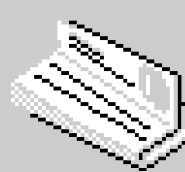
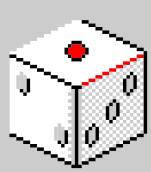
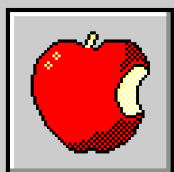


# گزارش پروژه دوم گروه ششم

اعضا:  
فاطمه خجسته  
مانا سامانی  
مهدی اسدی  
آناهیتا الیاسی

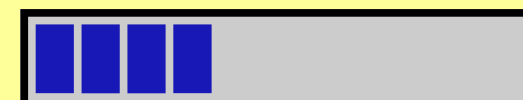


بوتکمپ هوش مصنوعی دلتا کوئرا



11:11PM

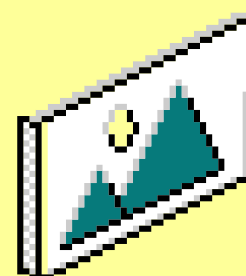




# پوشش پلاک خودرو

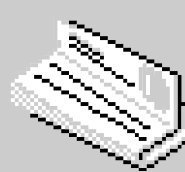
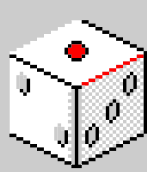
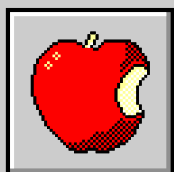
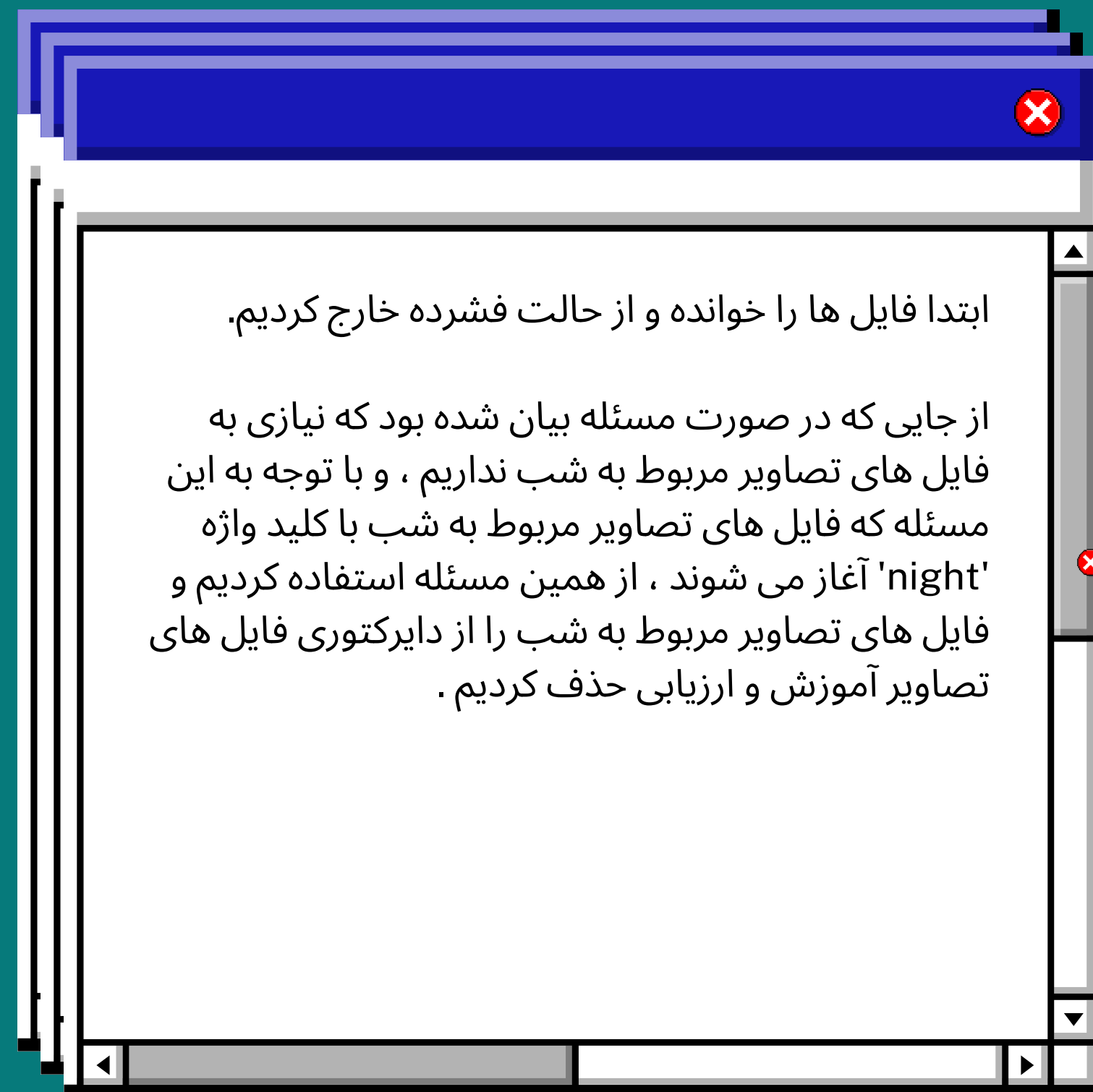
License plate cover

[Back to Agenda Page](#)



# پوشش پلاک خودرو

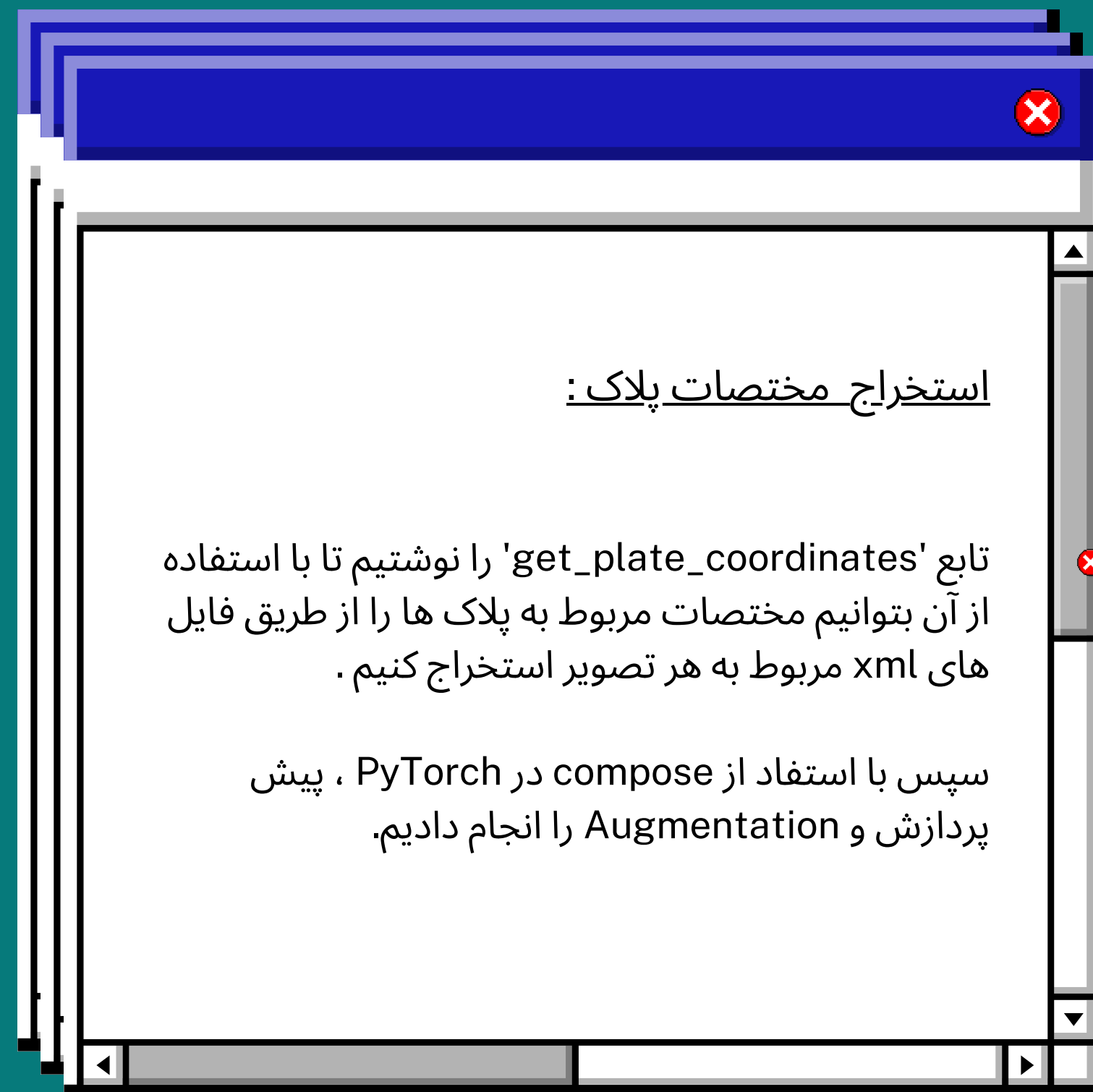
loading data & preprocessing



[Back to Agenda Page](#)

# پوشش پلاک خودرو

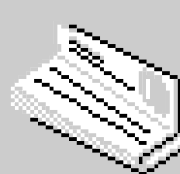
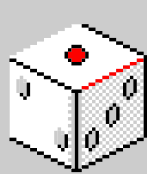
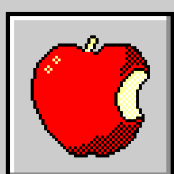
loading data & preprocessing



استخراج مختصات پلاک:

تابع 'get\_plate\_coordinates' را نوشتیم تا با استفاده از آن بتوانیم مختصات مربوط به پلاک ها را از طریق فایل های xml مربوط به هر تصویر استخراج کنیم.

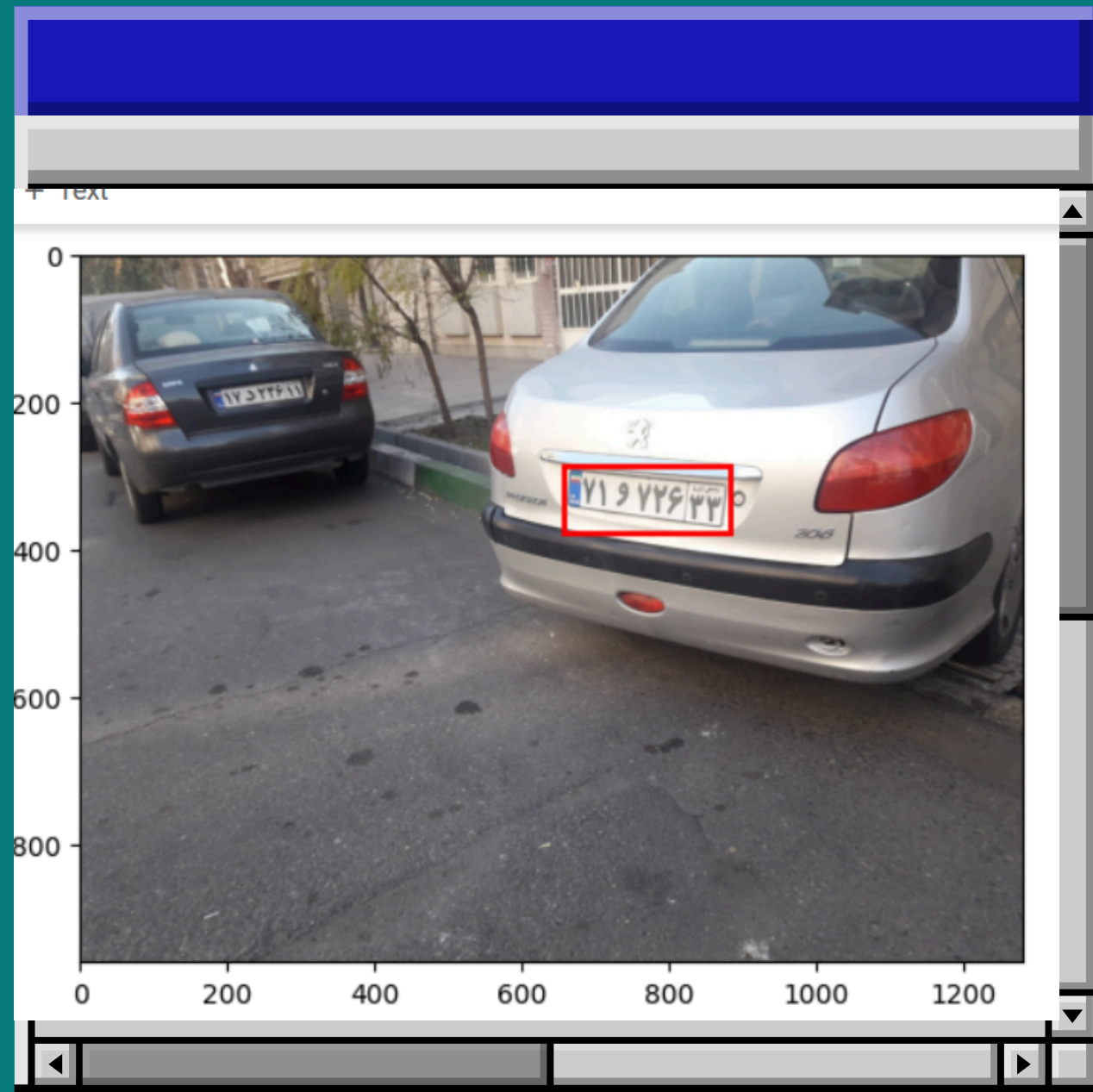
سپس با استفاده از compose در PyTorch ، پیش پردازش و Augmentation را انجام دادیم.



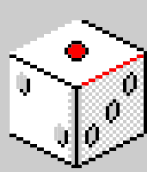
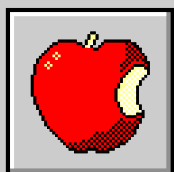
[Back to Agenda Page](#)

# پوشش پلاک خودرو

loading data & preprocessing



تابع 'show\_image\_with\_bounding\_box' را نوشتیم که با دریافت ورودی های آدرس دایرکتوری تصویر و مختصات پلاک خروجی از تابع 'get\_plate\_coordinates' بتواند پلاک را روی تصویر مربوطه نمایش بدهد.



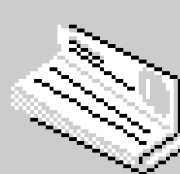
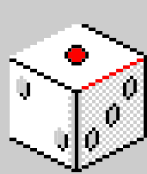
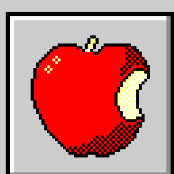
[Back to Agenda Page](#)

# پوشش پلاک خودرو

loading data & preprocessing



با کمک کلاس 'LicensePlateDataset' دیتاهای آموزشی و اعتبارسنجی را به مشخصات مورد نیاز تبدیل میکنیم و با استفاده از این کلاس برای داده های آموزش و ارزیابی دیتا لودر تعریف کردیم که دیتاست ایجاد شده از تصاویر توسط کلاس 'LicensePlateDataset' را آماده سازی و لود کند.



[Back to Agenda Page](#)

# پوشش پلاک خودرو

loading data & preprocessing

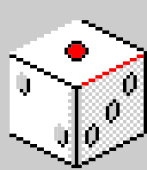
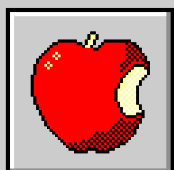


در ادامه مدل از پیش آموزش دیده YOLO موجود در کتابخانه Torch را بارگذاری کردیم و تنظیمات مربوط به آنرا انجام دادیم.

با استفاده از بهینه ساز 'Adam' و محاسبه loss از روش 'CrossEntropy' مدل را با کمک یک حلقه آموزش دادیم.

پس از اتمام آموزش مدل ، آنرا ذخیره می کنیم تا بتوانیم آنرا برای استفاده ی مجدد در دسترس داشته باشیم.

برای انجام ادامه مراحل مدل ذخیره شده را لود و بارگذاری می کنیم.



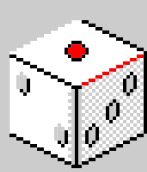
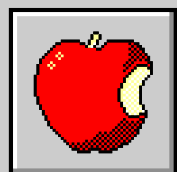
[Back to Agenda Page](#)



# پوشش پلاک خودرو



loading data & preprocessing



[Back to Agenda Page](#)

# پوشش پلاک خودرو

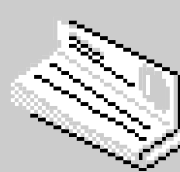
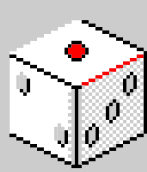
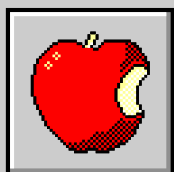


loading data & preprocessing

قابلیت‌ها: شناسایی اشیاء در سایزهای مختلف در تصویر با دقت نسبتاً بالا، حتی در شرایط مختلف نور و پس‌زمینه.

نحوه کار YOLO3 در Torch :  
YOLOv3. مرحله اول: تصویر به مدل داده می‌شود به جای تقسیم تصویر به بخش‌های کوچک، کل تصویر را هم‌زمان پردازش می‌کند.

مثلاً مرحله دوم: مدل به بخش‌های مختلف تصویر نگاه می‌کند تا بتواند اشیاء  $52 \times 52$  و  $26 \times 26$  و  $13 \times 13$  را در سایزهای مختلف پیدا کند.



[Back to Agenda Page](#)

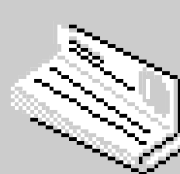
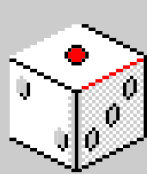
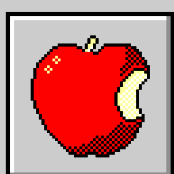
# پوشش پلاک خودرو

loading data & preprocessing



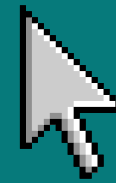
مرحله سوم: برای هر بخش، مدل تعدادی جعبه (Bounding Box) پیشنهاد می‌دهد که ممکن است محل اشیا باشند. هر جعبه شامل مختصات، نوع شی (مثل ماشین یا آدم) و میزان اطمینان مدل به شناسایی است.

مرحله آخر: YOLOv3 با حذف جعبه‌های تکراری، فقط جعبه‌های با احتمال بالا را نگه می‌دارد و مکان اشیا را در تصویر برمی‌گرداند.



[Back to Agenda Page](#)

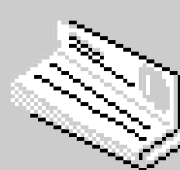
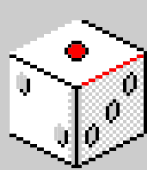
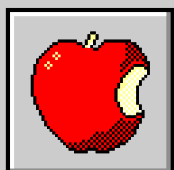
# پوشش پلاک خودرو



loading data & preprocessing

در ادامه تابع 'evaluate\_model' را برای مشاهده پیش بینی و ارزیابی مدل در زمینه تشخیص پلاک نوشتیم.

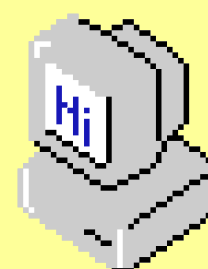
متأسفانه مدل پس از تست های متوالی و تغییرات مختلف و تلاش برای بهینه سازی و تغییر و تغییر در نحوه آموزش ، نمی تواند به درستی پلاک ها را تشخیص بدهد و دارای مشکل و خطا می باشد !



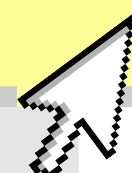
[Back to Agenda Page](#)



# NLP

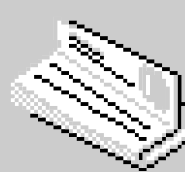
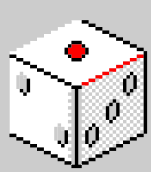
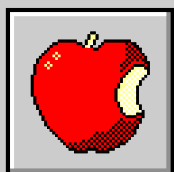


[Back to Agenda Page](#)



# NLP

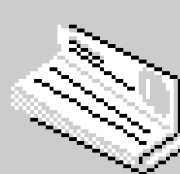
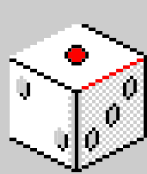
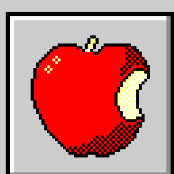
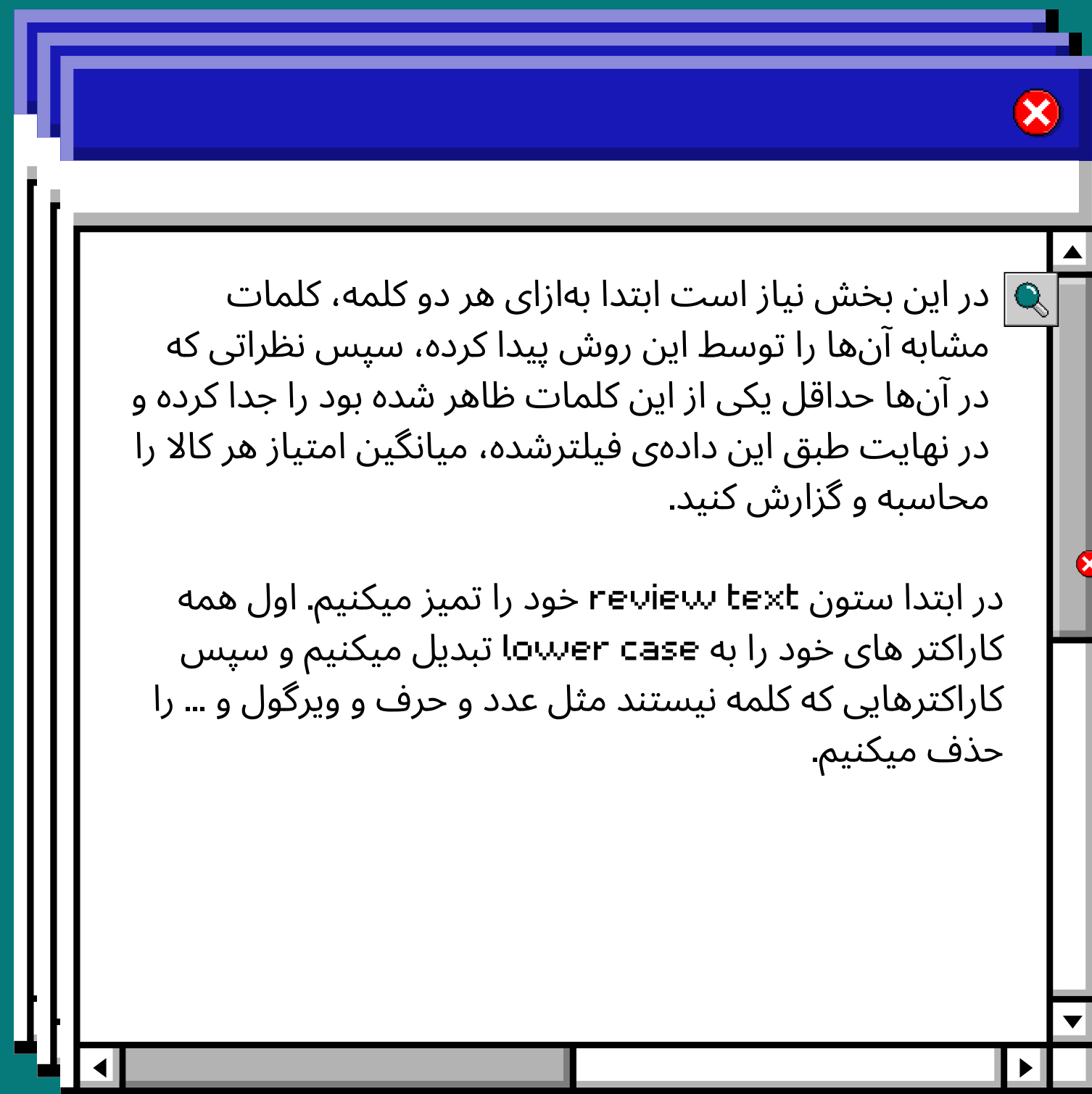
بخش اول  
تجزیه و تحلیل اولیه از داده‌ها



[Back to Agenda Page](#)

# NLP

## بخش دوم میزان رضایت از یک جنبه‌ی مشخص



[Back to Agenda Page](#)

## بخش دوم

# میزان رضایت از یک جنبه‌ی مشخص

**مدل سازی:**

در بخش مدل سازی از کتابخانه gensim که داری مدل Word2Vec است استفاده میکنیم.

میکنیم. یعنی هر جمله را tokenize خود را review text ابتدا ستون به کلمه های مختلف تقسیم میکنیم. یعنی یک آرایه از کلمات را ما داریم.

بعد پارامتر های مدلی خود را تعریف کرده ایم.

پارامتر vector size : تعداد ویژگی ها را در بردار های کلمات را مشخص می کند.

پارامتر window: اندازه فاصله را تنظیم می کند که چند کلمه اطراف هر کلمه هدف هنگام یادگیری روابط بین کلمات در نظر گرفته شود.

پارامتر min count: کلماتی که کمتر از این مقدار ظاهر شده اند نادیده می گیرد و نویز کلمات نادر را کاهش می دهد.

سپس مدل خود را با این هایپرپارامتر ها آموزش میدهیم و در لیست ذخیره میکنیم.

مدل word2vec روابط بین کلمات را بر اساس زمینه تعریف شده توسط پارامتر ها یاد می گیرد.



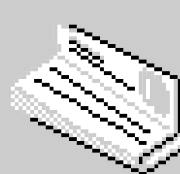
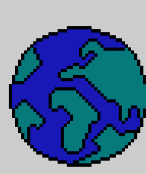
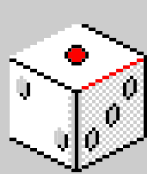
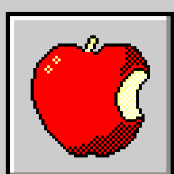
# NLP

## بخش دوم میزان رضایت از یک جنبه‌ی مشخص



سپس با استفاده از مدل هایی که آموزش داده ایم، کلماتی که مشابه با warranty و guarantee هستند را پیدا می‌کند. سپس همه این ها را در یک لیست ذخیره می‌کند. از set استفاده کرده ایم که اگر duplicate داشتیم، حذف کند.

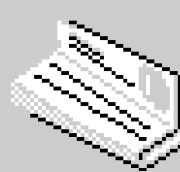
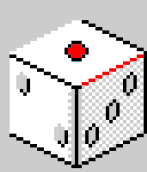
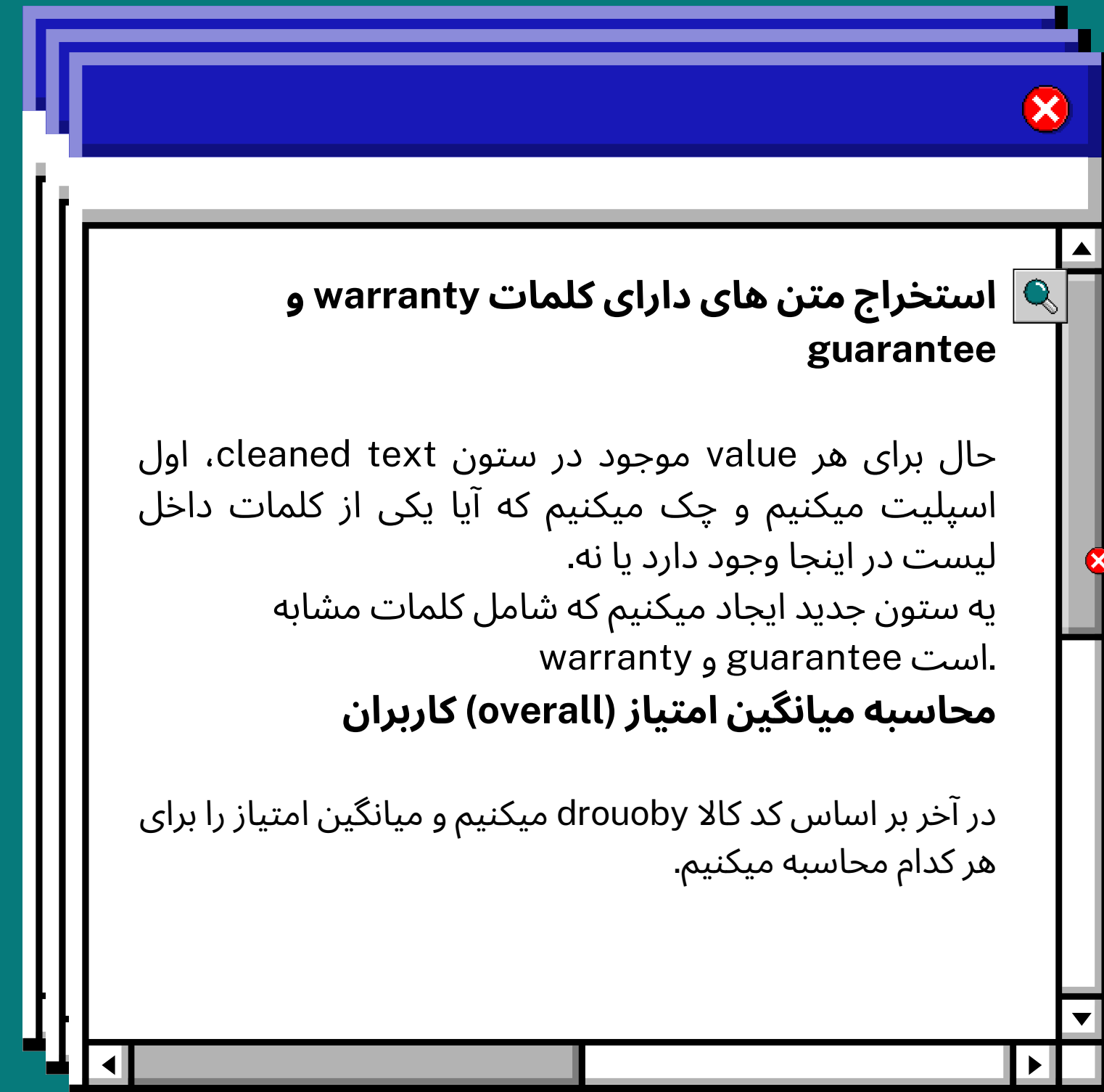
سپس کلمات داخلی لیست را چک میکنیم و بعضی از کلمه هارا که مشابهتی ندارند را حذف میکنیم. سپس یک دور دیگر بر روی تک تک کلماتی که داریم، حروف مشابه به آن ها را پیدا میکنیم. (در حدود ۱۳۳ تا کلمه) از بین این ۱۳۳ کلمه، کلماتی که معنی آن ها بای حرف هدف ما متفاوت است را حذف میکنیم.



[Back to Agenda Page](#)

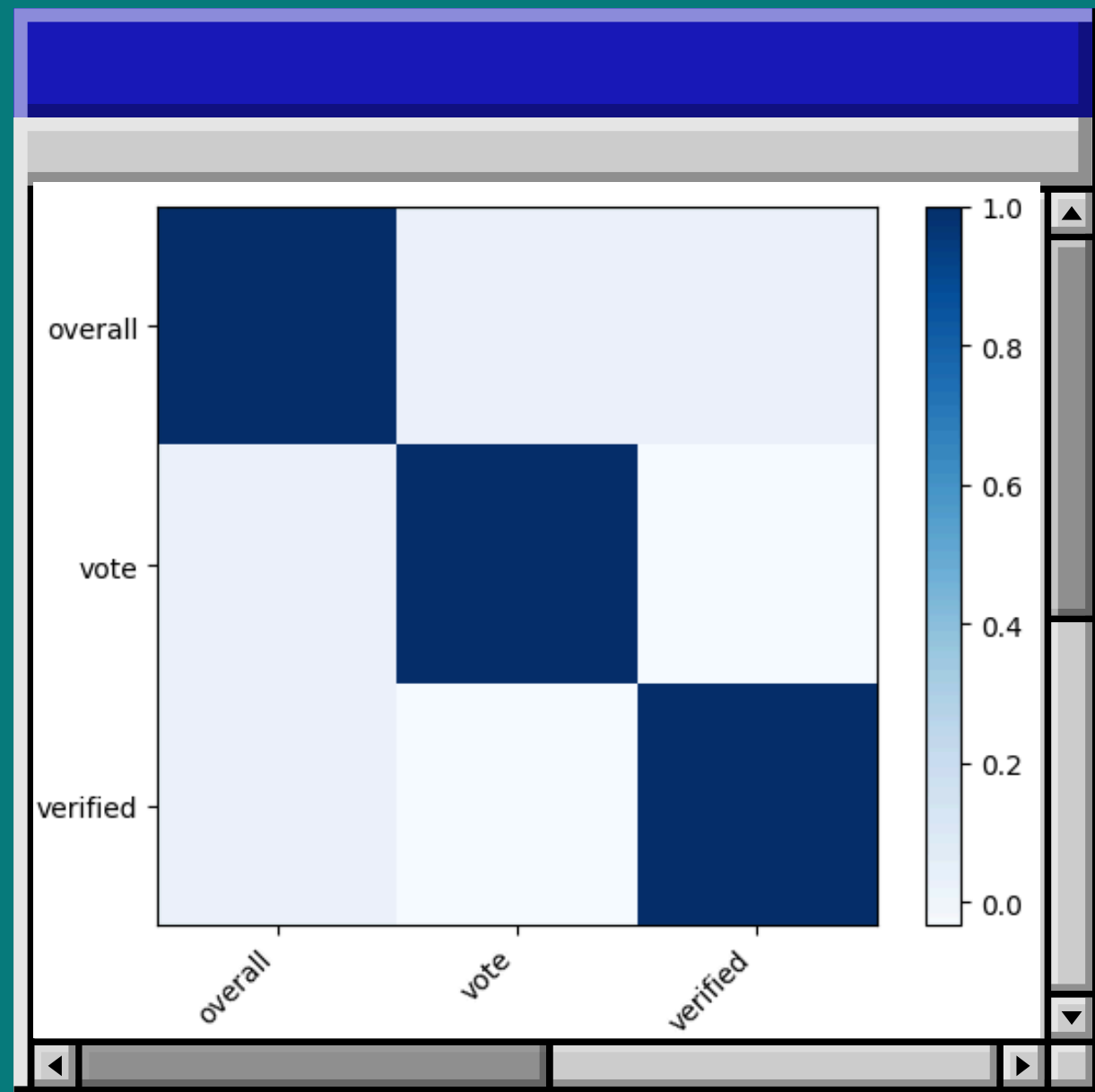
# NLP

## بخش دوم میزان رضایت از یک جنبه‌ی مشخص



# NLP

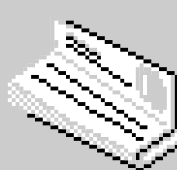
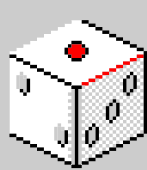
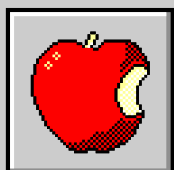
## بخش سوم مدل تحلیل احساس



نیاز است مدلی طراحی کنید که با دریافت متن نظر کاربر، احساس رضایت وی نسبت به کالا را بین عددی از ۱ تا ۵ تعیین کند.

مهندسی ویژگی:

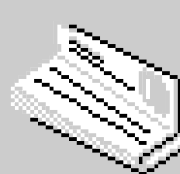
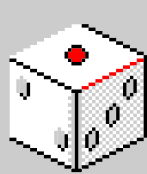
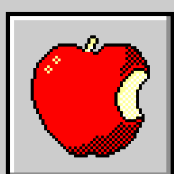
- ستون `style` به دلیل پیچیدگی و محتوای کم از دیتای آموزش حذف شد
- طبق تحقیقاتی که انجام شد، ستون های `asin`, `title`, `brand` مفید نبودند و حذف شدند
- بررسی ارتباط ستون های `verified`, `overall`، دو ستون دیگر `vote` از نظر همبستگی نشان داد به دلیل همبستگی بسیار پایین با `overall`، دو ستون دیگر باید حذف می شدند.



[Back to Agenda Page](#)

# NLP

## بخش سوم مدل تحلیل احساس



[Back to Agenda Page](#)

# NLP

## بخش سوم مدل تحلیل احساس

```
words_to_keep = {  
    "good", "bad", "perfect", 'great',  
    "poor", "nice", "amazing", "disappointing", ...
```

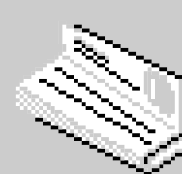
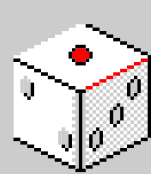
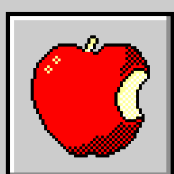
```
contractions = {  
    "can't": "cannot", "won't": "will not", "n't": "  
not", "re": " are", ...
```

```
important_stopwords = {  
    "not", "no", "never", "very", "all" ...
```

```
negation_words = ("not", "no", "never",  
    "none", "nothing", ...
```

### پردازش متن :

در این قسمت با استفاده از کتابخانه‌های `pandas` و `nlTK`، نظرات کاربران از فایل داده بارگذاری شده و در قالب یک ستون متنی ترکیب شدند. با اعمال یک تابع پیش‌پردازش، متن به حروف کوچک تبدیل شده، اختصارات به صورت کامل نوشته شدند، نشانه‌گذاری‌های غیرضروری حذف شد، و کلمات پرتکرار کاهش یافتند.



[Back to Agenda Page](#)

# NLP

## بخش سوم مدل تحلیل احساس

```
words_to_keep = {  
    "good", "bad", "perfect", 'great',  
    "poor", "nice", "amazing", "disappointing", ...
```

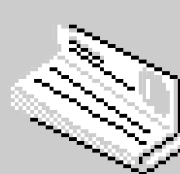
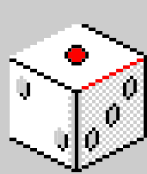
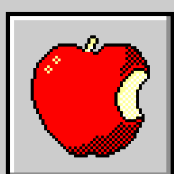
```
contractions = {  
    "can't": "cannot", "won't": "will not", "n't": "  
not", "re": " are", ...
```

```
important_stopwords = {  
    "not", "no", "never", "very", "all" ...
```

```
negation_words = ("not", "no", "never",  
    "none", "nothing", ...
```

پردازش متن :

سپس کلمات باقیمانده بدون توقف به جز کلمات مهم،  
حذف و پس از ریشه‌یابی، الگوی منفی‌سازی در جملات  
نیز اعمال شد تا عبارات منفی به شکل برجسته نشان  
داده شوند. در نهایت، این متن‌های پردازش‌شده ذخیره  
شدند.



[Back to Agenda Page](#)

# NLP

## بخش سوم مدل تحلیل احساس

```
pipeline = Pipeline([  
    ('tfidf',  
     TfidfVectorizer(max_features=20000,  
                     ngram_range=(1, 3))),  
    ('clf',  
     LogisticRegression(multi_class='multinomial',  
                        solver='lbfgs', max_iter=5000))  
])
```

### مدل سازی :

مدل 1:

در مدل یک ابتدا از logistic regression به همراه

tfidfvectorizer استفاده کردم.

مدل tfidf، تکست ما را به یک فیچر عددی تبدیل می‌کند که

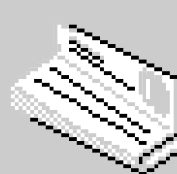
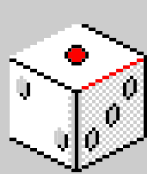
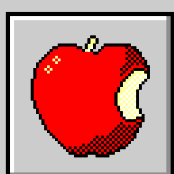
اهمیت هر کلمه را در تکست نشان می‌دهد.

این مدل از ngram استفاده می‌کند که phrase های مهم را  
دربرمی‌گیرد.

یک آپشن max feature دارد که تعداد کلمه های مهم ما را  
محدود می‌کند.

بعد این که مدل خود را فیت کردیم با ما f1 score پایین را داد:

F1 Score = 0.7105



[Back to Agenda Page](#)

# NLP

## بخش سوم مدل تحلیل احساس

```
pipeline4 = Pipeline([  
    'count_vec',  
    CountVectorizer(max_features=20000,  
                    ngram_range=(1, 3)),  
    'clf',  
    LogisticRegression(multi_class='multinomial',  
                       solver='lbfgs', max_iter=5000))  
])
```

### مدل سازی :

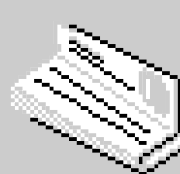
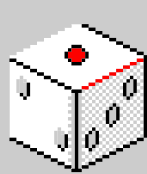
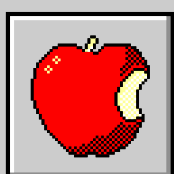
مدل 2:

در مدل دوم به جای `td idf` از `count vectorzise` استفاده کرده ایم.

این مدل تکست ها را تبدیل به یک فیچر عددی می کند که با استفاده شمردن تعداد `frequency` هر کلمه است.

بعد این که مدل خود را فیت کردیم با `f1 score` پایین را داد:

`F1 Score = 0.7014`



[Back to Agenda Page](#)



# NLP

## بخش سوم مدل تحلیل احساس

### Example:

Let's say we have two short documents:

- Document 1: "The cat sat on the mat."
- Document 2: "The cat lay on the rug."

Using CountVectorizer, we get a matrix like this:

	the	cat	sat	on	mat	lay	rug
Doc 1	2	1	1	1	1	0	0
Doc 2	2	1	0	1	0	1	1

Each cell represents how many times each word appears in each document.

### مدل CountVectorizer:

دیتا رو به یک ماتریس از توکن ها تبدیل می‌کند. با ایجاد یه ماتریس که هر سطر مربوط به یک تکست است و هر ستون مربوط به توکن های ما است.

برای مثلا برای جمله:

The cat sat on the mat

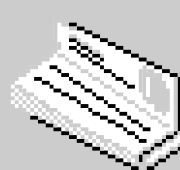
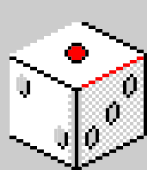
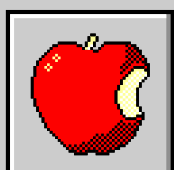
توکن های ما میشود

["the", "cat", "sat", "on", "the", "mat"]

حال به مثال روبرو توجه کنید

میبینیم که ستون های همان کلمات unique در متن هستند و سطر های ما هم تکست های ما.

به طور مثال برای سطر یک برای هر ستون تعداد دفعاتی که کلمه ظاهر شده است را قرار می‌دهد.



[Back to Agenda Page](#)

# NLP

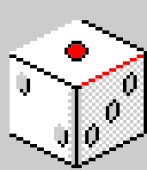
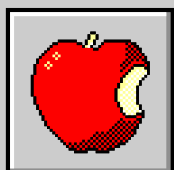
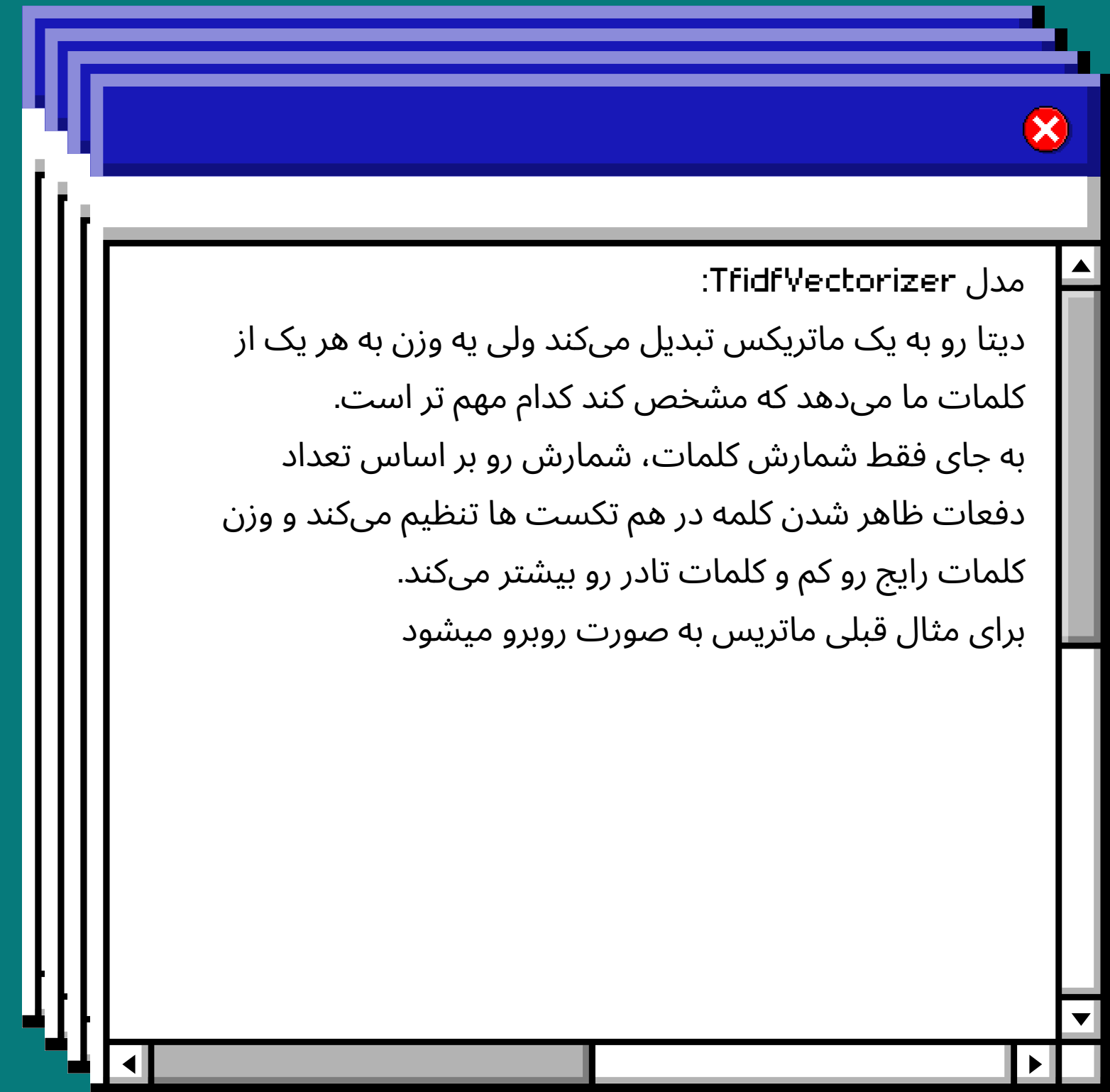
## بخش سوم مدل تحلیل احساس

### Example:

Using the same documents as before:

	the	cat	sat	on	mat	lay	rug
Doc 1	0.3	0.4	0.5	0.3	0.5	0	0
Doc 2	0.3	0.4	0	0.3	0	0.5	0.5

Here, words like "the" and "on" get lower scores because they appear in both documents, making them less unique. Words like "mat" and "rug," which are unique to each document, get higher scores.

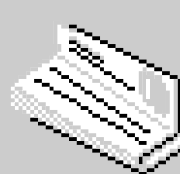
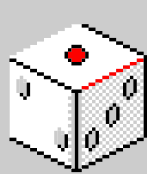
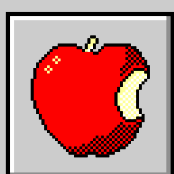


[Back to Agenda Page](#)

# NLP

## بخش سوم مدل تحلیل احساس

```
def sampling(df, num_target,  
            each_class_size):  
    min_class_size = min(each_class_size,  
                          df[num_target].value_counts().min())  
    balanced_df = pd.concat([  
        df[df[num_target] ==  
sentiment].sample(min_class_size,  
replace=True)  
        for sentiment in df[num_target].unique()  
    ])  
    return  
balanced_df.sample(frac=1).reset_index(drop  
=True)
```



# NLP

## بخش سوم مدل تحلیل احساس

```
undersampler =  
RandomUnderSampler(random_state=42)  
X_resampled, y_resampled =  
undersampler.fit_resample(X.to_frame(), y)
```

