

ایمپورت کردن کتاب خانه های لازم:

```
import pandas as pd
from collections import Counter
from nltk import ngrams, pos_tag
from hazm import sent_tokenize, word_tokenize, stopwords_list, Normalizer
from nltk.probability import LaplaceProbDist, SimpleGoodTuringProbDist, FreqDist
```

مرحله صفر:

باز کردن فایل و تبدیل ستون کامنت به str

```
# Step 0: Open the Excel file
file_name = "digikala_comment.xlsx"
data_file = pd.read_excel(file_name)
data_file['comment'] = data_file['comment'].astype(str)
```

مرحله ۱-۱:

توکن کردن کامنت ها:

این مرحله با استفاده از hazm انجام شده است.

```
# Step 1: Preprocessing
# Step 1-1: Tokenize comments into sentences
print('step 1-1')
data_file['comment_sentences'] = data_file['comment'].apply(sent_tokenize)
```

مرحله ۲-۱:

حذف فضاهای خالی اضافی با استفاده از strip

```
# Step 1-2: Remove extra spaces
print('step 1-2')
data_file['comment_sentences'] = data_file['comment_sentences'].apply(lambda sentences: [sentence.strip() for sentence in sentences])
```

مرحله ۳-۱:

حذف استاپ ورد، علایم نگارشی و سایر موارد اضافی:

```
# Step 1-3: Tokenize, remove stop words, non-alphabetic characters, HTML tags, emojis, and normalize text
print('step 1-3')
stopwords = set(stopwords_list())
normalizer = Normalizer()
data_file['tokenized_comment'] = data_file['comment_sentences'].apply(lambda sentences: [
    [normalizer.normalize(word.replace('\u200c', ' ')) for word in word_tokenize(sentence) if word.isalpha() and word not in stopwords]
    for sentence in sentences
])
```

مرحله ۱-۴:

تبدیل اعداد به آدرس ها به جایگزین گفته شده:

```
# Step 1-4: Replace numbers with <NUM> and URLs with <URL>
print('step 1-4')
data_file['tokenized_comment'] = data_file['tokenized_comment'].apply(lambda sentences: [
    ['<NUM>' if word.isdigit() else '<URL>' if word.startswith(('http', 'www')) else word for word in sentence]
    for sentence in sentences
])
```

مراحل ۲-۲ ۲-۳ ۳-۲

ساخت n-gram:

```
# Step 2: Create language models
# Step 2-1: Create n-grams
print('step 2-1')
Comment Code
def get_ngrams(data, n):
    ngram_list = [ngrams(sentence, n) for sentence in data]
    ngrams_flat = [ng for sent in ngram_list for ng in sent]
    return ngrams_flat

unigrams = get_ngrams(data_file['tokenized_comment'], 1)
bigrams = get_ngrams(data_file['tokenized_comment'], 2)
trigrams = get_ngrams(data_file['tokenized_comment'], 3)
```

چاپ موارد مربوط به این مرحله:

نتیجہ چاب:

روش Laplace smoothing برای n-gram ها به خوبی عمل نمی کند زیرا در این روش به هر کلمه یک تعداد ثابت از تعداد تکرارها اضافه می شود تا از احتمال صفر برای کلماتی که در داده ها دیده نشده اند، جلوگیری شود. این اضافه کردن تعداد ثابت به تمام کلمات ممکن است باعث افزایش تعداد تکرارها و در نتیجه افزایش احتمالات شود. این امر می تواند منجر به افزایش دقت نا همگن در تخمین احتمالات شود.

به طور مثال، اگر یک کلمه در داده ها دیده نشده باشد، روش Laplace smoothing باعث افزایش احتمال آن کلمه می شود، حتی اگر دیگر کلمات در متن با آن هماهنگ نباشند. این مشکل به خصوص برای n-gram هایی با تعداد زیادی از واژگان ممکن خودش را نشان می دهد.

محاسبه perplexity:

```
import math
Comment Code
def calculate_perplexity(model, test_data):
    N = sum(len(sentence) for sentence in test_data)
    cross_entropy = -sum(math.log2(model.prob(word)) for sentence in test_data for word in sentence) / N
    perplexity = 2 ** cross_entropy
    return perplexity

test_sentences = [
    "بوی تند ولی خوشبو داره",
    "بلوتوثش کار نمیکنه حالا تا بدستم رسیده باید برش گردونم",
    "بلند گوهاش بیس بالا و صدای زیادی بمیداره که بعد از مدتی باعث خسته شدن مغز آدم میشه",
    "لطفاً کلاهی مورد نظر رو در پیشنهاد ویژه قرار بدید"
]

for i, sentence in enumerate(test_sentences, 1):
    laplace_unigram = LaplaceProbDist(freq_unigrams, bins=len(freq_unigrams))
    perplexity_unigram = calculate_perplexity(laplace_unigram, [sentence])

    laplace_bigram = LaplaceProbDist(freq_bigrams, bins=len(freq_bigrams))
    perplexity_bigram = calculate_perplexity(laplace_bigram, [sentence])

    laplace_trigram = LaplaceProbDist(freq_trigrams, bins=len(freq_trigrams))
    perplexity_trigram = calculate_perplexity(laplace_trigram, [sentence])

    print(f"Sentence {i} - Laplace Unigram Perplexity: {perplexity_unigram}")
    print(f"Sentence {i} - Laplace Bigram Perplexity: {perplexity_bigram}")
    print(f"Sentence {i} - Laplace Trigram Perplexity: {perplexity_trigram}")
    print('-' * 50)
```

گزارش این مرحله:

```

Sentence 1 - Laplace Unigram Perplexity: 82120.99999999996
Sentence 1 - Laplace Bigram Perplexity: 63562.999999999796
Sentence 1 - Laplace Trigram Perplexity: 56783.000000000005
-----
Sentence 2 - Laplace Unigram Perplexity: 82120.999999999873
Sentence 2 - Laplace Bigram Perplexity: 63563.000000000735
Sentence 2 - Laplace Trigram Perplexity: 56783.00000000054
-----
Sentence 3 - Laplace Unigram Perplexity: 82120.99999999975
Sentence 3 - Laplace Bigram Perplexity: 63563.00000000112
Sentence 3 - Laplace Trigram Perplexity: 56783.000000000684
-----
Sentence 4 - Laplace Unigram Perplexity: 82120.999999999894
Sentence 4 - Laplace Bigram Perplexity: 63563.000000000495
Sentence 4 - Laplace Trigram Perplexity: 56783.00000000054
-----

```

مرحله ۲-۴:

```

print('step 2-4')
Comment Code
def predict_next_words(model, input_sequence, length=15):
    for i in range(length):
        next_word = model.generate()
        input_sequence.append(next_word)

    return input_sequence

for sentence in ["یک تن ماهی خوب", "رنگ قرمز کفش", "گوشی سامسونگ", "یکی از چراغهای وضعیت", "سرفه جویی در پودر ماشین"]:
    print(sentence)
    print(predict_next_words(laplace_unigram, sentence.split()))
    print('one \n')
    print(predict_next_words(laplace_bigram, sentence.split()))
    print('two \n')
    print(predict_next_words(laplace_trigram, sentence.split()))
    print('three \n')

```

خروجی این قسمت یک مقدار طولانی بود قسمتی از خروجی نمایش داده شده:


```
# Step 3: POS Tagging
# Step 3-1: Perform POS tagging on the preprocessed data
print('step 3-1')
data_file['pos_tags'] = data_file['tokenized_comment'].apply(lambda sentences: [pos_tag(sentence) for sentence in sentences])

with open('pos.txt', 'w', encoding='utf-8') as file:
    for sentence_tags in data_file['pos_tags']:
        for tags in sentence_tags:
            sentence = ' '.join([f"{word} --> {tag}" for word, tag in tags])
            file.write(sentence + '\n')
        file.write('\n')
```

خروجی کار(فقط بخشی از خروجی در اسکرین شات است، فایل خروجی نیز ارسال شده):

```
ردم --> NNP هستن --> NNP می‌کنیم --> NNP ثبت --> NNP کالها --> NNP نظراتی --> NNP اشاره --> NNP موضوع --> NNP می‌خواستم --> NNP بگم --> NNP نظرم --> JJ سلام
NN --> NNP ندیم --> NNP الکی --> NNP کارشماهی --> NNP ک --> JJ ک --> JJ بهتره
- گذشت --> NNP عالیته --> NNP واقعا --> NNP خریدم --> NNP پیشنهاد --> NNP خریدم --> NNP اشتنا --> NNP نفر --> NNP دارم --> NNP بانک --> NNP پاور --> JJ سال
. داره --> NNP شارژ --> NNP بازدهی --> NNP کمتر --> NNP افت --> NNP دروعد --> NNP اونم --> NNP خریدم --> NNP همسرم --> NN اdata --> NNP پاور --> JJ همزمان
- واقعا --> NNP داره --> NNP فیک --> NNP باشه --> NNP اصل --> NNP می‌خرید --> NNP ک --> NNP کالایی --> NNP شیاومی --> NNP ک --> NNP اطرافیانم --> JJ تجربه
--> NN طولانی --> NN بیخفید
NN --> NNP میشه --> NNP خوشبوتر --> NNP میگذره --> NNP چقدر --> NNP خوبه --> NNP ماندگاریش --> NNP داره --> NNP خوشبو --> NNP تند --> JJ بوی
NNF --> NNP نحوی --> NNP میکنه --> NNP گیر --> NNP نمیکه --> NNP کار --> NNP بیاطری --> NNP دیگه --> NNP سال --> NNP هست --> NNP ماشین --> NNP مفید --> JJ عمر
```

مرحله ۲-۳:

تعداد رخداد هر تگ:

```
# Step 3-2: Count occurrences of each POS tag
print('step 3-2')
pos_tags_flat = [tag[1] for sentence in data_file['pos_tags'].sum() for tag in sentence]
pos_tags_count = Counter(pos_tags_flat)
print(f"3-2: POS Tags Count: {pos_tags_count}")
```

خروجی:

```
step 3-2
3-2: POS Tags Count: Counter({'NNP': 426441, 'NN': 47401, 'JJ': 38203, 'FW': 684, 'VBZ': 402, 'VBD': 332, 'DT': 181, 'VBP': 117, 'VB': 110, 'NNS': 86, 'IN': 65, 'MD': 57, 'CC': 53, 'RB': 51, 'CD': 43, 'VBG': 39, 'PRP': 13, 'JJR': 12, 'RP': 10, 'TO': 8, 'VBN': 7, 'PRP$': 7, 'JJS': 2, 'SYM': 1, 'NNPS': 1, '``': 1, 'RBR': 1})
```

مرحله ۳-۳:

```
# Step 3-3: Extract and count proper nouns (NNP)
print('step 3-3')
proper_nouns = [word[0] for sentence in data_file['pos_tags'].sum() for word in sentence if word[1] == 'NNP']
proper_nouns_count = Counter(proper_nouns)

print(f"3-3: Top 15 Proper Nouns: {proper_nouns_count.most_common(15)}")
```

خروجی:

```
3-3: Top 15 Proper Nouns: [('داره', 5550), ('کوفی', 5038), ('هات', 4233), ('استفاده', 4123), ('هست', 4016), ('یه', 3677), ('میفه', 3364), ('کیفیت', 3335), ('خریدم', 2948), ('واقعا', 2697), ('اها', 2426), ('فیمه', 2405), ('کار', 2299), ('خرید', 2298), ('دیجی', 2203), ('PS D:\sku\ntp\2>|
```

چون کلمات فارسی هستن یکم خوانش پایین اومده این تو کنسول vscode

سوال ۱

- مرحله یک: تولید کاندید ها (Candidate Generation): در این مرحله، با استفاده از مدل های مختلف، یک لیست از کوئری های پیشنهادی برای یک پیشوند (prefix) خاص تولید می شود.
- مرحله دو: رتبه بندی کاندید ها (Candidate Ranking): کاندید ها بر اساس اهمیت و ارتباط با کوئری ورودی رتبه بندی می شوند. این مرحله با استفاده از مدل های یادگیری عمیق انجام می شود.

سوال ۲

- حافظه محدود: توانایی مدل در به خاطر سپاری کاندید ها و اطلاعات مرتبط با کاربر.
- کیفیت بالای پیش بینی: توانایی مدل در پیش بینی کاندید های با کیفیت بالا برای ورودی های مختلف.
- کارایی و سرعت: اجرای سریع و کارایی بالا در محیط های واقعی.
- بهبود بازبازی (Recall): افزایش تعداد کوئری های مناسب در لیست پیشنهادی.
- بهبود دقت (Precision): افزایش ارتباط کاندیدها با نیاز واقعی کاربر.
- مدیریت بازه زمانی (Latency): افزایش سرعت و کاهش تاخیر در تولید پیشنهاد ها.
- مدیریت منابع (Resource Management): بهره گیری بهینه از منابع محاسباتی برای اجرای بهتر سیستم.

سوال ۳

- توانایی در مدل سازی ارتباطات پیچیده: مدل های زبان عصبی قابلیت مدل سازی ارتباطات پیچیده بین کلمات را دارا هستند.
- قابلیت یادگیری از داده: مدل های زبان عصبی با توجه به حجم بالای داده ها، قابلیت یادگیری و به روز رسانی بهتری دارند.
- انعطاف پذیری در تطبیق با داده جدید: مدل های زبان عصبی قابلیت تطبیق با داده های جدید را دارا هستند و می توانند به تغییرات در زبان و استفاده کنندگان پاسخ دهند.

- مدل سازی همه جانبه (End-to-End Modeling): قابلیت یادگیری ارتباطات پیچیده درون کاندیدها و کوئری ورودی.

- قابلیت افزودن ویژگی های اضافی (Incorporating Additional Features): قابلیت اضافه کردن ویژگی های شخصی سازی مانند شناسه کاربر در مدل.

- مدیریت همسانی دنباله (Sequence Coherence Modeling): توانایی مدل سازی همسانی دنباله های کلمات.

سوال ۴

رویکرد MCG یا Maximum Context Generation که در این مقاله مطرح شده است، در فاز تولید کاندید ها بر اساس الگوریتم گریدی (Greedy Matching) بر روی کلمات پیشوندی ورودی عمل می کند. این الگوریتم از آخرین چندین کلمه پیشوندی به صورت گریدی (از طولانی ترین به کوتاه ترین) برای پیدا کردن مطلوب ترین کاندید ها برای تکمیل خودکار کوئری استفاده می کند.

سوال ۵

Unnormalized Language Model (LSTMEmb): یک مدل زبان عصبی بر اساس شبکه عصبی با ساختار LSTM برای محاسبه امتیاز رتبه بندی.

Frequency-Based Models (MPC, LWG, MCG): مدل های مبتنی بر فراوانی که بر اساس تکرار و فراوانی واژگان محاسبه می شوند.

Hybrid Models (NN+Frequency): مدل های ترکیبی که از اطلاعات فراوانی و مدل های عصبی برای رتبه بندی استفاده می کنند.

Convolutional Latent Semantic Model (CLSM): یک مدل عصبی مبتنی بر شبکه های عصبی کانولوشنی برای محاسبه امتیاز رتبه بندی.