

Generativní umělá inteligence – Díl první: příliš velká očekávání

EDUARD BARTL

Přírodovědecká fakulta UP, Olomouc

Umělá inteligence není v současné době námětem diskuzí pouze v odborné komunitě, jisté povědomí má o ní prakticky každý, kdo má přístup ke sdělovacím prostředkům, zejména k internetu. Vliv umělé inteligence na většinu oblastí lidského konání je skutečně nebývalý. Většina veřejnosti má velká očekávání. Někteří lidé vidí v umělé inteligenci lék na spoustu problémů tohoto světa, jiní se jí bojí a byli by nejraději, kdyby bylo její použití zakázáno. Série článků začínající tímto dílem se bude zabývat jednou z nejvíce diskutovaných podoblastí umělé inteligence, které se dnes říká generativní umělá inteligence a do které spadají zejména velké jazykové modely, jakým je například ChatGPT.

1. Cíl série a další informace na úvod

Hlavním cílem těchto článků je nahodit přílišný optimismus a současně rozehnat přílišný pesimismus, jenž panuje kolem generativní umělé inteligence. Myšlenka postupující celým seriálem je celkem prostá. Pokusíme se ponořit do historie vzniku generativní umělé inteligence a především do principů, na kterých je založena. Pochopení těchto principů by totiž mělo vést ke střízlivějšímu pohledu na věc a tedy k rozptýlení neopodstatněného očekávání.

Učitelům na základních a středních školách (zejména pak učitelům výpočetní techniky a informatiky) snad tato série článků napoví, jaké místo by mohla generativní umělá inteligence zaujmout ve vzdělávání. Jsem si vědom, že se jedná o cíl dosti ambiciózní, ale ze svojí pozice již delší dobu cítím, že bych se o něj měl alespoň pokusit.

Čtenáře musím na samotném začátku upozornit, že některé názory, které v průběhu série vyslovím, jsou založeny na neúplných informacích (vždy na to upozorním) a mohou tedy být nepřesné. Neúplnost informací je důsledkem toho, že OpenAI, Google a jiní velcí hráči na poli generativní umělé inteligence zdaleka nezveřejňují vše; OpenAI tím tak poněkud popírá svůj název. Navíc vývoj v oblasti generativní umělé inteligence chvátá

obrovským tempem kupředu, spousta věcí se teprve usazuje a dokonce i největší odborníci na problematiku nemají ve všem jasno a jejich pohled na věc se může významně lišit.¹⁾

Při psaní vycházím z nejrůznějších zdrojů. Počínaje svými zápisky z přednášek ze svých studentských let, přes odborné články vědců, kteří se zasloužili o současný rozmach umělé inteligence až po Wikipedii a jiné zdroje na internetu. Nezpochybnitelnou inspirací jsou pro mě také rozhovory a přednášky dvou významných českých odborníků na umělou inteligenci (a mimochodem velmi charismatických řečníků), dr. Jana Romportla a dr. Tomáše Mikolova. Některé z jejich rozhovorů nebo přednášek je možné najít na YouTube.

2. Kdy to vše začalo

Generativní umělá inteligence je schopna generovat smysluplný text, obrázky nebo data jiného typu (třeba video nebo audio) a je úzce spjata s takzvanými velkými jazykovými modely. Co přesně znamená *velký jazykový model* a jakým způsobem souvisí s pojmem generativní umělá inteligence se budeme postupně dozvídat později.

Generativní umělá inteligence vstoupila do povědomí široké veřejnosti náhle a nečekaně, a to v listopadu roku 2022. Tedy v okamžiku, kdy společnost OpenAI poskytla k volnému používání program ChatGPT, který byl založený na velkém jazykovém modelu GPT-3.5. Podle mého pozorování si spousta lidí spojuje podzim roku 2022 s okamžikem vzniku velkých jazykových modelů a za jejich tvůrce považuje právě společnost OpenAI. Skutečnost je ale taková, že se první jazykové modely objevily již roku 2018 a podle mého názoru větší přínos na vzniku velkých jazykových modelů v podobě, jak je známe dnes, měla firma Google; k tomu se dostaneme, až se začneme bavit o takzvaném self-attention mechanismu a transformerové architektuře.

Následující seznam uvádí názvy některých velkých jazykových modelů a jejich tvůrce. Seznam není zdaleka úplný, velkých jazykových modelů je výrazně více a přibývají takřka jako houby po dešti:

- říjen 2018 – BERT od Googlu,

¹⁾ Jako příklad uveďme nositele Turingovy ceny Yanna LeCuna a Geoffreyho E. Hintonu. První jmenovaný poukazuje na velmi omezené schopnosti umělé inteligence, druhý jmenovaný vidí věci opačně a předpovídá nástup takzvané *obecné umělé inteligence* (tedy inteligence, která – stručně řečeno – dosahuje kvalit inteligence lidské) již v následujících dvaceti letech.

- červen 2020 – GPT-3 od OpenAI,
- květen 2022 – LaMDA 2 od Googlu,
- březen 2022 – Chinchilla od DeepMind,
- listopad 2022 – GPT-3.5 od OpenAI,
- únor 2023 – LLaMA od Meta,
- prosinec 2023 – Gemini od Google,
- květen 2024 – GPT-4o,
- září 2024 – Gemini 1.5 Pro od Google.

Uvedený seznam napovídá, že podzim 2022, kdy byl zveřejněn Chat-GPT založený na modelu GPT 3.5, nepředstavoval žádný technologický zlom na poli velkých jazykových modelů. Šlo spíše o výborným způsobem zrealizovaný marketingový tah společnosti OpenAI. Za významný je možné spíše považovat rok 2018. Tomuto roku však předcházela dlouhý vývoj v oblasti umělé inteligence, který je možné vystopovat až do 40. let 20. století. Stručně lze tedy říci, že navzdory všeobecnému názoru nejsou velké jazykové modely novou technologií; staví na vědeckém bádání mnoha vědců (zejména matematiků a informatiků, ale také jazykovědců, neurovědců a filozofů) dokonce i z období předpočítačové éry. Historickému vývoji se budeme v následujících dílech také věnovat.

3. Jazykové modely

Začneme pozvolna; vysvětlíme si nejprve, co je to jazykový model.²⁾ *Jazykový model* je statistický model přirozeného jazyka, jakým je čeština, němčina nebo angličtina. Hodí se však také pro modelování umělých jazyků, jakými jsou například programovací jazyky.

Jazykové modely zjednodušeně řečeno fungují tak, že pro předem zadanou posloupnost slov (tedy pro určitý *kontext*³⁾) $w_0w_1 \cdots w_{t-1}$ umí stanovit podmíněnou pravděpodobnost

$$P(w_t | w_0w_1 \cdots w_{t-1}),$$

tedy pravděpodobnost, že bezprostředně po kontextu $w_0w_1 \cdots w_{t-1}$ následuje slovo w_t . Tuto pravděpodobnost umí stanovit pro libovolné slovo w_t daného přirozeného jazyka.

²⁾Pozorný čtenář si jistě všiml, že místo „velký jazykový model“ nyní píšou pouze „jazykový model“. Tento rozdíl bude objasněn vzápětí.

³⁾Kontext je posloupností slov

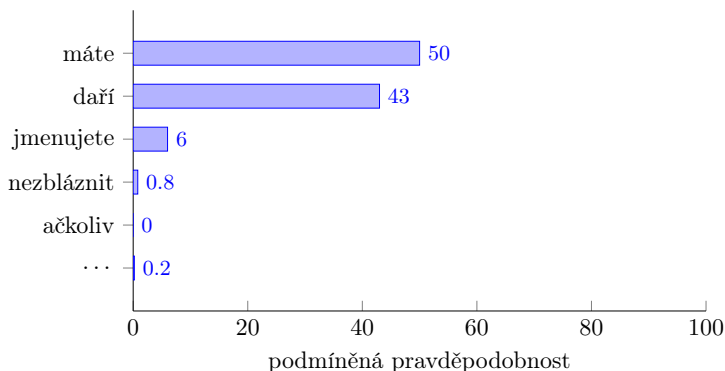
Uveďme si jednoduchý příklad. Bude-li kontextem posloupnost slov „Dobrý den, jak se“, každému z nás začnou automaticky naskakovat v hlavě různá slova, která by mohla být vhodným pokračováním této posloupnosti. Patrně to budou slova „máte“, „daří“ a podobně. I v naší hlavě totiž funguje jazykový model, který jsme si vytvořili každodenním kontaktem s česky psaným a mluveným textem. Jazykový model tedy vezme všechna česká slova a jistým způsobem jim přiřadí pravděpodobnost, že se tato slova vyskytnou po kontextu „Dobrý den, jak se“. Tyto pravděpodobnosti pak mohou vypadat třeba tak, jak je uvedeno na obr. 1 na konci článku.

Další činnost jazykového modelu pak vypadá tak, že jako pokračování kontextu je náhodně zvoleno jedno ze slov daného jazyka. Při tomto náhodném výběru jsou však zohledněny podmíněné pravděpodobnosti, o kterých byla řeč v předchozích odstavcích. Vrátime-li se k příkladu kontextu „Dobrý den, jak se“, tak můžeme říci, že s největší pravděpodobností bude jako pokračování vybráno slovo „máte“. Nikoliv však nutně! Zvoleno může být i slovo „nezbáznit“, pravděpodobnost, že se tak stane je ovšem dosti malá. S jistotou může říct pouze to, že nebude vybráno slovo „ačkoliv“ nebo jiné slovo, pro které vyšla podmíněná pravděpodobnost nulová. Z toho také plyne, při opakování tohoto procesu mohou být vybrána různá slova, jednou tedy dostaneme pokračování „Dobrý den, jak se máte“, jindy zas „Dobrý den, jak se daří“.

Velké jazykové modely, jako je třeba GPT, fungují na velmi podobném principu.⁴⁾ Tento princip je, jak jsme mohli vidět, velmi jednoduchý. Výrazně složitější je ovšem způsob, jakým velké jazykové modely stanovují podmíněné pravděpodobnosti. Slovo je v přirozeném jazyce velké množství a různých kontextů je ještě výrazně více. Velké jazykové modely tedy musí jistým způsobem vypočítat a poté uložit obrovské množství podmíněných pravděpodobností. Generování textu na základě těchto pravděpodobností je pak poměrně jednoduchou záležitostí.

O tom, jak velké jazykové modely počítají a ukládají podmíněné pravděpodobnosti se budeme bavit v následujících pokračováních série. Na závěr tohoto dílu pouze poznamenejme, že to nějakým způsobem souvisí s dalšími často diskutovanými termíny dnešní doby, konkrétně s *umělými neuronovými sítěmi* a *hlubokým učením*.

⁴⁾Pro úplnost dodejme, že jednou z významných odlišností oproti našemu výkladu je skutečnost, že velké jazykové modely nepracují se slovy, ale s *tokeny*, které si pro jednoduchost můžeme představit jako slabiky nebo interpunkční symboly. Abychom další výklad příliš nekomplikovali, budeme vždy pracovat na úrovni celých slov.



Obr. 1 Podmíněné pravděpodobnosti slov po kontextu „Dobry den, jak se“. Tři tečky na konci jsou uvedeny pro zjednodušení a nahrazují všechna ostatní slova; 0,2 je tedy součet pravděpodobností všech zbývajících slov, která nejsou v grafu uvedena.

* * * * *

POZVÁNKA

Komise pro vzdělávání učitelů matematiky a fyziky JČMF,
 Pobočný spolek JČMF Olomouc
 a Gymnázium Jevíčko

vás srdečně zvou na IV. ročník semináře

MATEMATIKA A FYZIKA VE ŠKOLE

Seminář se bude konat ve dnech 20.–22. srpna 2025 v prostorách Gymnázia v Jevíčku. Obsahem semináře jsou přednášky z pedagogiky, matematiky a fyziky. Vedle odborných přednášek bude věnován prostor na diskusi k otázkám současného vzdělávání v ČR. Seminář je určen učitelům matematiky středních a vysokých škol. Na seminář je možné se přihlásit na webových stránkách gymnázia: <https://www.gymjev.cz>.

Za organizační výbor Dag Hrubý