

技術ワードから読み解く ChatGPTのメカニズム

横山 裕樹

2022年11月にOpenAIがChatGPTを公開してから1年が経とうとしています。本稿の執筆時点でも、本体のバージョンが上がったりプラグインが増えたりしていますし、利用者の間でも、活用法やプロンプト・エンジニアリングなどでさまざまな提案がされています。

ChatGPTはまるで人間かのように答えを返してくるため、内部でどのような計算が行われているのか、意識することは少ないのではないのでしょうか。しかし、いろいろと試してみると、ChatGPTには特有のクセのようなものもあることが分かってきます。ChatGPTがどのような構造を持っていて、どのように学習しているのかを理解しておくことは、これを利用する際の手助けになると考えられます。

ChatGPTの詳しい仕様は公表されていませんが、元になっているGPT (Generative Pretrained Transformer) やInstructGPTについてはソースコード⁽²⁾や論文⁽³⁾が公開されています。ここでは、それらに基づいてChatGPTの仕組みに迫りたいと思います。

理解へのステップ1…次の単語を選ぶ

● 次の単語は確率で選ばれる

GPTがどのように文章を生成しているのかを解説します。GPTは書きかけの文章にトークンを1つずつ追加していくことで文章を生成します。トークンについては後ほど詳しく解説しますが、ここでは単語と考えてください。

図1の例では、GPTに“you”というトークンを入力しています。GPTは入力されたトークンに対して、その次にどのようなトークンが出現しそうか、その度合いを連続値で出力します(図の棒グラフ)。“you”が英文の主語だとすると、その後には“are”などの動詞が来そうです。こういった次に来そうなトークンについては大きい値を出力します。候補となるトークンは非常に多いのですが、図では模式的に5個だけ表示しています。

次に、これらの数値に後処理を行った後、softmaxという関数を用いて確率に変換します(図の円グラフ

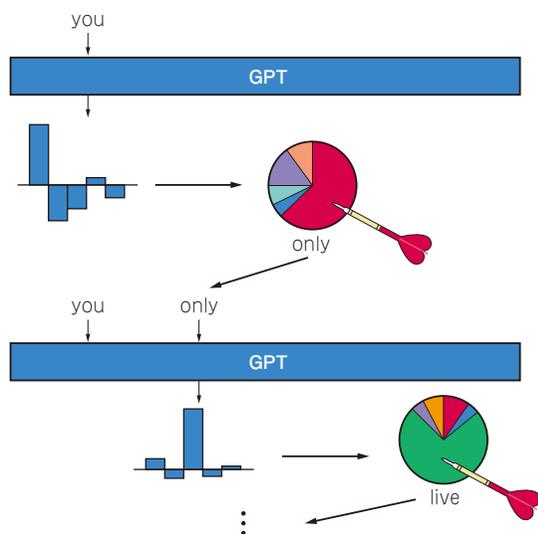


図1 GPTによる文章の生成

入力されたトークンに対して、次にどのようなトークンが出現しそうか、その度合いを連続値で出力する

フ)。そして、この確率に基づいて次のトークンを選択します。図の例では“only”が選択されています。このようにしてGPTによって1つのトークンが生成されます。

今度は生成したトークンを末尾に加えたトークン列“you/only”を入力すると、同じように“only”の次に来るトークンの確率が得られます。このプロセスを何度も繰り返すことで、長い文章を生成できます。

初めの入力“you”の1トークンだけなので、さまざまなトークンが次に出現する候補として挙げられるでしょう。一方で、“you/only”のように、2つのトークンの組み合わせが入力されると、次に出現し得るトークンは絞り込まれてきます。GPTは長いトークン列から必要な情報を取り出して、次にどのようなトークンを生成するべきかを判断できます。ChatGPTに同じプロンプト(コマンド)を繰り返し与えると、その都度異なる回答が得られることがありますが、それは上記のような確率的な生成過程によるためです。