

سوال اول

فایل های این سوال در پوشه ی Source Codes با عنوان P1 قرارداد شده است.

بخش ۱

از کتابخانه های ذیل برای انجام امور استفاده نمودیم. (مراحل این بخش در فایل P1_crawling.ipynb قرار داده شده است).

```
from bs4 import BeautifulSoup
import requests
import time
```

در ابتدا اطلاعات مربوط به صفحه ی خودروی ۲۰۶ تیپ دو را به صورت دستی از سایت دیوار بیرون کشیدیم، سپس دریافتیم که در هر صفحه صرفاً ۲۴ مورد فایل وجود دارد و برای لود شدن باقی موارد نیز به اسکرول است، پس کتابخانه های خزگر این قابلیت را نداشتند از این رو API که صدا زده می شد برای قدم بعدی را بررسی نمودیم و با استفاده از آن در ابتدا یک مجموعه دادگانی تحت عنوان Tokens.txt ساختیم که توکن مربوط به هر محصول در دیوار بود که این توکن در انتهای آدرس هر مورد قرار می گرفت و اجازه دسترسی را برای ما بوجود می آورد.

```
url = 'https://divar.ir/v/-/{token}'.format(token=token)
```

حال لازم است تا توکن های استخراج شده (۴۰۰۰ مورد توکن) را داخل لینک فوق قرار داده و اطلاعات مورد نیاز هر خودرو را بیرون بکشیم، از ابتدا می دانستیم که اجرای این دستورات و ارسال درخواست های زیاد به سایت دیوار موجب می شود تا بخشی از داده ها را برای ما نفرستد از این رو در برخی مقاطع جهت اعمال فعالیت ها از sleep استفاده نمودیم.

از ۴۰۰۰ توکنی که داشتیم توانستیم در نهایت اطلاعات ۱۰۰۰ خودرو را با توجه به وقت محدود از سایت دریافت کنیم که این تعداد نیز طی چندین مرحله پیاپی اتفاق افتاد و صرفاً به این تعداد بسنده کردیم چرا که مشخصاً انجام کار به صورت کلی حائز اهمیت است و صرف زمان بیشتر می تواند مشکلات را رفع نماید. (این فرایند در چند مرحله انجام شد که آخرین مرحله ۸۵ دقیقه زمان برد)

تصویری از دیتاست نهایی که ساخته شد در ادامه قرار داده شده است.

	شاسی عقب	تعداد اقساط	مبلغ هر قسط	مدت اعتبار پیش پرداخت	مبلغ جواز	شاسی جلو	مایل به معاوضه	نوع آگهی	قیمت فروش نقدی	گیربکس	مهلت بیمه شخص ثالث	وضعیت بدنه	وضعیت شاسی ها	وضعیت موتور	نوع سوخت	رنگ	مدل (سال تولید)	کارکرد	توکن
0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	۲۴۰,۰۰۰,۰۰۰ تومان	دنده ای	۶ ماه	خط و خش جزیی	سالم و پلمپ	سالم	بنزینی	نقره ای	۱۳۷۷	۲۴۰,۰۰۰	AZd00Ehe
1	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	۳۳۰,۰۰۰,۰۰۰ تومان	دنده ای	۱۰ ماه	رنگ شدگی	سالم و پلمپ	سالم	بنزینی	سفید	۱۳۹۰	۱۰۵,۰۰۰	AZdKkX0Z
2	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	۲۸۰,۰۰۰,۰۰۰ تومان	دنده ای	۶ ماه	سالم و خش	NaN	سالم	بنزینی	نوکمدادی	۱۴۰۰	۶۸,۰۰۰	AZdG0oTg
3	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	۲۴۲,۰۰۰,۰۰۰ تومان	دنده ای	۳ ماه	سالم و بی خط و خش	سالم و پلمپ	سالم	بنزینی	سفید	۱۳۹۹	۱۳۸,۰۰۰	AZ0IEkbl
4	NaN	NaN	NaN	NaN	NaN	NaN	NaN	فروشی	۲۶۹,۰۰۰,۰۰۰ تومان	دنده ای	۵ ماه	سالم و بی خط و خش	سالم و پلمپ	سالم	بنزینی	سفید	۱۴۰۱	۲۴,۰۰۰	AZROGMH3

بخش ۲

ما در بخش قبلی ۱۰۰۰ توکن از ۴۰۰۰ توکنی که داشتیم را اطلاعات اساسی اش را بررسی کردیم و در یک فایل با عنوان output.csv ذخیره نمودیم، برای اضافه کردن اطلاعات هر توکن مانند تصویر و توضیحات نیز قصد داریم تا بخشی از داده ها را استفاده کنیم چرا که کار زمان بر و با حجم بالایی می شود و از طرفی از سمت سایت دیوار محدودیت بیشتری قطعاً اعمال می گردد.

(مراحل این بخش در فایل P1_2.ipynb قرار داده شده است.)

در قدم اول اقدام به تمیز کردن فایل دیتاست نمودیم سپس دو قطعه کد توسعه داده شد تا توضیحات فارسی مربوط به اطلاعات هر ماشین و همچنین تصویر اصلی هر مورد به دیتاست افزوده شود.

دیتاست نهایی :

توضیحات	قیمت فروش نقدی	گیرکس	مهلت بیمه شخص ثالث	وضعیت بدنه	وضعیت شاسی ها	وضعیت موتور	نوع سوخت	رنگ	مدل (سال تولید)	کارکرد	توکن	تصویر
b ^{۰۰}	۲۴۰,۰۰۰,۰۰۰ تومان	دنده ای	۱۲ ماه	رنگ شنگی	سالم و پلمپ	سالم	بنزینی	نقره ای	۱۳۷۷	۱۴۷,۰۰۰	AZa66Zw	575
b ^{۰۰}	۲۶۰,۰۰۰,۰۰۰ تومان	دنده ای	۲ ماه	خط و خش جزئی	سالم و پلمپ	سالم	بنزینی	سفید صدفی	۱۳۷۷	۲۰۰,۰۰۰	AZbODON2	410
b ^{۰۰}	تومان ۱۱,۲۱,۱۱۱	دنده ای	۸ ماه	سالم و بی خط و خش	سالم و پلمپ	سالم	بنزینی	سفید	۱۴۰۱	۱۶۰۰	AZbSO_0P	330
b ^{۰۰}	۳۷۰,۰۰۰,۰۰۰ تومان	دنده ای	۸ ماه	رنگ شنگی	سالم و پلمپ	سالم	بنزینی	سفید صدفی	۱۳۹۸	۶۰,۰۰۰	AZayprhc	593
b ^{۰۰}	۲۴۰,۰۰۰,۰۰۰ تومان	دنده ای	۶ ماه	سالم و بی خط و خش	سالم و پلمپ	سالم	بنزینی	خاکستری	۱۳۹۸	۵۰,۰۰۰	AZZGpbta	795

یک نمونه تصویر خروجی از هر مورد :



تصویر ارائه شده مربوط به اولین ردیف از دیتاست نهایی می باشد که در بخش فوقانی قرار داده شده است.

بخش ۳

اقدامات زیادی روی داده ها در این مرحله انجام شد، به طوریکه تمامی مراحل پیش پردازش پیاده سازی شد.

- داده هایی که صفر بودند پر شدند.
- اعداد فارسی به اعداد انگلیسی تغییر کرد.
- پسوند های اضافه مانند تومان و ماه حذف شد
- داده های عددی به عدد تبدیل شدند.
- سر ستون ها تغییر نام داده شدند.
- ستون های بدون استفاده مانند توکن و نوع سوخت حذف شدند.
- عملیات لیبیل زدن داده های غیر عددی انجام شد و لیبیل ها نیز در فایل ذخیره شد.
- استاندارد سازی انجام شد و میانگین ها و انحراف از معیار ها در فایل ذخیره شد.
- مدل رگرسیون خطی ساخته شد و مدل نیز ذخیره شد.

	Usage	Year	Color	Engine	Chassis	Body	Insurance	Gearbox	Price
0	240000.0	1387	17	1	1	2	6	1	240.0
1	105000.0	1390	10	1	1	4	10	1	303.0
2	6800.0	1400	18	1	1	7	6	1	480.0
3	38800.0	1399	10	1	1	7	3	1	432.0
4	24000.0	1401	10	1	1	7	5	1	469.0
5	16000.0	1400	5	1	1	7	10	1	455.0
6	228000.0	1387	10	1	1	7	8	1	260.0
7	117000.0	1396	10	1	1	2	3	1	365.0
8	252000.0	1389	17	1	1	4	6	1	265.0
9	107000.0	1385	3	1	1	7	9	1	285.0

(مراحل این بخش و بخش ۴ در فایل P1_3.ipynb قرار داده شده است.)

بخش ۴

خروجی نهایی معیار های مدل:

R-squared: 0.6205064152622783
 Mean Squared Error: 4759.280698979096
 Mean Absolute Error: 38.36179018644281

خروجی برنامه :

price prediction : [226.80667661]

قیمت و نظر داخل دیوار :

۲۷۰ میلیون تومان

