

## キー入力の癖の学習に基づく入力間違い推定法の検討

西村 希樹<sup>†</sup> 寺澤 卓也<sup>†</sup><sup>†</sup>東京工科大学 メディア学部 メディア学科

## 1 はじめに

近年、PC を使用して、数千字の文書を書く機会は少なくない。特に PC 操作が不得意であると誤入力を起こす可能性が考えられる。そのため文書校正を行う必要性は高い。

ワードプロセッサには、辞書に基づいて単語の校正を行う機能がある。しかし、登録されていない単語は校正できない。そこで本研究では、辞書に依存せず、ユーザごとのキー入力の癖から間違いやすいキーを推測して校正する方法を考案した。

## 2 提案手法

## 2.1 概要

本手法では、まず、ユーザにサンプル文書を入力してもらい、まとまった量のキー入力をキーロガーで取得する。次に作成したプログラムで、取得したログデータから誤入力に関連したデータを抽出し、機械学習が行える形に変換する。そして機械学習で癖のパターンを学習させ、確率の計算を行う。最後に作成した学習データに基づき、校正を行う文書に対して、確率によって間違いやすいと判断された箇所を中心に、ハイライト表示を行い、校正の支援を行う。

## 2.2 ユーザの入力の取得

機械学習を行うためには、まとまった量のキー入力を修正動作も含めて取得する必要がある。キー入力を取得するため、キーロガーを使用する。これにより機械学習を行うために十分なログデータを蓄積させることができる。

## 2.3 ログデータ中の必要なデータの抽出と変換

次にログデータ中に不要なデータがあることによって学習結果に影響を及ぼさないよう、誤入力のみを抽出を行う。ログデータから誤入力部分を抽出するには、BackSpace、もしくはDelete キーで訂正された部分を手がかりとする。同じ文字を何度も訂正していればその文字の入

力をユーザが苦手としている可能性があるかと判断できる。図1の場合、mを打とうとして、aと誤って入力しmに訂正したことが考えられるため、aと[BackSpace]を抽出する。

```

K
A
[BackSpace]
M
A

```

図1：打ち間違いの例

また、penと入力するつもりが、誤ってoenと入力し、3文字消して訂正した場合、eとnは削除して入力し直しており、誤入力した文字と訂正した文字が同じとなる。間違いであるのはpと打とうとしてoと誤って入力した箇所であるため、この場合のeとnの修正は誤入力ではないと見なし、除外する。これにより、誤入力の箇所のみを抽出できる。

次に、使用するライブラリは、学習させるデータを読み込ませる際、二次元配列の形に変換する必要があるため、Pythonのプログラムを作成し、抽出した文字を二次元配列に変換する。

## 2.4 機械学習

本研究では機械学習環境として、Pythonのライブラリであり、トピックモデルと呼ばれる文章中の単語の出現確率を推定できるgensim[1]を使用する。gensimは他の機械学習ライブラリと比較すると図2の位置づけだと考えられる。

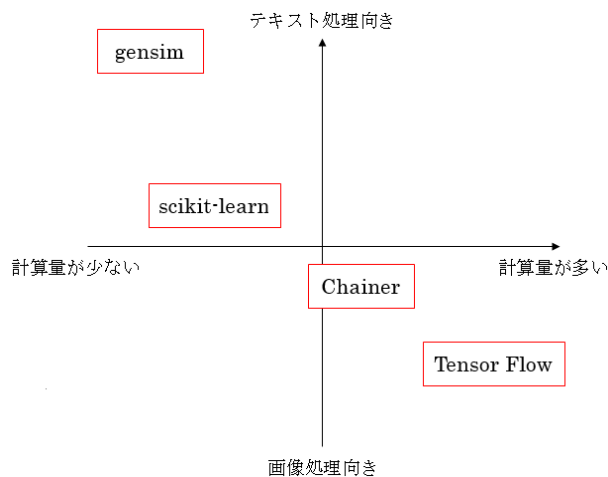


図2:機械学習ライブラリの比較

A presumption method to find typos based on machine-learning of typing habits

Kizuki NISHIMURA<sup>†</sup> and Takuya TERASAWA<sup>†</sup>

<sup>†</sup>School of Media Science, Tokyo University of Technology

本研究では、gensim のクラスの一つで、独立した単語同士が、設定回数以上出現したら、1つの単語として結合し扱うことができる Phrases[2] を使用する。これを、誤入力のパターンが一定回数出現した場合、誤入力の癖として認識させるために用いる。これにより、1回しか間違えなかったといった偶発的な誤入力を除去する。また、間違いのパターンの出現回数をカウントし、以下の数式で誤入力の確率を算出する。

$$\text{あるキーを誤入力する確率} = \frac{\text{そのキーの誤入力の回数}}{\text{そのキーの入力回数}}$$

## 2.5 文章構成の手法

校正自体は人間が行う。校正対象の文書のうち、以上の処理によって判明したチェックすべき箇所を示すために、複数の単語を一度にハイライトできるテキストエディタを使用する。その際、確率を参照し、確率毎に3レベルにハイライト表示を色分けして行うことで、特に注意深く見るべき箇所を把握しやすくした。

## 3 検証と実験結果

### 3.1 1500文字での検証

ユーザのキー入力の癖を取得するため、まずある程度の量のテキストを入力してもらいキーロガーで記録した。ユーザの使用環境を考慮し、キーボードは普段使用しているものを使用するなどを行い、癖以外で起こる誤入力を防いだ。また、キーロガーは日本語を記録できないため、ログデータを取得するためには、pykakasi[3]によって日本語をローマ字化したものを入力する形にした。

まず1500文字で行った結果、文字数が少ないため、偶発的に起こった間違いを除去することができなかつた。しかし、検証のためにユーザに文書校正を行わず作成させたテキストを解析すると、あらかじめ人間が調査し発見した誤入力のうち、75%の誤入力を本研究の手法で検出することができた。

### 3.2 10000文字での検証

次に行った実験は、前回の結果を踏まえ10000文字の入力で行った。検証を行うために入力をしたユーザから、PCを使用して実際に過去に講義で作成したレポートの文書を提供してもらい、学習結果を反映させた。目視で調査することによって検出した結果との比較から、10000回文字入力で3回以上間違えた文字を誤入力のパターンとすることで、偶発的に起きる間違いを取り除けると判断した。

### 3.3 ハイライト表示の検証

まず、偶発的な誤入力のみ取り除き、2.5の手法でハイライト表示を行った結果、ほとんどの文字がハイライトされてしまった。そのため、ある文字について、誤入力の回数に応じて設定したしきい値を、誤入力する確率の値が超えていけば、その文字をハイライトすることにした。今回は、10回以上間違えた文字であれば誤入力確率2%以上の文字をハイライト、10回以下であれば3%以上の文字をハイライトするようにした。その結果、図3のようにハイライトする箇所を減らすことができた。

```

15 | to|sanka|soli|ni|wakerarer|↓
16 | fi|hengcha|io|inoreben|ni|Ari|↓
17 | biteki|↓
18 | chiteki|↓
19 | rin|iteki|in|mosa|ikou|nomonowome|zast|Ononi|
20 | genchi|gohaka|inoreben|ni|Ari|↓
21 | sonoyou|namonowome|zast|onon|tha|Dwana|i|↓
22 | genchi|gowok|kokkanon|ar|ita|totomon|ibungaku|
23 | bukateki|kachi|iwotakame|tamonogakok|igode|
24 | nihongohame|ji|ishi|nochi|sama|zama|ji|in|

```

図3：ハイライトする文字を限定した具体例

## 4 結論

本研究の手法、すなわち、キーロガーによりキー入力を取得し、それに含まれる誤入力を元にした機械学習によって、校正対象文書中の誤入力をある程度検出できることが分かった。また、ハイライト表示の手法では、誤入力回数に応じて決めたしきい値を、算出した確率の値が超えていけば、その文字に対し、ハイライト表示を行うことで適切に校正支援できるめどをつけられた。

## 5 今後の課題

本手法によって校正する文章中の入力ミス発見の精度をどこまで上げることができるのかは明らかになっていない。そのため、文字をどの程度解析することが適切なかを明確にする必要がある。また、解析したユーザの他の文書に対してハイライト表示を行った場合、同じ精度で検出することが可能なかを明確化できなかったため、今後、検証する必要がある。

## 参考文献

- [1] Radim Řehůřek, “gensim”  
<https://radimrehurek.com/gensim/models/phrases.html> 最終閲覧日 2016/12/13
- [2] Radim Řehůřek, “Phrases”  
<https://radimrehurek.com/gensim/models/phrases.html> 最終閲覧日 2016/12/13
- [3] miurahr, “pykakasi”  
<https://github.com/miurahr/pykakasi>  
最終閲覧日 2016/12/19