

AIツールとAI人材の育成

AI Singaporeは、Adobe PDF Extract APIを使用して深層学習モデルを加速しています。



創業
2017年
従業員数 90名
Singapore
<https://aisingapore.org>

40%
機械学習（ML）モデルを提供
するために必要なスプリント
の削減

導入製品
[Adobe Acrobat Services](#) >
[Adobe PDF Extract API](#) >

■ 課題

- 自然言語のパイプライン処理によって取り込むデータの品質向上
- 企業の持続可能性報告書（サステナビリティレポート）に必要なコンテンツの提供
- 10スプリントかつ7か月以内に、最小限の実用的な機械学習モデル開発

■ 成果

- 異なる500のドキュメントのコンセプトを証明するため、**高精度**でPDFからデータを抽出
- PDF Extract APIを**わずか2週間**で実装
- より優れたモデル結果を得るために、**より良いコンテキスト**と構造を持つデータを提供
- 計画よりも**40%早く**展開可能なモデルを提供

シンガポールでは、ビジネスの課題解決のために人工知能（AI）を活用したいと考える企業は、AI Singaporeに相談をします。AIを活用して将来のシンガポールのデジタル経済を強化するために立ち上げられたこの国家プログラムは、シンガポールに拠点を置くすべての研究機関と、AI製品を開発するスタートアップや企業のエコシステムを結集し、AIへの取り組みを強化するために、実利用を想定した研究を行い、知識を深め、ツールを開発し、AI取り組みを支える人材を育成しています。

Siavash Sakhaviは、シンガポールのAIイノベーションを促進するために設立された「100 Experimentsプログラム」の副代表です。この組織は、既成のAIソリューションが存在しない領域において問題を提起し、シンガポールの研究者のエコシステムとAISGのエンジニアリングチームによって、9～18ヶ月以内に解決することを目指しています。

Sakhaviのチームは、大規模な多国籍金融サービスクライアントから持続可能性報告書プロジェクトに関する依頼を受けました。クライアントは、さまざまな情報源からPDF形式で提供される多様なレポートやパンフレットからテキストを抽出するのに苦労していました。プロジェクトチームは複数のAI・データ・プラットフォームエンジニア、およびAISGの実習生からなり、抽出された情報を自然言語のパイプライン処理に入力したいと考えていましたが、当時使用していたPDF抽出ツールによって大量の非構造化の意味不明なテキストが返されており、自然言語のパイプライン処理が期待通りに機能していないことに気付きました。

通常のプロジェクトでは、チームは1～2か月でモデルを開発することができます。しかし、このときはすでに10スプリントのうちの6スプリント目に入っています、結果を出すためのプレッシャーが高まっています。幸いなことに、彼らは当時ベータ版から一般提供に移行していたAdobe PDF Extract APIについて知りました。アドビの新しいWebサービスであるPDF Extract APIは、ネイティブおよびスキャンされたPDFファイルからデータとコンテンツを解析し、構造化されたJSONファイル内のテキスト、表、および画像要素を抽出できます。

「Adobe PDF Extract APIはまさに救世主でした。
これがなければ、さまざまなデータを使った自然言語処理ソリューションを
計画通りに構築することは難しかったでしょう。」

Siavash Sakhavi
Assistant Head, 100E, AI Singapore

より良い構造とコンテンツ抽出の実現

「1回目のデモで、私たちはそれまで使っていたPDF抽出ツールからはAdobe PDF Extract APIに完全に切り替えることを決めました。」とSakhaviは述べています。スプリントの終わりまでに、チームは大きな進歩を遂げていきました。自然言語のパイプライン処理への迅速なデータ取り込みを実現し、最終的にはプロジェクトのスポンサーに約束された作業を予定よりも早く提供することができました。

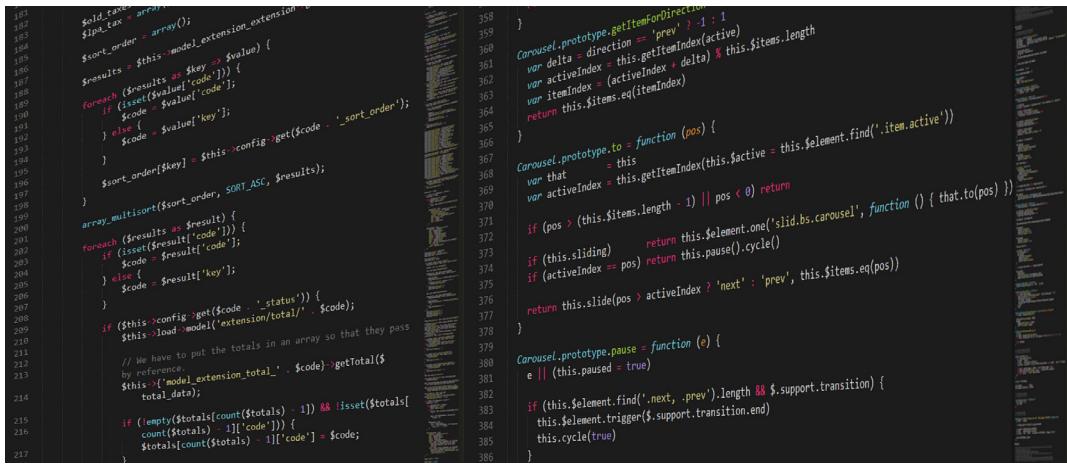
「Adobe PDF Extract APIはまさに救世主でした。これがなければ、さまざまなデータを使った自然言語処理ソリューションを計画通りに構築することは難しかったでしょう。」とSakhaviは述べています。「それまで、機械学習モデルがどのように機能すべきかを構築し、改良してきました。しかし、PDF抽出ツールでの抽出結果がボトルネックになっていたのです。」

チームがテストしていたオープンソースのPDF抽出ツールは、段落を正確に識別することができませんでした。大部分のテキストは文章の途中で途切れているなど活用が出来ず、また図表の数字やラベルが本文として誤って抽出されるなど、文章の構造情報が完全に欠落していました。

「これらのテキストから必要な情報を手動で抽出し、最終的にそれらを分類するのは簡単な仕事ではありません。自動化と効率性を提供するより良いPDF抽出ツールを見つける必要がありました。幸いなことに、私たちは自然言

語処理量を急速に増やし始めるタイミングでPDF Extract APIを組み込むことができました。」Sakhaviは述べています。

PDF Extract APIは、単なる文や断片ではなく、段落に基づいて出力できます。「段落をグループ化するコンテキストデータと機能は非常に貴重でした。これにより、自然言語のパイプライン処理が大幅に強化され、機械学習アルゴリズムでもより良い結果が得られるようになりました。」とSakhaviは述べています。



```
184     $old_text = $new_text;
185     $old_text = array();
186     $sort_order = array();
187     $results = $this->model_extension->extension->get('results');
188     foreach ($results as $key => $value) {
189         if (isset($value['code'])) {
190             $code = $value['code'];
191         } else {
192             $code = $value['key'];
193         }
194         $sort_order[$key] = $this->config->get($code . '_sort_order');
195     }
196     array_multisort($sort_order, SORT_ASC, $results);
197     foreach ($results as $result) {
198         if (isset($result['code'])) {
199             $code = $result['code'];
200         } else {
201             $code = $result['key'];
202         }
203         $code = $result['code'];
204         if ($this->config->get($code . '_status')) {
205             if ($this->config->get($code . '_extension_total')) {
206                 $this->load_model('extension/total/' . $code);
207             }
208             // We have to put the totals in an array so that they pass
209             // by reference.
210             $this->['model_extension_total'] . $code)->getTotal($
211             'total_data');
212             if (!empty($totals[count($totals) - 1])) && !isset($totals[
213             count($totals) - 1]['code']) {
214                 $totals[count($totals) - 1]['code'] = $code;
215             }
216         }
217     }
358     Carousel.prototype.getItemOrDirection = function (pos) {
359         var delta = direction == 'prev' ? -1 : 1;
360         var activeIndex = this.getItemIndex(active);
361         var itemIndex = (activeIndex + delta) % this.$items.length;
362         return this.$items.eq(itemIndex);
363     }
364     Carousel.prototype.to = function (pos) {
365         var that = this;
366         var activeIndex = this.getItemIndex(this.active = this.$element.find('.item.active'));
367         if (pos > (this.$items.length - 1) || pos < 0) return;
368         if (this.sliding) return this.$element.one('slide.bs.carousel', function () { that.to(pos) });
369         if (activeIndex == pos) return this.pause().cycle();
370         return this.slide(pos > activeIndex ? 'next' : 'prev', this.$items.eq(pos));
371     }
372     Carousel.prototype.pause = function (e) {
373         e || (this.paused = true);
374         if (this.$element.find('.next, .prev').length && $.support.transition) {
375             if (this.$element.trigger($.support.transition.end));
376             this.cycle(true);
377         }
378     }
379     Carousel.prototype.cycle = function (e) {
380         e || (this.paused = false);
381         if (this.$element.trigger($.support.transition.end));
382         this.$element.removeClass('carousel-fade').addClass('carousel-slide');
383         this.$element.trigger($.support.transition.end);
384         this.$element.removeClass('carousel-slide').addClass('carousel-fade');
385     }
386 }
```

上位のモデル結果のための関連性向上

Sakhaviは次のように述べています。「私たちの結果は驚くべきものでした。PDF Extract APIは非常にうまく組み込まれました。私たちは、必要なものだけを抽出するためのアダプタを作成しました。APIは驚くほど正確であり、プロジェクトスケジュールを加速させました。」

AISGチームは、Bidirectional Encoder Representations from Transformers (BERT) の深層学習モデルを開発することを計画していました。このモデルにデータを供給するために、プロジェクトのスポンサーは、特定の企業ESG（環境・社会・ガバナンス）イニシアチブに重要な要因を示すキーワード用語集の定義を提供しました。AISGチームの目標は、これらの定義と分析対象ドキュメントの間で類似性マッチングプロセスを実行し、テキストのどの部分がプロジェクトに関連するかを特定することでした。

「これらは環境持続可能性に固有のコンテキストだったので、用語集の情報と一致する関連テキストを抽出するまでモデルをトレーニングすることができませんでしたが、PDF Extract APIは環境持続可能性に関連するトピックに関する段落の文脈で重要な文レベルの情報を正しく識別することができ、高品質なデータの取り込みを実現しました。」とSakhaviは述べています。

「私たちの結果は驚くべきものでした。

PDF Extract APIは非常にうまく組み込まれました。APIは驚くほど正確であり、プロジェクトスケジュールを加速させました。」

Siavash Sakhavi

Assistant Head, 100E, AI Singapore

企業のスポンサーは、様々なソースからの500もの多様な文書を分析したがっており、Sakhaviのチームに最初の10文書を処理し、コンセプトを証明した後、残りの490文書を処理できる展開可能なモデルを提供するよう求めていました。「関係するレポートやパンフレットは、どれも簡単ではありませんでした。テキスト要素の他に、ページ

の各所に画像がたくさんあり、サイズや内容も様々でした」と彼は述べています。Sakhaviと彼のチームは、たった2週間でPDF Extract APIを実装し、データ品質の大幅な向上を実現しました。「学習曲線は非常に短かったです」と彼は述べています。チームはすぐに、データのリストをパイプラインに組み込み、パス属性に基づいてすべてを解析して、望ましい結果を得る方法を学びました。

「スプリント6の終わりには、私たちは有望な結果をクライアントに提示することができ、チームのさらなる作業によって本番環境のソリューションを実現することができました。」と、Sakhaviは述べています。

Innovative Singaporean AI ecosystemの育成

プロジェクトの範囲はその後、新しい領域の研究調査に拡張され、AISGのAIエンジニアリングチームと連携して作業しているシンガポールの大学に引き継がれました。企業のスポンサーにとっての良いニュースは、AISGプロジェクトチームで働いていた同じ人々が、引き続きこのイニシアチブを実現させるために大学研究チームに参加し、最終的にはラベル付きデータを使用して新しいモデルをトレーニングし、優れた結果を出していることです。

Sakhaviは次のように述べています。「AISGにはPDF解析が必要なプロジェクトが多数あります。私からすべてのチームへの推奨事項は、PDF Extract APIを活用することであり、スポンサー企業にも継続的な使用を提案することです。」



© 2024 Adobe. All rights reserved.

Adobe and the Adobe logo are either registered trademarks of Adobe in the United States and/or other countries. All other trademarks are the property of their respective owners.