

言語間比較によるWikipediaの 補完情報抽出手法の提案

☆ 藤原 裕也 (甲南大学)
鈴木 優 (名古屋大学)
小西 幸男 (甲南大学)
灘本 明代 (甲南大学)

背景



- **Wikipedia**

- 特徴

- 284以上の多言語版が存在
 - 誰でも記事を編集することが可能

- 問題

- ある話題に対しての情報が不足している記事が多く存在



情報補完

- 他のWebから情報補完する
 - **他の言語版から情報補完する** etc...

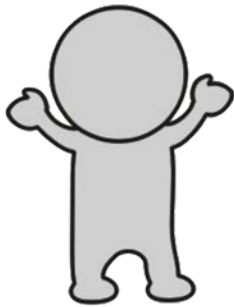
言語版によって書いてある内容が異なる

例:日本の文化

検索Query:たこ焼き

英語版

日本語版



外国人ユーザ

Contents

- 1 Takoyaki pan
- 2 See also
- 3 References

目次

- 1 概要
- 2 材料
 - 2.1 生地
 - 2.2 具
 - 2.3 味付け
- 3 食べ方
- 4 歴史
- 5 大阪のたこ焼き
- 6 器具
 - 6.1 業務用
 - 6.2 家庭用
- 7 日本国外におけるたこ焼き
 - 7.1 朝鮮文化圏
 - 7.2 中国文化圏
 - 7.3 東南アジア
- 8 関連事項

補完

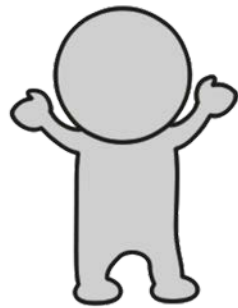
理由

- 編集者が少ない
- 十分な知識を持っていない



目的

閲覧しているWikipediaに不足している情報を
他言語Wikipediaから補完する



外国人ユーザ

理由

- 編集者が少ない
- 十分な知識を持っていない

英語版

Contents

- 1 Takoyaki pan
- 2 See also
- 3 References

日本語版

目次

- 1 概要
- 2 材料
 - 2.1 生地
 - 2.2 具
 - 2.3 味付け
- 3 食べ方
- 4 歴史
- 5 大阪のたこ焼き
- 6 器具
 - 6.1 業務用
 - 6.2 家庭用
- 7 日本国外におけるたこ焼き
 - 7.1 朝鮮文化圏
 - 7.2 中国文化圏
 - 7.3 東南アジア
- 8 関連事項

比較



補完情報を抽出

補完情報を抽出

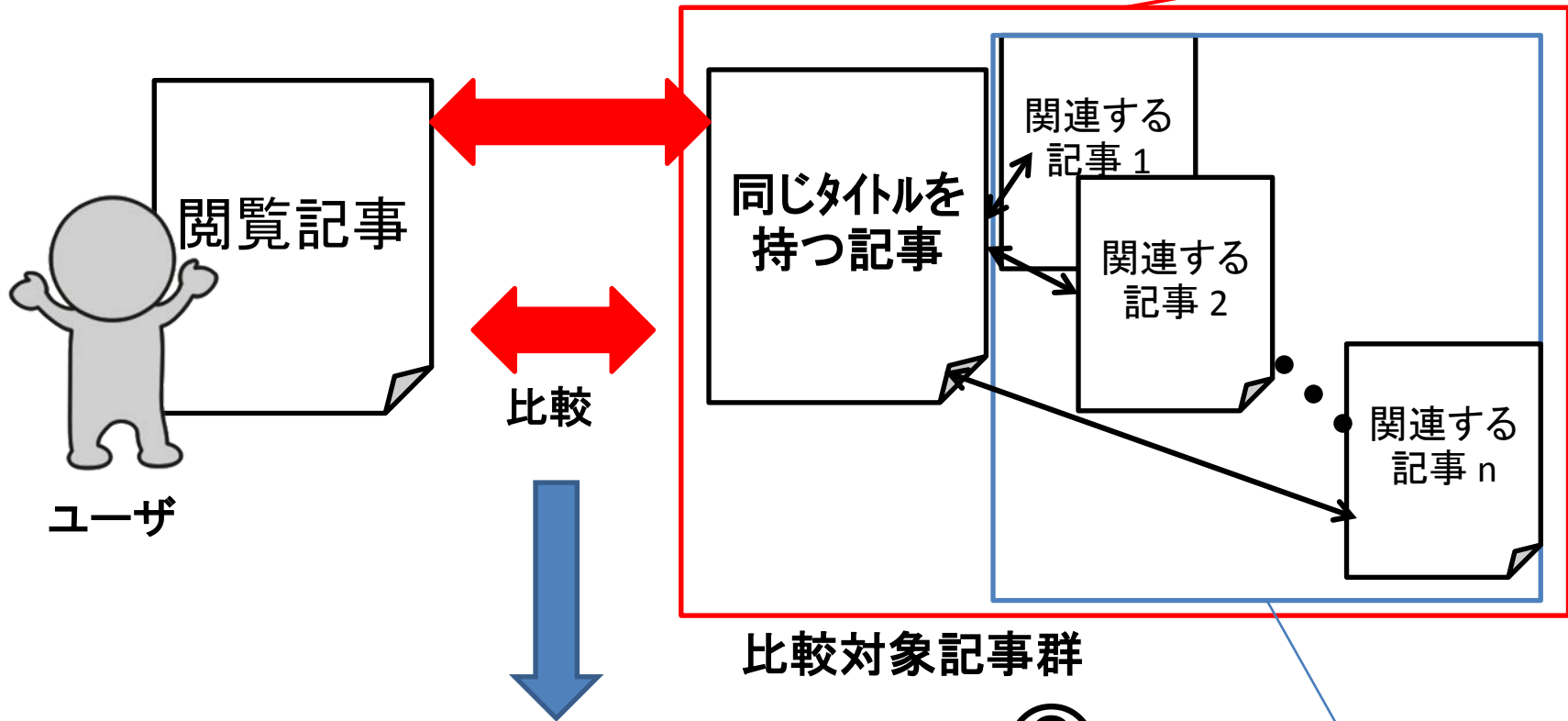
補完



全体の流れ

①

比較対象記事の決定



②

比較対象領域の決定

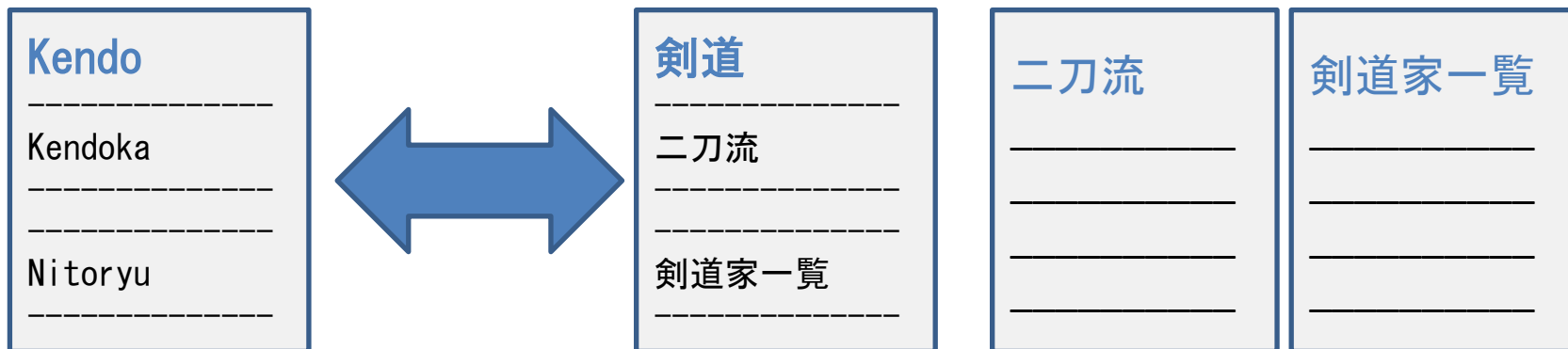
補完情報抽出

①比較対象Wikipediaの記事の決定

- 言語や文化の違いから情報の粒度が異なる
→対応する記事が複数にまたがる場合がある

– Ex:「剣道」

- 英語版:二刀流や剣道家の説明が含まれている
- 日本語版:二刀流, 剣道家一覧の記事が各々存在する



複数ページと比較して補完情報を
抽出する必要がある

①比較対象Wikipediaの記事の決定

リンク構造解析

関連している記事同士は
リンク関係にある

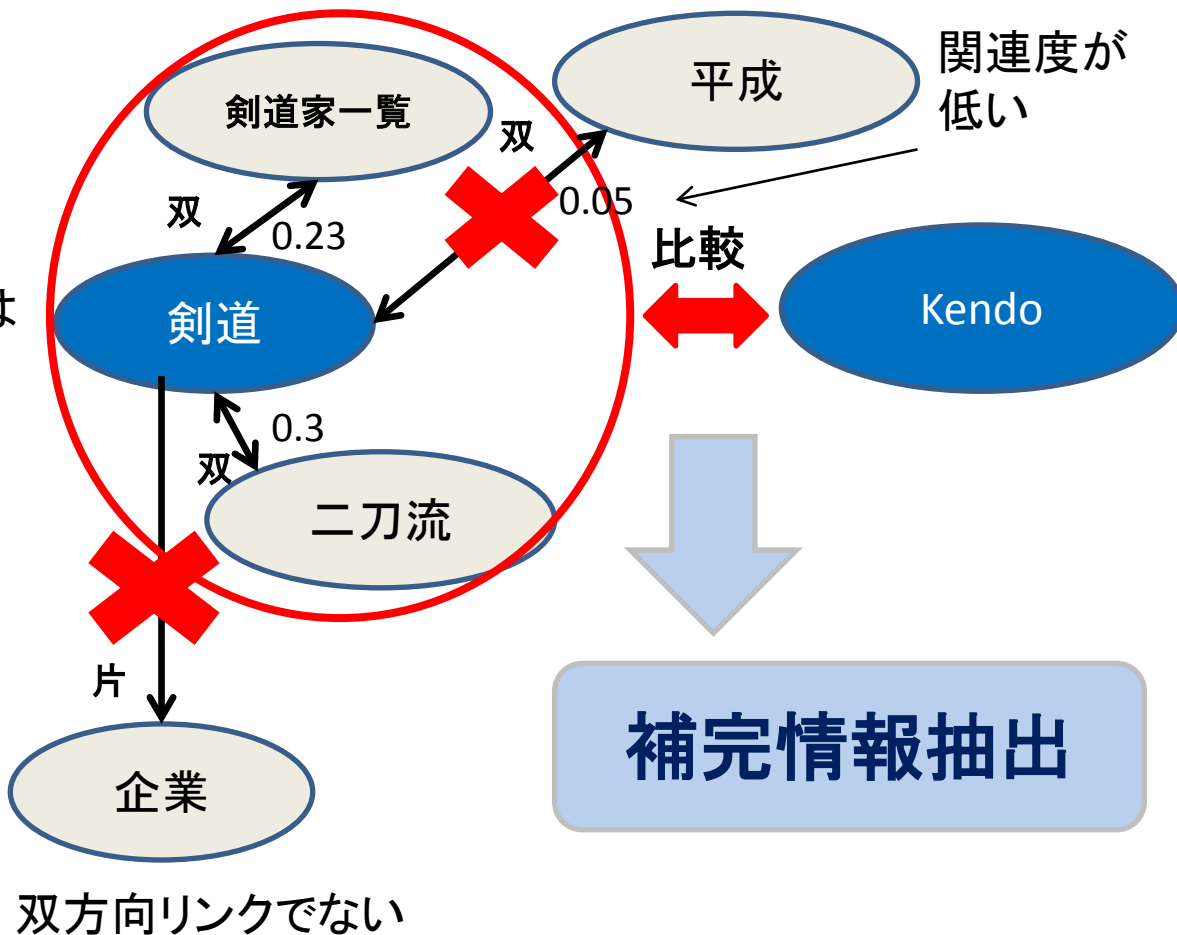
リンクグラフの生成

強連結: 関連性が強い

関連度計算

閾値以下を取り除く

比較対象ページの取得



補完情報抽出

①比較対象Wikipediaの記事の決定 関連度

双方向リンク注目した
記事と記事との関連する度合い

関連度

アンカー文字列の出現位置

サマりにリンクを張っている記事は
関連性が高い

アンカー文字列の出現回数

記事に何度も出現する
アンカー文字列は関連性が高い

コンテンツの類似性

関連する記事はある程度内容が
似ている

①比較対象Wikipediaの記事の決定 関連度

一番初めの説明部分

目次 [非表示]

- 1 歴史
 - 1.1 江戸時代
 - 1.2 明治・大正時代
 - 1.2.1 撃剣興行
 - 1.2.2 警視庁剣術
 - 1.2.3 大日本武徳会
 - 1.2.4 学校剣道
 - 1.2.5 剣道という名称について

剣道

サマリ

剣道(けんどう)は、剣術の竹刀稽古(撃剣)を競技化した武道。

歴史 [編集]

セグメント

「[剣術#歴史](#)」も参照

江戸時代 [編集]

剣道の直接の起源は防具と竹刀を使用する打ち込み稽古で、[手](#)を製作し、竹刀による打ち込み稽古法を確立した。[宝暦](#)年発達にともない袋竹刀より強固な割竹刀が作られるようになった。江戸時代後期から末期には、竹刀打ち中心の道場が興隆し、北辰一刀流玄武館、神道無念流練兵館や、幕府の設立した講武所頭取並の男谷信友は竹刀の全長を3尺8寸と定め、この時代の試合は審判規則や競技大会はなく、個人や道場

セグメント

明治・大正時代 [編集]

セグメント

撃剣興行 [編集]

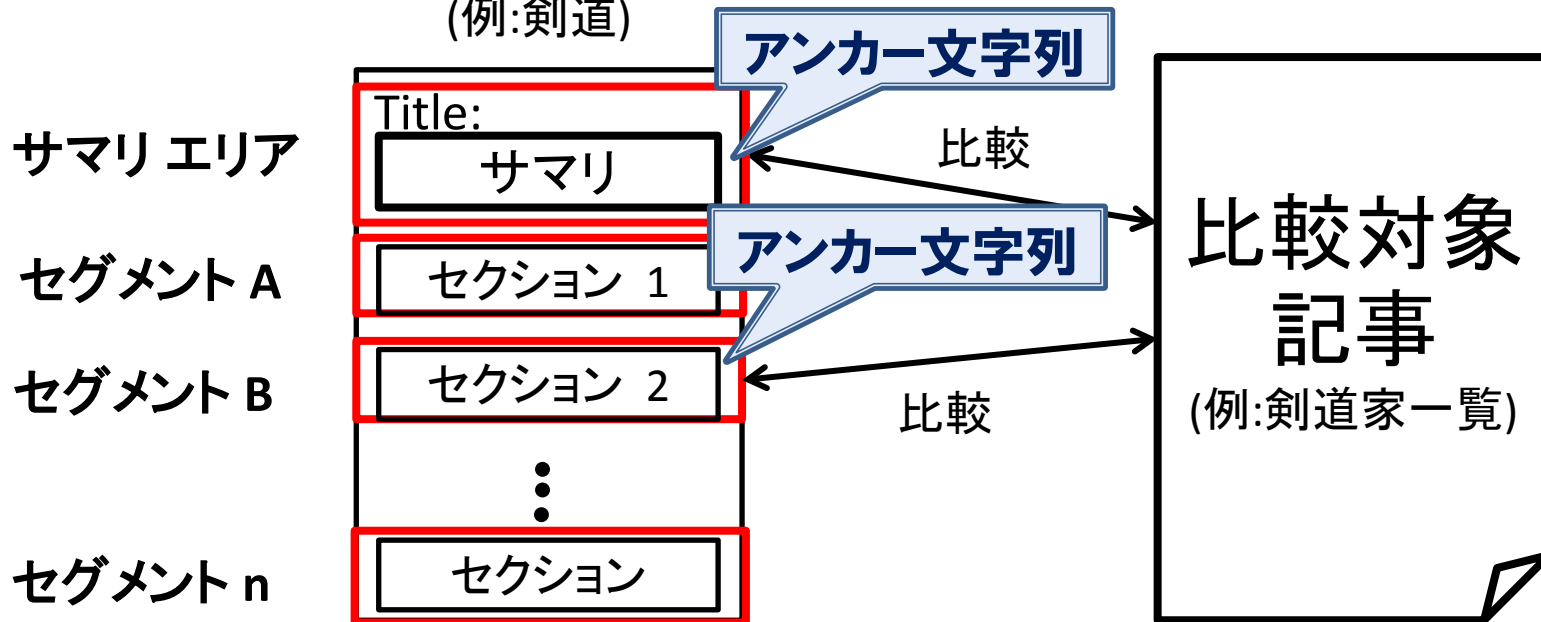
セグメント

明治維新によって武士の身分が廃止され、帯刀も禁止され、[榊原健吉](#)は明治6年(1873年)、撃剣興行という剣術見世物

記事をサマリとセグメントに分割

①比較対象Wikipediaの記事の決定 関連度

クエリがタイトルのページ(基準ノード)
(例:剣道)



比較対象記事のアンカー文字列が基準ノードのサマリ、どのセグメントに出現するかを求める

$$R_i = \{\alpha \cdot (TF_{sum_i} \cdot S_{sum_i}) + \sum_{k=1}^n (TF_{ik} \cdot S_{ik})\} / \max(R_{im})$$

i : 比較対象ノード

TF_{sum_i} のサマリのアンカー文字列の出現回数

S_{sum_i} のサマリと比較対象記事との類似度

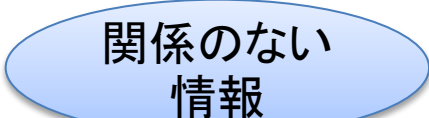
TF_{ik} のあるセグメントのアンカー文字列の出現回数

S_{ik} のあるセグメントと比較対象記事との類似度

n : ある l のリンクを張っているセグメントの数

$\max(R_{im})$: 比較対象記事群の R_i の最大値

②比較対象領域の決定

- 以前, 我々は関連度を用いて抽出した比較対象記事に対し補完情報の抽出を行った
 - 閲覧記事に関係ない情報が抽出される場合が存在した
 - 例: 剣道と二刀流の場合
 - 剣道の情報は一部分のみ
 - 西洋剣術の二刀流など
- 

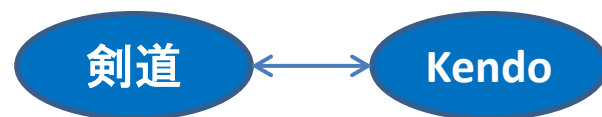
比較対象領域を決定する必要がある

②比較対象領域の決定

得られた比較対象記事に対し記事の分類を行う

- 比較基準記事

閲覧記事と同じタイトルを持つ記事
例:「Kendo」に対する「剣道」



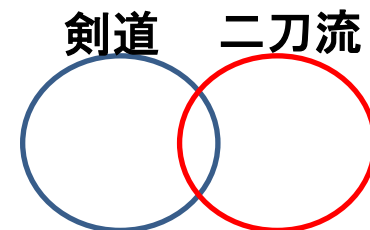
- 包含関係記事

比較基準記事と包含関係にある記事
例:「剣道」と「剣道家一覧」の関係



- 部分一致記事

記事の一部が比較基準記事と関係する記事
例:「剣道」と「二刀流」の関係

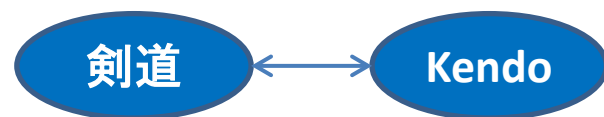


②比較対象領域の決定

得られた比較対象記事に対し記事の分類を行う

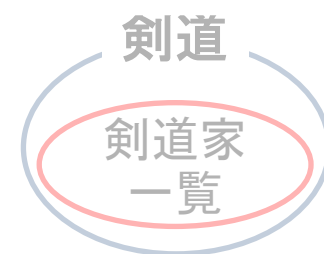
- 比較基準記事

閲覧記事と同じタイトルを持つ記事
例:「Kendo」に対する「剣道」



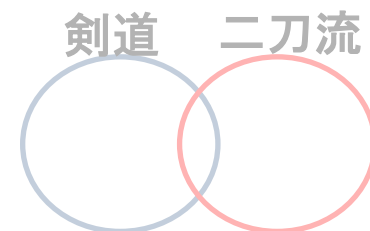
- 包含関係記事

比較基準記事と包含関係にある記事
例:「剣道」と「剣道家一覧の関係」



- 部分一致記事

記事の一部が比較基準記事と関係する記事
例:「剣道」と「二刀流」の関係



比較基準記事

- 閲覧記事と同じタイトルを持つ記事
– 例:「Kendo」に対する「剣道」

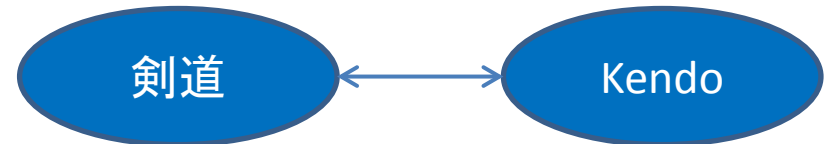
抽出方法

閲覧記事に存在する言語間リンクを用いて抽出

他言語へのリンク



英語版: Kendo

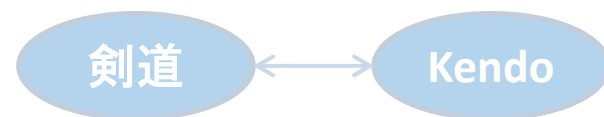


②比較対象領域の決定

得られた比較対象記事に対し記事の分類を行う

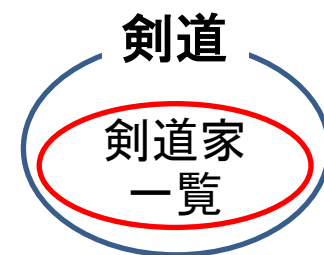
- 比較基準記事

閲覧記事と同じタイトルを持つ記事
例:「Kendo」に対する「剣道」



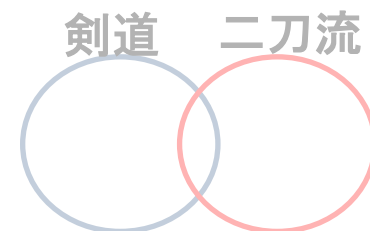
- 包含関係記事

比較基準記事と包含関係にある記事
例:「剣道」と「剣道家一覧」の関係



- 部分一致記事

記事の一部が比較基準記事と関係する記事
例:「剣道」と「二刀流」の関係



包含関係記事

- 比較基準記事と包含関係になっている記事
 - >比較基準記事とis-a関係になっている記事
 - 例:「剣道」と「剣道家一覧」の関係

抽出方法

中山らの提案するLSP法を用いて抽出

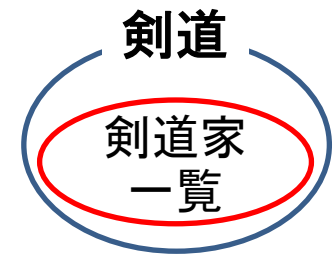


記事の冒頭文を重要文とし、その文に比較基準記事のアンカー文字列が存在する場合に包含関係記事として抽出

例: 剣道家一覧

比較基準記事への
アンカー文字列

剣道家一覧は、**剣道**で活躍した人の一覧

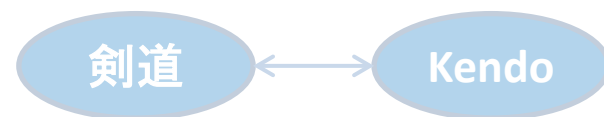


②比較対象領域の決定

得られた比較対象記事に対し記事の分類を行う

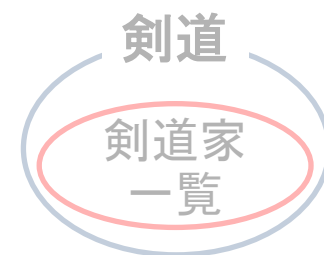
- 比較基準記事

閲覧記事と同じタイトルを持つ記事
例:「Kendo」に対する「剣道」



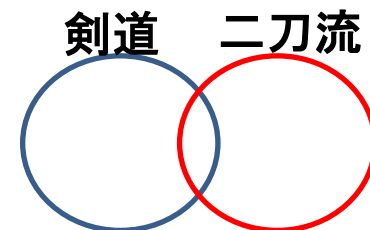
- 包含関係記事

比較基準記事と包含関係にある記事
例:「剣道」と「剣道家一覧」の関係



- 部分一致記事

記事の一部が比較基準記事と関係する記事
例:「剣道」と「二刀流」の関係

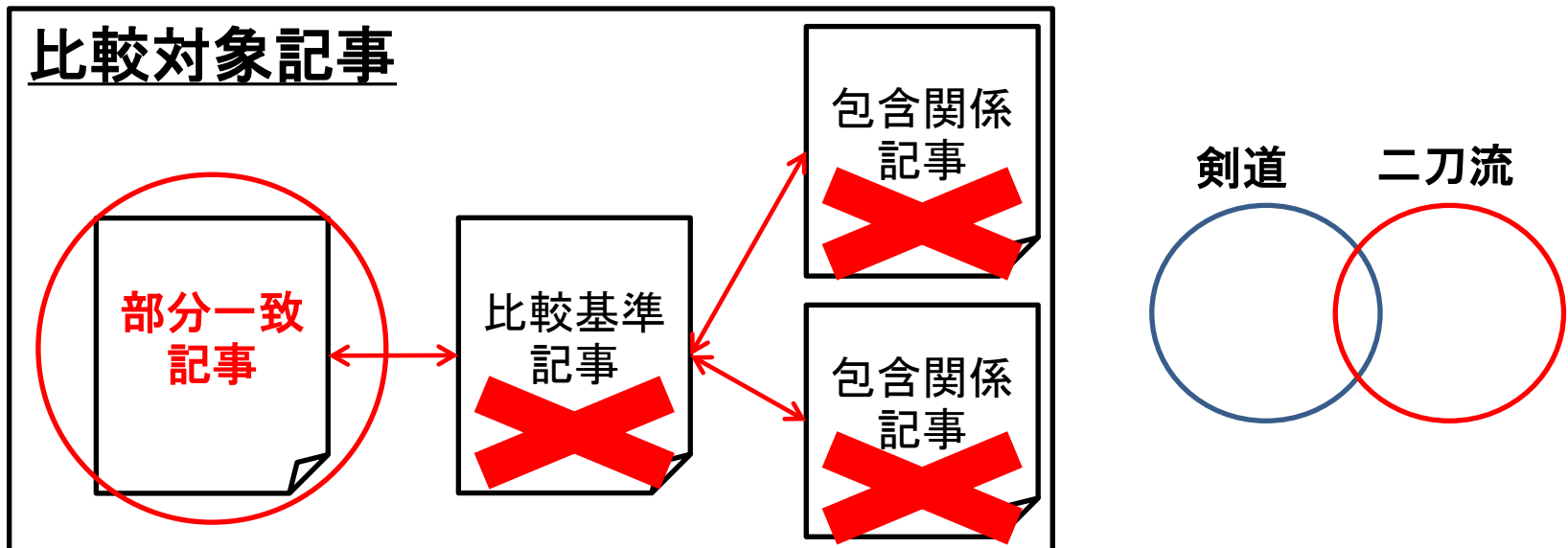


部分一致記事

- 記事の一部が比較基準記事と関係する記事
– 例:「剣道」と「二刀流」との関係

抽出方法

比較対象記事から比較基準記事と包含関係記事を
除いたすべての記事



②比較対象領域の決定

- 分類
 - － 比較基準記事
 - － 包含関係記事
 - － 部分一致記事
- 比較領域の決定
 - － 比較基準記事・包含関係記事
 - 記事全体を比較
 - － 部分一致記事
 - セクションのタイトルに比較基準記事のアンカー文字列を含む場合
 - サブセクションのタイトルに比較基準記事のアンカー文字列を含む場合
 - 記事本文中に比較基準記事のアンカー文字列を含む場合

比較対象領域の決定と補完情報抽出

- 分類ごとに補完情報を抽出するために、閲覧記事と比較対象となる領域を決定する

比較基準記事・包含関係記事

閲覧記事との関係が強いと考え記事全体を比較対象とする

概要 [編集]

タラやカレイ、オヒシなど白身魚の切り身に、小麦粉を卵や水または棒状に切ったで揚げたチップスと合わせて提供する。この場合の手で言うフライドポテト（アメリカで言うフレンチフライ）のイギリスでの呼称の切り身小一切れにジャガイモ中一個分）で450キロカロリー程。

歴史 [編集]

白身魚の切り身を揚げた料理は、少なくとも中世ヨーロッパに存在して、ヨーロッパ各地でジャガイモを揚げた料理も作られるようになった。両者はいつかは諸説入り乱れている。記録に残る限りでは、1860年にロンドンが最も古くのものである。19世紀後半に底引き網漁の技術革新が起こり、チップスは労働者階級の日常食になった。第二次世界大戦下のイギリスでフィッシュ・アンド・チップスであった。戦後もフィッシュ・アンド・チップス

食べ方 [編集]

モルトビネガー（麦芽を原料とする穀物酢）と食塩をかけてマッシュイビー一般的だが、マヨネーズやタルタルソースなどをつけて食べることもあなど好みにより、多様な味付けが行なわれてい。飲食店内では皿に芋のように、紙袋に入れるか円錐型に丸めた新聞紙に包まれて渡される店もある。ファストフード店では、フィッシュをパンズに挟み、チップス

セグメント
セグメント
セグメント

History

Main article: British cuisine

Fish and chips became a stock meal among the working classes in Great Britain as a cities during the second half of the 19th century.^[2] In 1860, the first fish and chip shop was opened in London. Deep-fried chips (slices or pieces of potato) as a dish may have first appeared in Britain in the earliest usage of "chips" in this sense, the mention in Dickens' *A Tale of Two Cities* in 1859.^[3] The word "chips" first occurred in southern England in 1680.^[4] The word "chips" first occurred in southern England in 1680.^[4] The word "chips" first occurred in southern England in 1680.^[4] In the United Kingdom, the word "chips" is particularly standard in the north and south.

England

The dish became popular in wider circles (Charles Dickens mentions a "fried fish with a side of deep-fried chipped potatoes" in *Temple Bar Market*).^[13] It remains unclear when the fish and-chip shop industry we know today. It was first mentioned in London in 1860 or in 1865, while a Mr Lee

閲覧記事との関係が強いと考え記事全体を比較対象とする

補完情報

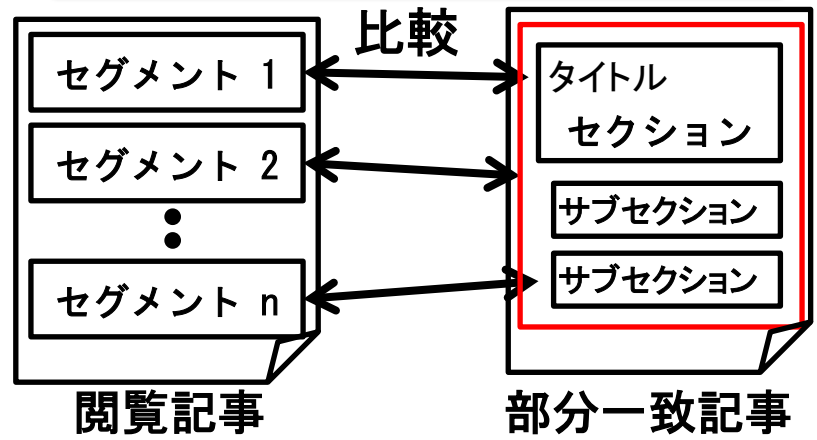
類似記事の閾値以下



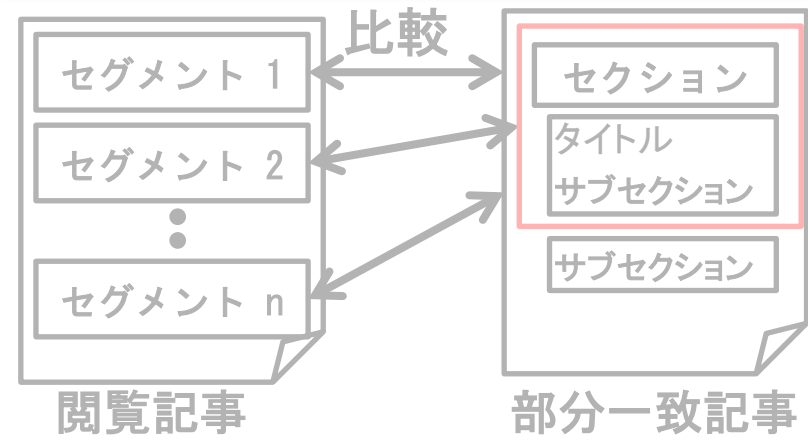
比較対象領域の決定と補完情報抽出

部分一致記事

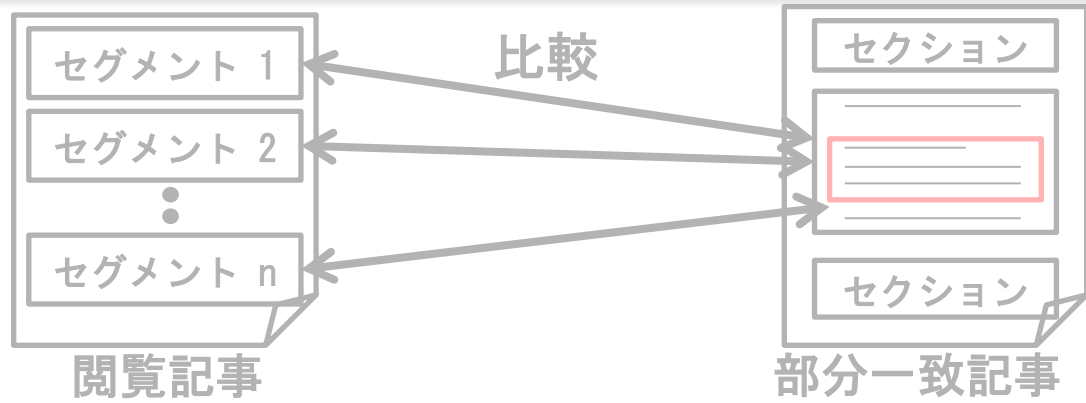
セクションのタイトルに比較基準記事のアンカー文字列を含む場合



サブセクションのタイトルに比較基準記事のアンカー文字列を含む場合



記事本文中に比較基準記事のアンカー文字列を含む場合

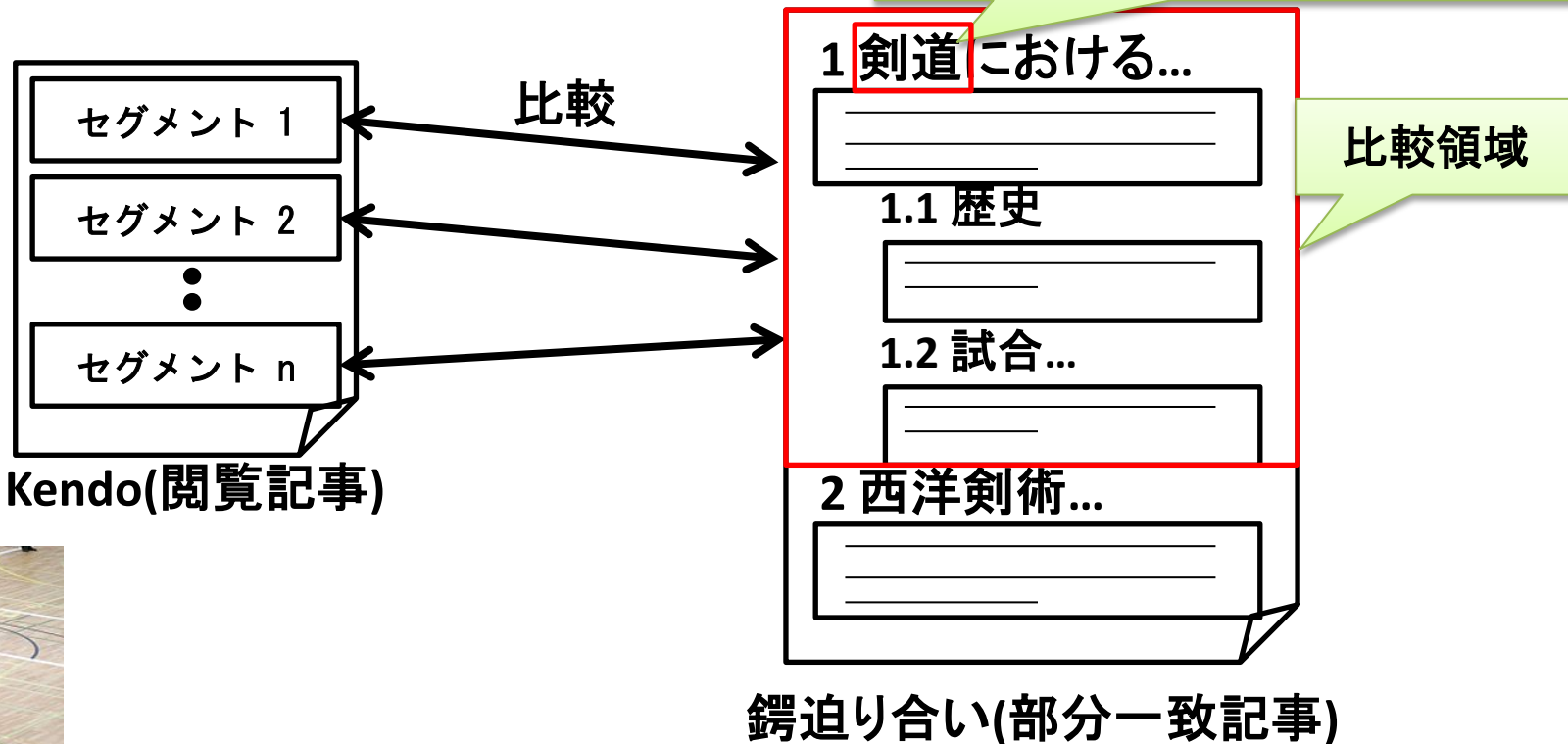


比較対象領域の決定と補完情報抽出

セクションのタイトルに比較基準記事のアンカー文字列を含む場合
サブセクションを含めそのセクション全体を比較対象とする

例：剣道(比較基準記事)と鰐迫り合い(部分一致記事)

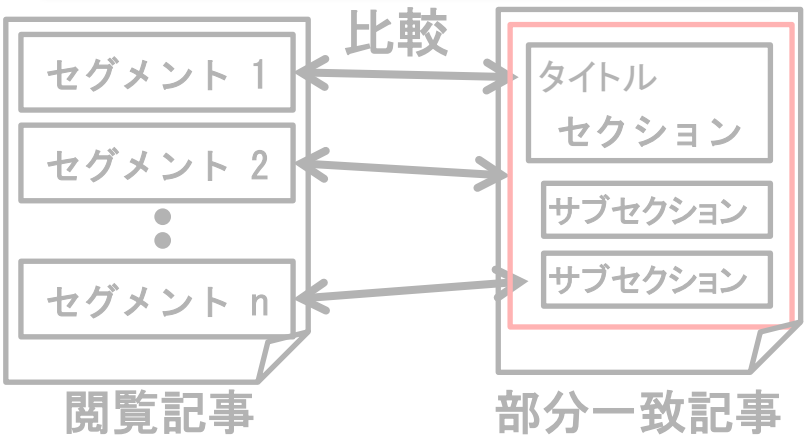
比較基準記事のアンカー文字列



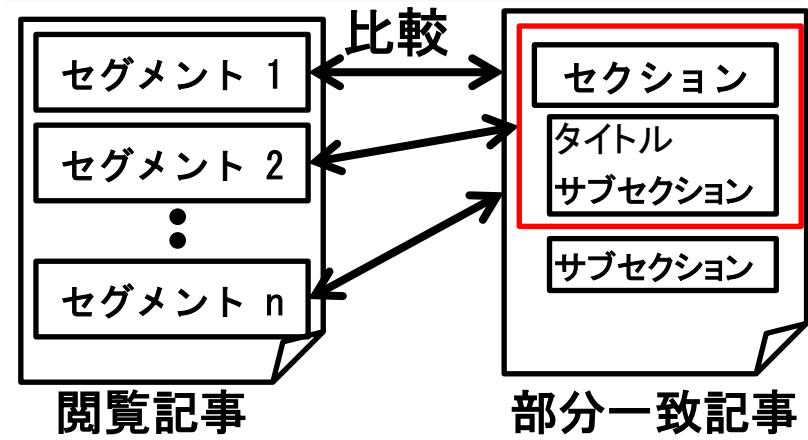
比較対象領域の決定と補完情報抽出

部分一致記事

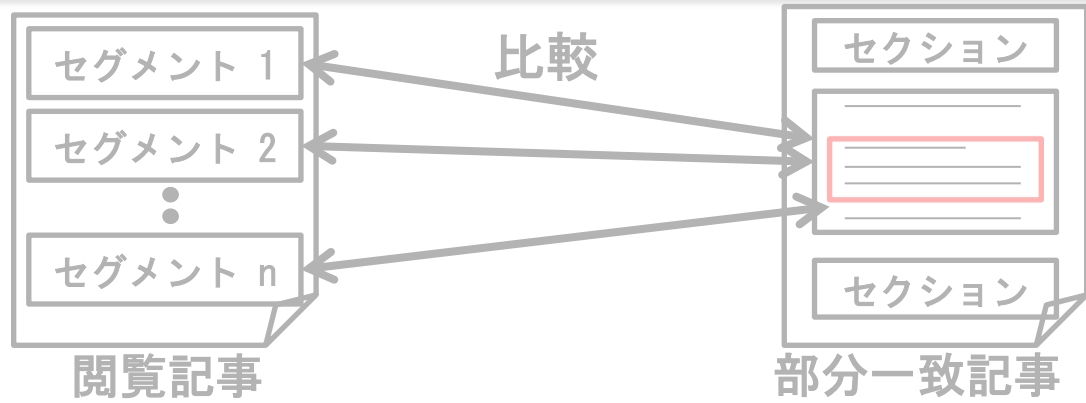
セクションのタイトルに比較基準記事のアンカー文字列を含む場合



サブセクションのタイトルに比較基準記事のアンカー文字列を含む場合



記事本文中に比較基準記事のアンカー文字列を含む場合

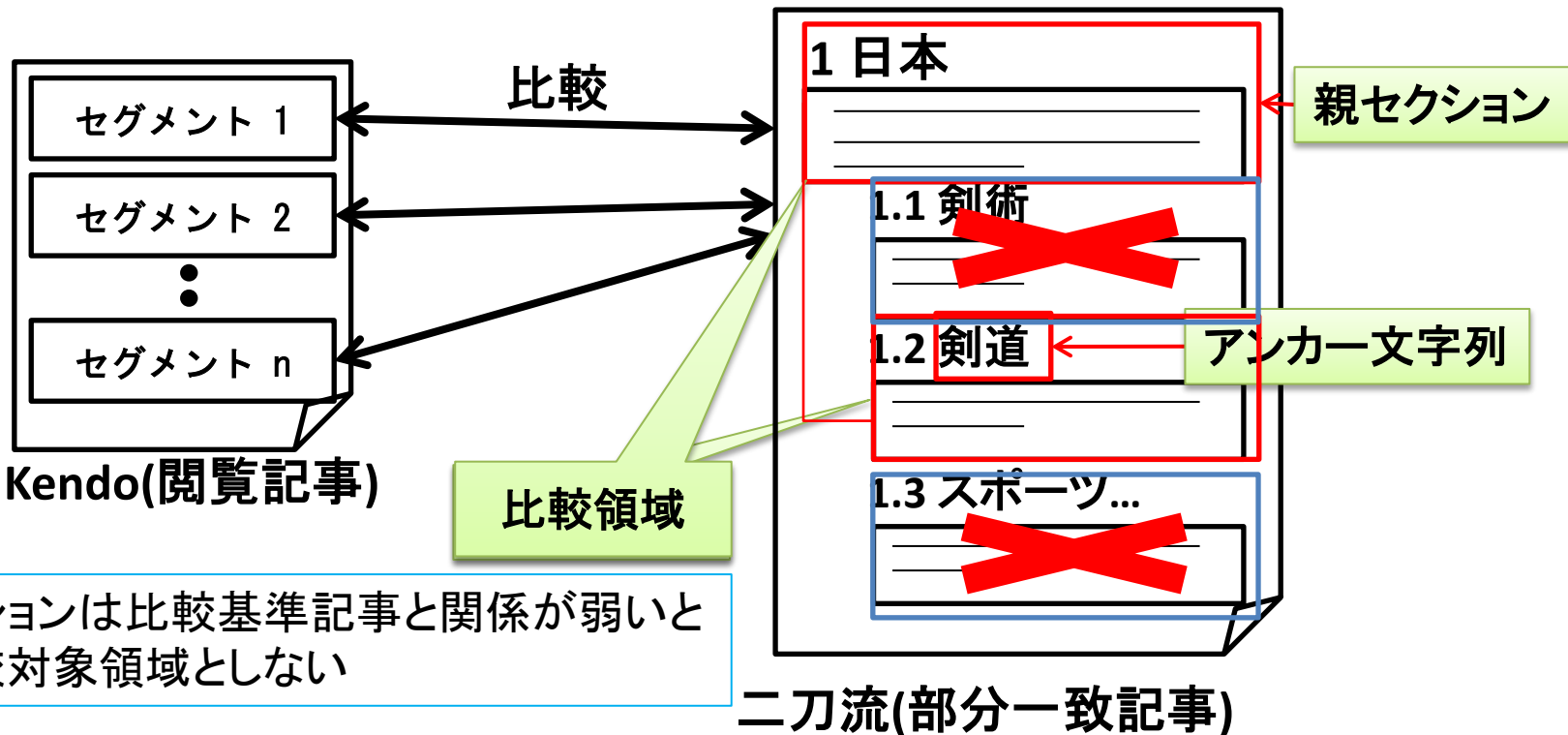


比較対象領域の決定と補完情報抽出

サブセクションのタイトルに比較基準記事のアンカー文字列を含む場合

親セクションとそのサブセクションを比較対象とする

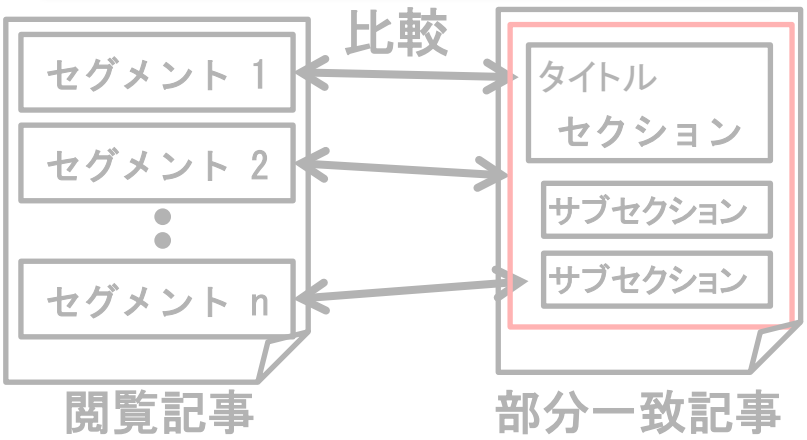
例: 剣道(比較基準記事)と二刀流(部分一致記事)



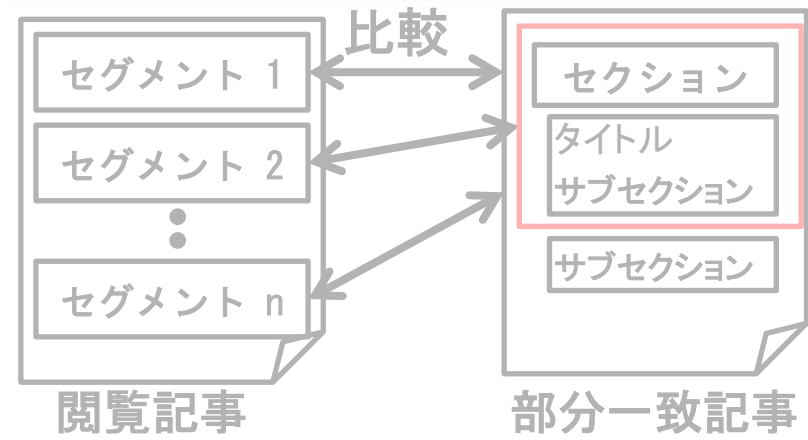
比較対象領域の決定と補完情報抽出

部分一致記事

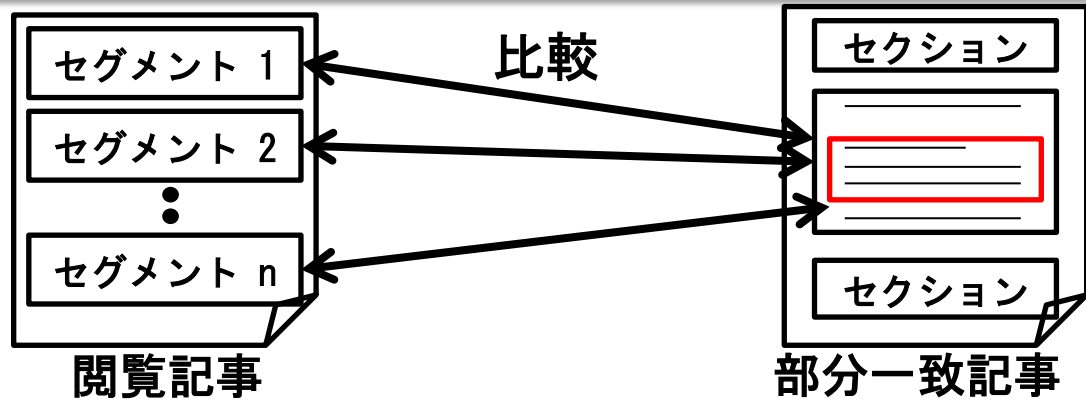
セクションのタイトルに比較基準記事のアンカー文字列を含む場合



サブセクションのタイトルに比較基準記事のアンカー文字列を含む場合



記事本文中に比較基準記事のアンカー文字列を含む場合

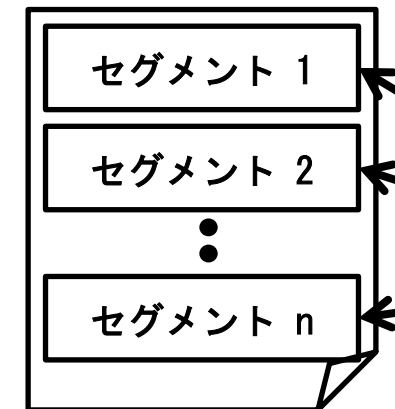


比較対象領域の決定と補完情報抽出

記事本文中に比較基準記事のアンカー文字列を含む場合

アンカー文字列の含まれている段落のみを比較対象とする

例: 剣道(比較基準記事)と道場(部分一致記事)



Kendo(閲覧記事)

概要 [編集]

各地にある道場の多くはいずれかの流派に属し、看板または施設名標にその流派名では、地域毎にそれらを束ねる組織が存在し、またその上位にも同様の組織が存在し、柔道の講道館などへと連なる。

道場の内装であるが、一般的な日本の道場では、正面には神座が設けられ、額縁は目によって異なり、例えば柔道では畳、剣道では板張りである。

現代における道場は、都市の過密化や地価の高騰などにより、それ専用の建物だけ教室などを開いている場合、自宅家屋の一角を道場に改造している場合もある。

道場(部分一致記事)

比較

比較基準記事のアンカー文字列

比較領域

プロトタイプシステム

プロトタイプシステム

Query
Yukata en

クエリと閲覧言語を入力



比較対象記事群

The screenshot shows a browser window with the URL localhost/pro_sys/src/Yukata_wiki.html. The page displays two side-by-side articles about 'Yukata'. The left article is the English version, and the right is the Japanese version. A blue arrow points from the search area to the browser. A blue circle highlights a navigation menu in the Japanese article with the label 'タブ' (Tab). A red box highlights a section of the Japanese article with the label '補完情報' (Supplemental Information).

閲覧記事(英語版)

比較対象記事(日本語版)

実験

- 提案手法の有用性を示す実験を行った
 - 実験内容
 - 提案手法とBaselineの比較
 - Baseline:比較対象領域の決定を行わない場合
 - 適合率, 再現率, F値の比較
 - 比較言語版
 - 閲覧言語:英語版
 - 比較対象言語:日本語版
 - 設定(前実験より)
 - 関連度の式の α :3.0
 - 関連度の閾値 β :0.2
 - コンテンツの比較の閾値 γ :0.2

$$R_i = \{ \alpha \cdot (TF_{sum_i} \cdot S_{sum_i}) + \sum_{k=1}^n (TF_{ik} \cdot S_{ik}) \} / \max(R_{im})$$

実験条件

Query	実験条件		
	閲覧記事の セクションの数	比較対象 記事の数	比較対象記事の セクションの数
My Neighbor Totoro (となりのトトロ)	9	2	25
Doraemon (ドラえもん)	9	5	38
Iaido (居合道)	10	4	24
Manzai (漫才)	3	4	22
Yukata (浴衣)	1	3	7
Urashima Taro (浦島太郎)	5	5	39
Pikachu (ピカチュウ)	6	2	25
Kinkaku-ji (鹿苑寺)	24	6	28
Hello Kitty (ハローキティ)	15	2	60
Kyudo (弓道)	16	3	47

$$\text{適合率} = \frac{\text{抽出した補完情報} \cap \text{正解データ}}{\text{抽出した補完情報}}$$

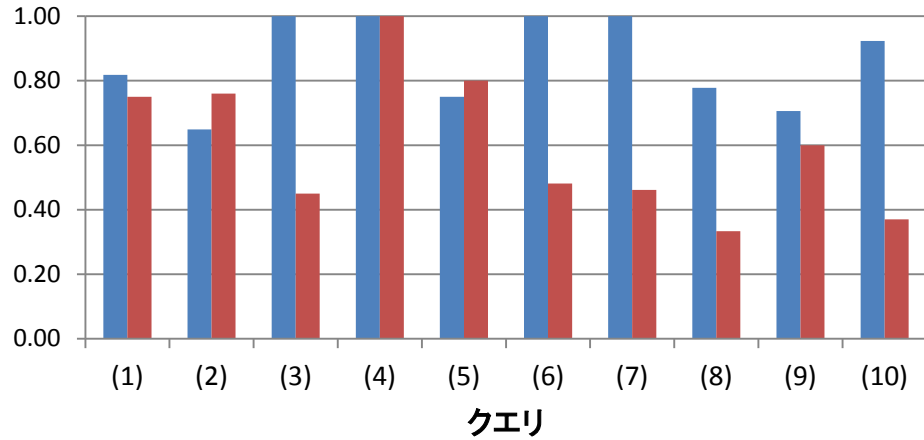
$$\text{再現率} = \frac{\text{抽出した補完情報} \cap \text{正解データ}}{\text{正解データ}}$$

正解データ: 閲覧記事に対し補完情報となる比較対象記事のセクションまたは段落

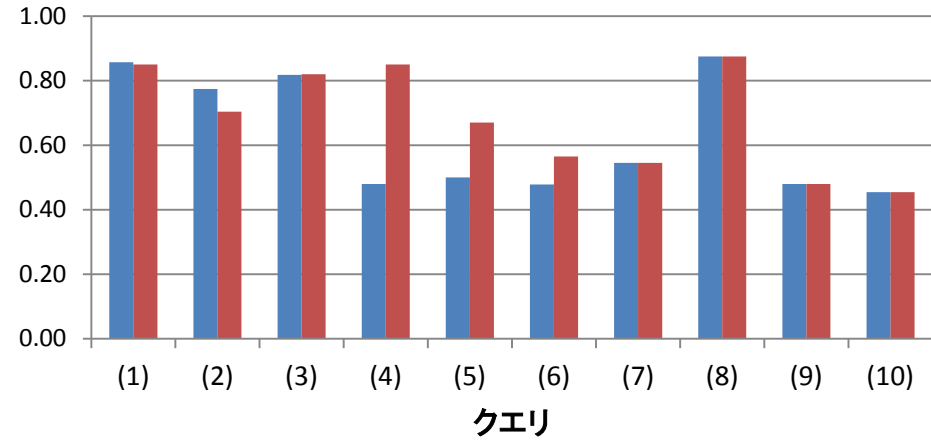
結果

■ 提案手法
■ baseline

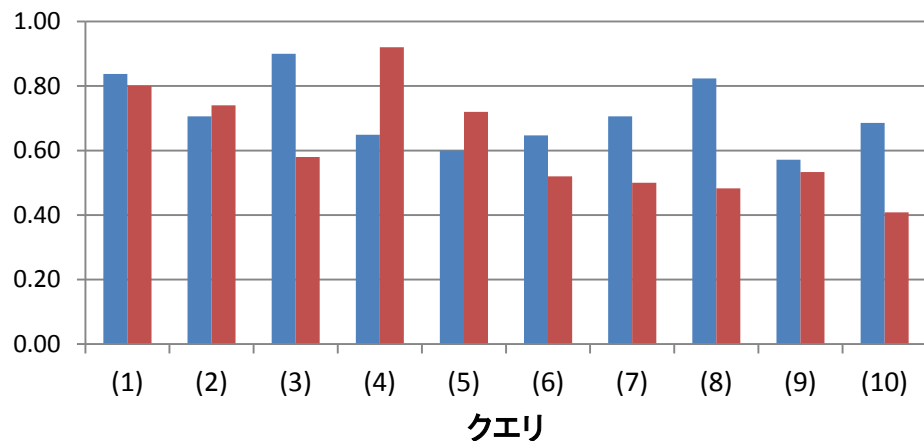
適合率



再現率



F値



Number	クエリ
(1)	My Neighbor Totoro(となりのトトロ)
(2)	Doraemon(ドラえもん)
(3)	Iaido(居合道)
(4)	Manzai(漫才)
(5)	Yukata(浴衣)
(6)	Urashima Taro(浦島太郎)
(7)	Pikachu(ピカチュウ)
(8)	Kinkaku-ji(鹿苑寺)
(9)	Hello_Kitty(ハローキティ)
(10)	Kyudo(弓道)

平均(適合率:0.60->0.86, 再現率:0.68->0.62, F値:0.62->0.71)

結果の良い例

- 居合道
 - Baseline
 - 部分一致記事である武道の称号の範士が抽出, しかし柔道や弓道の範士のように居合道に関係のない情報が抽出された
 - 提案手法
 - 範士, その中でも居合道の情報のみが抽出できた
- となりのトトロ
 - Baseline
 - 部分一致記事である狭山丘陵, 地理情報などとなりのトトロと関係のない情報が抽出された
 - 提案手法
 - 狭山丘陵, となりのトトロの舞台となったという情報が抽出できた

考察

- 部分一致記事において補完情報と成り得ない情報が抽出される場合が存在した
 - 例:ドラえもん
 - 部分一致記事として作者の藤子・F・不二雄が抽出
 - 藤子・F・不二雄の記事では多くのセクションでドラえもんのアンカー文字列が出現
 - ドラえもん以外に多くの漫画を描いており, ドラえもんの補完情報とならない情報が抽出された

考察

- 比較対象となる領域が正しく決定できない場合が存在
 - 例：浦島太郎
 - 部分一致記事として荘内半島が抽出

浦島伝説 [編集]

比較対象

浦島伝説は浦島太郎の同義語

荘内半島には、**浦島太郎**に関する地名が数多く残されている。

- 生里 - 太郎誕生の地。
- 箱浦 - 太郎が青年となり、移住し、かえって来た後、箱を開いた。
- 鴨ノ越、丸山島 - 亀を助けた海岸。
- 積 - 乙姫からもらった宝物を積んだ地。
- 金輪の鼻 - 乙姫との別れの際に姫の腕輪が落ちた。
- 紫雲出山 - 箱から出た煙から。
- 室浜 - 年をとった太郎が住んだ。
- 姫路(栗島) - 姫が太郎と別れた後、一時立ち寄ったところ。

荘内半島

まとめと今後の課題

- まとめ
 - Wikipediaの多言語性に着目し, 内容の充実していない記事に対し他言語のWikipediaを用いて, 情報の補完を行う手法を提案した
 - 提案手法
 - 比較対象記事の決定
 - 比較対象領域の決定
 - 比較基準記事, 包含関係記事, 部分一致記事
 - コンテンツの比較による補完情報抽出
- 今後の課題
 - 補完情報の分類(詳細な情報, 新しい情報 etc...)
 - 補完情報の提示のインタフェース
 - 比較領域の決定の際に同義語への対応