

Определение стоимости автомобилей (учебный проект)

Сервис по продаже автомобилей с пробегом «Не бит, не крашен» разрабатывает приложение для привлечения новых клиентов. В нём можно быстро узнать рыночную стоимость своего автомобиля. В нашем распоряжении исторические данные: технические характеристики, комплектации и цены автомобилей. Необходимо построить модель для определения стоимости.

Заказчику важны:

- качество предсказания;
- скорость предсказания;
- время обучения.

Для решения задачи заказчик предоставил датасет:

Признаки:

DateCrawled — дата скачивания анкеты из базы

VehicleType — тип автомобильного кузова

RegistrationYear — год регистрации автомобиля

Gearbox — тип коробки передач

Power — мощность (л. с.)

Model — модель автомобиля

Kilometer — пробег (км)

RegistrationMonth — месяц регистрации автомобиля

FuelType — тип топлива

Brand — марка автомобиля

NotRepaired — была машина в ремонте или нет

DateCreated — дата создания анкеты

NumberOfPictures — количество фотографий автомобиля

PostalCode — почтовый индекс владельца анкеты (пользователя)

LastSeen — дата последней активности пользователя

Целевой признак:

Price — цена (евро)

План работы

- Изучение и подготовка данных
- Обучение моделей
- Анализ моделей
- Тестирование лучшей модели

Подготовка данных

In [89]:

```
1 import pandas as pd
2 import numpy as np
3 import random
4 import seaborn as sea
5 import time
6
7 from sklearn.preprocessing import OrdinalEncoder, StandardScaler
8 from sklearn.model_selection import train_test_split, cross_val_score, GridSearchCV
9
10 from sklearn.linear_model import LinearRegression
11 from sklearn.metrics import mean_squared_error, make_scorer
12 from sklearn.ensemble import RandomForestRegressor
13 from lightgbm import LGBMRegressor
```

In [90]:

```
1 data_original = pd.read_csv('/datasets/autos.csv')
2 data = data_original
```

In [91]:

```
1 data.columns = ['date_crawled', 'price', 'vehicle_type', 'registration_year', 'gearbox', 'power', 'model', 'km', 'reg_month', 'fuel_type', 'brand', 'not_repaired', 'date_created', 'pictures', 'postal_code', 'last_seen']
```

In [92]:

```
1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 354369 entries, 0 to 354368
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  -
0   date_crawled          354369 non-null object
1   price                 354369 non-null int64
2   vehicle_type          316879 non-null object
3   registration_year     354369 non-null int64
4   gearbox               334536 non-null object
5   power                 354369 non-null int64
6   model                 334664 non-null object
7   km                    354369 non-null int64
8   reg_month             354369 non-null int64
9   fuel_type             321474 non-null object
10  brand                 354369 non-null object
11  not_repaired          283215 non-null object
12  date_created          354369 non-null object
13  pictures              354369 non-null int64
14  postal_code           354369 non-null int64
15  last_seen             354369 non-null object
dtypes: int64(7), object(9)
memory usage: 43.3+ MB
```

price

У 10772 цена 0, чаще встречается цена от 500 до сих 6 тысяч

Строки, где нет цены, удалю, так как это целевой признак

Есть цены 1 евро и др <10 евро. Или они ошибочны, или единица измерения - тысяча евро? Хотелось бы эту информацию уточнить у заказчика

In [93]:

```
1 data = data.drop(index = data[data.price==0].index, axis=0)
2 data.price.isna().sum()
```

Out[93]:

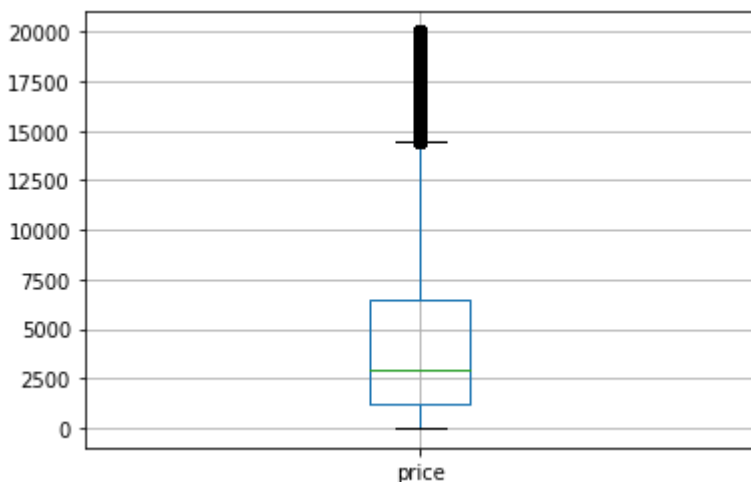
0

In [94]:

```
1 data[['price']].boxplot()
```

Out[94]:

<AxesSubplot:>



vehicle_type

Тип кузова: нет данных - 37490, other - 3288. Я предполагала, что для ряда моделей можно взять данные из аналогичных моделей с заполненным типом кузова, но у одной и той же модели может быть разный тип кузова и этот вариант не подошел. Поэтому выделим их отдельной группой

Выделяю пустые значения в отдельную группу "нет данных"

In [95]:

```
1 data.loc[data.vehicle_type.isna(), 'vehicle_type'] = 'no_data'
```

Отделю тестовую выборку

In [96]:

```
1 data.loc[data['price']<=1200, 'price_class'] = 0
2 data.loc[(data['price']>1200)&(data['price']<=2890), 'price_class'] = 1
3 data.loc[(data['price']>2890)&(data['price']<=6500), 'price_class'] = 2
4 data.loc[(data['price']>6500)&(data['price']<=20000), 'price_class'] = 3
5
6 train, test = train_test_split(
7     data,
8     test_size=0.25,
9     random_state=12345,
10    stratify=data['price_class']
11 )
12 print(train.shape, test.shape)
13
```

```
(257697, 17) (85900, 17)
```

registration_year

Есть неадекватные года - больше даты скачивания и меньше 1768. [статья на вики](#)

([https://translated.turbopages.org/proxy_u/en-ru.ru.a2922144-63446787-7d62ef71-](https://translated.turbopages.org/proxy_u/en-ru.ru.a2922144-63446787-7d62ef71-74722d776562/https/en.wikipedia.org/wiki/Vehicle_registration_plates_of_the_Soviet_Union#1931_to_1946)

[74722d776562/https/en.wikipedia.org/wiki/Vehicle_registration_plates_of_the_Soviet_Union#1931_to_1946](https://translated.turbopages.org/proxy_u/en-ru.ru.a2922144-63446787-7d62ef71-74722d776562/https/en.wikipedia.org/wiki/Vehicle_registration_plates_of_the_Soviet_Union#1931_to_1946)).

Думаю, началом учета могу условно считать 1931 год. Стандартизация учета произошла в СССР в 1975 году и в России в 1993 году, но как я понимаю и ранее учет был. При этом не исключаю, что в этом поле указывают год производства по техпаспорту. Автомобили начали производить в 1768 году. По этому вопросу я бы уточнила информацию у заказчика. Также данные можно проверить по дате выпуска модели. так у форда мондео стоит год регистрации 1000, хотя модель произведена в 1993. Всего таких 166 строк. 82 строк в периоде с 1768 по 1933. Их хочу проверить, действительно ли продаются древние модели

1. Моделям с некорректным годом регистрации установлю год случайно из диапазона по годам выпуска модели (меньше 1768 и больше текущего). Предполагаю, что возможно база выгружалась в 2016-м году, но возможно обновлялась после выгрузки? Так как не могу это проверить, использую текущий год.
2. Проверив вручную дату выхода моделей, у которых в данных стоит год до 1933 сделала вывод, что все эти модели вышли позже, а значит данные ошибочны. Заполню их также случайными значениями по годам выпуска
3. После выполнения первых двух пунктов осталась 41 строка с некорректным годом. Часть из них - sonstige_autos (прочие бренды в переводе с немецкого) и еще несколько строк без указания модели. Эти строки удалю, так как данные по наименованию модели определяющие для формирования цены.

In [97]:

```
1 data[(data.registration_year<=1933)|(data.registration_year>2022)][['brand','mod
```

Out[97]:

brand	model	
volkswagen	golf	17
opel	corisa	7
ford	other	6
bmw	3er	6
	other	5
opel	other	4
volkswagen	polo	4
opel	zafira	4
volkswagen	kaefer	4
fiat	other	4
mercedes_benz	other	3
peugeot	other	3
ford	mondeo	3
alfa_romeo	156	2
mercedes_benz	a_klasse	2
porsche	911	2
ford	ka	2

In [98]:

```
1  #данные выхода моделей брала из википедии
2  years = [['volkswagen', 'golf', 1974, 2022],
3           ['alfa_romeo', '156', 1997, 2007],
4           ['opel', 'corsa', 1982, 2022],
5           ['bmw', '3er', 1990, 1998],
6           ['opel', 'zafira', 1999, 2022],
7           ['volkswagen', 'kaefer', 1946, 2003],
8           ['volkswagen', 'polo', 1975, 2022],
9           ['mercedes_benz', 'a_klasse', 1997, 2022],
10          ['bmw', '5er', 2003, 2010],
11          ['volkswagen', 'passat', 1973, 2022],
12          ['subaru', 'impreza', 1992, 2022],
13          ['renault', 'twingo', 1993, 2022],
14          ['renault', 'laguna', 1994, 2015],
15          ['renault', 'clio', 1990, 2022],
16          ['porsche', '911', 1963, 2022],
17          ['skoda', 'octavia', 1996, 2022],
18          ['bmw', '1er', 2004, 2022],
19          #['opel', 'kadett', 1937-1940, 1962-1993]
20          ['mitsubishi', 'colt', 1962, 2012],
21          ['mercedes_benz', 'clk', 1996, 2009],
22          ['mazda', '6_reihe', 2002, 2022],
23          ['ford', 'fiesta', 1976, 2022],
24          ['ford', 'escort', 1968, 2003],
25          ['fiat', 'panda', 1980, 2022],
26          ['citroen', 'c4', 2004, 2022],
27          ['bmw', 'x_reihe', 1999, 2022],
28          ['volkswagen', 'transporter', 1949, 2022],
29          ['ford', 'ka', 1996, 2016],
30          ['ford', 'mondeo', 1992, 2022],
31          ['citroen', 'c3', 2002, 2022],
32          ['volkswagen', 'caddy', 1980, 2022],
33          ['volkswagen', 'beetle', 1998, 2010],
34          ['renault', 'espace', 1984, 2022],
35          ['audi', 'a2', 1999, 2005],
36          ['opel', 'agila', 2000, 2015],
37          ['opel', 'calibra', 1989, 1997],
38          ['fiat', 'punto', 1993, 2018]]
39
40  def registration(data):
41      #отдельно для опеля, так как был перерыв в годах выпуска
42      data.loc[(data.registration_year<=1933) | (data.registration_year>2022)
43              &(data.brand=='opel')
44              &(data.model=='kadett'),
45              'registration_year']=random.choice(list(range(1937, 1940+1))+list(ra
46
47      for one in years:
48          data.loc[(data.registration_year<=1933) | (data.registration_year>2022)
49                  &(data.brand==one[0])
50                  &(data.model==one[1]),
51                  'registration_year']=random.randint(one[2], one[3]+1)
52      data = data.drop(index = data[(data.registration_year<=1933) | (data.registrat
53      print(data[(data.registration_year<=1933) | (data.registration_year>2022)])
54      return data
55
56  test = registration(test)
57  train = registration(train)
58
```

```
59  
Empty DataFrame
```

```
Columns: [date_crawled, price, vehicle_type, registration_year, gearbo  
x, power, model, km, reg_month, fuel_type, brand, not_repaired, date_c  
reated, pictures, postal_code, last_seen, price_class]
```

```
Index: []
```

```
Empty DataFrame
```

```
Columns: [date_crawled, price, vehicle_type, registration_year, gearbo  
x, power, model, km, reg_month, fuel_type, brand, not_repaired, date_c  
reated, pictures, postal_code, last_seen, price_class]
```

```
Index: []
```

gearbox

Автомат или ручная. 19833 не заполнены. Модели, выпущенные до 1940 года - МКПП. [ссылка на вики \(https://ru.wikipedia.org/wiki/Автоматическая_коробка_передач#История\)](https://ru.wikipedia.org/wiki/Автоматическая_коробка_передач#История).

Для моделей до 1940 года заполняю manual. Остальные пропуски объединяю в отдельную группу

```
In [99]:
```

```
1 def gearbox(data):  
2     data.loc[data.registration_year<=1940, 'gearbox'] = 'manual'  
3     data.loc[data.gearbox.isna(), 'gearbox'] = 'no_data'  
4     return data  
5 train = gearbox(train)  
6 test = gearbox(test)
```

power

40225 нулей. Самая слабая машина на рынке имеет мощность 5 и 17 лс, однако в данных есть 83 машины с мощностью меньше 5. Я погуглила несколько моделей - данные некорректны. Данные меньше 5 удаляю.

Все пустые данные заполню -999

Для моделей с мощностью более 1600 установлю медианное значение мощности.

```
In [100]:
```

```
1 def power(data):  
2     data.power = data.power.drop(index = data[(data.power<5) & (data.power>0)].index)  
3     data.loc[data.power.isna(), 'power'] = -999  
4     data.loc[data.power>1600, 'power'] = data.power.median()  
5     return data  
6 test = power(test)  
7 train = power(train)  
8  
9
```

Данные по самым мощным автомобилям по годам выпуска взяты из [таблицы хронологии \(https://translated.turbopages.org/proxy_u/en-ru.ru.05545da8-63465441-a8a9547e-74722d776562/https/en.wikipedia.org/wiki/Highest_horsepower_engine\)](https://translated.turbopages.org/proxy_u/en-ru.ru.05545da8-63465441-a8a9547e-74722d776562/https/en.wikipedia.org/wiki/Highest_horsepower_engine). Для тех моделей, где мощность больше, чем выпускаемая в этот год самая мощная модель, устанавливаю максимально возможное для этого года значение

In [101]:

```
1 horses = [[1894, 1.5],
2 [1897, 8],
3 [1899, 23],
4 [1901, 35],
5 [1902, 45],
6 [1903, 60],
7 [1907, 75],
8 [1908, 91],
9 [1910, 121],
10 [1912, 200],
11 [1928, 269],
12 [1932, 324],
13 [1935, 406],
14 [1958, 405],
15 [1963, 431],
16 [1965, 492],
17 [1991, 560],
18 [1992, 627],
19 [1995, 680],
20 [2004, 817],
21 [2005, 1001],
22 [2009, 1305],
23 [2014, 1360],
24 [2015, 1521],
25 [2021, 1600]
26 ]
27 horses = pd.DataFrame(horses, columns = ['year', 'power']).sort_values('power', as
28
29
```

Данные для моделей с незаполненной мощностью заполню медианой по бренду

In [102]:

```
1 def horses_f(data):
2     for one in range(len(horses)):
3         data.loc[(data.registration_year == horses.loc[one, 'year']) &
4                 (data.power == horses.loc[one, 'power']), 'power'] = horses['power']
5
6     for one in data.brand.unique():
7         data.loc[(data.power == 0) & (data.brand == one), 'power'] = data[data.brand == one].power.median()
8     return data
9 train = horses_f(train)
10 test = horses_f(test)
```

Проверяю данные >500 лс и сверяю их с данными с сайта drom.ru. Если значение не входит в диапазон модели, то заполняю случайным числом из диапазона.

In [103]:

```
1 horses2 = [['volkswagen', 'golf', 50, 310],
2 ['bmw', '3er', 75, 387],
3 ['opel', 'astra', 60, 200],
4 ['bmw', 'm_reihe', 195, 560],
5 ['mercedes_benz', 'e_klasse', 72, 612],
6 ['volkswagen', 'passat', 68, 300],
7 ['volkswagen', 'polo', 40, 220],
8 ['opel', 'corssa', 58, 210],
9 ['porsche', 'cayenne', 240, 570],
10 ['renault', 'twingo', 55, 133],
11 ['opel', 'vectra', 57, 280],
12 ['mercedes_benz', 's_klasse', 112, 630],
13 ['mercedes_benz', 'a_klasse', 60, 421],
14 ['renault', 'scenic', 64, 163],
15 ['chevrolet', 'matiz', 51, 82],
16 ['fiat', '500', 69, 135],
17 ['audi', 'a4', 75, 286],
18 ['audi', 'a6', 90, 350],
19 ['nissan', 'micra', 50, 110],
20 ['jaguar', 's_type', 200, 395],
21 ['volkswagen', 'transporter', 84, 204],
22 ['bmw', '5er', 86, 530],
23 ['ford', 'mondeo', 88, 240],
24 ['ford', 'ka', 50, 136],
25 ['ford', 'fiesta', 40, 200],
26 ['peugeot', '3_reihe', 45, 225],
27 ['renault', 'espace', 88, 241],
28 ['audi', 'a3', 110, 250],
29 ['opel', 'zafira', 82, 240],
30 ['mercedes_benz', 'm_klasse', 150, 558],
31 ['volvo', 'v40', 102, 249],
32 ['ford', 'galaxy', 90, 240],
33 ['citroen', 'c4', 88, 180],
34 ['fiat', 'punto', 60, 135],
35 ['ford', 'focus', 75, 350],
36 ['volkswagen', 'lupo', 50, 125],
37 ['honda', 'civic', 45, 205],
38 ['volkswagen', 'touran', 90, 190],
39 ['seat', 'arosa', 50, 100],
40 ['lancia', 'ypsilon', 54, 95],
41 ['audi', 'a8', 150, 571],
42 ['mercedes_benz', 'c_klasse', 75, 517],
43 ['citroen', 'c2', 60, 122],
44 ['skoda', 'fabia', 50, 180],
45 ['seat', 'leon', 85, 265],
46 ['volkswagen', 'sharan', 90, 220],
47 ['toyota', 'corolla', 55, 192],
48 ['skoda', 'octavia', 68, 230],
49 ['smart', 'forfour', 60, 177],
50 ['volkswagen', 'kaefer', 25, 50],
51 ['smart', 'fortwo', 41, 109],
52 ['toyota', 'rav', 116, 269],
53 ['toyota', 'yaris', 65, 272],
54 ['audi', '80', 54, 172],
55 ['renault', 'laguna', 83, 241],
56 ['bmw', '1er', 115, 340],
57 ['citroen', 'berlingo', 71, 130],
58 ['dacia', 'sandero', 65, 101],
59 ['daewoo', 'nubira', 106, 136],
```

```
60 ['ford', 'mustang', 85, 760],
61 ['kia', 'rio', 75, 138],
62 ['mazda', '6_reihe', 120, 274],
63 ['mercedes_benz', 'cl', 279, 517],
64 ['mercedes_benz', 'clk', 136, 582],
65 ['mitsubishi', 'colt', 55, 163],
66 ['mitsubishi', 'outlander', 98, 230],
67 ['nissan', 'x_trail', 114, 280],
68 ['peugeot', '2_reihe', 100, 170],
69 ['renault', 'kangoo', 55, 115],
70 ['mini', 'cooper', 98, 130]]
71
```

In [104]:

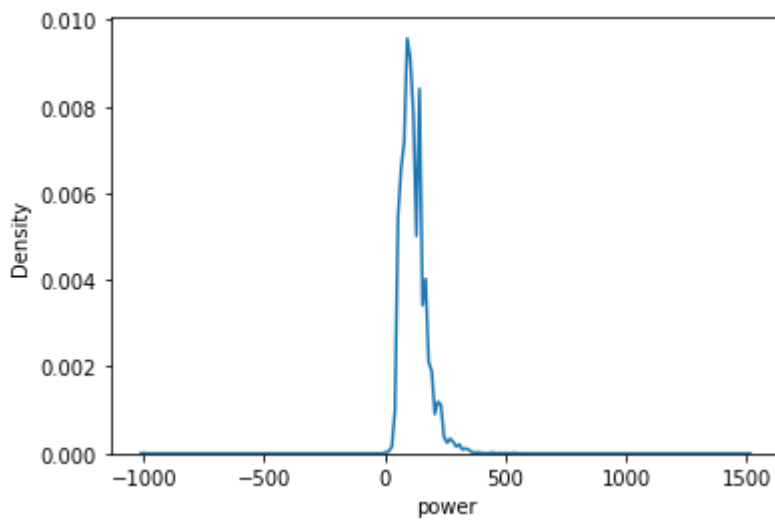
```
1 def horses_f2(data):
2     for one in horses2:
3         data.loc[(data.brand==one[0])
4                 &(data.model==one[1])
5                 &((data.power>one[3])|(data.power<one[2]))
6                 , 'power'] = random.randint(one[2], one[3]+1)
7     return data
8 train = horses_f2(train)
9 test = horses_f2(test)
```

In [105]:

```
1 sea.kdeplot(train['power'])
```

Out[105]:

<AxesSubplot:xlabel='power', ylabel='Density'>

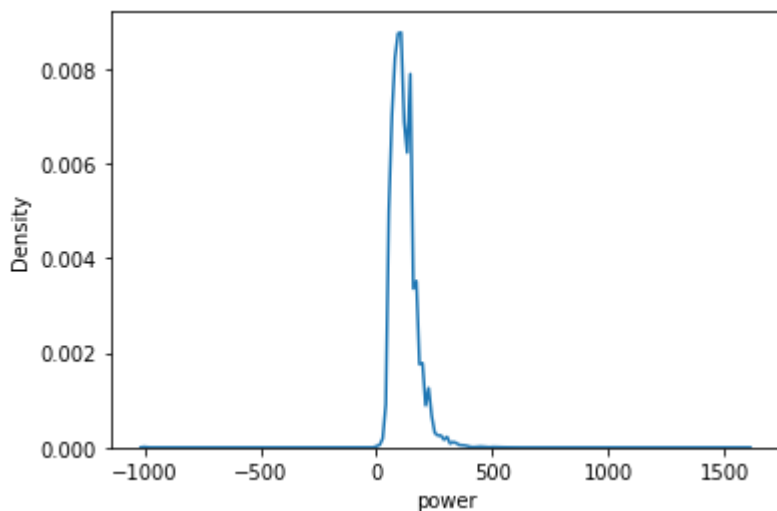


In [106]:

```
1 sea.kdeplot(test['power'])
```

Out[106]:

<AxesSubplot:xlabel='power', ylabel='Density'>



✓ **Успех:** Молодец 👍

model

Больше всего golf, other - 24421, пустые значения 19705. Им можно присвоить other. Некоторые модели есть у нескольких брендов (1_reihe). rangerover и range_rover - похожие названия модели, но выпускаются разными брендами rover и landrover

In [107]:

```
1 def model(data):
2     data.loc[data.model.isna(), 'model'] = 'other'
3     data.model.isna().sum()
4     return data
5 train = model(train)
6 test = model(test)
```

fuel_type

petrol (англ), gasoline (амер)- синонимы, lpg - пропан-бутан, cng - метан. hybrid - гибридные, electric - электрические, other - 204, нет данных - 32895.

Petrol и gasoline объединю. Пустые заполню как other

In [108]:

```
1 def fuel_type(data):
2     data.loc[data.fuel_type == 'petrol', 'fuel_type'] = 'gasoline'
3     data.loc[data.fuel_type.isna(), 'fuel_type'] = 'other'
4     return data
5 train = fuel_type(train)
6 test = fuel_type(test)
```

not_repaired

not_repaired - 71154 нет данных, yes/no. Нет данных - сделаю отдельной группой

In [109]:

```
1 def not_repaired(data):
2     data.loc[data.not_repaired.isna(), 'not_repaired'] = 'no_data'
3     data[data.not_repaired.isna()]
4     return data
5 train = not_repaired(train)
6 test = not_repaired(test)
```

last_seen

Довольно много строк, где год регистрации автомобиля больше, чем год последнего просмотра.
Возможно их стоит удалить

In [110]:

```
1 def last_seen(data):
2     data = data.drop(data[data.registration_year > 2022].index, axis = 0)
3     data.registration_year = pd.to_datetime(data.registration_year, format='%Y')
4     data = data.drop(data[data.registration_year > data.last_seen].index, axis = 0)
5     data.registration_year = data.registration_year.astype('int')
6
7     return data
8 train = last_seen(train)
9 test = last_seen(test)
```

In [111]:

```
1 train.power = train.power.astype('int')
2 test.power = test.power.astype('int')
```

Не требуют обработки

km - категориальный признак. Пустых значений нет. Больше всего с пробегом 150 тыс.

brand - пропусков и других проблем нет **postal_code** - пропусков нет, важен, так как по сути это регион продажи. По индексу видно, что это даже разные страны

Не нужны

date_crawled Дата скачивания анкеты из базы - если обнаружится какая-то ошибка, то будем изучать,

для модели не нужен. Анкеты скачаны в апреле и марте 2016 года, есть точное время скачивания

reg_month - не понятно, зачем нужен этот признаи. Нулей - 37352. Пропусков нет. Больше всего регистрировались в марте, июне и апреле. Меньше всего в феврале, августе и январе.

date_created - всего 1 дата за 2014 год, 18 дат за 2015, остальные - за 2016 (89 уникальных дат 2016-го года). Пустых нет.

pictures - 354369 с количеством 0, пропусков нет, столбец можно удалить

last_seen - пропусков нет, все пользователи были активны в 2016

In [112]:

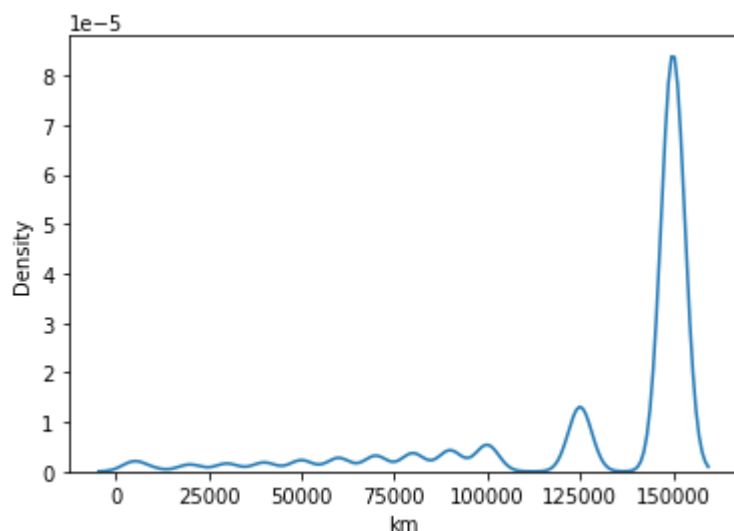
```
1 def drop(data):
2     data = data.drop(['date_crawled', 'reg_month', 'date_created', 'pictures', 'last_seen'])
3     data = data.drop_duplicates()
4     return data
5 train = drop(train)
6 test = drop(test)
```

In [113]:

```
1 sea.kdeplot(train['km'])
2
```

Out[113]:

<AxesSubplot:xlabel='km', ylabel='Density'>



In [114]:

```
1 train[train['power'] > 1500]
```

Out[114]:

	price	vehicle_type	registration_year	gearbox	power	model	km	fuel_type	brand
145154	1200	sedan	946684800000000000	manual	1503	other	150000	gasoline	bmw

In [115]:

```
1 test[test['power']>1500]
```

Out[115]:

	price	vehicle_type	registration_year	gearbox	power	model	km	fuel_type	bra
146238	1500	sedan	820454400000000000	manual	1596	other	150000	gasoline	br

Кодирование признаков

In [116]:

```
1 target_train = train['price']
2 features_train = train.drop(['price', 'price_class'], axis=1)
3 target_test = test['price']
4 features_test = test.drop(['price', 'price_class'], axis=1)
5
6 print(target_train.shape, features_train.shape, target_test.shape, features_test
```

(234310,) (234310, 10) (80779,) (80779, 10)

In [117]:

```
1 encoder = OrdinalEncoder(handle_unknown = 88888)
2 encoder.fit(features_train)
3 features_train_ordinal = pd.DataFrame(encoder.transform(features_train), columns =
4 features_test_ordinal = pd.DataFrame(encoder.transform(features_test), columns =
5 features_test_ordinal
```

Out[117]:

	vehicle_type	registration_year	gearbox	power	model	km	fuel_type	brand	not_repaired	postal_code
0	6.0	73.0	1.0	101.0	173.0	6.0	2.0	38.0	0.0	3162.0
1	6.0	63.0	1.0	86.0	33.0	12.0	2.0	20.0	0.0	796.0
2	6.0	69.0	1.0	56.0	83.0	9.0	2.0	24.0	2.0	3799.0
3	5.0	71.0	0.0	146.0	180.0	12.0	2.0	23.0	0.0	6118.0
4	5.0	61.0	1.0	86.0	59.0	12.0	2.0	20.0	0.0	7745.0
...
80774	4.0	72.0	1.0	64.0	166.0	5.0	2.0	23.0	0.0	4505.0
80775	2.0	69.0	1.0	260.0	248.0	11.0	2.0	2.0	0.0	4841.0
80776	6.0	68.0	1.0	61.0	245.0	11.0	2.0	36.0	0.0	3211.0
80777	8.0	66.0	1.0	105.0	166.0	12.0	2.0	25.0	0.0	5710.0

In [118]:

```
1 features_train_ohe = pd.get_dummies(features_train, drop_first = True)
2 features_test_ohe = pd.get_dummies(features_test, drop_first = True)
3 #for one in features_train_ohe.columns[~features_train_ohe.columns.isin(features_test_ohe.columns)]:
4 #     features_test_ohe[one] = 0
5 for one in features_train_ohe.columns:
6     if one not in features_test_ohe.columns:
7         features_test_ohe[one] = 0
8 for one in features_test_ohe.columns:
9     if one not in features_train_ohe.columns:
10         features_train_ohe[one] = 0
11 features_train_ohe.shape, features_test_ohe.shape
12
```

Out[118]:

```
((234310, 309), (80779, 309))
```

Корреляция признаков

In [119]:

```
1 features_train.corr()
```

Out[119]:

	registration_year	power	km	postal_code
registration_year	1.000000	0.108428	-0.214832	0.035511
power	0.108428	1.000000	0.115379	0.057420
km	-0.214832	0.115379	1.000000	-0.013770
postal_code	0.035511	0.057420	-0.013770	1.000000

Ярко выраженной линейной зависимости между признаками нет

Обучение моделей

Константная модель

In [120]:

```
1 def rmse(target, prediction):
2     rmse = mean_squared_error(target, prediction, squared = False)
3     return rmse
4
```

In [121]:

```
1 prediction = [data.price.median()]*len(data)
2 rmse(data.price, prediction)
3
```

Out[121]:

4812.3097981057535

Линейная регрессия

In [122]:

```
1 features_train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 234310 entries, 275374 to 39122
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  -
0   vehicle_type          234310 non-null object
1   registration_year      234310 non-null int64
2   gearbox               234310 non-null object
3   power                 234310 non-null int64
4   model                 234310 non-null object
5   km                    234310 non-null int64
6   fuel_type             234310 non-null object
7   brand                 234310 non-null object
8   not_repaired          234310 non-null object
9   postal_code           234310 non-null int64
dtypes: int64(4), object(6)
memory usage: 19.7+ MB
```

In [123]:

```
1 train_rmse = [0,0,0]
2 test_rmse = [0,0,0]
3 fit_time = [0,0,0]
4 p_time = [0,0,0]
5
6 pd.options.mode.chained_assignment = None
7 scaler = StandardScaler()
8 scaler.fit(features_train_ohe[['power','km','postal_code','registration_year']])
9 features_scaled_train = features_train_ohe
10 features_scaled_test = features_test_ohe
11 features_scaled_train[['power','km','postal_code','registration_year']] = scaler.
12 features_scaled_test[['power','km','postal_code','registration_year']] = scaler.
```

In [124]:

```
1 model_l = LinearRegression()
2 scor = make_scorer(rmse, greater_is_better=False)
3 scores = cross_val_score(model_l, features_scaled_train, target_train, cv=5, scor=
4 train_rmse[2] = scores.mean()
5 print(train_rmse[2])
6
7 start_time = time.time()
8 model_l.fit(features_scaled_train, target_train)
9 fit_time[2] = (time.time() - start_time)
10 print("Время обучения %s с" % (fit_time[2]))
11 start_time = time.time()
12 predict_l = model_l.predict(features_scaled_train)
13 p_time[2] = time.time() - start_time
14 print("Время предсказания %s с" % p_time[2])
```

-2755.2868212521785

Время обучения 21.788912057876587 с

Время предсказания 0.40083742141723633 с

LightGBM

In [125]:

```
1
2 trees = [10, 50, 100]
3 best_scores = 999999
4 params = {'n_estimators': range(10, 40, 10),
5           'max_depth': range(10, 40, 10),
6           'num_leaves': range(10, 40, 10)}
7 model_lgbm = LGBMRegressor()
8 grid = GridSearchCV(model_lgbm, params, cv=5, scoring = 'neg_root_mean_square
9
10 start_time = time.time()
11 grid.fit(features_train_ordinal, target_train)
12 fit_time[0] = (time.time() - start_time)
13
14 train_rmse[0] = grid.best_score_
15 print(grid.best_estimator_, train_rmse[0])
16 print("Время обучения %s с" % (fit_time[0]))
17
18 start_time = time.time()
19 predict_lgbm = grid.predict(features_train_ordinal)
20 p_time[0] = time.time() - start_time
21 print("Время предсказания %s с" % p_time[0])
22
```

LGBMRegressor(max_depth=10, n_estimators=30, num_leaves=30) -1975.0848
555962239

Время обучения 782.7031135559082 с

Время предсказания 0.5144298076629639 с

RandomForest

In [126]:

```
1 trees = [10, 50, 100]
2 best_scores = 999999
3 parametr = {'n_estimators': range(10, 40, 20),
4             'max_depth': range(10, 40, 10),
5             'min_samples_split': [0.1, 0.2, 0.3],
6             'min_samples_leaf': range(10, 40, 10)}
7
8 model_rf = RandomForestRegressor()
9 grid_rf = GridSearchCV(model_rf, parametr, cv=5, scoring = 'neg_root_mean_squa
10 start_time = time.time()
11 grid_rf.fit(features_train_ordinal, target_train)
12 fit_time[1]=(time.time() - start_time)
13 train_rmse[1] = grid_rf.best_score_
14 print(grid_rf.best_estimator_, train_rmse[1])
15 print("Время обучения %s с" % (fit_time[1]))
16 start_time = time.time()
17 predict_rf= grid_rf.predict(features_train_ordinal)
18 p_time[1] = time.time() - start_time
19 print("Время предсказания %s с" % p_time[1])
20
```

RandomForestRegressor(max_depth=30, min_samples_leaf=20, min_samples_s
plit=0.1,

n_estimators=10) -2944.008969945214

Время обучения 760.2367017269135 с

Время предсказания 0.08147072792053223 с

Анализ моделей

In [127]:

```
1 work_time = pd.DataFrame(columns = ['name', 'time_predict', 'time_fit', 'train_rms
2 work_time.time_fit = fit_time
3 work_time.time_predict = p_time
4 work_time.name = ['LGBM', 'RandomForest', 'LinearRegression']
5 work_time.train_rmse = np.array(train_rmse).astype('int')
6 work_time
```

Out[127]:

	name	time_predict	time_fit	train_rmse
0	LGBM	0.514430	782.703114	-1975
1	RandomForest	0.081471	760.236702	-2944
2	LinearRegression	0.400837	21.788912	-2755

Лучший результат rmse - у градиентного бустинга. Быстрее работает линейная регрессия.

Тестирование лучшей модели

In [131]:

```
1 start_time = time.time()
2 best_prediction = grid.predict(features_test_ordinal)
3 fit_time_best = (time.time() - start_time)
4 best_rmse = rmse(best_prediction, target_test)
5 print('RMSE лучшей модели', best_rmse)
6 print("Время обучения %s с" % work_time.time_fit[0])
7 print("Время предсказания %s с" % fit_time_best)
8
```

RMSE лучшей модели 1988.5825310269688

Время обучения 782.7031135559082 с

Время предсказания 0.19648075103759766 с

Рекомендую заказчику использовать модель LightGBM, так как она показывает лучшее качество предсказаний и единственная выполняет поставленную задачу по метрике качества. При этом она показывает худший результат по времени. В случае неудовлетворенности результатом предложу продолжить подбор гиперпараметров для RandomForest.

Вывод

In []:

```
1
2 # Не забывай про финальный вывод)
```

В рамках проекта проделана работа по предобработке данных, обучению 3 моделей для предсказания рыночной цены автомобиля. Модель LightGBM работает дольше других моделей, но показывает лучшее качество в текущей задаче регрессии и является рекомендованной к применению.