

Semantic Loss を用いた意図解釈モデルの 半教師ありマルチタスク学習の試み

土田 正明

株式会社コトバデザイン

masa.tsucchi@cotobadesign.com

1 はじめに

本稿では、対話インターフェースの基本機能である意図解釈について、深層学習による半教師ありマルチタスク学習を考える。意図解釈としては、ユーザの発話から_intentとスロットを認識するタスクを扱う。

intent 発話意図を扱うためのフレームの種別。

スロット intentのフレームに定義される属性。

例えば、intent「乗り換え案内」のフレームが「乗り換え案内 (始点駅, 終点駅)」の場合、「始点駅」「終点駅」がスロットとなる。intent認識とスロット認識のマルチタスク学習により、独立に学習するよりも高精度になるという結果が報告されている [6, 5, 2]。

本稿では、訓練中にintent認識とスロット認識の出力間の制約を満たすようにラベル無しデータを活用する半教師ありマルチタスク学習法を提案する。これは、ラベル付きデータでは、ラベルに従って学習することで自然に制約を満たすように学習されるが、ラベル無しデータは制約を満たすように学習されるとは限らない問題に対処することで、精度の向上を狙った手法である。本稿の貢献は以下の通りである。

1. intent認識とスロット認識の出力間の制約を用いた半教師ありマルチタスク学習として、Semantic Loss Function [9] を用いた一手法を示す。
2. 2つのデータセットによる評価実験を通し、提案手法の有効性を示す。
3. 我々の知る限り、intent認識とスロット認識の出力間の制約を用いた半教師ありマルチタスク学習は、本稿が初である。

2 提案手法

Semantic Loss Function 命題論理の文を用いたロス関数であり、文 α に含まれる原始命題を $\mathbf{X} =$

$\{X_1, \dots, X_n\}$ とし、各原始命題が true の確率を表す $\mathbf{p} = \{p_1, \dots, p_n\}$ を用いて、以下で定義される。

$$L^s(\alpha, \mathbf{p}) \propto -\log \sum_{\mathbf{x} \models \alpha} \prod_{\mathbf{x} \models X_i} p_i \prod_{\mathbf{x} \models \neg X_i} (1 - p_i)$$

$\mathbf{x} \models \alpha$ の \mathbf{x} は、文 α を充足する真理値割り当てである。直感的には、 $L^s(\alpha, \mathbf{p})$ は \mathbf{p} から \mathbf{X} の真理値をサンプリングした場合に、文 α が充足される真理値割り当てが生成される確率の負の対数尤度に比例する関数と考えられる。Semantic Loss Function は、文が確率 1 で充足される場合に 0 になる。

intent認識とスロット認識の出力間制約 スロットは、intentのフレームに定義される。そのため、スロットは、それが定義されたintentと同時に認識される必要がある。例えば「渋谷からの電車の乗り継ぎを知りたいんだけど」の「渋谷」を「始点駅」として扱う場合は、「始点駅」が定義された「乗り換え案内」等のintentも同時に認識されるべきである。一方で、スロットが定義されたintentの発話であっても、スロットの情報が存在するとは限らない。例えば「乗り換えを調べたいんだけど」のintentは「乗り換え案内」ではあるが、スロットの情報は存在しない。また、「渋谷からの乗り換えはいつも混んでるね」が「始点駅」のスロットがないintentである場合、「渋谷」は文脈的には類似するが「始点駅」のスロットには該当しないと認識できるであろう。

上記を考慮し、発話がいずれかのintentとして認識され、かつ、発話内の全 token が、認識されたintentのスロットか、どのスロットでもない場合に true になる文を導入する。

$$\text{Intent-slot-rel} : \forall_i (\text{intent}_i(u) \wedge (\bigwedge_{t \in u} (\bigvee_{j \in S_i} \text{slot}_j(u_t))))$$

\vee は論理和、 \wedge は論理積、 $\text{intent}_i(u)$ は発話 u が i 番目のintentの場合に true となる原始命題、 S_i は

i 番目のインテントに定義されたスロットのインデックスの集合, $\text{slot}_j(u_t)$ は発話 u の t 番目の token が j 番目のスロットの場合に true となる原始命題である. 各 token が必ずスロットに属するわけではないため, 全ての S_i に対して, どのスロットでもないことを表す特殊なスロットを含める.

Intent-slot-rel を用いた半教師ありマルチタスク学習のロス関数は以下となる.

$$L(\mathbf{L}, \mathbf{U}) = \frac{1}{|\mathbf{L}|} \sum_{u_i \in \mathbf{L}} \text{MultiTaskLoss}(u_i) + \frac{w_c}{|\mathbf{U}|} \sum_{u_j \in \mathbf{U}} L^s(\text{Intent-slot-rel}, \mathbf{p}_{u_j}^I, \mathbf{p}_{u_j}^S, u_j)$$

\mathbf{L} はラベル付き発話集合 (以降, ラベル付きデータと呼ぶ), \mathbf{U} はラベル無し発話集合 (以降, ラベル無しデータと呼ぶ), $\text{MultiTaskLoss}(u)$ はラベル付き発話 u に対するマルチタスクの教師あり学習のロス関数である. $L^s(\text{Intent-slot-rel}, \mathbf{p}_{u_j}^I, \mathbf{p}_{u_j}^S, u_j)$ は, Semantic Loss Function の定義と文 Intent-slot-rel に従って, 発話 u のインテントの確率 p_u^I , 発話 u 内の各 token のスロットの確率 p_u^S から計算する. w_c は重みである.

3 評価実験

本節では, 下記 2 点の評価について述べる.

1. インテントとスロットの出力間の制約を用いた半教師ありマルチタスク学習の効果
2. ラベル付きデータサイズと精度向上効果の関係

3.1 意図解釈モデル

Encoder-Decoder に Attention を入れ, インテントとスロットの認識で Encoder と embedding のパラメータを共有したモデルを用いる.

Encoder は Bi-directional LSTM, Decoder は LSTM で, それぞれ一層である. 入力の Token の ID 列 $t = (t_1, \dots, t_T)$ を, Encoder, Decoder 共通の Embedding 層で $e = (e_1, \dots, e_T)$ に変換し, Encoder と Decoder の各時刻の入力とする. Embedding 層と LSTM の出力には dropout を適用する.

インテント認識は Encoder の順方向と逆方向のそれぞれの最後の出力を連結したベクトル $h_T^{enc} = [\overrightarrow{h_T^{enc}}, \overleftarrow{h_1^{enc}}]$ を入力に, 2 層のパーセプトロン (活性化関数は ReLU) でインテントの種類数次元のベクトルに写像して, softmax を適用する.

スロット認識は, Decoder の時刻 t の LSTM の出力と Attention を用いて計算したコンテキストベクトル c_t^{attn} を連結したベクトル $h_t^{dec} = [h_t^{dec}, c_t^{attn}]$ を入力に, 2 層のパーセプトロン (活性化関数は ReLU) でスロットの種類数次元 (IOB タグ変換後) のベクトルに写像して, softmax を適用する. コンテキストベクトル c_t は, Encoder の出力 $h_t^{enc} = [\overrightarrow{h_t^{enc}}, \overleftarrow{h_t^{enc}}]$ と Decoder の時刻 t の出力 h_t^{dec} を用いて, 以下の通りに計算する.

$$c_t^{attn} = \sum_{j=1}^T \alpha_{t,j} h_j^{enc}$$

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{k=1}^T \exp(e_{t,k})}$$

$$e_{t,k} = W_{conv} h_k^{enc} \cdot h_t^{dec}$$

W_{conv} は, 次元を揃えるための線形変換で, 学習パラメータである.

訓練では, ラベルの有無の区別をせずにミニバッチを作り, 内部で分けてロスを計算する. $\text{MultiTaskLoss}(u)$ には, 発話 u のインテント認識の交差エントロピーと, 発話 u 内の各 token のスロット認識の交差エントロピーの平均との和を用いる.

3.2 実験設定

3.2.1 データセット

In-house 我々が開発した日本語のデータセットである. 38 種類のインテント (一つは「その他」) と 54 種類 (IOB フォーマット変換前) のスロットが定義されている. インテントとスロットは, 天気予報, 乗り換え検索, グルメ検索など, ウェブサービスの API を参考に設計されている. 「新宿から渋谷の乗り換えはいつも混んでるね」などの紛らわしい例や, 「今日友達と飲むから, 新宿から渋谷の乗り換えを教えて」などの端的ではないがインテントを含む例も含まれている. 約 40% が「その他」である. モデル選択のための開発用を含む訓練に 4 万例, テストに 4 千例を用いた.

Snips[1] 一般に公開¹ されている英語のデータセットである. 7 種類のインテント (AddToPlaylist, BookRestaurant, GetWeather, PlayMusic, RateBook, SearchCreativeWork, SearchScreeningEvent) と 39 種類のスロット (IOB フォーマット変換前) が定義されている. validate から始まるファイルにま

¹<https://github.com/snipsco/nlu-benchmark/tree/master/2017-06-custom-intent-engines>

700 例をテスト用とし、それ以外の 15,884 例を全てモデル選択のための開発用を含む訓練に用いた。

3.2.2 実験方法

比較手法 1)Intent-slot-rel を用いた半教師ありマルチタスク学習 (提案), 2) 出力間制約を用いない半教師ありマルチタスク学習, 3)Intent-slot-rel を用いた教師ありマルチタスク学習, 4) 教師ありマルチタスク学習を比較する。2 には, Semantic loss function を用いた半教師あり学習である One-hot ロス [9] を用いる。

$$L^s(\text{One-hot}, \mathbf{p}) \propto -\log \sum_{i=1}^n p_i \prod_{j=1, i \neq j}^n (1 - p_j)$$

One-hot ロスは, いずれかの原始命題が true で, それ以外が全て false であれば 0 となる。ロス関数は w_o は重み, u_j は発話 u_j の token の集合として, 以下の通りとなる。 p_t^S は t のスロットの確率である。

$$\begin{aligned} L(\mathbf{L}, \mathbf{U}) &= \frac{1}{|\mathbf{L}|} \sum_{u_i \in \mathbf{L}} \text{MultiTaskLoss}(u_i) \\ &+ \frac{w_o}{|\mathbf{U}|} \sum_{u_j \in \mathbf{U}} L^s(\text{One-hot}, \mathbf{p}_{u_j}^I, u_j) \\ &+ \frac{w_o}{|\mathbf{U}|} \sum_{u_j \in \mathbf{U}} \frac{\sum_{t \in u_j} L^s(\text{One-hot}, \mathbf{p}_t^S, t)}{|u_j|} \end{aligned}$$

実験設定 ロス関数以外を同じ条件するために, 各手法で使用しない重み w_c, w_o を 0 に設定した。Intent-slot-rel と One-hot のロスを用いる場合の重みは, 予備実験から定め, In-house では, $w_c = 0.05, w_o = 0.01$, Snips では, $w_c = 0.1, w_o = 0.005$ を用いた。ラベル付きデータのサイズは 1 千, 5 千, 1 万に設定し, 残りをラベル無しデータとして利用した。

評価指標には, テストデータの F 値のマイクロ平均 (以降, マイクロ平均 F 値と呼ぶ) を用いた。スロット認識は, token 毎のラベルの一致ではなく, ラベル列からデコードして, token 列を抽出し, 完全一致する場合に正解とした。各手法で同じ乱数シード 20 個を用いて, マイクロ平均 F 値を 20 サンプル取得した。モデル選択では, 訓練データから 1 千例を開発用にサンプリングし, 30 epoch まで学習し, 開発用データの Intent 認識とスロット認識のマイクロ平均 F 値の調和平均が最大となる epoch のモデルを採用した。

In-house では, 日本語 Wikipedia を用いて word2vec(skip-gram, negative-sampling) で訓練した 128 次元の Embedding を用いた。token は mecab で取得した。Snips では, Glove で訓練された 200 次元

表 1: In-house: Intent のマイクロ平均 F 値の平均

| | 1,000 | 5,000 | 10,000 |
|-----------------------|-------|-------|--------|
| 半教師 (Intent-slot-rel) | 0.873 | 0.967 | 0.982 |
| 半教師 (One-hot) | 0.872 | 0.966 | 0.982 |
| 教師 (Intent-slot-rel) | 0.870 | 0.967 | 0.981 |
| 教師 | 0.869 | 0.966 | 0.980 |

表 2: In-house: スロットのマイクロ平均 F 値の平均

| | 1,000 | 5,000 | 10,000 |
|-----------------------|--------------------|--------------------|--------------------|
| 半教師 (Intent-slot-rel) | 0.599 \diamond † | 0.807 \diamond † | 0.863 \diamond † |
| 半教師 (One-hot) | 0.579 | 0.796 | 0.858 |
| 教師 (Intent-slot-rel) | 0.572 | 0.796 | 0.856 |
| 教師 | 0.573 | 0.798 | 0.857 |

表 3: Snips: Intent のマイクロ平均 F 値の平均

| | 1,000 | 5,000 | 10,000 |
|-----------------------|-------|-------|--------|
| 半教師 (Intent-slot-rel) | 0.966 | 0.980 | 0.983 |
| 半教師 (One-hot) | 0.966 | 0.980 | 0.983 |
| 教師 (Intent-slot-rel) | 0.965 | 0.980 | 0.983 |
| 教師 | 0.968 | 0.979 | 0.981 |

表 4: Snips: スロットのマイクロ平均 F 値の平均

| | 1,000 | 5,000 | 10,000 |
|-----------------------|--------------------|--------------------|--------|
| 半教師 (Intent-slot-rel) | 0.769 \diamond † | 0.852 \diamond † | 0.880 |
| 半教師 (One-hot) | 0.759 | 0.848 | 0.879 |
| 教師 (Intent-slot-rel) | 0.756 | 0.847 | 0.881 |
| 教師 | 0.755 | 0.847 | 0.877 |

表 5: スロット認識の提案手法の効果量 (Hedges の g)

| データ | 比較手法 | 1,000 | 5,000 | 10,000 |
|----------|----------------------|-------|-------|--------|
| In-house | 半教師 (One-hot) | 1.118 | 1.051 | 0.740 |
| | 教師 (Intent-slot-rel) | 1.305 | 1.124 | 0.911 |
| | 教師 | 1.633 | 0.891 | 0.973 |
| Snips | 半教師 (One-hot) | 0.950 | 0.612 | - |
| | 教師 (Intent-slot-rel) | 1.126 | 0.713 | - |
| | 教師 | 1.325 | 0.904 | - |

の Embedding² を用いた。token はスペース区切りで取得した。いずれも, Encoder の各方向の LSTM, Decoder の LSTM の隠れ状態は 256 次元とした。

3.3 結果

In-house の結果を表 1, 2 に, Snips の結果を表 3, 4 に示す³。効果の分析のために, 4 つの手法の全ペアに対し, マイクロ平均 F 値 20 サンプルを用いて対応のあるデータの t 検定で p 値を取得し, 多重比較による False Discovery Rate の制御のために Benjamini-Hochberg 法で補正した q 値を算出した。表 1 から 4 の” \diamond ”は, 提案手法と 2 つの教師ありとの比較で全て q 値が 0.05 以下で, かつ, 値が大きい方に付与している。”†”は, 提案手法と One-hot による半教師ありとの比較についての同様を表す。表 2, 4 に示す通り, スロット認識

²<http://nlp.stanford.edu/data/glove.6B.zip>

³表の値は全て小数点第 4 位を四捨五入している。

で有意な差が見られたため、表5に、有意差がある場合の効果量 (Hedges の g^4) を示した。表1から5の結果の傾向は下記の通りである。

- 提案手法は、_intent認識に対する有意な変化無しで、スロットの認識精度向上が見られた。
- ラベル付きデータのサイズが小さい場合に効果が大きい傾向にあった。
- Intent-slot-rel ロスは、ラベルありデータのみでは効果が見られなかった。
- One-hot ロスは効果が見られなかった。

4 関連研究

スロット認識に関連する系列ラベリングの半教師ありマルチタスク学習の先行研究として、文献 [8, 4] が存在する。これらは、ラベリング対象の token の周辺単語を予測する補助タスクを用いた半教師ありマルチタスク学習で、複数の主タスクの出力間の制約を用いる方法ではないため、本稿とは相補的である。

シンボリックな制約や知識を活用するための先行研究には、文献 [3, 7] などが存在する。文献 [3] は、論理式で制約を表現でき、ラベルなしデータの活用も可能な点で本稿と類似している。評価極性判定における「A but B なら文全体と B の評価極性は同じ」や固有表現認識における「B-ORG の次に I-PER は取れない」など、同じタスク内の制約のみ検討されているが、枠組みとしての一般性は高いので、本稿の着眼点や提案した制約は適用可能と考えられる。文献 [7] は、intent とスロットを認識するための正規表現ルールを利用し、ルールがマッチしたことを表す情報を深層学習に取り込む方法を提案している。intent とスロットのそれぞれの認識精度を向上するための方式であるため、本稿とは相補的である。

意図解釈のマルチタスク学習の先行研究として、文献 [6, 5, 2] などが存在する。本稿の提案は、特定のモデルに依存しないため、これらにも応用可能である。

5 おわりに

本稿では、intent 認識とスロット認識の出力間の制約を用いた半教師ありマルチタスク学習法について述べた。評価実験により、提案手法は、スロットの認

識の精度向上効果が見られ、また、ラベル付きデータのサイズが少ない場合に効果が大きい傾向にあることを示した。今後は、提案手法の汎用性の評価のために、他のモデル、データセットを用いた評価や、他の半教師あり学習との併用等で実験を行う予定である。

謝辞

本研究の一部は、平成30年度総務省委託研究「高度対話エージェント技術の研究開発・実証」の成果によるものである。

参考文献

- [1] Alice Coucke, Alaa Saade, Adrien Ball, héodore Bluche, Alexandre Caulier, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, Maël Primet, and Joseph Dureau. Snips voice platform: an embedded spoken language understanding system for private-by-design voice interfaces. *CoRR*, 1805.10190v3, 2018.
- [2] Chih-Wen Goo, Guang Gao, Yun-Kai Hsu, Chih-Li Huo, Tsung-Chieh Chen, Keng-Wei Hsu, and Yun-Nung Chen. Slot-gated modeling for joint slot filling and intent prediction. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)*, pp. 753–757, 2018.
- [3] Zhiting Hu, Xuezhe Ma, Zhengzhong Liu, Eduard Hovy, and Eric Xing. Harnessing deep neural networks with logic rules. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 2410–2420, 2016.
- [4] Ouyu Lan, Su Zhu, and Kai Yu. Semi-supervised training using adversarial multi-task learning for spoken language understanding. In *Proceedings of 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6049–6053, 2018.
- [5] Changliang Li, Cunliang Kong, and Yan Zhao. A joint multi-task learning framework for spoken language understanding. In *Proceedings of 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6054–6058, 2018.
- [6] Bing Liu and Ian Lane. Attention-based recurrent neural network models for joint intent detection and slot filling. In *Proceedings of Interspeech 2016*, pp. 685–689, 2016.
- [7] Bingfeng Luo, Yansong Feng, Zheng Wang, Songfang Huang, Rui Yan, and Dongyan Zhao. Marrying up regular expressions with neural networks: A case study for spoken language understanding. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 2083–2093, 2018.
- [8] Marek Rei. Semi-supervised multitask learning for sequence labeling. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 2121–2130, 2017.
- [9] Jingyi Xu, Zilu Zhang, Tal Friedman, Yitao Liang, and Guy Van den Broeck. A semantic loss function for deep learning with symbolic knowledge. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, pp. 5502–5511, 2018.

⁴平均値の差が各群のサンプルサイズを考慮した標準偏差の平均値の何倍かを表す指標であり、0.5 以上で中程度、0.8 以上で大きな効果と言われている。