

丹羽健太 特別研究員

データが分散蓄積される時代にも機械学習モデルを最適化。「非同期分散型の深層学習技術」の研究

現在の深層学習ではデータを1カ所に集約したうえでモデルを学習する場合が一般的ですが、近い将来には取り扱うデータ量の増加やプライバシー保護の観点から、データが分散蓄積されることが予想されます。今回はそうした時代にあってもあたかも1カ所にデータを集約して学習したかのような機械学習モデルを最適化できる、「非同期分散型の深層学習技術」の研究に取り組む丹羽健太特別研究員にお話を伺いました。

◆PROFILE：2008年日本電信電話株式会社入社、NTTメディアインテリジェンス研究所にて收音処理に関する研究開発に従事。遠くの音をクリアに收音するズームアップマイク、マイク技術の商用化貢献などの成果を挙げる。その後、2017年～2018年のVictoria University of Wellington留学を機に、NTTコミュニケーション科学基礎研究所にて分散最適化など機械学習の研究を開始。非同期分散型の深層学習技術の研究に尽力。現在、NTTメディアインテリジェンス研究所兼務。



NTTの基礎研究— 非同期分散型の深層学習技術とは

◆研究されている内容を教えてください。

機械学習による画像認識、音声認識などのモデルを構築する際にはよく「深層学習」(ディープラーニング)が利用されます。

現在の深層学習では、1つの巨大なデータセンタを設置してすべてのデータをそこに集約し、そのデータを使用してモデルを学習させる言わば「一極集中型」が一般的です。しかし、今後、自動運転、ファクトリーオートメーション、分散電源、個人ごとにオートチューン化されたモデルなどが普及することを考えた場合には、取り扱うデータの量が増加するため、1カ所にすべてのデータを集めて処理し展開することは困難となっていくことでしょう。また、GDPR(EU一般データ保護規則)などの影響もあり、プライバシー保護の観点からもデータの集約は難しくなります。これらの要因から、近い将来にはデータの蓄積や推論処理は分散化されるのではないかと考えています。

そこで、各サーバで分散蓄積されたデータを使って高度な知を創造するための研究をしており、例えば、あたかもデータを1カ所に集めて学習を行ったかのような機械学習モデルを得ることができる技術を最近成果として挙げました。

◆具体的にはどのような方法がとられているのでしょうか。

現在の分散深層学習では、各サーバ間でモデルを交換し、平均化するという手段(平均化合意形成)がよく用いられています。

平均化合意形成は非常に簡単な演算で、かつ効果のあるものです。ただし、各サーバが統計的に同じようなデータを持っている場合にはこの方式でもうまくいきますが、各サーバのデータの統計的な偏りが大きい場合には学習がうまく進まないことが多いです。また、ネットワークを構成するサーバの数が多くなると同期して通信を行うことも困難となります。

そこで、ネットワークに接続された各サーバが非同期的に通信を行い、協調してモデルを学習する深層学習のアルゴリズムを構築しました。詳しい数式やアルゴリズムは割愛しますが、簡単に言うと、この「協調する」という概念を数学的に表現する研究、と思っていただければと思います。

例えば、仲の良い人たちが集まって相談し、ある結論を導くことは簡単でしょう。しかし、仲の悪い、個性の強い人たちが集まっても皆がバラバラな方向に向かってしまいます。こうした場合に「平均化」はあまり大きな意味を持ちません。今回のアルゴリズムは、各自を協調させるような1つに纏まる力がうまく表現されていると考えてください。

図1は、CIFAR-10と呼ばれる画像分類器のテストで標準的に用いられるデータセットを用いて、8カ所にデータを統計的な偏りを持たせて分配したときの、提案方式と従来方式との比較を示しています。縦軸は分類誤差率を表していて、値が小さくなるほど性能が良いと言えます。青い点線はデータを1カ所に集めて学習した「グローバルモデル」の性能を示し、緑色の実線は従来方式(Gossip法)、赤色の実線は今回提案した新方式の性能を示しています。

従来方式では途中で性能向上が止まってしまっている一方、新方式では学習が進むにつれてグローバルモデルの性能に近づいて

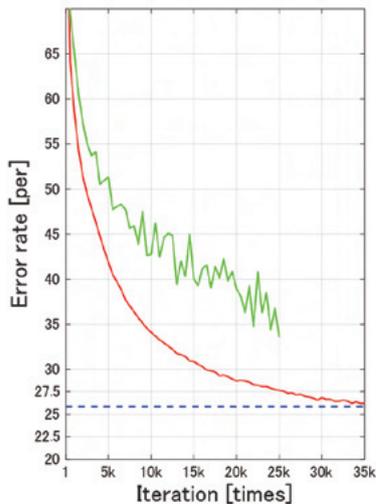


図1 従来方式およびグローバルモデルとの比較

ており、非同期通信下であっても、よりデータ全体に適合するモデルが得られていることが読み取れると思います。

◆この研究により、どのようなことが可能となるのでしょうか。

図2は複数拠点の医療画像データを用いて画像解析モデルの学習を行うデモンストレーションです。

医療データはプライバシー保護の観点において第一に挙がるもので、基本的に病院外に持ち出すことはできません。ましてや国外に持ち出すことはほぼ不可能です。そこで、N1~N8の8つの病院をネットワークでつなぎ、各病院から医療画像データを出さずにモデルを学習させることを考えます。具体的には胸部レントゲン写真から疾病の有無および約14種類の疾病を判断する医療画像診断補助のモデルです。現実社会において病院により取り扱う疾病に違いがあったり、地域性などにより状況が異なったりすることを加味し、N1~N8の持つデータの件数および収録疾病について人工的に偏った状態をつくり出しています。

N1~N8それぞれの青色のグラフは新方式、橙色のグラフは従来方式の学習の進捗を示しています。8カ所に分散した医療画像データがうまく協調し、「スーパードクター」のような高度な知の創造がなされていることが分かります。

非同期分散型の深層学習技術により、プライバシーを犠牲にすることなく、機械学習の恩恵を受けることができるわけです。こ

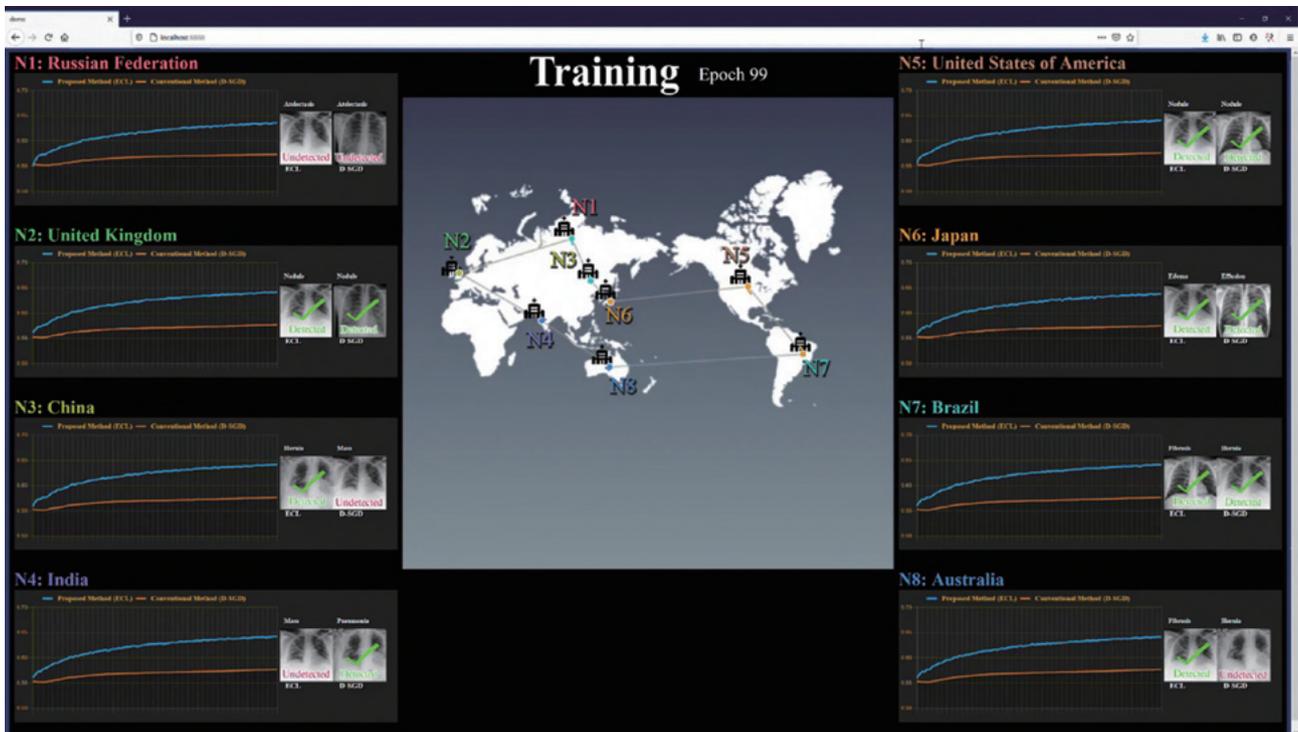


図2 複数拠点の医療画像解析モデル学習の融合



(今回はリモートにてインタビューを実施しました)

の技術は、特にアプリケーションに依存した技術ではありませんので、他にも、スマートフォンのデータを使ったテキスト変換、自動運転モデルの作成、コールセンタの音声認識、工場機器の故障検知など、非常に応用範囲の広いアルゴリズムだと思います。

NTTの地の利を活かしたサービスの提供を

◆NTTの強みはどのような点にあるとお考えでしょうか。

NTTは全国各地に局舎を所有しており、それぞれの局舎はネットワーク構造により接続されています。そうした地の利を活かし、各局舎にサーバを設置し、データの蓄積・処理してみるのはいかがでしょうか。当然、各局舎には全く違うデータが蓄積され、全体としてはとても巨大なシステムとなるでしょう。

そうしたデータを活用して、よりユーザの近くでデータ処理などのサービスを提供することは今後のNTTのビジネスとして有力な選択肢ではないかと思えます。いわゆる「エッジコンピューティング」の考え方ですね。

また、そうしたことを考えたとき、非同期分散型の機械学習のサービスを提供していくことは有用なのでは、と考えています。

◆今後の方向性についてはどのようにお考えでしょうか。

もともと私は音響、特にマイクロホンやスピーカーを使ったコミュニケーションに関連した研究をしていました。その後、入社から10年ほど経ち、ニュージーランドの大学に留学したことをきっかけとして「分散」というものをテーマに機械学習、特に最適化と呼ばれる分野の研究を始めました。

現在のシステムの多くは一極集中型です。データが1カ所でさばききれぬ量であること、求められているものも万人共通ととらえられ、同じサービスが分配されていることなどに起因しているのでしょう。しかし、今後はサービスの個人化が進むのではないかと考えています。個人に合わせたモデルの学習・推論の提供なども登場するのではないのでしょうか。そうした世の中になったとき、「このまで一極集中型でいいのか」という疑問があります。分散化すると考えたほうが自然なのではないのでしょうか。

考えてみれば、私たちの脳も「分散型」です。何かを話した

がらも、全然別のことを考えていたりする。私の脳が焼き切れていないところを見ると、現在の深層学習・推論のような重い計算をしているとは思えません(笑)。頭の中では何か分散した計算処理群が非同期的につながっていて、非常に軽い演算の膨大な組み合わせで高度なことを成し遂げているという気がしています。そう考えると、次世代のコミュニケーション、社会全体を支えるシステムというものも、高度な処理ができ、低消費電力で、柔軟性があり、少しくらい傷がついても簡単には壊れないものになるのではないかと考えています。

例えばIOWN(Innovative Optical and Wireless Network)構想では社会全体を最適化しようという文脈で、街における自動運転車の整列なども議論されています。それぞれの車は当然行先も違いますし、運転の仕方も違うべきだと思います。しかし、自分の利益だけを考えて動いてもうまくはいきませんし、また平均化も好ましくない気がします。

それぞれの車を協調させて取りまとめ、制御するようなシステムはまだないと思いますが、「分散」をテーマにしながら、そうしたものに寄与するシステム、ソフトウェアの根幹みたいなものをつくっていきたいと思っています。

◆これから基礎研究に取り組みたいと思っている方へメッセージがあればお願いします。

機械学習という分野は競争過多で、数カ月で状況が変わります。毎日とは言いませんがそれに近い頻度でいろいろな論文が投稿されています。では分散型のシステムに関しては競争過多かという点、分散学習に関しては競争過多ですが、私の考えている「柔軟に高度な知を与えられるような非同期型の分散システムを考える」ということ自体はそうではないように思います。「一極集中型」「End-to-End」という流れが主流な中で、あまり開拓されていない領域であり、そこに可能性を見出したいと思っています。

新しい分野を創造するにはエネルギーが必要です。それなりに知名度のある人が「こうだ」と言えば簡単ですが、そうでない人が99.9%を占める場合には、共通の課題を見つけ、議論したり参照し合って比較したりしながら競争することが非常に重要となります。

一方で、自分のポリシーのある独自の部分をつくり上げていくことも重要です。本当に自分が表現したいもの、つくりたいものを少しずつ表現していくことが大事なのではないかと思えます。

現在、データを集めることができれば、誰にでも(高校生でも)任意のアプリケーションのモデルをつくれるような時代です。そうした状況で、基礎研究と実用開発のどちらに軸足を置くべきか、私の周りの研究者たちも悩んでいる人は多いです。私は、たまたま留学させていただく機会があり、そこで分散に関する高度な知識を習得したので、基礎研究に軸足を置くという道を選びました。もちろん、その逆を選択するという道もあったでしょう。どちらのポジションにつけば良いかというのは、人によって向き不向きもあるので結論はありませんが、どちらか極端なほうに賭けたほうが良い時代なのかな、という気がしています。