



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

تحلیل داده های متنی گزارش 1 هوش مصنوعی دکتر قطعی

امیررضا رادجو
دانشجوی علوم کامپیوتر دانشگاه صنعتی امیرکبیر

11 اسفند 99

تحلیل داده‌ی متنی

متن کاوی و داده کاوی

داده کاوی روشی بسیار کارا برای کشف اطلاعات از داده های ساخت یافته است. متن کاوی مشابه داده کاوی است، اما ابزارهای داده کاوی طراحی شده اند تا داده های ساخت یافته از پایگاه داده را به کار ببرند. میتوان گفت، متن کاوی راه حل بهتری برای شرکتها است. پس تفاوت متن کاوی و داده کاوی این است که داده کاوی بر روی داده های ساخت یافته پایگاه داده کار می کند و متن کاوی، بر روی داده های غیر ساخت یافته و نیم ساخت یافته مانند Email و مستندات تمام متنی کار می کند. در متن کاوی سعی می گردد از همان تکنیکهای داده کاوی استفاده گردد. برای این منظور به تکنولوژیهای دیگری مانند پردازش زبان طبیعی، یادگیری ماشین و ... نیاز است تا به صورت اتوماتیک آمارهایی را جمع آوری نموده و ساختار و معنای مناسبی از متن استخراج گردد. در این موارد، دیدگاه عمومی استخراج ویژگیهای کلیدی از متن است. ویژگیهای استخراج شده بعنوان داده برای تحلیل استفاده می گردد.

متن کاوی و بازیابی اطلاعات

معمولاً در بازیابی اطلاعات با توجه به نیاز مطرح شده از سوی کاربر، مرتبط ترین متون و مستندات و یا در واقع «کیسه کلمه» از میان دیگر مستندات یک مجموعه بیرون کشیده میشود. بازیابی اطلاعات یافتن دانش نیست بلکه تنها آن مستنداتی را که مرتبط تر به نیاز اطلاعاتی جستجوگر تشخیص داده به او تحویل میدهد. این روش در واقع هیچ دانش و حتی هیچ اطلاعاتی را به ارمغان نمی آورد. متن کاوی ربطی به جستجوی کلمات کلیدی در وب ندارد. این عمل در حوزه بازیابی اطلاعات گنجانده می شود. به عبارتی بازیابی اطلاعات جستجو، کاوش، طبقه بندی و فیلتر نمودن اطلاعاتی است که در حال حاضر شناخته شده اند و در متن قرار داده شده است. ولی در متن کاوی مجموعه ای از مستندات بررسی شده و اطلاعاتی که در هیچیک از مستندات به صورت مجرد یا صریح وجود ندارد، استخراج می گردد.

پردازش زبان طبیعی یا NLP

هدف کلی آن رسیدن به یک درک بهتر از زبان طبیعی توسط کامپیوترهاست. تکنیک های مستحکم و ساده ای را برای پردازش سریع متن به کار می برد. همچنین از تکنیکهای آنالیز زبان شناسی نیز برای پردازش متن استفاده می کند. نقش NLP در متن کاوی فراهم کردن یک سیستم در مرحله استخراج اطلاعات با داده های زبانی است.

در اینجا یک آشنایی مروری با مباحث موجود در پردازش زبان های طبیعی می پردازیم

به طور کلی می توان گفت که زبان از لحاظ مفهومی نسبت به تصویر دارای پیچیدگی های متعددی است که فهم و پردازش آن توسط هوش مصنوعی را سخت تر می کند. از جمله این ویژگی ها می توان به پیچیدگی های حافظه ای و شبکه های معنایی اشاره کرد که در پردازش تصویر معمولاً چنین پیچیدگی هایی یافت نمی شود

با ظهور یادگیری ژرف در فضای یادگیری ماشین، ایده های جدیدی در زمینه ی مسائل حوزه پردازش زبان های طبیعی نیز به وجود آمد که به مرور اجمالی آن ها می پردازیم

(word embeddings) واژه‌نهفت‌ها

پس از برخورد با داده‌های متنی برای استفاده از تکنیک‌های یادگیری ژرف ابتدا لازم است تا نمایشی از آن‌ها در فضای برداری ایجاد کنیم. به همین منظور در سال ۲۰۱۳ توسط آقای میکولوف و همکاران روشی برای نشان دادن واژگان در فضای چند بعدی ارائه کردند.

ایده‌ی اصلی این روش برای هر کلمه استفاده از کلماتی است که در نزدیکی آن قرار دارند. برای مثال در شعر زیر:

بلبلِ برکِ گلِ خوش‌رنگ در مزار داشت

واندرانِ برکِ و نواخوش ناله‌های زار داشت

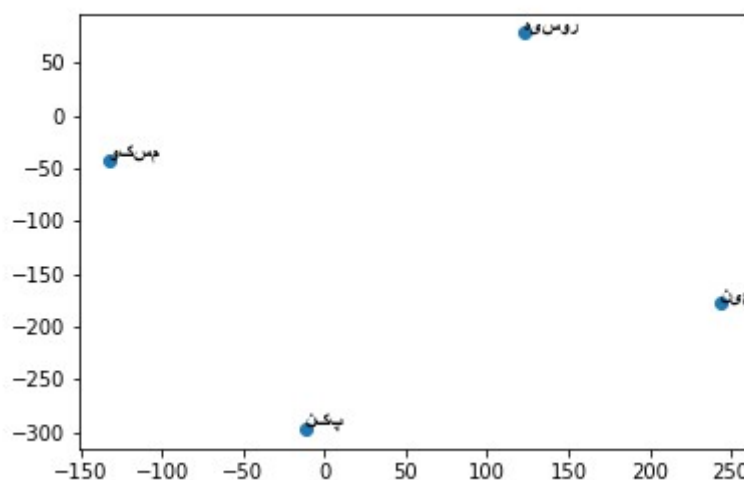
واژه «گل» در همسایگی واژگان «برگ» و «بلبل» و «خوش‌رنگ» آمده‌است. در این روش به نوعی سعی بر آن است تا با استفاده از فراوانی کلمات موجود در همسایگی یک کلمه شبکه‌ی معنایی آن را تخمین بزنند. این نیز معروف است (word2vec) روش به کلمه به بردار.

چند ویژگی جالب در این روش مشاهده شده‌است.

واژگانی که معنای نزدیک به یکدیگر دارند به جهت آن‌که در همسایگی آن‌ها واژگان مشترک با احتمال خوبی به تعداد زیادی یافت می‌شوند، نزدیک به یکدیگر قرار می‌گیرند.

ویژگی دیگر این است که اختلاف بردارهایی که ارتباط معنایی مشخصی با یکدیگر دارند تا حد قابل توجهی به یکدیگر شبیه هستند. برای مثال مشاهده شده که بردارهای کشورها به پایتخت آن‌ها تا حد قابل توجهی به یکدیگر شبیه اند.

برای مثال به شکل زیر توجه کنید



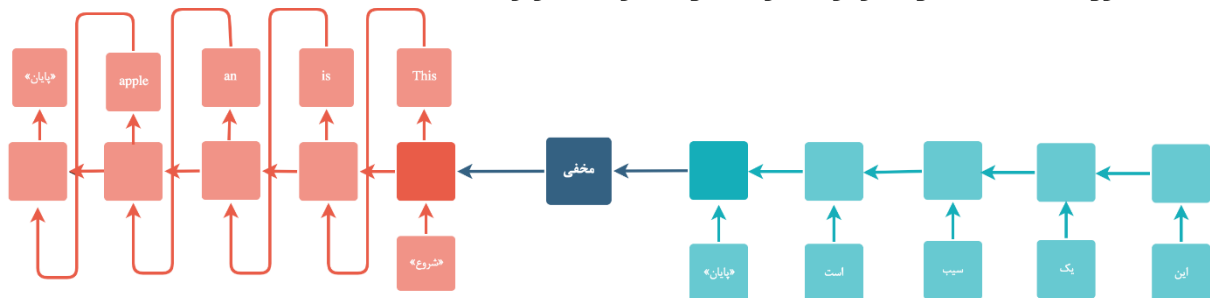
این ایده سنگ بنای بسیاری از ایده‌های آتی در زمینه یادگیری ماشین حوزه پردازش زبان‌های طبیعی است.

ترجمه ماشینی و دنباله به دنباله

ترجمه ماشینی تبدیل خودکار جملات از یک زبان به زبان دیگر است که تا پیش از به کارگیری یادگیری ژرف در زمینه پردازش زبان‌های طبیعی به صورت آماری انجام می‌گرفت. اما با به وجود آمدن روش‌های نوین یادگیری ماشین (Neural Machine Translation) و استفاده از واژه‌نهفت‌ها راه برای به وجود آمدن ایده‌ی ترجمه ماشینی عصبی هموار شد (Sequence to Sequence) و به طور کلی‌تر دنباله به دنباله (Sequence to Sequence Translation).

در روش دنباله به دنباله از دو شبکه‌ی عصبی بازگشتی (شبکه عصبی بازگشتی نوعی شبکه عصبی است که یکی برای کدگذاری و دیگری برای کدگشایی داده‌ها استفاده LSTM حاوی گره‌هایی با یال به خود است) مانند می‌شود.

در این معماری ابتدا متن توسط واژه‌نهفت‌ها تبدیل به تعدادی بردار می‌شود. سپس این بردار به ترتیب وارد کدگذار شده و پس از پردازش، وضعیت مخفی (وضعیت فعلی شبکه عصبی بازگشتی پس از ورود واژه‌نهفت‌ها به کدگذار) به عنوان ورودی به کدگشا داده می‌شود و کدگشا پس از دریافت وضعیت مخفی شروع به تولید دنباله خروجی می‌کند. نحوه کارکرد آن را می‌توانید در شکل زیر مشاهده کنید



مدل‌های از پیش‌آموزش‌دیده

اخیراً استفاده از مدل‌های از پیش‌آموزش‌دیده (مدل‌هایی که وزن‌های آن‌ها را طی فرایند آموزش ذخیره می‌کنند و در زمان اجرا از آن وزن‌ها استفاده می‌کنند) به شکل گسترده‌ای مورد استفاده قرار می‌گیرند.

یکی از معماری‌هایی است که از کدگذارهای **تبدیل‌گر** برای ایجاد یک مدل زبانی BERT برای مثال معماری دوطرفه استفاده کرده‌است. این مدل از پیش‌آموزش‌دیده بر روی مسأله‌های مختلف به خوبی جواب داده‌است. نیاز دارد. تدقیق به آموزش مدل موجود با (fine-tuning) همچنین برای مسائل بسیار دیگری تنها به تدقیق تعداد محدودی داده در فضای مسأله جدید گفته می‌شود که سبب می‌شود تا مدل از پیش‌آموزش‌دیده با اندکی تغییر در وزن‌هایش برای مسأله جدید قابل استفاده باشد.

:References

<https://quera.ir/college>

<https://www.coursera.org/specializations/deep-learning>

<https://virgool.io/@sadeqhasan>

<https://www.kaggle.com/learn/natural-language-processing>