

D-72 統計調査の約600分類の 符号付与システムについて

床 裕佳子, 下野 寿之, 和田 かず美, 坂下 佳一郎
(独立行政法人 統計センター)
2015/11/26 IBIS2015 つくば国際会議場

目的 文字列 (家計簿記入内容) から
3桁の符号を高精度で正解を与える!

性能評価

信頼度スコアのパーセンタイルごとの一致率

自然言語処理を活用して、
約570通りの分類を**高精度**かつ**高速**に行えるような
システムを構築する。

データ 総務省統計局が実施する家計調査のデータ

データの件数

学習用: 約187万件 (約190MB) ... 1年分

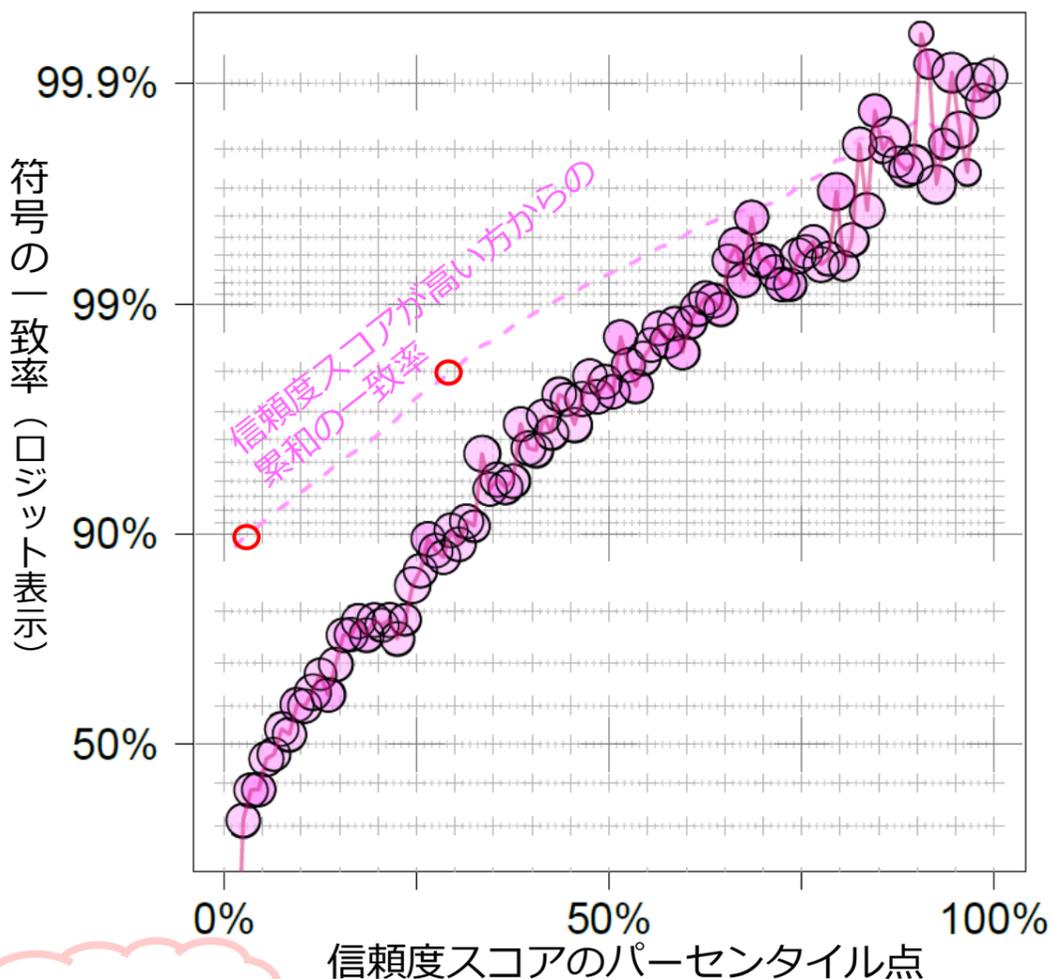
符号付与用: 約32万件 (約30MB) ... 1ヶ月分

正解の
分類符号

↓データイメージ

001	パン↓
548	カニカマ↓
313	給与 (3月分) ↓
281	にんじん2本↓
299	バナナ↓
411	〇〇バーガーにて外食↓
002	チョコレートパン↓
328	チョコ↓
001	ロールぱん6ヶ入↓
088	キャベツ入りメンチカツ↓
315	所得税 (6ヶ月分) ↓
572	カニ↓
412	外食 (うどん) ↓
331	豆乳入りどら焼き↓

分類対象の
文字列



全体の結果

一致率 (一致件数 / 符号付与された件数) **89.2%**
符号付与率 (符号付与された件数 / 全件数) **98.5%**

さらに **一致率90%保持ライン**
97.1%

一致率98%保持ライン
69.6%

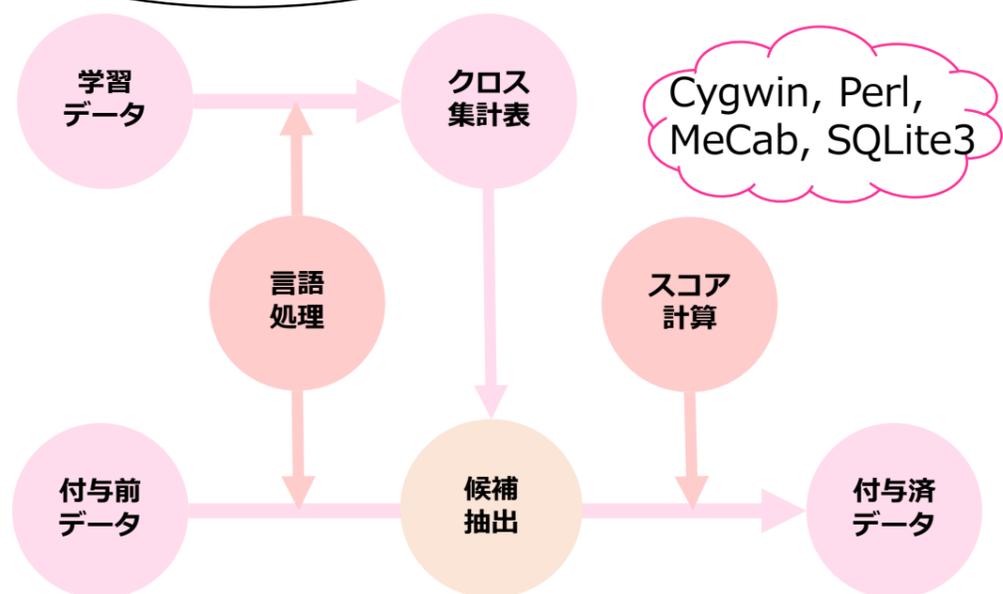
処理時間

(Xeon, 3GHz)

学習フェーズ
約**2秒**/1万データ

符号付与フェーズ
約**5秒**/1万データ

システム構成



考察・課題

- 単純**...クロス集計表のみに依存
- 高精度**...専門的な知識を持った職員に近い精度
- 速い**...処理スピードが**実用的**

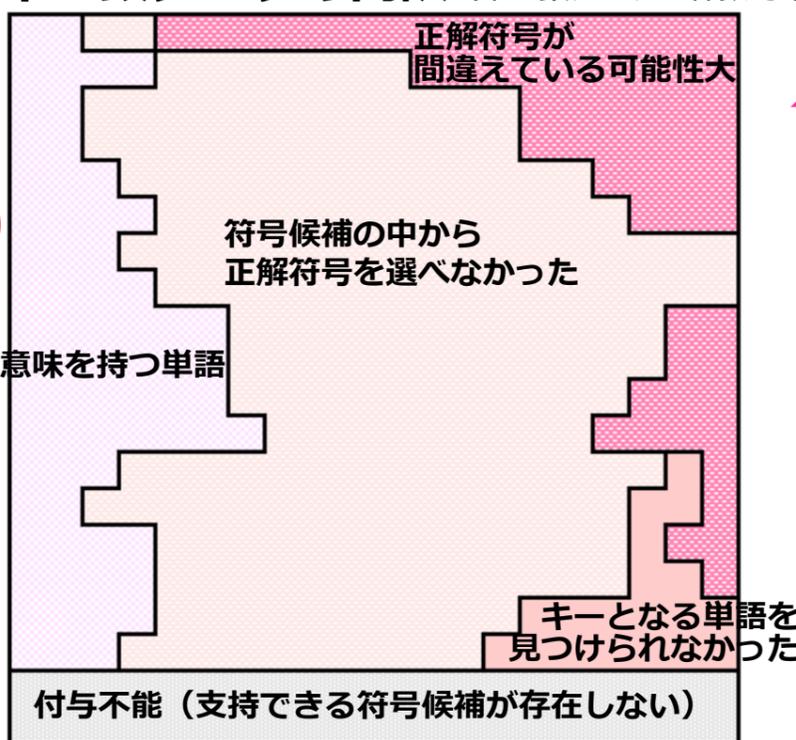
要改善ポイント

- 複数の意味をもつ単語
- "強い"キーワードの副作用
- 表記ゆれ
- カタカナ連続記入 など

上記を踏まえ... **今後の課題**

- **他の属性情報の活用**⇒同音異義語/地域間表記ゆれ/収入か支出か
- **MeCabユーザ辞書の拡充**⇒カタカナ/表記ゆれ/記入ミスをカバー

不一致データの内訳 (不一致データの件数で等積表示)



注: 本ポスターの内容は、報告者本人の見解であり、必ずしも所属組織を代表するものではありません。