



表1 Yahoo!知恵袋に投稿されたすべての回答および3個以上連続する不読符号からなる不読符号列が末尾にある回答について、回答者数、回答件数、ベストアンサーおよびファーストアンサーの件数

	回答者数	回答件数	ベスト アンサー数	ファースト アンサー数
すべての回答	183,242	13,477,785	3,116,009	3,116,009
末尾に不読符号列がある回答	89,133	3,242,694	477,462	927,296

最もよいと判定された回答1件がベストアンサーに必ず選ばれる。また、それぞれの質問に対して最初に投稿された回答を本研究ではファーストアンサーとよぶ。それぞれの質問についてベストアンサーとファーストアンサーは必ず1件ずつあるので、それらの数は質問の数と等しい。なお、1つのアカウントから1件の質問に投稿できる回答は1件だけである。この公開データには、以下に示す時間に関する情報も記録されている。

質問投稿日 質問が最初に投稿された時刻、あるいは最後に更新された時刻

回答投稿日 回答が投稿された時刻

質問最終更新日 回答の受付が締め切れ、ベストアンサーが決定された時刻

「質問投稿日」「回答投稿日」「質問最終更新日」はいずれも日付だけでなく、時刻も秒単位で記録されている。回答が投稿されてから質問が更新されることがあるので、「回答投稿日」の時刻が「質問投稿日」の時刻よりも前である回答が1,706,325件あった。この1,706,325件の回答をYahoo!知恵袋に投稿されたすべての回答13,477,785件から取り除いた回答11,771,460件を対象に調査したところ、質問が投稿されてから回答が投稿されるまでの時間の平均値は9,706秒、中央値は600秒であった。

### 3 末尾に不読符号列がある回答の調査

テキストには、読み上げる時には発音されない記号や符号がある。また、読み上げる時には本来発音される文字であっても、用いられている位置や状況によって、発音されないことがある。このような文字や記号や符号をわれわれは不読符号とよぶことにする。例えばYahoo!知恵袋では、その末尾に3個以上の不読符号が連続して用いられている回答がおよそ4件に1件の割合で投稿されている。これほどさかんに用いられている表現であるのに、この不読符号列についての研究はほとんど行われていない。そこで本研究では、Yahoo!知恵袋に投稿された回答を対象に、その末尾で連続して用いられる場合に発音されないことが多い

- 英数字以外のアスキー文字 (!#\$%&.:;?@{} など)
- 記号 (、 。 , . . : ; など)
- ギリシャ文字
- キリル文字
- 罫線

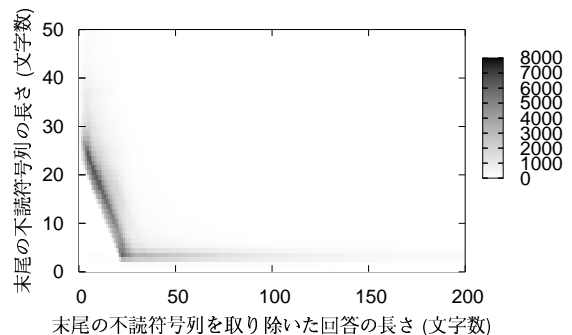


図2 3個以上連続する不読符号からなる不読符号列が末尾にある回答について、回答末尾の不読符号列の長さとその不読符号列を除いた回答の長さの関連性を示すヒートマップ

が3個以上連続して回答末尾で用いられている不読符号列について調査を行う。メールやチャットなどでのコミュニケーションでは、文末の不読符号列は一般にコミュニケーションを円滑に行うために用いられる。例えば(例文3)の末尾の不読符号列は、回答者の意見がおだやかに伝わるようにするため用いられている。

(例文3) サウンドレコーダーでもある程度は出来るけど、やっぱり SoundEngine がお勧めかな。。。

一方、Yahoo!知恵袋では他のユーザとのコミュニケーションの円滑化以外の目的、具体的には、投稿文字数制限を回避するために不読符号列が用いられることがある。投稿文字数制限とは2004年5月27日からYahoo!知恵袋に導入された規則である。この規則によって、全角25文字未満の回答の投稿が禁止された<sup>3</sup>。さらに、全角スペースのみの投稿、一定数以上の連続改行、絵文字(アスキーアートなど)での投稿も禁止された。この規則の適用を回避するため、(例文4)では回答の末尾に13個の「！」が用いられている。

(例文4) アルミホイールに包んで火の中にボン!!!!!!!!!!!!!!

表1は、Yahoo!知恵袋に投稿されたすべての回答と3個以上連続する不読符号からなる不読符号列が末尾にある回答について、回答者数、回答件数、そしてそのうちファーストアンサーであるものとベストアンサーであるものの件数を示す。図2は、3個以上連続する不読符

<sup>3</sup><http://chiebukuro.yahoo.co.jp/docs/whats2004.html>

表 2 3 個以上連続する不読符号からなる不読符号列が末尾にあり、それ以外の長さが 25 文字未満および 25 文字以上の回答について、回答者数、回答件数、ベストアンサーおよびファーストアンサーの件数

末尾の不読符号列以外の回答の長さ	回答者数	回答件数	ベストアンサー数	ファーストアンサー数
25 文字未満	52,998	1,745,797	191,791	616,702
25 文字以上	77,299	1,496,897	285,671	310,594

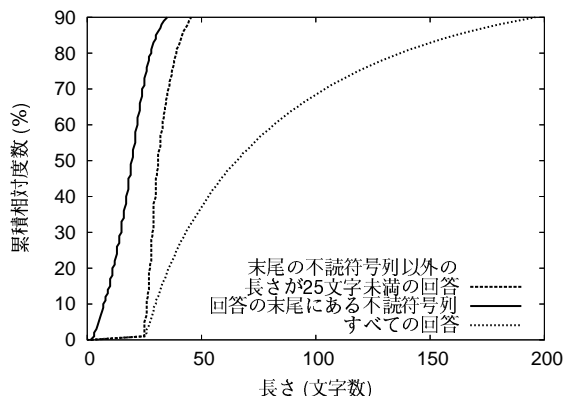


図 3 3 個以上連続する不読符号からなる不読符号列が末尾にあり、それ以外の長さが 25 文字未満の回答について、回答およびその末尾の不読符号列の長さの累積相対度数分布

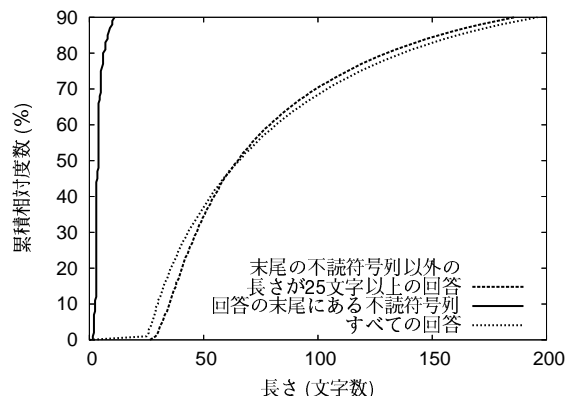


図 4 3 個以上連続する不読符号からなる不読符号列が末尾にあり、それ以外の長さが 25 文字以上の回答について、回答およびその末尾の不読符号列の長さの累積相対度数分布

号からなる不読符号列が末尾にある回答について、回答末尾の不読符号列の長さとその不読符号列を除いた回答の長さの関連性を示すヒートマップである。図 2 で濃い色で示されているのは、回答末尾の不読符号列とそれ以外の回答の長さの組み合わせで出現頻度の高いものである。なお本研究では、全角文字 1 個は 1 文字、半角文字 1 個は 0.5 文字として回答および不読符号列の長さを測定した。

回答末尾の不読符号列の長さの中央値は 10 文字で、(例文 3) の末尾にある不読符号列の長さ比べると 2 倍以上長い。図 2 のヒートマップを見ると、10 文字よりも長い不読符号列は末尾の不読符号列以外の長さが 25 文字未満の回答に集中している。そこで末尾の不読符号列以外の長さが 25 文字未満の回答に注目すると、末尾の不読符号列とそれ以外の回答の長さの和、すなわち、回答の長さが 25 ~ 30 文字である場合の頻度が高い。これは、投稿文字数制限を回避するために回答の末尾に不読符号列がさかんに用いられていることが考えられる。一方、末尾の不読符号列以外の長さが 25 文字以上の回答では、末尾の不読符号列の長さが 3 ~ 4 文字の場合が多く、不読符号列以外の回答の長さはさまざまである。このように図 2 のヒートマップから、末尾の不読符号列以外の長さが 25 文字未満の回答と 25 文字以上の回答とで、不読符号列の用いられ方が異なることが考えられる。そこで、末尾に不読符号列がある回答を

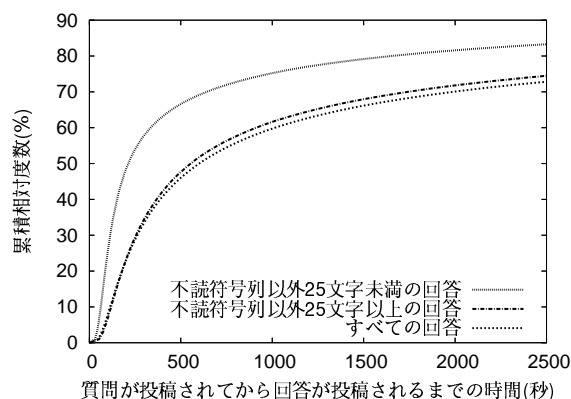


図 5 3 個以上連続する不読符号からなる不読符号列が末尾にあり、それ以外の長さが 25 文字未満および 25 文字以上の回答について、質問が投稿されてから回答が投稿されるまでの時間の累積相対度数分布

- 末尾の不読符号列以外の長さが 25 文字未満の回答
  - 末尾の不読符号列以外の長さが 25 文字以上の回答
- に分けて以下の調査を行った。

- 回答者数、回答件数とそのうちファーストアンサーおよびベストアンサーであるものの件数 (表 2)
- 末尾の不読符号列と回答の長さ (図 3 と図 4)
- 質問が投稿されてから回答が投稿されるまでの時間 (図 5)

末尾の不読符号列と回答の長さの調査 (図 3 と図 4) は、末尾に不読符号列がある回答 3,242,694 件すべてを対象にしている。一方、質問が投稿されてから回答が投稿されるまでの時間の調査 (図 5) は、「質問投稿日」以後に投稿された回答 2,790,352 件のみを対象にしている。

末尾の不読符号列以外の長さが 25 文字以上の回答では、不読符号列は投稿文字数制限の回避のために用いられることはない。したがって、回答の末尾で不読符号列はコミュニケーションを円滑に行うために用いられていると考えられる。そして図 4 から、末尾の不読符号列以外の長さが 25 文字以上の回答の長さの分布が Yahoo!知恵袋に投稿されたすべての回答の長さの分布に近いことがわかる。また図 5 から、質問が投稿されてから回答が投稿されるまでの時間の分布も、Yahoo!知恵袋に投稿されたすべての回答の場合の分布に近いことがわかる。これらのことから、末尾の不読符号列以外の長さが 25 文字以上の回答の場合、末尾の不読符号列がその回答の長さや投稿するタイミングに及ぼす影響は少ないと考えられる。また、末尾の不読符号列以外の長さが 25 文字以上の回答のファーストアンサー率は 20.75%であり、Yahoo!知恵袋に投稿されたすべての回答を対象にした場合のファーストアンサー率 23.12% より低い。したがって、末尾の不読符号列以外の長さが 25 文字以上の回答を投稿する場合、ファーストアンサーを投稿することにこだわるユーザは少ないと考えられる。

一方、末尾の不読符号列以外の長さが 25 文字未満の回答では、不読符号列は投稿文字数制限の回避のために用いられている。ただし、投稿文字数制限の回避のためだけでなく、コミュニケーションを円滑に行うためにも用いられることがある。末尾の不読符号列以外の長さが 25 文字未満の回答の場合、図 5 から、質問が投稿されてから回答が投稿されるまでの時間が短いことがわかる。この原因の 1 つに、一番最初に回答を投稿すること、すなわち、ファーストアンサーを投稿することを楽しむことが目的のユーザがいることが考えられる。ファーストアンサーを投稿するためには、できるだけ早く回答を投稿しなければならない。このため、短い回答に、投稿文字数制限を回避できる長さの不読符号列をつけ、すばやく投稿することが考えられる。(例文 2) や (例文 4) のように、同じ不読符号を連続して用いれば、長い不読符号列の入力も簡単である。そこで 4 章では、末尾の不読符号列以外の長さが 25 文字未満の回答を対象に、ファーストアンサーを投稿することを楽しむことが目的と考えられるユーザが不読符号列を回答の末尾に用いるかどうかについて検討する。

#### 4 末尾に不読符号列がある回答をくりかえし投稿するユーザの回答の分析

ファーストアンサーを投稿することを楽しむことが目的のユーザがいるのなら、そのユーザはファーストアンサーをくりかえし投稿していると考えられる。したがって、ファーストアンサーを投稿するために回答の末尾で不読符号列を用いるユーザがいるなら、そのユーザは末尾に不読符号列があるファーストアンサーをくりかえし投稿していると考えられる。そこで、末尾に不読符号列があるファーストアンサーをくりかえし投稿しているユーザを検出するため、以下の 2 つの仮説を用いる。

仮説 1 ユーザ  $i$  が末尾の不読符号列以外の長さが 25 文字未満の回答を異常に多く投稿していないならば、ユーザ  $i$  は末尾の不読符号列以外の長さが 25 文字未満の回答を  $N_1(i)$  件投稿していると期待できる。

$$N_1(i) = P_1 \times ans(i)$$

ここで  $ans(i)$  とはユーザ  $i$  が投稿した回答の件数で、 $P_1$  とは投稿されたすべての回答の中から無作為に 1 つ選んだ回答が末尾の不読符号列以外の長さが 25 文字未満の回答である確率である。したがって  $P_1$  は

$$P_1 = \frac{N_{\text{不読符号列以外 25 文字未満}}}{N_{ans}}$$

となる。ここで  $N_{\text{不読符号列以外 25 文字未満}}$  とは末尾の不読符号列以外の長さが 25 文字未満の回答の数で、 $N_{ans}$  とは投稿されたすべての回答の数である。本研究では、表 1 と表 2 に示すように、 $N_{\text{不読符号列以外 25 文字未満}}$  は 1,745,797 件、 $N_{ans}$  は 13,477,785 件であるので、 $P_1$  は 0.129531 となる。もしこの仮説が片側二項検定で棄却されれば、ユーザ  $i$  は末尾の不読符号列以外の長さが 25 文字未満の回答を異常に多く投稿していると判定する。

仮説 2 ユーザ  $i$  がファーストアンサーを異常に多く投稿していないならば、ユーザ  $i$  はファーストアンサーを  $N_2(i)$  件投稿していると期待できる。

$$N_2(i) = P_2 \times ans(i)$$

ここで  $ans(i)$  とはユーザ  $i$  が投稿した回答の件数で、 $P_2$  とは投稿されたすべての回答の中から無作為に 1 つ選んだ回答がファーストアンサーである確率である。したがって  $P_2$  は

$$P_2 = \frac{N_{\text{ファーストアンサー}}}{N_{ans}}$$

となる。ここで  $N_{\text{ファーストアンサー}}$  とはファーストアンサーである回答の数で、 $N_{ans}$  とは投稿されたすべての回答の数である。本研究では、表 1 と表 2 に示すように、 $N_{\text{ファーストアンサー}}$  は 3,116,009 件、 $N_{ans}$  は 13,477,785 件であるので、 $P_2$  は 0.231196 となる。もしこの仮説が片側

	回答者数	回答件数	ベスト アンサー数	ファースト アンサー数
グループ 1	2,302	691,357	75,883	334,247
グループ 2	4,442	535,879	52,540	115,654

表 3 グループ 1 とグループ 2 のユーザたちが投稿した末尾の不読符号列以外の長さが 25 文字未満の回答について、回答者数，回答件数，ベストアンサーとファーストアンサーの件数

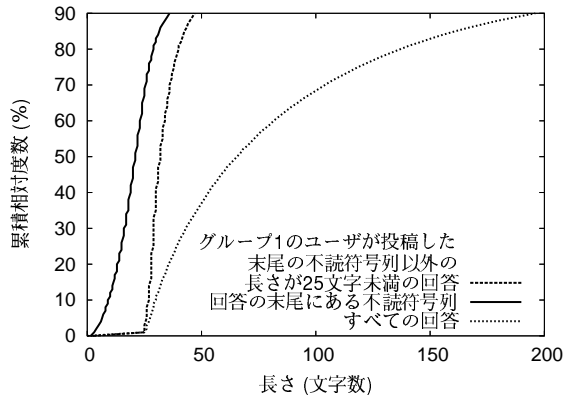


図 6 グループ 1 のユーザたちが投稿した末尾の不読符号列以外の長さが 25 文字未満の回答について、回答とその末尾の不読符号列の長さの累積相対度数分布

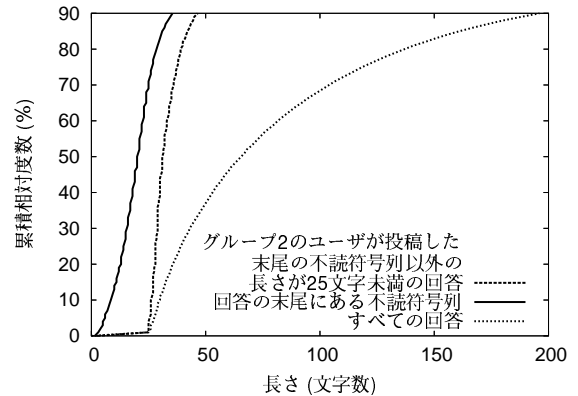


図 7 グループ 2 のユーザたちが投稿した末尾の不読符号列以外の長さが 25 文字未満の回答について、回答とその末尾の不読符号列の長さの累積相対度数分布

二項検定で棄却されれば、ユーザ  $i$  はファーストアンサーを異常に多く投稿していると判定する。

Yahoo! 知恵袋に回答を投稿したすべてのユーザ (183,242 ユーザ) を対象に、まず、仮説 1 を用いて末尾の不読符号列以外の長さが 25 文字未満の回答をくりかえし投稿したユーザを検出する。次に、仮説 1 で検出されたユーザを仮説 2 を用いて以下の 2 つのグループに分類する。

グループ 1 仮説 1 も仮説 2 も棄却されたユーザのグループ。すなわち、末尾の不読符号列以外の長さが 25 文字未満の回答およびファーストアンサーをくりかえし投稿したと判定されたユーザのグループ。ファーストアンサーを投稿するために回答の末尾で不読符号列を用いるユーザがいるなら、そのユーザはこのグループに分類される。

グループ 2 仮説 1 は棄却されたが、仮説 2 は棄却されなかったユーザのグループ。すなわち、ファーストアンサーの投稿はふつうかそれ以下だが、末尾の不読符号列以外の長さが 25 文字未満の回答はくりかえし投稿したと判定されたユーザのグループ。

実験で用いた有意水準は、仮説 1 も仮説 2 も 0.005 である。グループ 1 とグループ 2 のユーザたちを検出した後、それらのユーザたちが投稿した末尾の不読符号列以外の長さが 25 文字未満の回答について、以下の調査を行った。

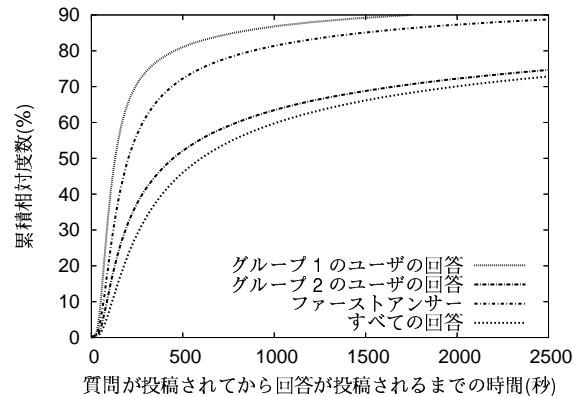


図 8 グループ 1 とグループ 2 のユーザたちが投稿した末尾の不読符号列以外の長さが 25 文字未満の回答について、質問が投稿されてから回答が投稿されるまでの時間の累積相対度数分布

- 回答者数，回答件数とそのうちファーストアンサーおよびベストアンサーであるものの件数 (表 3)
- 回答および末尾の不読符号列の長さ (図 6 と図 7)
- 質問が投稿されてから回答が投稿されるまでの時間 (図 8)

回答の長さや末尾の不読符号列の長さの調査 (図 6 と図 7) は、グループ 1 およびグループ 2 のユーザたちが投稿した末尾の不読符号列以外の長さが 25 文字未満の回答 1,227,236 件すべてを対象にしている。一方、質問が投稿されてから回答が投稿されるまでの時間の調査 (図 8)

は、グループ1およびグループ2のユーザたちが投稿した末尾の不読符号列以外の長さが25文字未満の回答のうち、「質問投稿日」以後に投稿された回答1,049,242件のみを対象にしている。同様に、ファーストアンサーおよびすべての回答もそれぞれ「質問投稿日」以後に投稿されたもののみを対象にしている。なお、(例文2)の回答投稿者と(例文4)の回答投稿者はいずれもグループ1に分類された。

最初に、回答の長さで末尾の不読符号列の長さについて検討する。図6と図7に示すように、グループ1とグループ2のユーザたちが投稿した末尾の不読符号列以外の長さが25文字未満の回答を比べると、その回答の長さで末尾の不読符号列の長さの分布はそれぞれ非常に近い。したがって、回答の長さおよび末尾の不読符号列の長さから、その回答がグループ1のユーザによって投稿されたのか、それともグループ2のユーザによって投稿されたのか、区別することはむずかしい。

次に、質問が投稿されてから回答が投稿されるまでの時間について検討する。図8に示すように、グループ1とグループ2のユーザたちが投稿した末尾の不読符号列以外の長さが25文字未満の回答を比べると、質問が投稿されてから回答が投稿されるまでの時間の分布はグループ1とグループ2では大きく異なる。グループ1のユーザたちが投稿した末尾の不読符号列以外の長さが25文字未満の回答のファーストアンサー率は48.35%であるのに、図8に示すように、質問が投稿されてから回答が投稿されるまでの時間は、ファーストアンサーの場合よりも短い時間に分布している。このことから、グループ1のユーザたちは、末尾の不読符号列以外の長さが25文字未満の回答をすばやく投稿することに強いこだわりがあると考えられる。一方、グループ2のユーザたちの回答の投稿時間の分布は、Yahoo!知恵袋に投稿されたすべての回答の投稿時間の分布に近い。したがって、グループ2のユーザたちが投稿した末尾の不読符号列以外の長さが25文字未満の回答の場合、末尾の不読符号列がその回答を投稿するタイミングに及ぼす影響は少ないと考えられる。この点は、末尾の不読符号列以外の長さが25文字以上の回答の場合と同じである。

最後に、ファーストアンサー率について検討する。表3に示すように、グループ1のユーザたちが投稿した末尾の不読符号列以外の長さが25文字未満の回答のファーストアンサー率は48.35%である。このファーストアンサー率は、末尾の不読符号列がなく、回答の長さが47文字以下のすべての回答2,590,836件のファーストアンサー率27.14%よりもはるかに高い。回答の長さが47文字以下の回答とファーストアンサー率について比較したのは、図6に示すように、グループ1のユーザたちが

投稿した末尾の不読符号列以外の長さが25文字未満の回答の90%は回答の長さが47文字以下だからである。これらのことから、末尾の不読符号列以外の長さが25文字未満の回答を投稿する場合、グループ1のユーザたちはファーストアンサーを投稿することについてのこだわりが強いと考えられる。一方、グループ2のユーザたちが投稿した末尾の不読符号列以外の長さが25文字未満の回答のファーストアンサー率は21.58%である。これは、Yahoo!知恵袋に投稿されたすべての回答を対象にした場合のファーストアンサー率23.12%よりも低い。したがって、末尾の不読符号列以外の長さが25文字未満の回答を投稿する場合、グループ2のユーザたちはファーストアンサーを投稿することについてのこだわりが弱いと考えられる。

## 5 おわりに

Yahoo!知恵袋には、ファーストアンサーを投稿することに強いこだわりがあり、その目的のために短い回答の末尾に長い不読符号列を用いるユーザがいることがわかった。このユーザたちは、ファーストアンサーを投稿することを楽しむことが目的と考えられる。一方、ファーストアンサーを投稿することにはこだわりがなく、投稿文字数制限を回避するために短い回答の末尾に長い不読符号列を用いるユーザがいることもわかった。これらのユーザを検出し区別するには、2つの仮説を用いた提案手法が有効である。

## 参考文献

- [1] 松村真宏, 三浦麻子, 柴内康文, 大澤幸生, 石塚満: 2ちゃんねるが盛り上がるダイナミズム, 情報処理学会論文誌, Vol.45, No.3, pp.1053-1061, 2004.
- [2] 野島久雄: ネットワークにおける感情伝達的手段としての:-)(smily face), 情報処理学会夏のプログラミング・シンポジウム報告集, pp.41-48, 1989.
- [3] Walther, J. B. and Burgoon, J. K.: Relational Communication in Computer-Mediated Interaction, Human Communication Research, Vol.19 Issue 1, pp.50-88, 1992.
- [4] Witmer, D. F. and Katzman, S. L.: On-line Smiles: Does Gender make a Difference in the Use of Graphic Accents?, Journal of Computer-Mediated Communication, Vol.2 No.4, 1997.
- [5] 井上みづほ, 藤巻美菜子, 石崎俊: 電子メール文における感情表現の解析システムについて:感情表現の収集・分類・解析, 電子情報通信学会技術研究報告, TL 思考と言語 96(608), pp.1-8, 1997.
- [6] 原田登美: 「顔文字」による日本語の円滑なコミュニケーション: 「配慮」と「ボライトネス」の表現機能, 言語と文化, Vol.8, pp.205-224, 2004.
- [7] 加藤尚吾, 加藤由樹, 小林まゆ, 柳沢昌義: 電子メールで 사용되는顔文字から解釈される感情の種類に関する分析, 教育情報研究, Vol.22, No.4, pp.31-39, 2007.
- [8] 加藤尚吾, 加藤由樹, 島峯ゆり, 柳沢昌義: 携帯メールコミュニケーションにおける顔文字の機能に関する分析: 相手との親しさの程度による影響の検討, 教育情報研究, Vol.24, No.2, pp.47-55, 2008.