

複層意味フレーム分析を用いた 意味役割タグつきコーパス

評価版の公開

黒田 航 井佐原 均

(独) 情報通信研究機構
けいはんな情報通信融合研究センター

NLP 11 [03/17/2005]

研究の背景と目的

- 黒田・井佐原 (2004) [NLP10] は Berkeley FrameNet (BFN) を拡張した意味役割タグつきコーパス開発を開始
- BFNと別に複層意味フレーム分析 (MSFA) と呼ぶ独自の意味タグづけ体系=手法を定義
- 今回の研究の目的: MSFA が技能として習得可能か確かめ, 可能ならば文書化する

概要

- MSFA の目的, これまでの成果
- MSFA の理論と実例
- 現状と今後の方針

MSEFA の目的

MSFA は何の記述か？

- 文 s の MSFA はヒト x が s を読んだり/聞いたりしたときに x が理解する内容 $m(s)$, すなわち x による s の理解内容 $F(x, m(s))$ の可能な限り明示的で体系的な記述
- ただし最適でないし, 今は(まだ)完全でもない
 - $m(s)$ ではなく $F(x, m(s))$ の最適性は x に依存する
 - これは“読み”の“深さ”, “観点”のパラメータ化

意味役割タグづけの動機

- 解析用の辞書とは別に文脈化された意味のデータベース化が必要
 - 意味(内容)分析は(当分)自動化できない
 - ヒトの意味直観は恐ろしく微妙，かつ正確である
- 言語学者だって“使える”(意味)分析を提供できることの証明
 - 役に立たない統語派生の分析や語の分類しかできない訳ではない

ちよいと大袈裟な話

- (自然言語)情報処理の十八番の「知識の自動獲得」もいいのだけれど
- イキナリ自動獲得に走る前に自然言語データから手動で、どんな知識が、どの程度まで獲得できるのかを見極めて/見積もっておく必要はないのでしょうか?
 - Text/Data Miningで「掘る」前に「どこ」に「何」があるか事前調査しなくていいのでしょうか?

これまでの成果

- 作業者四人による人手コーディングと結果の編集 (数値的評価はまだ)
- 京大コーパスから三記事文 (合計64文)
 - 950103083-001,018 (ラグビー)
 - 950101075-001,036 (将棋の名人戦)
 - 950107210-002,010 (マラドーナ逮捕)
 - フレームの延べ数にして500個ぐらい
- <http://61.115.230.87/~mutiyama/cgi-bin/hiki/hiki.cgi?FrontPage> で部分的に公開
 - 内山将夫 (NiCT) 氏の好意による

なぜ「評価」版？

- 記述の自由度の過剰を收拾するための情報収集
 - MSFA はキリがない = 終了条件が不明確
- 分野ごとの需要は尊重するが、特定の利用目的に解析を特化させるつもりはない
 - 私たちの考える意味役割タグつきコーパスは NLP の特殊な用途の他にも言語学/認知科学/日本語教育を含め、可能な限り広い研究分野への利用価値をもつもの
 - 言語学者の(無謀な)突っ走りへの「保険」

お断り

- MSFA は言語学者/言語の認知科学者の仕事なので、NLP での応用は視野に入っているとは言え、それ自体は目的ではない
- 「何のためにやっているの?」と「これは何の役にたつの?」という質問の答えは、別
 - 私が自信をもって答えられるのは前者のみ
 - 最も簡潔な後者への答えは「あなたの想像力次第」

MSEFA の理論と実例

MSEFA の設計思想

- 格フレーム辞書より細粒度の意味記述
- 読みのポテンシャル空間を特定する必要
 - 解釈は必ずしも“正しい読み”の特定ではない
- 読みの深さをパラメータ化する必要
 - 最適な読み(の深さ)は課題依存的である
- 意味解析の統語情報への依存性を軽減
 - 統語情報は意味解析に不可欠か? – No

MSEFA の異例な特徴

- 理解内容の正確で詳細な記述に特化
- 格(助詞)パターンなどの意味の統語的実現
状況は記述の対象外
 - 具現化のための最適な形式の選択の問題は捨象
 - この問題は, 語彙概念構造分析 (LCS), FrameNet,
格フレーム辞書にお任せします

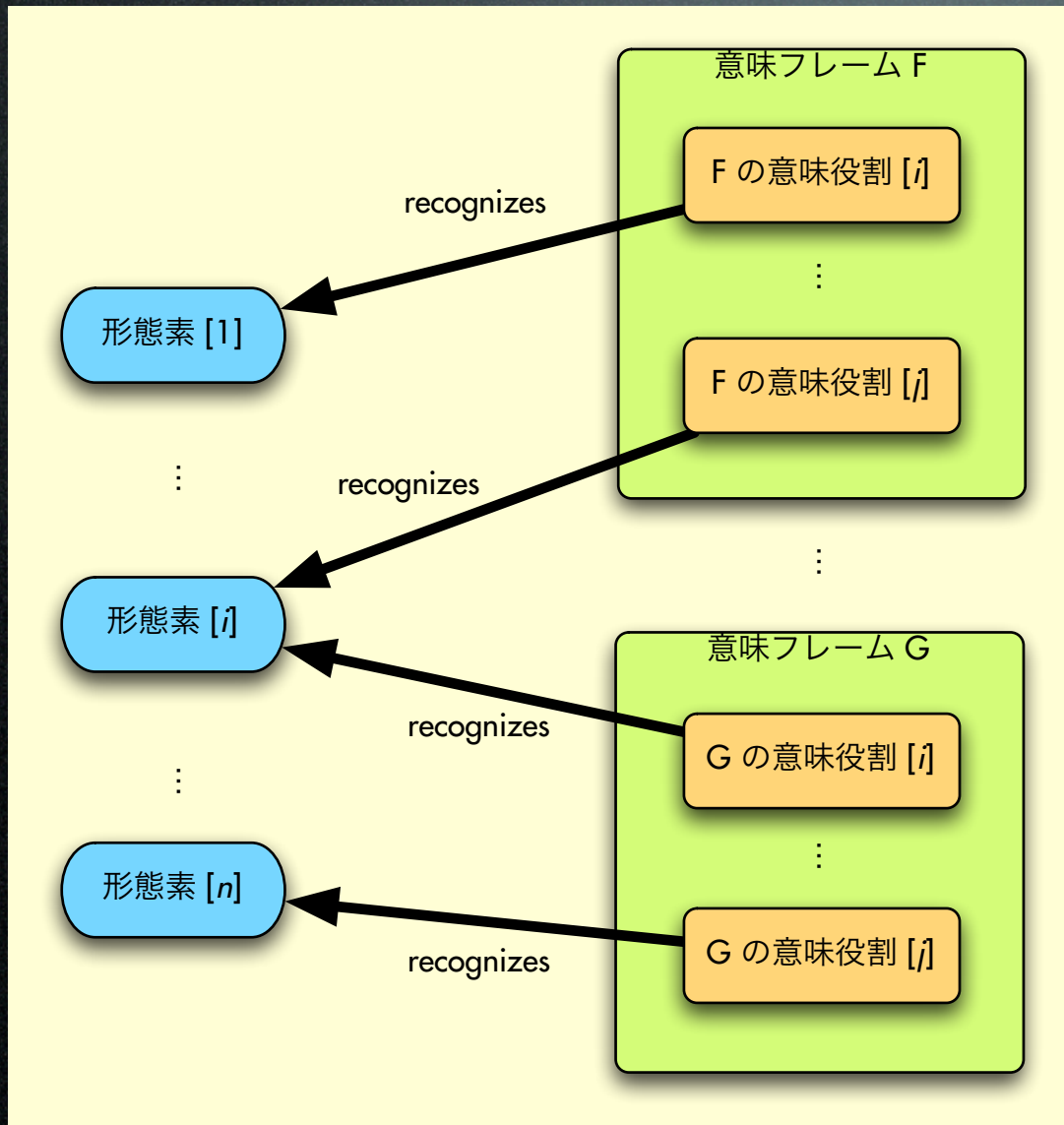
意味役割タグづけとは？

- 状況に依存しない意味型 (e.g., ヒト, 地震) の特定ではなく
- 状況内の意味役割 (e.g., 加害体, 被害者, 犠牲者) の特定
- MSFA は形態素列と複数の意味フレーム列 (= 意味役割のグループ化の列) との対応づけシステム

望ましい意味分析の条件

1. 機械学習可能なぐらい十分な一貫性
2. 文書分野を選ばない網羅性
3. (半)自動化可能な解析 (e.g., 格フレーム辞書) の精度を越える具体性, 特定性
4. 用途 (e.g., 機械翻訳) を限定しない汎用性
5. 文脈内の語の意味の多元性を記述する柔軟性
 - MSFA は [3,4,5] を最重要視
 - シソーラス型 (e.g., EDR, IPAL, WN) は [5] で失格

形態素と意味役割の対応

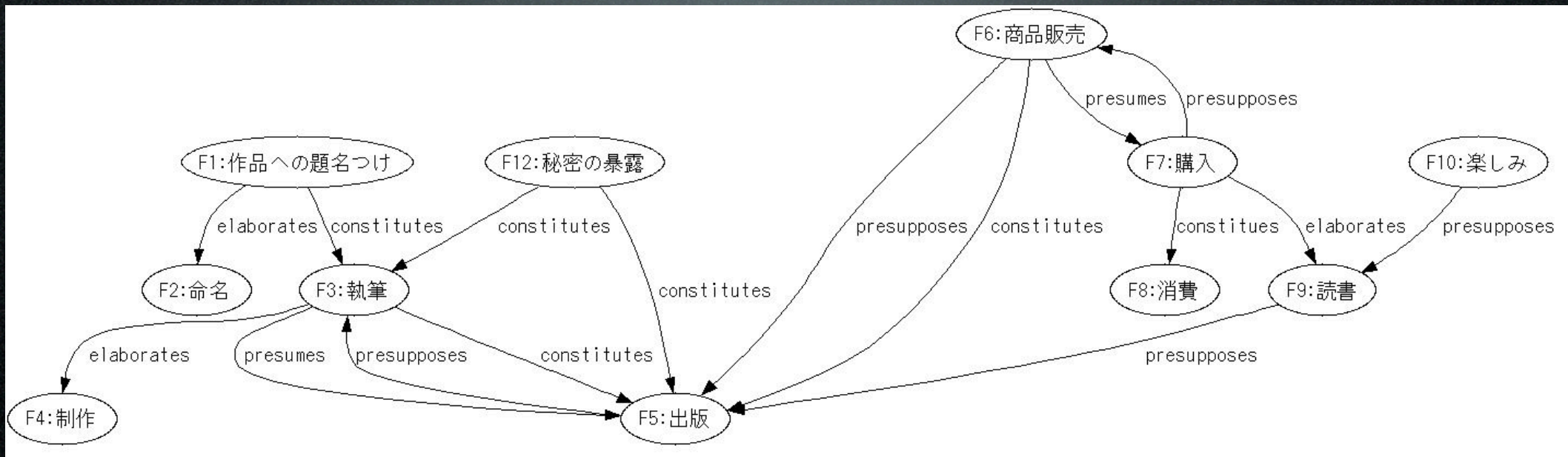


- 意味解釈は並列パターン認識の一種
- 文中での形態素 (= 語) と意味役割 (= フレーム要素) の対応は一對多
- 矛盾のない限り, 形態素は幾つの意味役割を実現しても構わない

意味フレーム分析の要点

- 「先日二人組の男に襲われたXY銀行の支店は、その先日も不渡りに見舞われたばかりだった」
 - [二人組の男, XY銀行]: --linked to--> <銀行強盗>フレームの<強盗, 銀行>役割
 - [不渡り, XY銀行]: --linked to--> <打撃発生>フレームの<打撃, (経済)活動体>役割

MSFA を基にした意味フレームの組織化の記述



- “「ホワイトハウスの内側」と題する本が十四日、米国で発売される” (S-ID: 950112062-001) の MSFA を基に FOCAL Wiki ページ (<http://61.115.230.87/~mutiyama/cgi-bin/hiki/hiki.cgi?FrontPage>) で自動生成 (Graphviz 使用)

文中の語の意味の多元性

- 任意の文 s 中で任意の形態素 m/s に割り当てられる解釈役割 (= 意味役割) $r(m/s)$ はただ一つだとは限らない
 - 意味役割の文脈内相対化による語義の曖昧性解消への非一義化的アプローチ
- (BFN と異なり)形態素に付与される意味役割の数に理論上の上限なし
 - 理論言語学の“常識”からは意識的に逸脱

メタファーやメトニミー

- 生成辞書理論 (Pustejovsky 1995) と同じぐらいはうまくメトニミーを扱える
 - GL と異なり, Qualia 構造の記述はフレーム動詞の相互作用の結果として動的に与えられる
- メタファー (“圧力をかける”, “手を封じる”) がちゃんと扱える
 - メタファーは動的に生成されると考えるので, 概念メタファー理論 (Lakoff & Johnson 1980, 1999) とは違い, メタファーの閉じた目録があるとは考えない

現状と今後の方針

現状

- 基礎が固まる前に開発してもしょうがない
 - 評価版を段階的に開発，公開し，軌道修正しながら小規模だが良質の(教師)データの提供
- 難しいのはタグづけ仕様の決定だけでなくタグづけ作業者の確保 (育成を含めて)
 - MSFA は教授可能，習得可能な技能だと判明したが，習得/体得には集中的な訓練が必要
 - 適当な作業者を必要に応じて雇うより，適切な大学の言語学研究室と連携したほうが得なはず

究極の未来?

- 「英辞郎」のように有志の (Web 上での) Open Development に委ねるのがイチバン
 - 何しろ開発費はタダ
 - 作業者(言語学者)も納期に縛られないで嬉しい
- そのために必要なのは
 - MSFA の体系化, 作業の明示化
 - 解析結果の品質管理 (FOCAL Wiki で部分的に実現)

今後の方針と課題

- タグづけ結果から有意味な意味フレームを選
定し明示的な定義を与える
 - 当面は一部の比較的優先順位の高いフレームのみに定
義を与える予定
 - 評価版公開の目的は選定の条件の特定
- 次の段階で意味フレームのデータベース化
- タグづけ仕様の文書化と公開

お断りとお願い

- 仮に意味役割タグづけの企画が正しいゴールをめざしているとしても、先はまだまだ長い
- 従って
 - 過度の期待はもたないでください
 - フィードバックがあれば、その分だけ皆さんの期待に答えられる可能性が高くなります
 - “こういう文はどうするの?!”のような挑戦も歓迎

謝辞

内山将夫 (NICT)

金丸 敏幸, 中本 敬子, 野澤 元, 龍岡昌弘

(FOCAL 研究グループ)

黒宮 公彦 (大阪学院大学)

竹内孔一 (岡山大学)

石山 昌代, 大谷 直輝, 鬼頭 修, 横森大輔

(京大山梨研究室)

