

تحليل متن فارسی

چالش‌ها و تجربه‌ها

دانشگاه اصفهان

گروه پژوهشی تجارت الکترونیکی شناختی

محمد مهدی بخشی

مهر ۱۴۰۲



پایان نامه‌ی من چه بود؟

- بخش بندی فروشندگان با استفاده از روند احساسات نظرات برخط در بازارهای الکترونیکی
- تحلیل احساسات
- پیش پردازش متن
- مدل سازی تحلیلگر احساسات
- خوشه بندی روندهای احساسات

دیتاست

دیتا از کجا بیارم؟

جمع کردن دستی -> روش من

- <https://nlpdataset.ir/>
- <https://github.com/MEgooneh/awesome-Iran-datasets>
- <http://dataheart.ir/>
- <https://github.com/persiannlp/parsinlu>

پیش‌پردازش متن

زودن متون از عبارت‌ها، کلمه‌ها و سایر مواردی که باعث ایجاد کج‌فهمی در مدل تشخیص احساسات می‌شود یا تغییر آن‌ها به نحوی که متن مطلوب به دست آید.

چه روشی برای پیش‌پردازش متن‌ها خوبه؟

- PERSIAN SENTIMENT ANALYSIS: FEATURE ENGINEERING, DATASETS, AND CHALLENGES (Asgarnezhad, Monadjemi)

چه جوری بفهمم خوبن یا نه؟



پیش‌پردازش متن

■ نرمالسازی و عملیات‌های مشابه مثل تبدیل متن محاوره‌ای به رسمی

■ با چه ابزاری؟

■ <https://github.com/ICTRC/Parsivar>

■ <https://www.roshan-ai.ir/hazm/docs/index.html>

■ <https://text-mining.ir>

■ تبدیل متن محاوره‌ای به رسمی



پیش‌پردازش متن

پاکسازی

چه چیزهایی باید حذف بشن؟

شماره‌های تماس، آدرس پست الکترونیکی، پیوند به سایت‌ها، نماد ارزها، اعداد و شماره‌ها، علائم نگارشی

به جز علامت تعجب، تگ‌های HTML و هشتک‌ها

کلمات بی‌اثر -> از کجا پیدا کنیم؟

ایموجی‌ها -> کدوم ایموجی‌ها؟

پیش پردازش متن

ریشه یابی |

Stemmer |

Lemmatizer |

پیش‌پردازش متن

نرمالسازی

پاکسازی تگ‌ها، هشتگ‌ها ...

کلمات بی‌اثر

ریشه‌یابی



مدل سازی و آموزش تحلیلگر

استخراج ویژگی‌ها -> کدوم روش؟

- TF-IDF
- Word embedding

یادگیرنده‌ها

الگوریتم‌های یادگیری ماشینی مرسوم

الگوریتم‌های جدیدتر مثل XLNET، BERT، GPT و T5 -> نیاز به سرورهای گرافیکی

<https://cloud.jarvislabs.ai/> همراه با شارژ رایگان در صورت عضویت در <https://forums.fast.ai/>

<https://ferdowsi.cloud/> به همراه تخفیفات <https://labsnet.ir/>

<https://lambdalabs.com/service/gpu-cloud>



پایان

ممنون که مرا شنیدید!