

(۱) الف)

NER:

در پردازش زبان طبیعی، NER یا Named Entity Recognition یا «تشخیص انتساب موجودات نامدار» است. این یکی از وظایف مهم در NLP است که به ما کمک می‌کند تا موجودات مختلفی که در یک متن ذکر شده‌اند را شناسایی و دسته‌بندی کنیم، مانند افراد، مکان‌ها، شرکت‌ها، تاریخ‌ها و غیره. به عبارت دیگر، NER به ما کمک می‌کند تا بفهمیم که یک کلمه یا یک مجموعه از کلمات در یک متن به چه معنی‌ای اشاره دارد و این اطلاعات را استخراج کنیم.

چالش‌ها:

- **Ambiguity.** کلمات ممکن است فریبنده باشند. یک اصطلاح مانند "amazon" ممکن است به رودخانه یا شرکت، بسته به زمینه، اشاره داشته باشد که شناسایی موجودیت را به یک تلاش دشوار تبدیل می‌کند.
- **Context dependency.** کلمات اغلب معنای خود را از متن اطراف می‌گیرند. کلمه "Apple" در یک مقاله فنی احتمالاً به شرکت اشاره دارد، در حالی که در یک دستور غذا، احتمالاً میوه است. درک چنین تفاوت‌های ظریف برای تشخیص دقیق موجودیت بسیار مهم است.
- **Language variations.** زبان انسان، با زبان عامیانه، گویش‌ها و تفاوت‌های منطقه‌ای آن، می‌تواند چالش‌هایی ایجاد کند. آنچه در یک منطقه رایج است ممکن است در منطقه دیگر بیگانه باشد و فرآیند NER را پیچیده کند.
- **Data sparsity.** برای روش‌های NER مبتنی بر یادگیری ماشین، در دسترس بودن داده‌های برچسب‌گذاری شده جامع بسیار مهم است. با این حال، به دست آوردن چنین داده‌هایی، به ویژه برای زبان‌های کمتر رایج یا حوزه‌های تخصصی، می‌تواند چالش برانگیز باشد.
- **Model generalization.** در حالی که یک مدل ممکن است در تشخیص موجودیت‌ها در یک حوزه برتر باشد، ممکن است در حوزه دیگر دچار تزلزل شود. اطمینان از تعمیم مدل‌های NER به خوبی در حوزه‌های مختلف یک چالش دائمی است.

(۱) ب)

- **پیچیدگی:** متون با ساختار واضح و پیچیدگی کم، اغلب با دقت بالاتری در سیستم‌های NER شناسایی می‌شوند. اما متونی که ساختار پیچیده‌تری دارند یا به عنوان مثال از اصطلاحات، اختصارات، یا جملات غیرمعمول استفاده می‌کنند، ممکن است باعث کاهش دقت سیستم‌های NER شود.
- **مبهم بودن:** متونی که شامل عبارات مبهم یا دارای ابهامات معنایی هستند، ممکن است برای سیستم‌های NER چالش برانگیز باشند. این ابهامات می‌توانند از نظر معنایی یا تعیین محدوده دسته‌بندی نام‌ها و اجزای دیگر متن ایجاد شوند.
- **انحراف از الگوهای استاندارد:** برخی از متون، ممکن است از الگوها یا قوانین استاندارد خارج شوند که سیستم‌های NER بر اساس آن‌ها آموزش داده شده‌اند. در این صورت، دقت سیستم ممکن است کاهش یابد زیرا الگوهای غیرمعمول برای آنها چالش بیشتری ایجاد می‌کند.
- **تنوع فرهنگی و زبانی:** تفاوت‌های فرهنگی و زبانی در بین زبانها و مناطق می‌تواند چالش‌هایی را برای دقت در NER ایجاد کند. زبانهای مختلف ممکن است دارای کنوانسیون‌های نامگذاری و ساختارهای گرامری متمایز باشند، که نیاز به رویکردهای مدل‌سازی خاص زبان و داده‌های آموزش دارند. علاوه بر این، تغییرات در منابع فرهنگی و هنجارهای متنی ممکن است بر تفسیر و طبقه‌بندی موجودات نامگذاری شده در متون چند فرهنگی یا چند زبانه تأثیر بگذارد.

- وابستگی به ساختار و محتوای جملات: ساختار و محتوای جملات متن می‌تواند نقش بسیار مهمی در تشخیص موجودیت‌ها داشته باشد. عناصر متن مانند کلمات مشابه، واژگان اضافی، و فعل‌ها ممکن است معنای موجودیت‌ها را متغیر کنند. برای مثال، جمله "مریم کتاب خواند" ممکن است به این معنا باشد که مریم یک کتاب را خوانده است یا کتابی با عنوان "مریم" وجود دارد.

(ج)

- برخورد با انعطاف پذیری بیشتر: یکی از محدودیت‌های اصلی HMM‌ها این است که فقط به ویژگی‌های محلی یا لحظه‌ای توجه می‌کنند و وابستگی به سایر ویژگی‌های موجود در داده را نادیده می‌گیرند. در عوض، CRF‌ها قادرند انعطاف بیشتری را در مدل‌سازی این وابستگی‌ها ارائه دهند، به خصوص با استفاده از ویژگی‌های چند متغیره که اطلاعات مرتبط با تمام دنباله را در نظر می‌گیرند.
- استفاده از ویژگی‌های بیشتر و متغیرهای بیشتر: در CRF‌ها، می‌توان از ویژگی‌های متعدد و متغیرهای بیشتری برای مدل‌سازی استفاده کرد. این امر به افزایش اطلاعات وارد شده به مدل و بهبود دقت در پیش‌بینی دنباله‌ها کمک می‌کند و به CRF‌ها امکان مطابقت بیشتر با واقعیت را می‌دهد، زیرا آنها می‌توانند ویژگی‌های مرتبط را به طور مستقیم با یکدیگر مدل کنند.
- استنتاج کلی: CRF‌ها هنگام پیش‌بینی، استنتاج سراسری را در کل دنباله انجام می‌دهند. این به آن‌ها اجازه می‌دهد تا وابستگی‌های بین برچسب‌ها را در کل دنباله در نظر بگیرند، نه صرفاً در نظر گرفتن انتقال‌های محلی مانند HMM. در نتیجه، CRF‌ها می‌توانند وابستگی‌های دوربرد را بهتر دریافت کنند و پیش‌بینی‌های آگاهانه‌تری انجام دهند.
- دسته‌بندی بهتر با داده‌های نامتوازن: CRF‌ها می‌توانند با داده‌های نامتوازن بهتر کار کنند، یعنی داده‌هایی که دارای نسبت نمونه‌های مثبت به منفی نامتوازن هستند. با استفاده از تابع‌های هزینه متناسب با اهمیت نمونه‌ها و با افزودن ویژگی‌های مرتبط، CRF‌ها می‌توانند عملکرد بهتری در مواجهه با این نوع داده‌ها ارائه دهند.
- کارایی آموزش: CRF‌ها را می‌توان با استفاده از تکنیک‌هایی مانند نزول گرادینت تصادفی (SGD) یا روش‌های گرادینت شرطی، که امکان همگرایی سریع‌تر و مقیاس‌پذیری را برای مجموعه داده‌های بزرگ فراهم می‌کند، به طور موثر آموزش داد. در مقابل، آموزش HMM‌ها اغلب شامل الگوریتم‌های تکراری مانند الگوریتم Baum-Welch است که ممکن است کندتر همگرا شوند و برای مجموعه داده‌های بزرگ کمتر مقیاس‌پذیر باشند.

(د)

1. I/PRP need/VBP a/DT flight/NN from/IN Atlanta/NN

Atlanta/NN درست نیست زیرا یک مکان خاص (a specific location) است و باید "NNP" باشد.

2. Does/VBZ this/DT flight/NN serve/VB dinner/NNS

"NNS" stands for plural noun, in this context "dinner" refers to a meal rather than individual items, so "NN" is a better label for that.

3. I/PRP have/VB a/DT friend/NN living/VBG in/IN Denver/NNP

Have should be labeled as "have/VBP" because it's a verb in the form of the base form.

4. Can/VBP you/PRP list/VB the/DT nonstop/JJ afternoon/NN flights/NNS

Can should be labelled as "MD" because it is a modal.

• روش برچسب‌گذاری (BIO (Begin, Inside, Outside)

یک روش متداول برای برچسب‌گذاری دنباله‌های متنی است، به ویژه برای شناسایی موجودیت‌های نام‌دار (Named Entities یا NE) در متن. در این روش، هر کلمه یا توکن در متن با یکی از برچسب‌های "I"، "B"، یا "O" برچسب‌گذاری می‌شود.

- "B" برای شروع یک NE استفاده می‌شود.

- "I" برای ادامه NE استفاده می‌شود.

- "O" برای هر کلمه‌ای که به NE مرتبط نیست، به عنوان خارج از موجودیت نام‌دار برچسب‌گذاری می‌شود.

• برچسب‌گذاری (IO (Inside, Outside)

در این روش، تنها دو برچسب "I" و "O" وجود دارد. "I" برای تمامی کلماتی که داخل یک موجودیت نام‌دار قرار دارند و "O" برای تمامی کلماتی که خارج از هر موجودیت نام‌دار هستند. از این روش برای مواردی که تنها به شناسایی موجودیت‌های نام‌دار می‌پردازیم و جزییات داخلی هر موجودیت برای ما مهم نیست استفاده می‌شود.

• برچسب‌گذاری (BIOES (Begin, Inside, Outside, End, Single)

در این روش، به علاوه‌ی برچسب‌های "I"، "B"، و "O"، دو برچسب جدید "E" (End) و "S" (Single) نیز وجود دارد. "E" برای آخرین کلمه یک NE استفاده می‌شود و "S" برای مواردی که یک NE تنها از یک کلمه تشکیل شده است. این روش به ما اجازه می‌دهد تا دقیقاً ابتدا و انتهای هر موجودیت را مشخص کنیم، که زمانی که نیاز به شناسایی دقیق محدوده‌های موجودیت‌های نام‌دار داریم می‌تواند مفید باشد.

(۳) الف)

• **ابهام در نام‌ها:** از آنجایی که در مجموعه داده نام فیلم‌ها وجود دارند، ممکن است نام افراد و دیگر نام‌های خاص شناخته نشوند.

• **نام فیلم (Movie Name):** بعضی از فیلم‌ها ممکن است نام‌های شبیه به هم داشته باشند که باعث ایجاد اشتباه در تفکیک و تشخیص موجودیت‌ها شوند، به عنوان مثال فیلم‌هایی که عنوان‌شان شامل کلمات مشابهی مثل "The A" و غیره باشد. همچنین عناوین فیلم‌ها اغلب حاوی کلمات یا عبارات رایجی هستند که می‌توانند تعبیر متعددی داشته باشند.

• **کاراکترهای خاص:** برخی از عنوان‌های فیلم ممکن است شامل کاراکترهای خاص، عبارات غیر استاندارد، یا ترکیباتی از اعداد و حروف باشند که باعث ایجاد اشتباه در تشخیص نام فیلم‌ها توسط الگوریتم‌های NER شود. همچنین ممکن است برخی از نام‌های کارگردان‌ها شبیه به نام‌های عمومی یا دیگر افراد باشند که باعث ایجاد اشتباه در تشخیص موجودیت‌ها شود.

• **همپوشانی موجودات نامگذاری شده:** در برخی موارد، عناوین فیلم ممکن است کلمات یا عباراتی را با موجودیت‌های نامگذاری شده دیگر به اشتراک بگذارند. به عنوان مثال، فیلم "The Godfather" کلمه

"Godfather" با نام شخصیت "Corleone Don" از همان فیلم مشترک است. این زمینه همپوشانی می‌تواند برای سیستم NER چالش برانگیز باشد تا بین موجودیت‌های نامگذاری شده مختلف به درستی تمایز قائل شود.

• **تناقضات یا تغییرات:** عناوین فیلم‌ها می‌توانند نسخه‌های زبانی متفاوت یا تغییراتی به دلیل محلی‌سازی یا تفاوت‌های فرهنگی داشته باشند. وجود چنین ناسازگاری‌هایی می‌تواند شناسایی و طبقه‌بندی دقیق موجودیت‌های نام‌دار را برای سیستم NER دشوارتر کند.

• **سال (Year) و امتیاز (Rating):** این فیلدها شامل اعداد و فرمت‌های عددی هستند که ممکن است باعث ایجاد اشتباه در تشخیص موجودیت‌ها توسط الگوریتم‌های NER شوند، به خصوص اگر الگوریتم به طور خودکار هر چیزی که به شکل عددی باشد، موجودیت مربوط به سال یا امتیاز تلقی کند.

General:

<https://www.youtube.com/watch?v=2XUhKpH0p4M>

https://www.youtube.com/watch?v=HKLbI_i7kuw

<https://www.youtube.com/playlist?list=PLvcbyUQ5t0UEK2KAGyUP7JO9K-Arct8OM>

Q1 theory:

<https://www.datacamp.com/blog/what-is-named-entity-recognition-ner>

<https://www.sketchengine.eu/penn-treebank-tagset>

https://www.ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.html

Q2 practical:

<https://www.youtube.com/watch?v=IqXdjdOgXPM>

<https://www.youtube.com/watch?v=fX5bYmnHqqE>