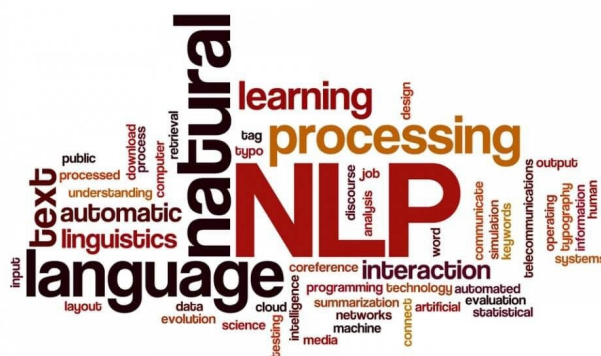


بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



# دانشکده مهندسی برق و کامپیوتر

## پردازش زبان‌های طبیعی - تمرین دوم

سید مهدی رضوی

استاد : آقای دکتر فیلی

اسفند - فرودین ماه ۱۴۰۲

## فهرست مطالب

۳	۱ تمرین اول
۳	۱.۱ پیش پردازش متن
۳	۲.۱ document term matrix
۵	۳.۱ TF-IDF matrix
۶	۴.۱ PPMI matrix
۷	۵.۱ نتیجه گیری
۸	۲ تمرین دوم
۹	۳ تمرین سوم

## فهرست تصاویر

۱۰	۱ سوال شماره ۲ که باز هم این سوال به هایپرپارامترهای ذکر شده در نتیجه گیری وابسته خواهد بود.
۱۰	۲ کاهش ابعاد مطلوب سوال ۳ (epochs = 10)
۱۱	۳ کاهش ابعاد مطلوب سوال ۳ (epochs = 20)
۱۱	۴ کاهش ابعاد مطلوب سوال ۳ (epochs = 100)
۱۲	۵ کاهش ابعاد مطلوب سوال ۳ (epochs = 40)

## ۱ تمرین اول

## ۱.۱ پیش پردازش متن

برای پیش پردازش بر روی متن ابتدا بعضی عبارات نگارشی را حذف می‌کنیم تا کلمات هم‌خانواده به راحتی قابل شناسایی باشند. سپس کلمات را به فرم کوچک آن‌ها تبدیل می‌کنیم و توکن گذاری خواهیم کرد. در مرحله بعد با استفاده از کتابخانه Porter Stemmer کلمات را هم‌ریشه خواهیم کرد.

## ۲.۱ document term matrix

برای ساختن این ماتریس به ازای هر سند ما یک بردار به طول اندازه مجموعه واژگان خود خواهیم داشت. بعد به ازای هر سند یک پیمایش بر روی واژگان سند خواهیم داشت و به ازای هر کلمه یک واحد درایه متناظر را اضافه خواهیم کرد. دقت داشته باشید چون که دو کلاس مثبت و منفی در میان داده‌ها داریم، به ازای هر کدام یک مجموعه واژگان جدا تشکیل می‌دهیم و برای هر کدام نیز یک ماتریس جداگانه تشکیل خواهیم داد. علت این امر محاسبه احتمالات likelihood جدا به ازای هر کلاس خواهد بود. بردار embedding یا همان احتمالات likelihood برای ساختن مدل Naive Bayes به صورت زیر محاسبه خواهد شد:

$$P(\text{term}|\text{class}) = \frac{\text{count}(\text{term}, \text{class}) + 1}{\sum_{\text{vocabulary}} (\text{count}(\text{term}, \text{class}) + 1)}$$

برای ساختن مدل Naive Bayes ما بایستی چند پارامتر مشهور این مدل را تخمین بزنیم. سپس با تخمین این پارامترها قادر به پیش‌بینی نمونه‌های جدید که مدل ندیده است، خواهیم بود. این پارامترها عبارتند از:

۱. احتمال پسین Prior Probability (چه میزان احتمال دارد که یک نمونه متعلق به یک کلاس باشد).

۲. احتمال likelihood (به شرط بودن در یک کلاس خاص، چه میزان احتمال دارد یک پیشامد مانند وقوع یک کلمه اتفاق بیوفتد).

در مدلی که برای این تمرین پیاده‌سازی کرده‌ام، ابتدا در تابع init به مقداردهی اولیه این متغیرها پرداخته‌ام. سپس در تابع fit به ازای نمونه‌های آموزشی به محاسبه هر یک از پارامترهای فوق پرداخته‌ام. چالش ما در این بخش محاسبه احتمال وقوع هر ترم در هر کلاس است. ما برای این کار دو مجموعه واژگان و دو ماتریس document term matrix تشکیل دادیم که در هر کدام به محاسبه  $P(\text{term} | \text{class})$  خواهیم پرداخت. در نهایت در تابع predict به ازای نمونه‌های تستی به محاسبه احتمال تعلق هر نمونه به هر کلاس را محاسبه خواهیم کرد و برچسب نمونه را تعیین خواهیم کرد.

$$P(\text{positive}|\text{sample}) <> P(\text{negative}|\text{sample})$$

نتایج ارزیابی ما به ازای این مدل تقریباً به صورت زیر خواهد بود :

precision : 0.6042345276872965

recall : 0.366600790513834

f1Score : 0.45633456334563344

## ۳.۱ TF-IDF matrix

ماتریس TF-IDF را به صورت حاصل ضرب Term Frequency در Inverse Document Frequency تشکیل خواهیم داد. بردار embedding یا همان احتمالات likelihood برای ساختن مدل Naive Bayes به صورت زیر محاسبه خواهد شد :

$$P(\text{term}|\text{class}) = \frac{tfidf(\text{term}, \text{class}) + 1}{\sum_{\text{vocabulary}} (tfidf(\text{term}, \text{class}) + 1)}$$

در واقع مانند سوال قبل بردار embedding را با استفاده از میانگین تشکیل خواهیم داد. البته با استفاده از هموارساز add-1-smoothing میانگین احتمالات را برای تشکیل مدل naive bayes را تشکیل می دهیم.

برای مدل Naive Bayes این Embedding نیز به مانند مدل قبلی توابع متناظر را پیاده سازی خواهیم کرد. فقط تابع احتمال likelihood ما تغییر خواهد کرد که برای محاسبه آن می بایستی از فرمول بالا استفاده کنیم.

نتایج ارزیابی این مدل به صورت زیر خواهد بود : نتایج ارزیابی ما به ازای این مدل تقریباً به صورت زیر خواهد بود :

precision : 0.7142857142857143

recall : 0.004940711462450593

f1Score : 0.009813542688910697

## ۴.۱ PPMI matrix

برای محاسبه این مرحله ، ابتدا نیاز به تشکیل ماتریس وقوع همزمان کلمات خواهیم داشت. (co-occurrence matrix) پس از ساختن این ماتریس ما قادر به ساخت ماتریس PPMI matrix خواهیم بود.

برای برچسب دادن به نمونه‌ها در مدل Naive Bayes در این قسمت ، بایستی که در هر جمله ، میانگین بردار PPMI کلمات آن جمله را محاسبه کنیم.  
در نهایت این بردار ماتریس ، بردار Likelihood ما را تشکیل خواهد داد.  
در نهایت نیز به پیش‌بینی نمونه‌های جدید خواهیم پرداخت.

**precision : 0.63468524567539514**

**recall : 0.046890927624872576**

**f1Score : 0.087329505**

## ۵.۱ نتیجه‌گیری

با توجه به این که ارزیابی ما با مدل Naive Bayes صورت می‌گیرد و این مدل نیز تا حدی به صورت زیادی به داده‌های آموزش معطوف می‌شود . در کل معیارهای ذکر شده در این تمرین تقریباً هر بار با تغییر random state تا حدود یک تا دو درصد میزان متغیرهای ارزیابی بالا یا پایین می‌شود. در کل همانطور که از نام این طبقه‌بند نیز مشخص است ، یک طبقه‌بند احمق است و در اینجا همانطور که از اعداد مشخص است ، با تغییر ماتریس document term به ماتریس TF-IDF که شاخص مناسب‌تری برای ارزیابی میزان مرتبط بودن کلمات در یک جمله است ، شاهد رشد ۱۰ درصدی در دقت بودیم ، اما متأسفانه در ماتریس PPMI که شاخص مناسب‌تری برای ارزیابی میزان وقوع کلمه در یک کلاس است ، ما شاهد بهبود نیستیم. برای شاخص recall نیز با توجه به این اعداد روش اول یعنی ماتریس document term . از همه روش‌ها مناسب‌تر است اختلاف فاحشی بین این روش با روش‌های دیگر embedding در این تمرین موجود است. بیشتر به نظر می‌آید با توجه به ماهیت این ماتریس که متناظر با فرکانس وقوع کلمه است و موارد دیگر خیلی تأثیر کمی دارد ، می‌تواند recall را به خوبی محاسبه کند.

## ۲ تمرین دوم

برای این تمرین ما به این صورت پیش رفتیم که تابع پیش پردازش را برای مجموعه تیتراهای اخبار به چند صورت پیاده سازی نمودیم . در نهایت به این نتیجه رسیدیم که پیاده سازی مرحله پیش پردازش باید تنها به صورت توکن کردن متن باشد و حذف کلمات به اصطلاح stop words و هم ریشه کردن کلمات هم خانواده بهتر است صورت نگیرد.

به نظرمی رسد که برای بررسی وجود و یا عدم وجود کنایه در متن ، این کلمات تاثیرگذار خواهند بود.

همچنین در میزان دقت و به خاطر سپاری مدل تاثیرگذار خواهند بود این کلمات.

سپس بردارها را از مخزن دانشگاه استنفورد دانلود ، و از آن ها به عنوان بردار embedding و در مدل logistic regression به عنوان بردار ویژگی استفاده خواهیم کرد.

ماتریس embeddings دوبعدی را با ابعاد اندازه واژگان در ابعاد هم نشینی که با توجه به فایل Glove صد می باشد ، خواهیم ساخت.

توزیع پراکندگی میزان متغیر هدف ما که همانطور که مشاهده می کنید میزان نمونه هایی که دارای کنایه هستند ۱۰۰۰ واحد بیشتر هستند.

has	sarcastics	:	13634
non	sarcastics	:	14985

نتایج حاصل از اجرای مدل ما به شرح زیر است : در جعبه متنی زیر قرار گرفته است.

استنتاجی که ما از این نتایج داریم به این صورت خواهد بود که ، نرخ میزان نمونه ای که به اشتباه بدون کنایه تشخیص داده شده است بسیار کم است. یعنی مدل تا حد بسیار خوبی توانسته است که نمونه هایی که کنایه دار است را تشخیص دهد.

اما نرخ میزان نمونه ای که به اشتباه دارای کنایه تشخیص داده شده است ، تا حد قابل توجهی بالا می باشد .

به نظرمی رسد کمی این موضوع به توزیع تعداد داده ها برمی گردد .

تعداد نمونه هایی که دارای کنایه است ۱۰۰۰ واحد کمتر است .

در نتیجه مدل کمتر بر روی داده هایی که دارای کنایه است آموزش دیده است تا داده هایی که بدون کنایه است.

Accuracy	:	0.556953179594689
precision	:	0.5202773691952091
recall	:	0.907292048369366
f1	:	0.6613247863247863



## ۳ تمرین سوم

برای این تمرین ، ما برای آموزش دادن بردارهای جانمایی به مدل negative sampling از تکنیک نمونه برداری منفی استفاده کرده ایم. بدین معنا که به ازای هر بردار کلمه با برچسب مثبت ، ۵ کلمه رندوم را انتخاب می کنیم . استفاده از دیکشنری word to index و بالعکس برای تبدیل کلمه به اندیس و بالعکس بسیار برای ما کارآمد بود. سپس ما بردارهای جانمایی و زمینه را به صورت رندوم مقداردهی اولیه خواهیم کرد. سپس با توجه به تعداد numEpochs به آموزش مدل Skip Gram خواهیم پرداخت. در شکل های زیر به تاثیر بسیار زیاد هایپر پارامتر numEpochs در آموزش مدل برخورد خواهیم کرد. همچنین با توجه به اجراهای متفاوت ، با توجه به بررسی حالات مختلف این ۴ کلمه ، به تاثیر بردارهای رندوم اولیه نیز پی خواهیم برد. جایی که در دو اجرای متفاوت ، با یک میزان numEpochs به بردارهایی کاملاً متفاوت خواهیم رسید.

اما با توجه به یک سری هایپر پارامتر و همچنین تاثیر میزان متن بر روی شکل گیری بردارهای جانمایی ، بردارهای حالت ایده آل کمی با موازی بودن فاصله دارد.

(به عنوان مثال بررسی رابطه بین عمه و عمو در متن شرلوک هلمز به مراتب ممکن است کمتر از بردار و خواهر باشد. جایی که حضور بیشتر معادل مذکر در یک بردار در دیگر حضور به مراتب کمتری دارد. در نتیجه بردارهای جانمایی آنها نیز کمتر فرصت حضور در ایپاک ها برای آموزش خواهند داشت. ) بنابر موضوع و محوریت مقاله ذکر شده در تمرین ، پیش بینی ما باید بر موازی بودن دو دوبردار ذکر شده در ایده آل ترین حالت ممکنه باشد.

اما با توجه به یک سری هایپر پارامتر و همچنین تاثیر میزان متن بر روی شکل گیری بردارهای جانمایی ، بردارهای حالت ایده آل کمی با موازی بودن فاصله دارد.

(به عنوان مثال بررسی رابطه بین عمه و عمو در متن شرلوک هلمز به مراتب ممکن است کمتر از بردار و خواهر باشد. جایی که حضور بیشتر معادل مذکر در یک بردار در دیگر حضور به مراتب کمتری دارد. در نتیجه بردارهای جانمایی آنها نیز کمتر فرصت حضور در ایپاک ها برای آموزش خواهند داشت. )

uncle	:	19
aunt	:	3
brother	:	8
sister	:	34

با توجه به تعداد epochs مشخص است هر چه مدل بیشتر فرصت آموزش داشته باشد ، بردارها به منطق مورد نظر ما نزدیک تر خواهند بود.

همچنین بردارهای رندوم اولیه نیز تاثیرگذار خواهند بود.

میزان تکرار کلمه مدنظر برای آموزش دیدن مدل skip gram برای تشکیل بردارهای جانمایی نیز تاثیر بسزایی خواهند داشت.

```
NLP_CA2_PART2.ipynb - Colab - Google Chrome
colabresearch.google.com/drive/12KvafE3yijFioghrnzH_kv5KwC7_ClytscrollTo=DXH9ASJHijJ
NLP_CA2_PART2.ipynb
File Edit View Insert Runtime Tools Help
+ Code + Text
# Get the indices of the words
queen_idx = word_to_index["queen"]
king_idx = word_to_index["king"]
man_idx = word_to_index["man"]
woman_idx = word_to_index["woman"]

# Calculate the new vector: king - man + woman
result_vector = embedding_vectors[king_idx] - embedding_vectors[man_idx] + embedding_vectors[woman_idx]

# Calculate the similarity using inner product
similarity = np.dot(embedding_vectors[queen_idx], result_vector)

print("Similarity between 'queen' and (king - man + woman) vector:", similarity)
Similarity between 'queen' and (king - man + woman) vector: 6.973933725831802

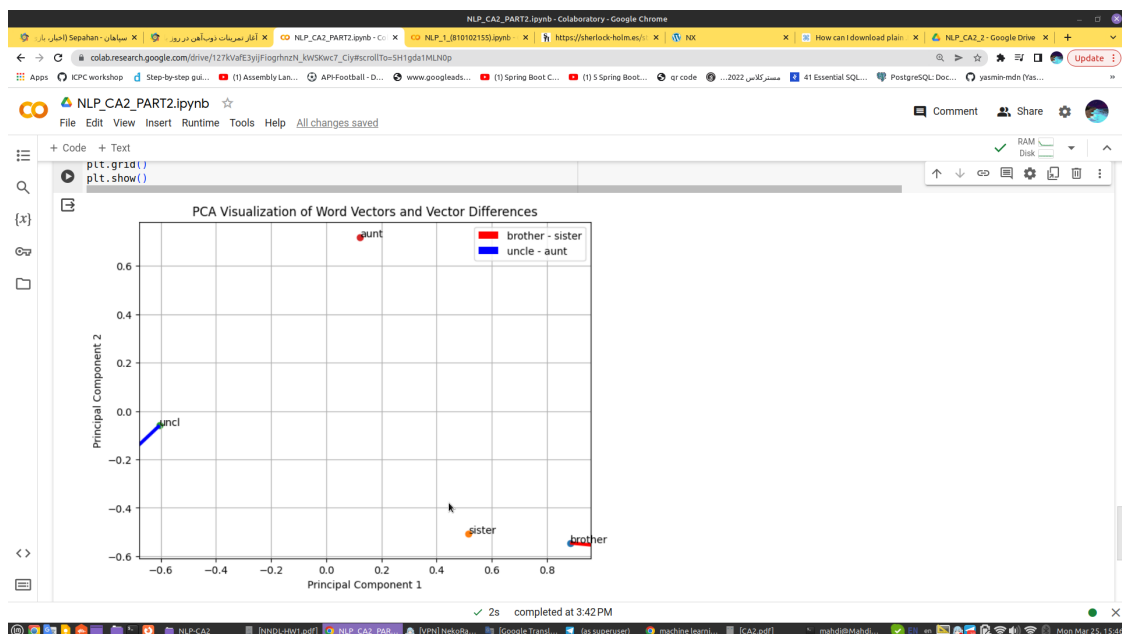
section 3 (PCA Transformation)

[19] import numpy as np
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt

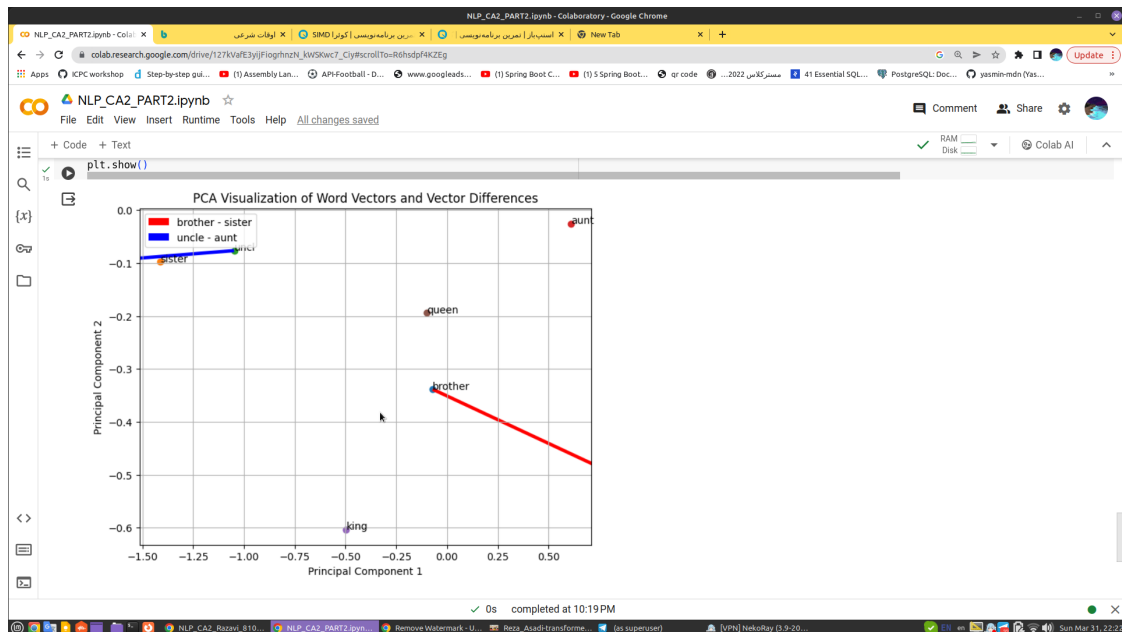
# Perform PCA transformation
pca = PCA(n_components=2)
embedding_vectors_2d = pca.fit_transform(embedding_vectors)

# Visualize the feature vectors of words in two dimensions
words_to_plot = ["brother", "sister", "uncle", "aunt", "king", "queen"]
```

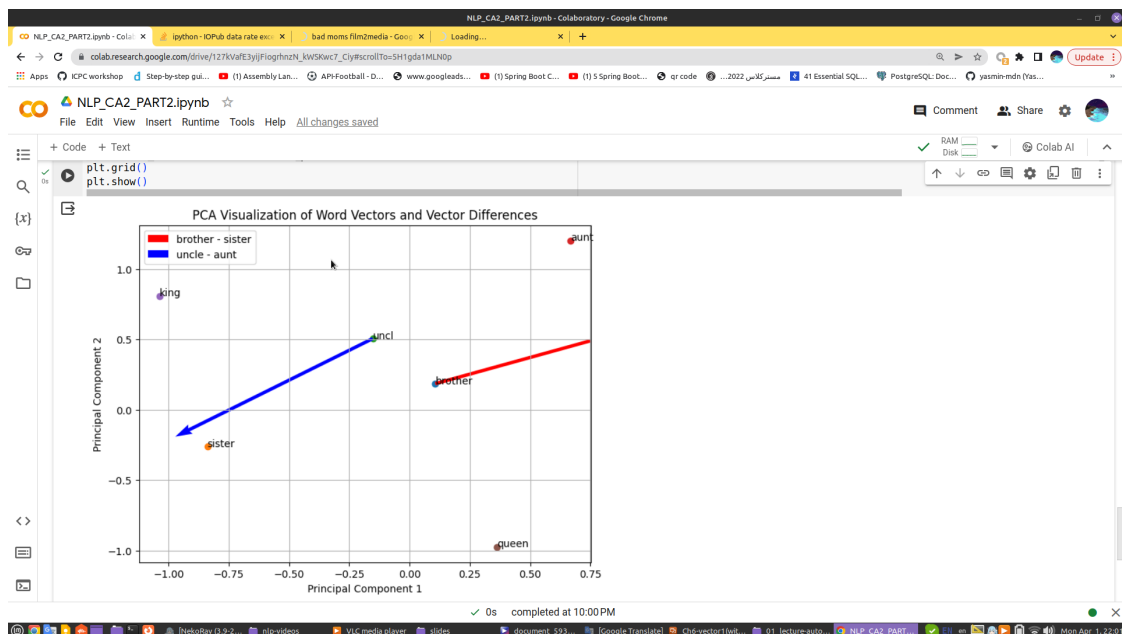
شکل ۱: سوال شماره ۲ که باز هم این سوال به هاپیرپارامترهای ذکر شده در نتیجه‌گیری وابسته خواهد بود.



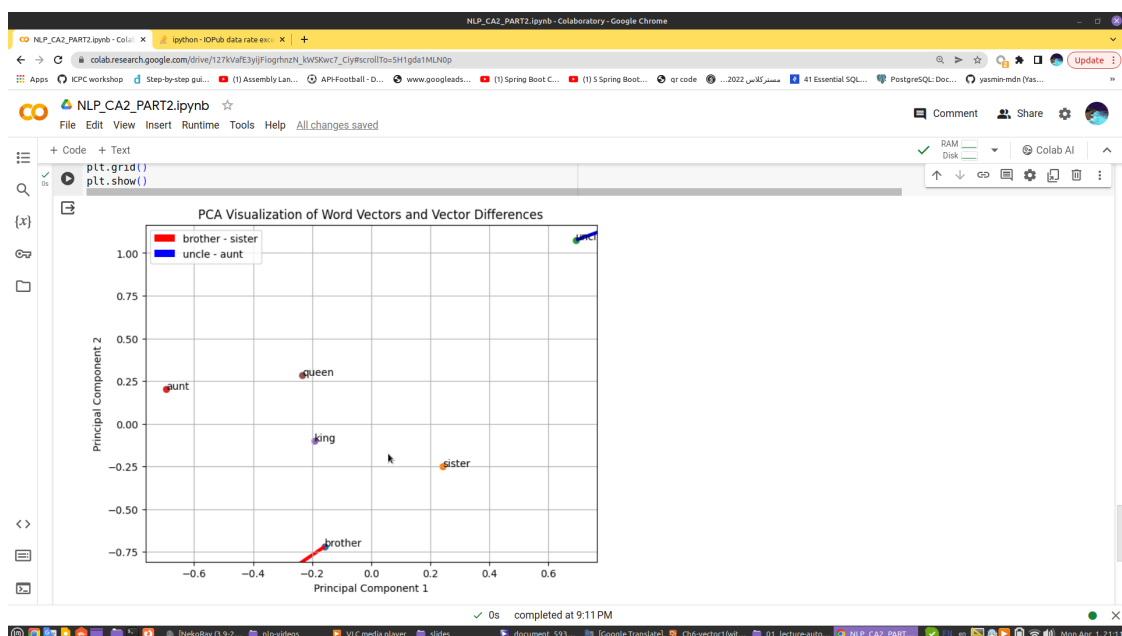
شکل ۲: کاهش ابعاد مطلوب سوال ۳ (epochs = 10)



شکل ۳: کاهش ابعاد مطلوب سوال ۳ (epochs = 20)



شکل ۴: کاهش ابعاد مطلوب سوال ۳ (epochs = 100)



شکل ۵: کاهش ابعاد مطلوب سوال ۳ (epochs = 40)