

ChatGPT と KM の技術/機能比較

—ChatGPT 解体新書—

■ ChatGPT とは、
GPT-3 を Alignment された InstructGPT を
対話型に特化したシステムです。〈右図上〉

■ GPT-3 とは、
・ Transformer の Decoder のみで構成された言語モデルで、これを N=96 層の深層にすることで GPT-3 になります。〈右図中〉

・ MaskedSelf-Attention とは、ひとつの単語をマスキングして、その前の文章からマスキングされた単語を学習するモジュールのこと。例えば、

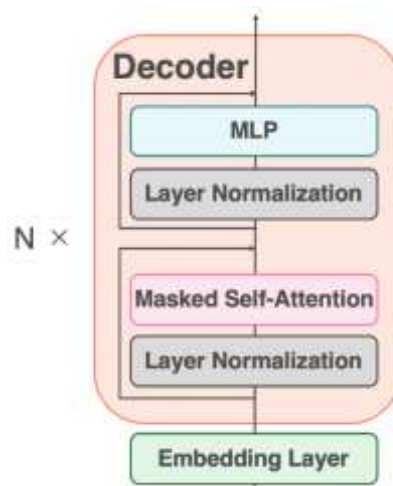
「日本の首都は東京です。」という文ならば、「東京」をマスキングして、入力が「日本の首都は」なら「東京」が出力となり、「日本の首都は東京」なら「です。」が出力になります。(意味的汎化性は無い)

〈右図下〉

・ MLP(Multi-LayerPerceptron)とは、多層パーセプトロンのことで、上記の N=96 層の深層にすることで、ディープラーニングのことです。

・ 特筆すべき点は、1,750 億個のパラメータで、570GB の文章コーパスを学習したことです。このコーパスは、主に CommonCrawl と呼ばれるネット上の文章を 45TB 以上集めて、これから学習に不適切な汚い文章を取り除いたのが 570GB でした。すなわち、Wikipedia などから採集された文章では、1%程度の文章しか学習に使えないことが判ります。(だから新聞記事は大変有効な学習データです。)

・ しかし、GPT-3 の出力は、不正確で非道徳的な出力をすることで問題になり、その問題の対処法として InstructGPT が誕生しました。



InstructGPT とは、人間のフィードバックを基にモデルを学習させたものです。(教師有り学習)

・これは RLHF(Reinforcement Learning from Human Feedback)という強化学習を使った教師あり学習データで、精度を上げようとしているものです。

・ InstructGPT の流れは、

- ①教師ありファインチューニング
- ②RewardModel の獲得
- ③RLHF

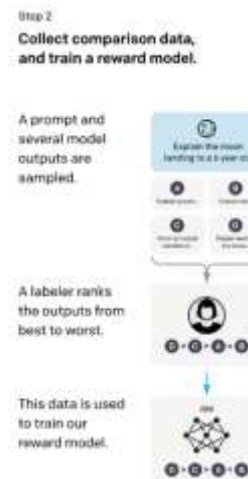
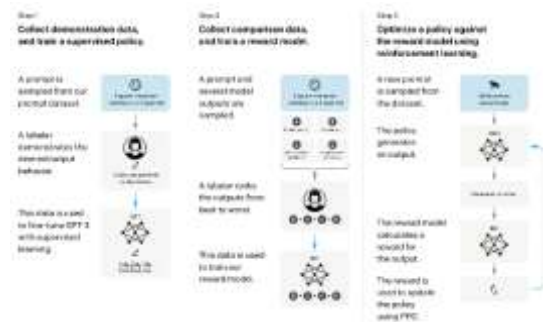
となっており、②～③は、ループになっています。

①教師ありファインチューニング(Supervised Fine-Tuning:SFT)とは、GPT-3 の出力データを元に、人間が入力プロンプトと出力文にラベルを付けて学習させることです。13,000 文程度の教師データでファインチューニングをしており、これを **SFT モデル**と呼んでいる。

②RewardModel とは、出力文の良し悪しを評価をするものです。評価はスカラー値 $r_{\theta}(x,y)$ で出力します。(θ は重みです)

- ・評価軸は、3 軸あります。
 - Truthfulness (真実性) :
デマやミスリードの情報ではないか
 - Harmlessness (無害性) :
人や環境の物理的・精神的に傷つけていないか
 - Helpfulness (有益性) :
ユーザのタスクを解決しているか

・学習パラメータ数は 60 億あり、いろいろなタスクからファインチューニングされた学習データを使い、文の良し悪しに基づいたランキングを学習させることで、曖昧性と矛盾性を解消している。

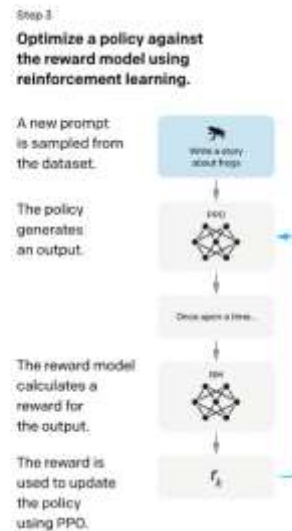


③Reinforcement Learning from Human Feedback:RLHF とは、
強化学習を用いて SFT モデルの精度を上げるものです。

・仕組みは簡単で、前ステップで獲得した RewardModel を最大化するように SFT モデルの方策 Policy を学習させるものです。(文章塊生成)

・ OpenAI の独自技術である PPO(Proximal Policy Optimization)によって、ポリシーの大きな更新を抑えながら最適化していく手法を使っている。

・これだけでは Policy が RewardModel をハックするような文を出してしまうという問題が生じるので、KL 正規化項を追加している。KL ダイバージェンスでは、31,000 文のプロンプトが使用され、対数尤度と共に期待値の制御をしている。(無意味な技巧的手法)



ChatGPT とは、
上記説明で判った通り、InstructGPT とほとんど同じものです。

・ ChatGPT と InstructGPT の違いは、
①モデル：GPT-3.5
②データ：会話データ
だけです。



・ GPT-3.5 とは、テキストだけでなく、コードも学習したものです。
コードを会話で検索/生成できることは、100 万エンジニア達を驚かせました。
(実際、使ってみると JavaScript 以外は精度が悪い)

・ 会話データは、
①会話データによる SFT モデルの学習
②SFT モデルの出力を人がランク付けし RewardModel (報酬) を学習
③RewardModel を最大化するように SFT モデルを PPO でファインチューニング
以上の 3 つのステップで学習されている。 以上□

テキストの意味的機能性項目

□テキストの意味的機能性項目は、

- ①上位概念/上位概念語が”付与”されているか
- ②属性 Attribute が”紐づけ”されているか
- ③共起性が”関連付け”されているか
- ④汎化性のある知識が”生成”されているか
- ⑤知識の構造化が”構築”されているか
- ⑥知識構造から新知識の”推論”ができていますか

の6つの機能になります。

①上位概念とは、

- ・各単語や知識の上位の概念クラスであり、上位概念語はその上位概念クラスに存在する単語や知識である。包含関係のこと。

例：上位概念（ポチ）＝動物、ペット 上概念語（ポチ）＝犬

（解説）

- ・上位概念の抽出は、意味解析で重要なファクターです。
- ・未知語「ポチ」が犬と判れば、年齢や尾があるという継承 Inheritance が付加できる。

②属性とは、

- ・上位/下位概念には該当せず、付帯されるもの。

例：属性（ラーメン）＝メンマ、チャーシュー、🍣

（解説）

- ・属性も上位概念と同様に意味解析では重要なファクターです。
- ・メンマの上位概念が食物だけでなく、ラーメンの概念を構成する役割 Role になる。

③共起性とは、

- ・短縮語や省略語並びに別名同義語などを指す。

例：「博士」→（ドクター、ドク、Dr、Ph.D.、最高学位、はかせ、…）

（解説）

- ・専門分野に於ける共起語の特定化は、文書の意味解析精度を左右するファクターです。
- ・博士と Dr が違う単語として処理されれば意味理解の精度劣化につながる。

④知識生成とは、

- ・平文から必要な知識を生成すること。

例：科学文書→「地球は太陽の周りを傾斜自転しながら公転している。」

(解説)

- ・文書の要点やトピックなどを抽出して、汎化性のある知識にまとめておけば、周辺の文書や知識との相関関係や因果関係に使えて、新発見などの可能性につながる。

⑤知識構造化とは、

- ・知識の包含関係や属性を基に「体系化」すること。(Ontology)

例：学校⊃生徒（包含関係）

生徒⇒年齢、性別、成績、身長、体重、…（属性）

(解説)

- ・知識の構造化は、知識の性質や成分など上位知識からの継承 Inheritance が明示化
- ・例：「生徒」から（学年、クラス、性別、成績、…）などの属性を継承項目から生成

⑥新知識推論とは、

- ・知識構造化を基に組合せ等で新知識を推論すること。(新知識発見)

例：ガーリックトースト⊕ラーメン＝「とろとろピリ辛ラーメン」🍜

(解説)

- ・知識構造の近傍や性質、成分などから組合せで合成ができる
- ・上記例「地球は太陽の周りを傾斜自転しながら公転している。」の知識があれば、「日出/日没」や「四季」などの因果関係の知識の発見につながる。

ChatGPT と KM の技術/機能比較

意味的機能性項目//商品名	ChatGPT	KnowledgeMiner
①上位概念/上位概念語の抽出と付与	✖Transformer/Decoder 内の MaskedSelf-Attention では、次の語彙や文の予測だけなので、上位概念や上位概念語の抽出は不可	◎自律学習で概念構造スケルトンが自動更新され、上位概念が概念タグ CT として自動付与される
②属性 Attribute の紐づけ	▲InstructGPT の RLHF では、関連属性の紐づけはできるが、上位概念の概念が無いので、汎化性が無い	◎自律学習の概念構造で関連属性は紐づけされ、概念タグ CT で自動付与される
③共起性の関連付け	▲MaskedSelf-Attention では、語彙や文の予測で共起性は期待できるが、やはり上位概念の概念が無いので、汎化性が期待できない	◎自律学習の照応共起テーブルから照応解析（代名詞、ゼロ代名詞、共起語）が自動で付与される
④知識生成の汎化性	▲SFT では、人間によるファインチューニングで政策 Policy として知識らしい文章が生成できる	◎KE という知識式への代入で、明示的に知識が自動生成される
⑤知識構造化の構築	✖RLHF では、Policy で文章生成はできるが、知識構造化はできない	◎知識式 KE の包含関係と属性を概念タグ CT と意味タグ ST で構造化できる
⑥新知識推論の仮説	✖InstructGPT では、汎化性のある知識生成も構造化もできないので、知識の組合せに依る新知識発見はできない	◎知識式 KE とその構造上の近傍の性質や成分などによって、組合せができて、新知識仮説/発見ができる
Comments	巷の SNS では、「おもちゃ」としては驚きだが、製品化としては障壁がまだまだ高いと云われている	汎用性のある解析エンジンの XML 出力は、タスクの開発工数の省力化が期待されている