

hispagestyleempty

修士論文

推論における重要度単語を転移させる知識蒸留法

武田 遥暉

主指導教員 白井 清昭

北陸先端科学技術大学院大学
先端科学技術研究科
(情報科学)

令和8年3月

Abstract

English abstract (1200 words)

A horizontal row of fifteen empty circles, evenly spaced, used as a visual element in a document.

概要

日本語の概要

A horizontal row of 20 empty circles, each with a thin black outline, arranged in a single horizontal line.

目 次

第1章 はじめに	1
1.1 背景	1
1.2 目的	1
1.3 本論文の構成	2
第2章 関連研究	3
2.1 知識蒸留	3
2.2 事前学習済み言語モデル	5
2.2.1 BERT の事前学習	5
2.2.2 教師モデルとしての BERT	5
2.3 知識蒸留によって構築されたモデル	5
2.3.1 TinyBERT	5
2.4 モデルの推論根拠の解釈に関する研究	6
2.4.1 Integrated Gradient	6
2.5 本研究の特徴	7
第3章 提案手法	8
第4章 実験・評価	9
第5章 おわりに	10

目 次

3.1 図のキャプション 8

表 目 次

4.1 表のキャプション	9
------------------------	---

第1章 はじめに

1.1 背景

大規模言語モデル（Large Language Model; LLM）は、大規模テキストで事前学習された言語モデルであり、文脈情報を高次元表現として獲得することできまざまな言語タスクにおいて高い性能を示している。特に分類モデルとして用いる場合、入力文全体の意味情報を統合した表現を利用することで、文書分類や感情分析などにおいて高い性能を発揮する。

このような能力を実現するため近年の LLM は大規模なモデルサイズを有しており、その運用には多数の GPU をはじめとする計算機資源を必要とする。その結果、十分な資金力を持つ一部の組織のみが利用可能である点や、クラウド経由で GPU 資源にアクセスできない環境では導入が困難である点が、実社会への幅広い展開を妨げる課題となっている。

そのため、LLM を教師モデル、より少ないパラメータを持つモデルを生徒モデルとし、教師モデルの性能を可能な限り維持しながら、タスク遂行能力を生徒モデルへ転移する手法として知識蒸留が提案されている。特に分類タスクにおいては、性能をほぼ維持したままモデルサイズを半分程度まで削減可能であるなど、高い有効性が示されており知識蒸留はモデル軽量化の代表的な手法として確立している。

しかしながら、知識蒸留は、計算資源の制約下において教師モデルの代替として利用可能な生徒モデルを構築することを目的とした技術であるにもかかわらず、蒸留によって学習された生徒モデルは分類精度を維持できている一方、推論時に重要と判断する単語が教師モデルと一致しないことが報告されている。このような差異は、医療分野や金融分野など判断根拠の信頼性が重視される領域において致命的な問題となり得るため、生徒モデルを教師モデルの単純な代替として用いることには依然として課題が残っている。

1.2 目的

本研究では、教師モデルと同じ単語に注目して生徒モデルが推論を行うよう学習させることを目的とする。対象タスクに対して十分に学習された教師モデルを用意し、以下の手順に従って知識蒸留を適用する。

まず、教師モデルよりも小規模な未学習の生徒モデルを準備し、両モデルに同一の入力インスタンスを与えて推論を行う。次に、得られた推論結果に対して Integrated Gradients を用いて単語重要度を算出し、それらを重要度分布として表現する。最後に、教師モデルと生徒モデルの重要度分布が近づくよう損失関数を設計し、先行研究で提案されている知識蒸留手法に当該損失項を導入することで、生徒モデルの学習を行う。本手法により、分類精度を先行研究と同程度に維持したまま、推論時に重要と判断される単語に関して教師モデルとの整合性を向上させることを狙う。

1.3 本論文の構成

本論文の構成は以下の通りである。2章では先行研究について述べる。3章では提案手法について詳細に説明する。4章では提案手法の評価実験について述べる。最後に5章では、本研究の手法についてのまとめと今後の課題について述べる。

第2章 関連研究

本章では本研究の関連研究について述べる。本章で述べる研究トピックはKnowledge Distillation（以下、知識蒸留）[3]である。知識蒸留は、性能が高くパラメータ数も大きいモデルを教師モデルとし、パラメータ数が小さいモデルへとタスクを解くための能力を転移されるための代表的な技術である。これまでに様々な知識蒸留手法が提案されており、教師モデルの出力層の知識を転移するものや、隠れ状態からの知識も転移させるもの、両方から知識を転移させるもの等、様々なパターンの知識蒸留のアプローチが提案してきた。

また、分類モデルの知識蒸留における教師モデルとしてBERT[1]を用いられることが多い。BERTを教師モデルにして、知識蒸留されたモデルの代表的なものとしてTinyBERT[4]が存在する。

これらの言語モデルの推論根拠を提供する技術としてIntegrated Gradient[5]がある。Integrated Gradientsは画像処理分野で提案されたものではあるが本研究ではモデルが重要であると認識した単語を抽出するために利用する。

以下、2.1節では、知識蒸留の原理的な設計について述べる。2.2.1節では、教師モデルとして用いるBERTについて紹介する。2.3.1節では、知識蒸留によって構築された代表的なモデルであるTinyBERTについて述べる。??節ではモデルが推論する際の根拠を提示する技術としてIntegrated Gradientsを紹介する。最後に、2.5では、本研究の特色について述べる。

2.1 知識蒸留

知識蒸留（Knowledge Distillation）は、高性能である一方、パラメータ数が多い教師モデルから、より小規模な生徒モデルへ特定のタスクを解く能力を転移させるための代表的な技術である。一般に、教師モデルは対象タスクに対して十分に学習（ファインチューニング）されており、生徒モデルは教師モデルの振る舞いを模倣することで効率的な学習を行う。

分類タスクでは、通常、あるドメインに対して数百から数万のラベル付きインスタンスが与えられる。代表的な自然言語理解ベンチマークであるGLUE[7]は、複数の分類タスクから構成されるデータセット群であり、知識蒸留の評価においても標準的に用いられる。GLUEに含まれるMRPC（Microsoft Research Paraphrase Corpus）は、2つの英文が意味的に等価であるかどうかを判定する二値分類タス

クである。各インスタンスにはラベル $y \in \{0, 1\}$ が付与されており、意味的に等価な場合を 1、そうでない場合を 0 と定義する。

例

sentence1: Amrozi accused his brother, whom he called “the witness”, of deliberately distorting his evidence.

sentence2: Referring to him as only “the witness”, Amrozi accused his brother of deliberately distorting his evidence.

label: equivalent (1)

通常の分類モデルの学習では、生徒モデル f_s が入力 x に対して予測するクラス確率分布 $p_s(y | x)$ と、正解ラベル y との間のクロスエントロピー損失を最小化する。この損失関数は、次式で定義される。

$$\mathcal{L}_{\text{CE}} = - \sum_c y_c \log p_s(c | x) \quad (2.1)$$

ここで、 c はクラスを表し、 y_c は正解クラスに対応するワンホット表現である。

知識蒸留では、この正解ラベルに基づく学習に加えて、同一の入力 x に対する教師モデル f_t の出力分布 $p_t(y | x)$ を、生徒モデルが模倣することを目的とする。教師モデルと生徒モデルの出力分布の乖離を測る指標として、Kullback–Leibler (KL) ダイバージェンスが一般に用いられる。このとき、知識蒸留損失は、次式で表される。

$$\mathcal{L}_{\text{KD}} = \text{KL}(p_t(\cdot | x) \| p_s(\cdot | x)) \quad (2.2)$$

ここで、 \cdot はクラス全体を表し、教師モデルが出力する確率分布全体を、生徒モデルが近似することを意味する。

最終的な学習では、正解ラベルに基づくクロスエントロピー損失と、教師モデルの振る舞いを模倣する知識蒸留損失を組み合わせた、次の損失関数を最小化する。

$$\mathcal{L} = (1 - \lambda) \mathcal{L}_{\text{CE}} + \lambda \mathcal{L}_{\text{KD}} \quad (2.3)$$

ここで、 $\lambda \in [0, 1]$ はハイパーパラメータであり、正解ラベルと教師モデルの知識のどちらを重視するかを制御する。

このように知識蒸留では、データセットに明示的に含まれるラベル情報だけでなく、教師モデルが獲得した予測分布という暗黙的な知識を生徒モデルに転移することで、高い性能を維持したままモデルの軽量化を実現することを目的としている。

2.2 事前学習済み言語モデル

2.2.1 BERT の事前学習

BERT[2] は、Transformer[6] を基盤とした双方向言語モデルであり、大規模コーパスを用いた事前学習によって汎用的な言語表現を獲得する。BERT の事前学習では、Masked Language Model (MLM) と Next Sentence Prediction (NSP) の 2 つのタスクが用いられる。

MLM では、入力文中の一部の単語をマスクし、その単語を周辺文脈から予測することで単語表現を学習する。MLM の損失関数は以下のように定義される。

$$\mathcal{L}_{\text{MLM}} = - \sum_{i=1}^N \log p(t_i | T_i) \quad (2.4)$$

ここで、 t_i はマスクされた単語、 T_i は Ft_i 以外の単語からなる入力文を表す。

次に、NSP は 2 つの文が元の文書中で連続しているか否かを判定するタスクであり、文間の関係性を学習することを目的としている。NSP の損失関数は次式で表される。

$$\mathcal{L}_{\text{NSP}} = - [y \log P_{\text{NSP}} + (1 - y) \log(1 - P_{\text{NSP}})] \quad (2.5)$$

ここで、 y は文の連続性ラベルであり、2 つの文が連続している場合は 1、そうでない場合は 0 を取る。また、 P_{NSP} は 2 文が連続しているとモデルが予測する確率を表す。BERT の事前学習では、式 (2.4) および式 (2.5) で定義される損失を同時に最小化することで、単語レベルおよび文レベルの意味関係を考慮した言語表現を獲得する。

2.2.2 教師モデルとしての BERT

本研究では教師モデルとして、BERT を用いる。BERT は事前学習によって大規模コーパスでの事前学習により汎用的な言語表現を獲得しており、このモデルに対して本研究で対象とするタスクでファインチューニングを行うことを行い、教師モデルとして採用する。

2.3 知識蒸留によって構築されたモデル

2.3.1 TinyBERT

TinyBERT[4] は、BERT を教師モデルとして知識蒸留を適用することで構築された軽量な言語モデルである。TinyBERT は、教師モデルの出力層だけでなく、中間層の隠れ状態や Attention 機構の重みも模倣することで、モデルサイズを大幅に削減しながら高い性能を維持することに成功している。

データ拡張 (data augmentation)

学習において TinyBERT は下流タスクでの蒸留効果を高めるためにデータ拡張を併用することが知られている。具体的には、元のデータセットに対してランダムに単語を挿入・削除・置換することで多様な入力インスタンスを生成し、頑健性を向上させる。

学習の流れ

TinyBERT の蒸留は大きく 2 段階で行われる。

事前学習段階 (pre-training) 事前学習段階では、BERT の事前学習で用いられた Masked Language Model (MLM) や Next Sentence Prediction (NSP) に加えて、中間層の隠れ状態や Attention 重みの模倣を目的とした損失を導入する。これにより、生徒モデルは教師の内部表現を効率的に学習し、下流タスクへの基盤を構築する。

下流タスクのファインチューニング (downstream fine-tuning) 下流タスクへの適応はさらに 2 段階（段階的蒸留）で実施される。

第 1 段階 — 中間層蒸留 (intermediate layer distillation) まず中間層の隠れ状態や Attention 重みの模倣損失を導入し、生徒の内部表現を教師に近づけることで下流タスクへの適応準備を行う。

第 2 段階 — 出力層蒸留 (prediction layer distillation) 続いて教師の出力（ソフトラベル）とデータのハードラベルを併用して生徒の出力層を学習させ、教師モデルの推論能力を模倣させる。

TinyBERT はこの学習の流れを通じて BERT の

2.4 モデルの推論根拠の解釈に関する研究

2.4.1 Integrated Gradient

Integrated Gradients（統合勾配法）は、深層学習モデルの予測に対する各入力特徴量の寄与度を評価するための手法である。具体的には、モデルの出力に対する入力特徴量の勾配を積分することで、各特徴量が予測にどれだけ影響を与えたかを定量化する。本研究では、Integrated Gradients を用いて、言語モデルが推論時に重要と判断した単語を特定する。具体的な計算式は以下によって定義される。

2.5 本研究の特徴

第3章 提案手法

(図を貼る)

図 3.1: 図のキャプション

第4章 実験・評価

表 4.1: 表のキャプション

	a	b
1	0.25	0.33
2	0.75	0.66

第5章 おわりに

参考文献

- [1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [3] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015.
- [4] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. TinyBERT: Distilling BERT for natural language understanding. In Trevor Cohn, Yulan He, and Yang Liu, editors, *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 4163–4174, Online, November 2020. Association for Computational Linguistics.
- [5] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks, 2017.
- [6] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.
- [7] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In Tal Linzen, Grzegorz Chrupała, and Afra Alishahi, editors, *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 353–355, Brussels, Belgium, November 2018. Association for Computational Linguistics.