



生物医学ドメインコーパスである PubMed<sup>2</sup>, PMC<sup>3</sup> を用いて事前学習を行った BioBERT モデルを提案した。BioBERT を CHEMPROT でファインチューニングした結果, CHEMPROT を用いた関係抽出において現時点での最高精度を達成している。

大規模なコーパスからエンティティ間の関係を抽出する技術として, Open Information Extraction (Open IE) [1]がある。Open IE では, 主語・目的語となるエンティティのペアと, エンティティ間の関係を合わせた3つ組を抽出する。例えば, “Shakespeare is author of Hamlet.” という文に対して, Open IE を用いると (Shakespeare, is, author), (Shakespeare, is author of, Hamlet) といった3つ組が抽出できる。2つのエンティティとその関係を抽出するという点において, 関係抽出タスクとの関連性があり, 後述するマルチタスク学習に利用できると思われる。

マルチタスク学習とは, 主タスクに関連する補助タスクを同時に学習することにより, タスク間で共通する情報を学習させる手法である。Ruder ら [9]によると, 複数タスクによる「潜在的なデータ拡張」, タスク間に共通する特徴に着目する「Attention focusing」, タスク間の難易度の差を利用する「盗み聞き」といった要素を, マルチタスク学習の精度向上の根拠としている。本研究では補助タスクにラベル付きデータを使用せず, Open IE によって自動でラベル付きデータを作成する。

### 3 提案手法

本研究では, 後述する主タスクと補助タスクを同時に学習するマルチタスク学習を提案する。提案手法の概念図を図2に示す。

なお, CHEMPROT には文をまたいだ関係も存在するが, 全体の1%未満であるため, 本研究では同一文内に出現するタンパク質と化学物質のペアのみを対象とする。したがって, 入力には事前にアノテーションされたタンパク質と化学物質のペアを含む1文を用いる。各入力文に対して主タスクのラベルと補助タスクのラベルが1つずつ付与されている。この入力文を共通のエンコーダーに入力し, エンコーダーの出力を各タスクの分類器の入力に用いる。

訓練時には主タスクと補助タスクの学習を同時に行い, 各タスクの損失の重み付き和を最終的な損失関数

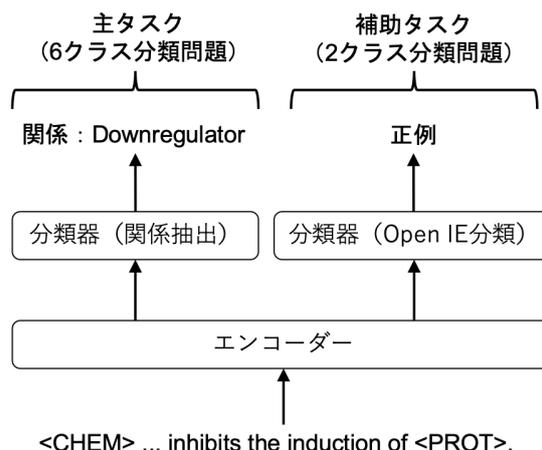


図 2: 提案手法の概念図

とする。主タスクの損失関数  $L_m$ , 補助タスクの損失関数  $L_a$ , 補助タスクの重みを表すパラメータ  $w$  を用いて, 最終的な損失関数は式 (1) で表せる。

$$L = (1 - w)L_m + wL_a \quad (1)$$

$w = 0$  のとき, 主タスクのみの学習が行われるので, これをベースラインとして用い, マルチタスク学習を行なった場合と比較する。ラベルを予測する際には補助タスクは行わず, 主タスクの出力のみを用いる。以下に主タスク, 補助タスクの詳細について述べる。

#### 3.1 主タスク

CHEMPROT では, タンパク質と化学物質のペアについて, 関係を持つペアに対しては5種類の関係のうちどの関係にあてはまるかがアノテーションされている。そのため主タスクでは, タンパク質と化学物質を含む文を入力とし, 2つのエンティティの関係を(正例5クラス+負例1クラス)の6クラスに分類するタスクを取り扱う。ここで, 正例の5クラスは「エンティティ間に関係があり, CHEMPROT で定義される5種類の関係のうちどれに当てはまるか」を示し, 負例は「CHEMPROT で定義される5種類の関係のどれにも当てはまらない」ということを示す。主タスクの損失関数には Cross Entropy Loss を用いる。

#### 3.2 補助タスク

補助タスクでは入力文に Open IE で抽出できる関係が含まれているかどうかを分類する。まず, 入力文

<sup>2</sup><http://www.ncbi.nlm.nih.gov/pubmed/>

<sup>3</sup><http://www.ncbi.nlm.nih.gov/pmc/>

に対して Open IE ツール Stanford CoreNLP[8] を用い、(エンティティ1, エンティティ間の関係, エンティティ2) の3つ組を抽出する。このとき、エンティティ1, 2 にそれぞれタンパク質と化学物質が含まれている3つ組が抽出できた場合、入力文に対する補助タスクのラベルは正例ラベルとなる。それ以外の場合、すなわち Open IE で抽出した3つ組の中にタンパク質と化学物質をそれぞれエンティティ1, 2 に含むものがない場合は負例となる。表1に例を示す。

以上で定義したラベルを用い、補助タスクではタンパク質と化学物質を含む文を入力とし、入力文が正例か負例かの2クラスに分類するタスクを取り扱う。補助タスクの損失関数には、主タスクと同様に Cross Entropy Loss を用いる。

## 4 実験

### 4.1 データセット

データセットとして CHEMPROT を用いる。CHEMPROT の統計量を表2に示す。

前処理は Lim ら [7] と同様のものを用いる。具体的には、以下の手順で前処理を行う。

1. アブストラクト中のタンパク質と化学物質を含む文を抽出する。
2. タンパク質, 化学物質をそれぞれ “<PROT>”, “<CHEM>” に置換する。
3. 関係判定対象のタンパク質と化学物質のペアを決め, 判定対象外のエンティティを “<OTHER>” に置換する。

図3の例では、文中にタンパク質が2つ、化学物質が1つ存在するため、関係判定対象の文が2つになる。

### 4.2 評価指標

評価には、主タスクの出力結果を用いる。評価指標としては、CHEMPROT で用いられており、一般的な関係抽出タスクの評価指標でもある精度、再現率、F値を用いる。なお、評価の際は CHEMPROT で用いられている評価スクリプトを使用して値を算出した。

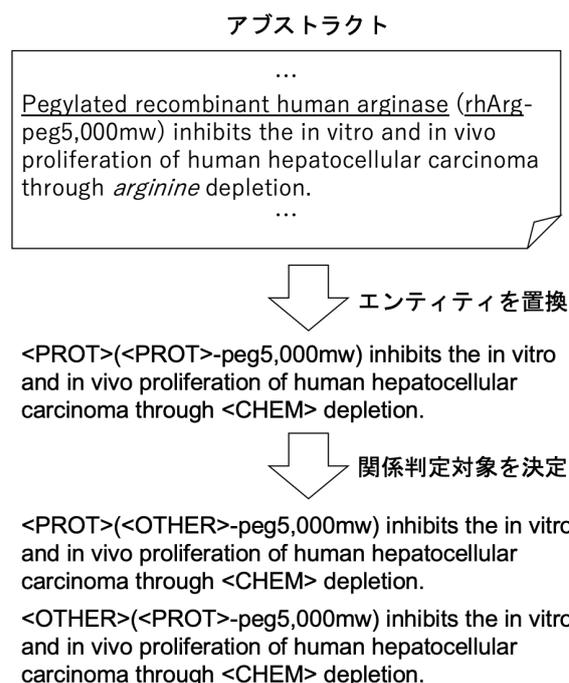


図3: 前処理の例。下線がタンパク質, 斜体が化学物質を表す。

### 4.3 実験設定

エンコーダーは BioBERT[6] を使用し、事前学習済みのモデルを用いて CHEMPROT によってファインチューニングを行う。式(1)における補助タスクの重み  $w$  は、ベースラインでは主タスクのみでの関係抽出とするため  $w = 0$ 、提案手法では  $w = 0.05, 0.10, \dots, 1.00$  のうち、開発データで最も F 値の高かった  $w = 0.05$  を用い、バッチサイズは 16 とした。モデルの最適化には Adam[4] を用い、Adam の各パラメータに関しては学習率  $\alpha = 0.00005$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , Weight decay = 0.01 とした。

### 4.4 結果

ベースラインと提案手法に加え、比較のために SPINN [7] と SciBERT [2]<sup>4</sup> と BioBERT の元論文 [6] で報告されている v1.1 による結果を表3に示す。ベースラインと比較すると、マルチタスク学習によって F 値が向上していることがわかる。

また、主タスクと補助タスクの相関を調べるため、学習データ中の各タスクの正例・負例数を表4に示す。

<sup>4</sup><https://paperswithcode.com/paper/scibert-pretrained-contextualized-embeddings>

表 1: 補助タスクのラベル付けの例

入力文	Open IE 出力	ラベル
<CHEM> inhibits the induction of <PROT>.	(<CHEM>, inhibits, induction of <PROT>) (<CHEM>, inhibits, induction)	正例
It is derived from <CHEM> and <PROT>.	(It, is derived from, <CHEM> and <PROT>)	負例

表 2: CHEMPROT の統計量

	訓練	開発	テスト
アブストラクト数	1,020	612	800
化学物質数	13,017	8,004	10,810
タンパク質数	12,735	7,563	10,018
関係判定対象文数	16,522	10,362	14,397
関係抽出正例文数	4,133	2,412	3,441
Open IE 正例文数	1,761	1,072	1,521

表 3: 実験結果

	適合率	再現率	F 値
SPINN	0.7480	0.5600	0.6410
SciBERT	-	-	0.7612
BioBERT v1.1	0.7702	0.7590	0.7646
ベースライン	0.7766	0.7510	0.7636
提案手法	0.7732	0.7730	0.7731

なお主タスクでの正例とは、CHEMPROT で定義されている 5 種類の関係に当てはまるものであり、逆にそれらの関係に当てはまらないものは負例となる。

表 4: 主タスクと補助タスクの相関

		主タスク		合計
		正例	負例	
補助タスク	正例	805	956	1,761
	負例	3,328	11,433	14,761
合計		4,133	12,389	16,522

主タスクでの正例発生確率と補助タスクでの正例発生確率でカイ二乗検定を用いて検定した結果、有意差が認められた ( $p < 0.01$ )。すなわち、主タスクで「関係あり」と判定されるエンティティのペアと、補助タスクにて Open IE によって関係が定義づけられるエンティティのペアの間には相関があり、マルチタスク学習によってこれらに共通の特徴量を学習したと考えられる。

## 5 おわりに

本研究では、タンパク質と化学物質間の関係抽出課題において、Open IE によって作成したラベル付き

データを用いてマルチタスク学習を行うことを提案した。CHEMPROT コーパスを用いた実験により、提案手法によって関係抽出の F 値を 0.773 に改善することができた。今後の課題としては、CHEMPROT 以外のデータでの提案手法の実験・評価などが挙げられる。

## 参考文献

- [1] Michele Banko, Michael J Cafarella, Stephen Soderland, Matthew Broadhead, and Oren Etzioni. Open information extraction from the web. In *Ijcai*, Vol. 7, pp. 2670–2676, 2007.
- [2] Iz Beltagy, Arman Cohan, and Kyle Lo. Scibert: Pretrained contextualized embeddings for scientific text. *arXiv preprint arXiv:1903.10676*, 2019.
- [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [4] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [5] Martin Krallinger, Obdulia Rabal, Saber A Akhondi, et al. Overview of the biocreative vi chemical-protein interaction track. In *Proceedings of the sixth BioCreative challenge evaluation workshop*, Vol. 1, pp. 141–146, 2017.
- [6] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. Biobert: pre-trained biomedical language representation model for biomedical text mining. *arXiv preprint arXiv:1901.08746*, 2019.
- [7] Sangrak Lim and Jaewoo Kang. Chemical-gene relation extraction using recursive neural network. *Database*, Vol. 2018, , 2018.
- [8] Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. The stanford corenlp natural language processing toolkit. In *Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations*, pp. 55–60, 2014.
- [9] Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*, 2017.