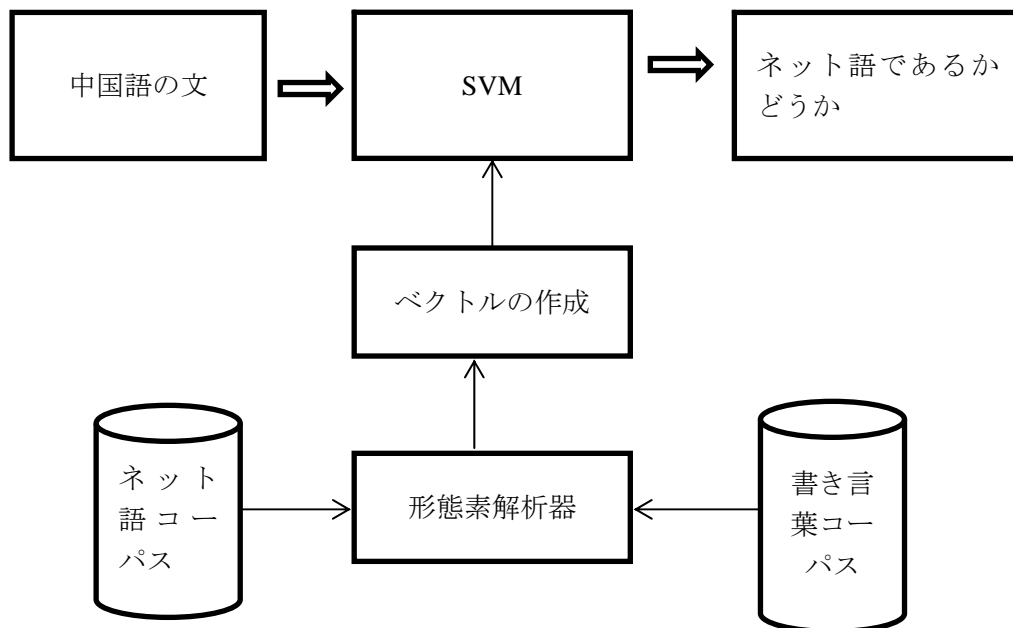


図1 中国語ネット語特徴の検出システムの構成

システムの入力は中国語文であり、出力は、ネット語であるかどうかである。まず、ネット語コーパスと書き言葉コーパスの形態素解析を行い、それを特徴としてベクトルを作成し、サポートベクターマシン (SVM) ([http://otndnld.oracle.co.jp/document/products/oracle11g/111/doc\\_dvd/datamine.111/E05704-02/ algo\\_svm.htm](http://otndnld.oracle.co.jp/document/products/oracle11g/111/doc_dvd/datamine.111/E05704-02/algo_svm.htm)) を用いて機械学習を行う。次に、入力 of 中国語文に対して同様に形態素解析を行い、ベクトルを作成し、先ほど学習した学習器を利用して、入力 of 中国語文がネット語かどうかを判定する。以下に各部分について順次述べる。

## 2. 1 形態素解析

形態素解析の部分は ICTCLAS という中国科学院計算技術研究所の無料形態素解析プログラム (<http://ictclas.org>) を使用した。このプログラムは HMM を基本解析方法として構成されている。プログラムは C 言語版と C++ 言語版が存在するが、本システムで使われているのはその C++ 言語版である。なお、この形態素解析プログラムの基本機能は単語分割、品詞つけ、未知語表記で、ユーザ自身でも辞書に単語を入れることが出来る。

この ICTCLAS プログラムの役割は、書き言葉コーパスの文、ネット語コーパスの文及び入力 of 文の形態素解析を行い、品詞タグを付けることである。また、ここで使用される品詞タグには、計算所漢語詞性標注集 (<http://icl.pku.edu.cn>) が使われている。

## 2. 2 ベクトルの作成と SVM

SVM には SVM-light という無料プログラムを使用した。SVM の入力はベクトルであるため、前処理としてベクトルを作成した。まず、ネット語コーパスと書き言葉コーパスの形態素解析の結果を合わせ、出現した全ての単語を統計し、単語毎に番号を付ける。ただし、ここでは、単語が同じでも、品詞が違う場合には、異なる単語として違う番号を付与した。特にここで統計するとき、アルファベットを除くため、タグ /x が付いている単語を全部除いた。ベクトル作成は、コーパスの文ごとに行った。素性は単語であり、素性値はコーパス中の頻度である。

## 3. 実験と結果

### 3. 1 データ

実験用のデータは、全てインターネットから収集したものである。

このうち、ネット語コーパスは、新浪微博 (<http://www.weibo.com>) から収集した。この新浪微博は、現在 2012 年 1 月までに、2.5 億を超えるユーザを持つ、中国最大のミニブログである。このミニブログはツイッターと同様、一発言として入力できるのは 140 文字までという制限がある。政府などの公式機関のユーザも多数あるが、ほとんどの発言はインターネット利用者の日常的な呟きなので、典型的なネット語があると考えられる。

これらのインターネットユーザ発言を収集するとき、火車採集器 (<http://www.locoy.com>) という無料ウェブデータ収集プログラムを利用し、無作為にユーザを選択してデータを取得した。最初に収集されたデータは既に HTML タグを全部除いた文である。このような文が一行一文という形で、テキストの中で記録されている。新浪微博から収集したネット語コーパスの文の数は 5000 文である。このほかに、百度貼バ (BBS サイト) (<http://tieba.baidu.com>) から 100 文を、ネット語のコーパスとしてテストに利用した。

書き言葉コーパスは前述のように中国国家文字委員会の現代中国語コーパス (<http://www.cncorpus.org/>) の中の新聞と社論というカテゴリのコーパスから取ったも

のである。このコーパスは、ネット語コーパスと同じく一行一文でテキストの中に記録されている。現代中国語コーパスから利用する文の数は2000文である。このほかに、SOHU ニュース (<http://www.sohu.com/>) から100文を、書き言葉のコーパスとしてテストに利用した。

### 3. 2 実験設計及び結果

ネット語コーパスと書き言葉コーパスのどちらの文かを判定する制度を見るために、実験1～実験4の四種類の実験を行った。以下にそれぞれについて述べる。

#### (1) 実験1 CLOSED テスト

まず、ネット語コーパスと書き言葉コーパスのどちらの文かを判定する制度を見るためのCLOSEDテストを行った。CLOSEDテストでは、テストデータは訓練データとして利用したものである。また、CLOSEDテストでは、顔文字などの符号を削除していない(略: 符号あり)実験を行った。以下に、CLOSEDテストにおける、訓練データ、テストデータの種類と数、また結果の求め方を説明する。さらに実験の結果を表1に示す。表1から、CLOSEDテストでは、全ての文がネット語であるかどうか、正しく判定できていることが分かる。

訓練データ：ネット語コーパス 5000文、書き言葉コーパス 2000文

テストデータ：ネット語コーパスからの1000文と書き言葉コーパスからの400文

結果：ネット語の率=SVMがネット語として認識した文/1000

書き言葉の率=SVMが書き言葉語として認識した文/400

ここで、ネット語におけるネット語の率はネット語の再現率であり、書き言葉における書き言葉の率は書き言葉の再現率となる。

表1 CLOSED テスト 結果

	ネット語の率	書き言葉の率
ネット語コーパスからの1000文 (符号あり)	100%	0%
書き言葉コーパスからの400文 (符号あり)	0%	100%

#### (2) 実験2 OPEN テスト (符号あり及び符号なし)

次に、訓練データとテストデータの重複を許さないOPENテストの二つを行った。OPENテストでは、符号と英文字によって構成され顔文字や略語の結果への影響を実証するため、符号および英文字がある実験と符号および英文字がない(略: 符号なし)実験を行って比較した。以下に、OPENテストにおける、訓練データ、テストデータの種類と数、また結果の求め方を説明する。さらに実験の結果を表2に示す。

訓練データ：ネット語コーパス 4000文、書き言葉 1600文

テストデータ：ネット語コーパスの上記4000文を除いて残った1000文、書き言葉

コーパスの上記1600文を除いて残った400文  
 結果：ネット語の率=SVMがネット語として認識した文/1000  
 書き言葉の率=SVMが書き言葉語として認識した文/400

表2 OPENテスト 結果

	ネット語の率	書き言葉の率	正解率
ネット語コーパスからの1000文(符号あり)	98.4%	1.6%	92.6%
書き言葉コーパスからの400文(符号あり)	22.0%	78.0%	
ネット語コーパスからの1000文(符号なし)	98.9%	1.1%	84.9%
書き言葉コーパスからの400文(符号なし)	50.2%	49.8%	

(3) 実験3 ネット語100文と書き言葉100文 テスト(符号あり及び符号なし)  
 次に、訓練データとして使用したコーパスとは異なるコーパスとして、ネット語コーパスにBBS、書き言葉コーパスに新聞を利用した際のOPENテストを行った。以下に、その際の訓練データ、テストデータの種類と数、また結果の求め方を説明する。さらに実験の結果を表3に示す。この実験の際にも実験2と同様、符号ありと符号なしの実験を行って比較した。

訓練データ：ネット語コーパス 5000文、書き言葉コーパス 2000文  
 テストデータ： 百度貼バ (BBSサイト) から取った100文をネット語とし、SOHUニュース から取った100文を書き言葉とし、テストを行う  
 結果：ネット語の率=SVMがネット語として認識した文/100  
 書き言葉の率=SVMが書き言葉語として認識した文/100

表3 ネット語100文と書き言葉100文 テスト 結果

	ネット語の率	書き言葉の率	正解率
ネット語100文(符号あり)	87%	13%	71.5%
書き言葉100文(符号あり)	44%	56%	
ネット語100文(符号なし)	69%	31%	83.5%
書き言葉100文(符号なし)	22%	78%	

#### (4) 実験4 アンケート (符号ありと符号なし)

最後に、比較対象として、人間に符号ありと符号なしの際、どの程度ネット語を判定できるかのアンケートを行った。以下に、アンケート実験における、訓練データ、テストデータの種類と数、また結果の求め方を説明する。さらに実験の結果を表4に示す。この実験の際にも実験2、実験3と同様、符号ありと符号なしの実験を行って比較した。

テスト方法：数人の中国人インターネット利用者にアンケート

テスト内容：上記の 百度貼バ (BBS サイト) から取ったネット語50文を、SOHU ニュース から取った50文を符号ありと符号なし二回、被験者に判断してもらった。顔文字で簡単に人間がネット語を判断し、その結果を覚えてしまう可能性を排除するため、先に符号なしのアンケートを行い、その後符号ありのアンケートを行った。また、

結果：ネット語の率=SVM がネット語として認識した文/50

書き言葉の率=SVM が書き言葉語として認識した文/50  
として計算した。

表4 ネット語50文と書き言葉50文の人工判定テスト結果

	ネット語の率 (符号あり)	書き言葉の率 (符号あり)	正解率 (符号あり)	ネット語の率 (符号なし)	書き言葉の率 (符号なし)	正解率 (符号なし)
20代中国人男性学生	44%	100%	72%	44%	100%	72%
20代中国人女性学生	20%	100%	60%	16%	100%	58%
平均	32%	100%	66%	30%	100%	65%

#### 4. 考察

まず、表2のOPENテストの結果が示すように、本研究のネット語コーパスおよび書き言葉コーパスは、確実に区別が存在する。特にネット語コーパスは、符号ありと符号なしの場合、それぞれ、98.4%と98.9%の再現率となった。しかし、書き言葉コーパスに対するOpenテストは、符号ありのとき、再現率が78.0%で、符号なしの場合の再現率は49.8%まで下がった。これは、訓練に使用したネット語コーパスの量が多いためと、ネット語コーパスの中でも、書き言葉のような文が多数存在するためであると思われる。また、符号がある場合とない場合の再現率の差から、機械学習において、符号の影響が大きいことが分かる。

続いて、実際の文に対するテストの結果(表3)を分析する。まず、符号がある場合とない場合と比べると、ネット語100文に対する認識正解率は18ポイント上がり、87%まで達成した。これに対し、書き言葉は符号なしの場合のほうが22ポイント上回り、78%まで達成した。これによって、符号がある場合、文がネット語として認識される傾向が強まり、符号がない場合には、書き言葉として認識される傾向が強まることが分かる。

アンケート(表4)の結果を見ると、全部書き言葉の判定は100%正解したが、ネット語の判定はいずれも44%と20%まで止まったことが分かる。また、符号がある場合の、符号がない場合と比べた正解率の上昇はわずかであった。実験4と実験3の正解率を

比べると、機械のほうが、正解率が上回ることがわかった。とこれは、ネット語といっても、BBS では、書き言葉的な表現も多数存在することが、判定の結果に大きく影響したためだと思われる。それに対し、書き言葉は、100%の正解率で、人間が書き言葉を認識するのは簡単だったことが分かる。ネット語の定義を人間が判断できるものとするれば、再現率にも変化があるだろう。

最後に表3と表4から機械学習と人間の正解率を比べる。表3から、機械学習は最高83.5%、表4から、人間の判断は最高66%であるため、作成したシステムの性能が人間に上回ることがわかる。

## 5. 結論

本論文では、文の入力に対してネット語かどうかの判定を行うシステムを作成した。入力文は形態素解析を行い、ベクトル化したあと、SVMを使ってネット語かどうかを判定した。

実験結果から、本システムは、符号がある場合、ネット語に対する判定の正解率が上がり、符号がない場合、書き言葉に対する判定再現率が上がるということが分かった。また、人間に対するアンケートの結果から、人間でもネット語かどうかの判定が難しいことが分かる。特に、BBS やミニブログなどの情報には、ネット語的な特徴がない言葉も多数存在するので、難しかったようである。また、機械学習と人間の正解率を比べると、機械学習は最高83.5%、人間の判断は最高66%であるため、作成したシステムの性能が人間に上回ることがわかった。

## 文 献

佐藤敏紀 「Perl で自然言語処理」 東京工業大学奥村研究室

<http://www.slideshare.net/overlast/perl-5460697>

谷岡 広樹、丸山 稔(2005)「形態素解析に基づく SVM を用いたアスキーアートの識別」  
電子情報通信学会技術研究報告. PRMU, パターン認識・メディア理解 「104:670」,  
pp.25-30

黒橋禎夫 機械学習に基づく自然言語処理 京都大学情報学研究科

<http://nlp.ist.i.kyoto-u.ac.jp/member/kuro/lecture/LIP10/LIP09.pdf>

Jin'ichi Murakami, HMM(Hidden Markov Model, 隠れマルコフモデル)

<http://unicorn.ike.tottori-u.ac.jp/murakami/doctor/node7.html>

語料庫在線 <http://www.cncorpus.org/>

SVM-light [http://www.cs.cornell.edu/People/tj/svm\\_light/](http://www.cs.cornell.edu/People/tj/svm_light/)

ICTCLAS <http://ictclas.org/>

情報と通信のハイパーテキスト <http://www.yobology.info/text/index.htm>

Shogo Computing Laboratory <http://sora-blue.net/~shogo82148/memo/algorithm/svm/>