Итоговая работа

Тема: Использование машинного обучения для прогнозирования отправления железнодорожных вагонов в ремонт.

Введение:

Отправка вагона в плановый ремонт может происходить по разным причинам, как по регламенту(срок/пробег), так и из-за накопления мелких дефектов:

- было много текущих ремонтов;
- не было вариантов на погрузку и т.д.

Этих причин много, и все они влияют на возможность осуществления ремонта.

Вагон отправляют в ремонт после получения уведомления об их неисправности. К сожалению, текущий процесс не позволяет распределять нагрузку на ремонт депо, управлять последней заявкой на погрузку пред ремонтом и многое другое.

Цель: снизить нагрузку на железнодорожную систему путем своевременного обслуживания вагонов.

Задача: создать модель прогнозирования даты отправления вагона в плановый ремонт.

Постановка задачи:

- 1. найти закономерности и оценить значимые признаки;
- 2. спрогнозировать, что вагон отправится в ПР в течение месяца;
- 3. спрогнозировать, что вагон отправится в ПР в течение 10 дней.

Обзор моделей:

Наивная модель построена на правилах с использованием минимального набора данных, без применения машинного обучения.

Реальный процесс выглядит следующим образом - в начале месяца берется срез по парку по всем вагонам, за ремонт которых несёт ответственность ПГК. Для выбранных вагонов требуется установить, какие из них будут отремонтированы в текущем месяце. Данная информация помогает планировать нагрузку на вагоноремонтное предприятие (ВРП). Вторая модель определяет критичные вагоны, которые будут отправлены в ремонт

в первую очередь (в ближайшие 10 дней). Это помогает фокусировать внимание диспетчеров.

Основными критериями, по которым вагон отправляется в плановый ремонт - является его остаточный пробег и срок до планового ремонта.

В регламентах РЖД используется следующее правило - если ресурс по пробегу не превышает 500 км, срок службы не превышает 500 дней, число текущих ремонтов больше 5 и/или плановый ремонт должен наступить через 15 дней (или меньше), то вагон может ехать только на ВРП.

Из этого регламента вытекают две особенности:

- 1. Диспетчер старается отправить вагон раньше положенных значений. Это позволяет выбрать предприятия, на которых ремонтироваться дешевле, а не ближайшее.
- 2. Компания-оператор может выбирать какому из нормативов нужно следовать ремонтировать вагон по сроку, или по пробегу, или по обоим критериям сразу. Поэтому встречаются вагоны, у которых пробег может не отслеживаться.

Вагон может быть отправлен в плановый ремонт и раньше положенного. На это может влиять, например, история грузовых операций и количество текущих(мелких) ремонтов.

Описание процесса решения (модель, тест, тренировка, результат):

Изначально производится импорт необходимых библиотек (pandas, numpy, os) и загрузка данных.

✓ Импорт библиотек

```
[1] from google.colab import drive drive.mount('/content/drive')

Mounted at /content/drive

import os # библиотеки работы с папками и системой

# и другие уже известные нам библиотеки import pandas as pd import numpy as np

from sklearn.neighbors import KNeighborsClassifier from sklearn.model_selection import train_test_split from sklearn.metrics import *
from sklearn.preprocessing import MinMaxScaler
```

Загрузка данных

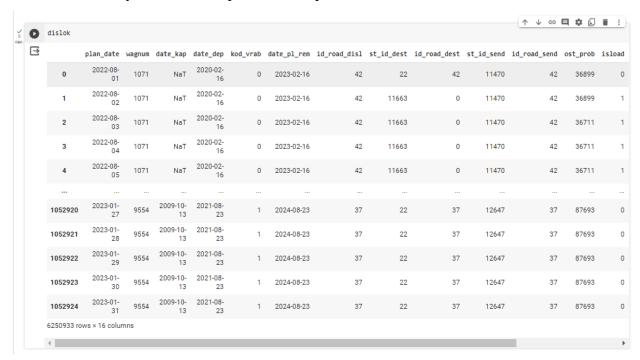
```
| unrar x /content/drive/MyDrive/DataWagon/train_2.rar /content/drive/MyDrive/DataWagon/train_2/
   \supseteq
       UNRAR 6.11 beta 1 freeware
                                       Copyright (c) 1993-2022 Alexander Roshal
       Extracting from /content/drive/MyDrive/DataWagon/train_2.rar
       Would you like to replace the existing file /content/drive/MyDrive/DataWagon/train_2/test/dislok_wagons.parquet
       11392957 bytes, modified on 2023-11-08 12:28
       with a new one
       11392957 bytes, modified on 2023-11-08 12:28
       [Y]es, [N]o, [A]ll, n[E]ver, [R]ename, [Q]uit q
       Program aborted
  [ ] !cp -r /content/train/content/имя_папки /content/drive/MyDrive/DataWagon
       cp: cannot stat '/content/train/content/имя_папки': No such file or directory
  [ ] !ls /content/drive/MyDrive/DataWagon/train
       dislok_wagons.parquet metrics_f1.py stations.parquet tr_rems.parquet freight_info.parquet prediction target wagons_probeg_ownersip.parquet
       freight_info.parquet prediction
       kti_izm.parquet
                           pr_rems.parquet task2_base_model.ipynb wag_params.parquet
5 [5] path = '/content/drive/MyDrive/DataWagon/train'
🟏 🚺 # список вагонов с остаточным пробегом на момент прогноза
        wag_prob = pd.read_parquet(path + '/wagons_probeg_ownersip.parquet').convert_dtypes()
       # данные по дислокации
       dislok = pd.read_parquet(path + '/dislok_wagons.parquet').convert_dtypes()
       wag_param = pd.read_parquet(path + '/wag_params.parquet').convert_dtypes()
       # данные по текущим ремонтам
       pr_rem = pd.read_parquet(path + '/pr_rems.parquet').convert_dtypes()
        # текущие ремонты вагонов
       tr_rem = pd.read_parquet(path + '/tr_rems.parquet').convert_dtypes()
        # данные по КТИ
       kti izm = pd.read parquet(path + '/kti izm.parquet').convert dtypes()
        # справочник грузов
        freight_info = pd.read_parquet(path + '/freight_info.parquet').convert_dtypes()
        # справочник станций
       stations = pd.read_parquet(path + '/stations.parquet').convert_dtypes()
        # таргет по прогноза выбытия вагонов в ПР на месяц и на 10 дней
        target = pd.read_csv('/content/drive/MyDrive/DataWagon/train_2/test/target/y_test.csv').convert_dtypes()
```

Исходные данные содержат 9 таблиц wag_prob, dislok, wag_param, pr_rem, tr_rem, kti_izm, freight_info, stations, target. Таблицы включают в себя следующее:

• Таблица wag_prob. Список вагонов, по которым известен пробег и тип владения на дату среза: Дата среза, Номер вагона, Остаточный пробег, Вид управления по договору (1-собственноые, 2-лизинговые, 11-сервисные), Тип РПС, Состояние в реестре (1-слеженение, 3-вкдючение, 4-готовится к искл.), Тип собственности (1-собственные, 2-принятые в аренду, 4-привелеченный парк, 6-сервис), Месяц.

[7]	wag_prob								
		repdate	wagnum	ost_prob	manage_type	rod_id	reestr_state	ownership_type	month
	0	2022-08-01	33361	7541	0	1	1	0	8
	1	2022-08-02	33361	7243	0	1	1	0	8
	2	2022-08-03	33361	6990	0	1	1	0	8
	3	2022-08-04	33361	6347	0	1	1	0	8
	4	2022-08-05	33361	6027	0	1	1	0	8
	•••							•••	
	9249584	2022-12-20	33350	35062	0	1	0	0	12
	9249585	2022-12-21	33350	35062	0	1	0	0	12
	9249586	2022-12-22	33350	35062	0	1	0	0	12
	9249587	2022-12-23	33350	35062	0	1	0	0	12
	9249588	2022-12-24	33350	35062	0	1	0	0	12
	6249857 ro	ws × 8 colum	ns						

• Таблица dislok. Информация по дислокации: Дата отчета, Номер вагона, Дата последнего капитального ремонта, Дата последнего деповского ремонта, Код предстоящего планового ремонта, Дата предстоящего планового ремонта, Код дороги дислокации, Ид. станции назначения, Код дороги назначения, Индекс станции отправления, Код дороги отправления, Остаточный пробег, Груженый/порожний, Код груза, Ид. последнего груза, Расстояние между станцией отправления и прибытия.



• Таблица wag_param. Данные по характеристикам вагона: Номер вагона, Модель вагона, Тип РПС, Грузоподъемность (в центнерах), Предельная грузоподъёмность, Объем кузова, Масса тары в центнерах, Дата постройки, Дата окончания срока службы, Завод

постройки, Норма пробега после ДР в тыс. км, Норма пробега после КР в тыс. км, Кузов, Код модели тележки, Тип тормозов, Тип воздухораспределителя, Межремонтный норматив пробега (<>0, если вагон на пробеге), Признак передачи вагона в аренду (1-собст, 2-арендованный, 3-инвентарный).

	wagnum	mode1	rod_id	gruz	cnsi_gruz_capacity	cnsi_volumek	tara	date_build	srok_sl	zavod_build	date_iskl	cnsi_probeg_dr	cnsi_probeg_kr
3218	26318	12- 600- 04	1		682	85.0	240	1992-12-25	2022- 04-27	5	2023-02-	160	160
19128	28344	12- 132	1	700	700	88.0	240	2003-08-12	2024- 12-24	0	2022-12- 14	110	160
21526	8099	11- 286	0	670	670	138.0	270	1995-08-31	2027- 10-01	1	NaT	110	160
32353	33350	12- 9850- 02	1	750	750	90.0	248	2014-10-27	2047- 02-05	19	NaT	250	500
81	5308	11- 276	0	680	680	122.0	260	1995-09-17	2027- 09-28	1	NaT	110	160
33703	18766	11- 280	0	680	680	138.0	259	2013-01-07	2046- 03-20	1	NaT	110	160
33704	18769	11- 280	0	680	680	138.0	259	2013-01-18	2046- 02-14	1	NaT	110	160
33705	18899	11- 280	0	680	680	138.0	260	2013-07-24	2044- 03-03	1	NaT	110	160
33706	18912	11- 280	0	680	680	138.0	259	2015-02-28	2044- 08-04	1	NaT	110	160
33707	18914	11- 280	0	680	680	138.0	259	2013-06-16	2046- 07-02	1	NaT	110	160

• Таблица pr_rem. Данные по плановым ремонтам: Номер вагона, Месяц ремонта, Тип РПС, модель, Дорога выбытия в ПР, Дорога прибытия на ВРП, Код ремонта, Станция выгрузки, код станции ВРП, Расстояние со станции до ВРП, Месяц.



• Таблица tr_rem. Данные по текущим ремонтам вагона: Номер вагона, Дата ремонта, Код работ, Код неисправности 1 (Заполняется по с.5353), Код неисправности 2 (Заполняется по с.5353), Код неисправности 3 (Заполняется по с.5353), Код модернизации 1 (Заполняется по с.5354), Код модернизации 2 (Заполняется по с.5354), Код модернизации 3 (Заполняется

по с.5354), Код модернизации 4 (Заполняется по с.5354), Код модернизации 5 (Заполняется по с.5354), Код модернизации 7 (Заполняется по с.5354), Код модернизации 7 (Заполняется по с.5354), ID дороги перевода вагона в неисправные(nsi_db.railway.RW_ID), Пробег вагона в груженом состоянии, Код станции перевода вагона в неисправные (nsi_db.station.ST_CODE).

	wagnum	rem_month	kod_vrab	neis1_kod	neis2_kod	neis3_kod	mod1_kod	mod2_kod	mod3_kod	mod4_kod	mod5_kod	mod6_kod	mod7_kod	road_id_send
0	29938	2022-08- 01	3	0	98	54	7	4	2	0	0	0	0	3
1	29938	2022-08- 01	3	14	98	54	7	4	2	0	0	0	0	3
2	29852	2022-08- 01	2	0	98	54	7	4	2	0	0	0	0	2
3	29852	2022-08- 01	2	36	98	54	7	4	2	0	0	0	0	2
4	13674	2022-08- 01	2	95	98	54	7	4	2	0	0	0	0	2
7695	31868	2023-01- 01	3	14	98	54	7	4	2	0	0	0	0	4
7696	22921	2023-01- 01	5	0	98	54	7	4	2	0	0	0	0	2
7697	22921	2023-01- 01	3	34	95	54	7	4	2	0	0	0	0	2
7698	4978	2023-01- 01	3	131	68	54	7	4	2	0	0	0	0	
7699	16226	2023-01- 01	3	6	98	54	7	4	2	0	0	0	0	4

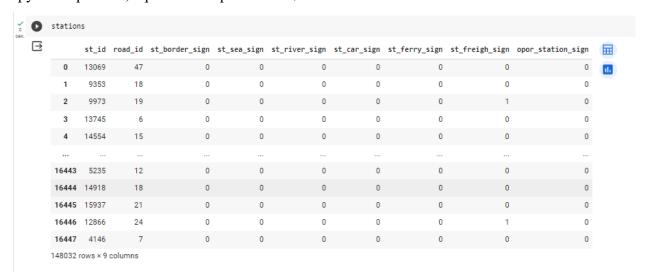
• Таблица kti_izm. Данные по КТИ: Номер вагона, Дата и время измерения КТИ, Пробег общий, Толщина гребня 1 ось слева, Толщина гребня 1 ось справа, Толщина гребня 2 ось слева, Толщина гребня 3 ось слева, Толщина гребня 4 ось справа, Толщина гребня 4 ось справа, Толщина обода 1 ось слева, Толщина обода 1 ось слева, Толщина обода 2 ось слева, Толщина обода 2 ось справа, Толщина обода 2 ось слева, Толщина обода 4 ось справа.

	wagnum	operation_date_dttm	mileage_all	axl1_l_w_flange	axl1_r_w_flange	axl2_l_w_flange	axl2_r_w_flange	ax13_1_w_flange	ax13_r_w_flang
0	325	2022-08-01	112091	30.5	31.0	26.7	27.4	28.8	28.
1	325	2022-08-03	112471	30.4	31.4	27.5	27.6	28.9	29.
2	325	2022-08-05	113938	30.0	31.1	27.1	27.9	28.5	28.
3	325	2022-09-15	121071	30.5	31.4	26.4	27.3	28.9	28.
4	325	2022-09-09	117341	30.7	31.5	26.9	27.1	28.9	28.
194723	28221	2023-01-31	27116	28.4	31.4	29.0	29.2	27.0	28
194752	29023	2023-01-31	1911	29.7	28.6	29.6	30.2	27.9	30
194780	33124	2023-01-31	118504	28.7	29.3	32.8	28.2	29.7	29
194825	28089	2023-01-31	22183	28.4	28.3	29.3	28.6	28.8	28
194859	28090	2023-01-31	158013	28.7	26.6	28.0	28.7	29.8	30

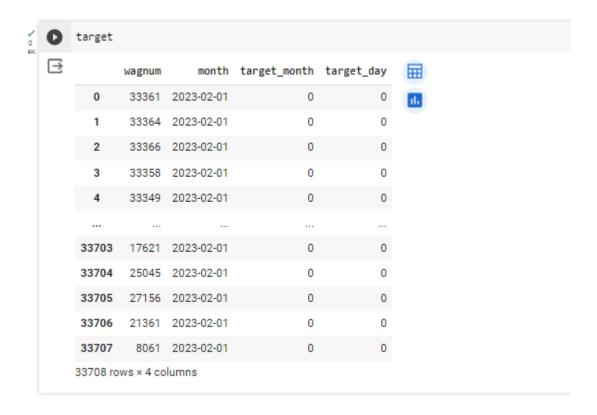
• Таблица freight_info. Справочник грузов: Ид груза, Класс груза, Признак скоропортящегося груза, Признак погрузки навалом, Признак погрузки насыпи, Признак погрузки наливом, Перевозится в открытом подвижном составе, Груз требует сопровождения, Признак смерзающегося груза.

0	freigh	t_info									
∃		fr_id	fr_class	skoroport	naval	nasip	naliv	openvagons	soprovod	smerz	E
	0	4989	2	0	0	0	0	0	0	0	t
	1	4990	2	0	0	0	0	0	0	0	
	2	4991	0	0	0	0	0	0	0	0	
	3	4992	2	0	0	0	0	0	0	0	
	4	4993	2	0	0	0	0	0	0	0	
	5074	989	0	0	1	1	0	0	0	0	
	5075	990	0	0	1	1	0	0	0	0	
	5076	991	0	0	1	1	0	0	0	0	
	5077	992	0	0	0	0	0	0	0	0	
	5078	993	0	0	1	1	0	0	0	0	
	5079 rd	ws×9c	olumns								

• Таблица stations. Справочник станции: Ид станции, Ид дороги, Признак граничной станции, Признак перевалки на море, Признак перевалки на реку, Признак перевалки на автотранспорт, Признак перевалки на паромную переправу, Признак станции, открытой для грузовой работы, Признак опорной станции.



• Таблица target: Номер вагона, Месяц ремонта, Ремонт в течение месяца (1-был ремонт, 0-не было ремонта), Ремонт в первые 10 дней (1-был ремонт, 0-не было ремонта). Данная таблица содержит две целевые переменные, так как стоит задача независимо друг от друга определить выбытие вагона в ремонт в первые 10 дней и в течении месяца.



Подготовка данных:

▼ Подготовка данных

Перед обучением модели необходимо подготовить данные. Для этого:

• Произвелась оценка среднесуточного пробега по пробегу вагона, на тот случай если данных по нормативу нет.

```
[16] target_data= pd.to_datetime('2023-01-01')
pred_data = pd.to_datetime('2023-02-01')

[17] wag_prob = wag_prob[(wag_prob.repdate == pd.to_datetime('2023-01-31')) | (wag_prob.repdate == wag_prob.repdate.min())]

[18] # оценим среднесуточный пробег из данных по пробегу вагона, на тот случай, если данных по нормативу нет
```

• Удалены дубликаты номеров вагонов из всех таблиц.

```
| [19] #yganaem gy6nu homepob barohob us bcex ta6nuq | wag_prob = wag_prob.drop_duplicates(subset='wagnum', keep='last') | dislok = dislok.drop_duplicates(subset='wagnum', keep='last') | wag_param = wag_param.drop_duplicates(subset='wagnum', keep='last') | pr_rem = pr_rem.drop_duplicates(subset='wagnum', keep='last') | kti_izm = kti_izm.drop_duplicates(subset='wagnum', keep='last')
```

• Из таблиц по фильтру отобраны все значения на дату обучения.

```
wag_prob['month'] = wag_prob['repdate'].apply(lambda x: str(pd.to_datetime(x, format='%Y-%m-%d').year) + '-' + str(pd.to_datetime(x, format='\color=xy-\color=xm-\color=xd').mo
dislok['month'] = dislok['plan_date'].apply(lambda x: str(pd.to_datetime(x, format='\color=xy-\color=xm-\color=xd').year) + '-' + str(pd.to_datetime(x, format='\color=xy-\color=xm-\color=xd').mo
```

• Посчитано сколько текущих ремонтов было за прошедший период.

```
у [23] # посчитаем сколько текущих ремонтов было за прошедший период 
скс tr_rem = tr_rem.groupby('wagnum', as_index= False).kod_vrab.count()
```

• Выбраны необходимые колонки из таблиц wag_prod_train, dislok_train, wag_param_pred, pr_rem_train, kti_izm_train, tr_term_train и объединены в один общий датафрейм для обучения модели.

```
| [24] #выбор необходимых колонок из таблиц
| wag_prod_train = wag_prob[wag_prob.month==(str(pd.to_datetime(target_data, format='%Y-%m-%d').year)+'-'+str(pd.to_datetime(target_data, format='%Y-%m-%d').year)+'-'+str(pd.to_datetime(target_data, format='%Y-%m-%d').won)
| wag_param_pred = wag_param[['wagnum', 'gruz', 'cnsi_gruz_capacity', 'cnsi_volumek', 'tara', 'date_build', 'srok_sl', 'zavod_build', 'date_iskl', 'cnsi_p
| pr_rem_train = pr_rem[['wagnum', 'road_id_send', 'road_id_rem', 'st_id_send', 'st_id_rem', 'distance']]
| kti_izm_train = kti_izm[['wagnum', 'mileage_all', 'axll_l_w_flange', 'axll_r_w_flange', 'axl2_l_w_flange', 'axl2_r_w_flange', 'axl3_l_w_flange', 'axl3_l_w_rim', 'axl4_l_w_rim', 'axl4_r_w_rim']]
| tr_term_train = tr_rem
| tr_term_train = tr_rem | tr_term_train, how = 'left')\|
| .merge(wag_param_pred, how = 'left')\|
| .merge(wag_param_pred, how = 'left')\|
| .merge(wag_param_pred, how = 'left')\|
| .merge(ti_izm_train, how = 'left')\|
| .merge(ti_izm_train, how = 'left')\|
| .merge(kti_izm_train, how = 'left')
```

• Отсеяны вагоны, которые проехали 0 км в день, т.к. по ним невозможно обучить модель.

```
Y [26] df_train[['cnsi_probeg_dr','cnsi_probeg_kr','mean_run']] = df_train[['cnsi_probeg_dr','cnsi_probeg_kr','mean_run']].fillna(0)#заполняем нулями

df_train['day_run'] = df_train.apply(lambda x : [val for val in [x.cnsi_probeg_kr, x.cnsi_probeg_dr, x.mean_run] if val != 0], axis = 1 )#заменяем нул

df_train['day_run'] = df_train.apply(lambda x : np.mean(x.day_run) if len(x.day_run) > 0 else 0, axis = 1 )#отсеиваем вагоны, которые проехали меньше 0 км
```

• Также рассмотрено количество дней, которое осталось до истечений срока службы вагона, дней до ближайшего ПР и какой остаточный ресурс будет на момент окончания месяца.

• Найдены параметры, которые больше всего коррелируют с target_month (будет ли в этом месяце ремонт).

```
í O
              #смотрим корреляции
                max_corr = df_train.corr()['target_month'].apply(lambda x:abs(x)).sort_values(ascending=False,axis=0)
max_corr.iloc[:50]
      xipython-input-32-08fa8b8a8f51>:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default max_corr = df_train.corr()['target_month'].apply(lambda x:abs(x)).sort_values(ascending=False,axis=0)
              max_corr = df_
target_day
mileage_all
prob_end_month
ost_prob
day_run
mean_run
axli_r_w_flange
axl2_l_w_flange
axl2_r_w_flange
axl2_r_w_flange
axl3_r_w_flange
axl3_l_w_flange
axl1_l_w_flange
axl3_l_w_flange
                                                              1.000000
                                                              0.312712
                                                              0.274162
0.271934
                                                              0.127527
                                                              0.124532
0.099007
0.089090
                                                              0.085205
0.082267
                                                              0.080600
                                                              0.080570
0.070061
0.069558
                                                              0.058047
0.050458
0.047280
0.045641
                wagnum
cnsi_probeg_dr
               telega
cnsi_probeg_kr
kuzov
tara
                                                              0.036369
0.035516
               rod_id
kod_vrab
cnsi_volumek
                                                              0.034066
0.030483
                                                              0.029897
                norma_km
zavod_build
                                                              0.027939
0.018505
                tippogl
isload
                                                              0.013471
               ax14_1_w_rim
ax14_r_w_rim
ax12_1_w_rim
cnsi_gruz_capacity
                                                              0.008525
0.008262
                                                              0.007614
                                                              0.006952
                gruz
                axl2_r_w_rim
axl3_r_w_rim
axl1_l_w_rim
                                                              0.006822
0.004735
0.002883
0.002289
              tormor 0.002289
distance 0.001744
ax13_1_w_rim 0.001706
ax11_rw_rim 0.000353
road_id_send NaN
road_id_rem NaN
st_id_send NaN
Name: target_month, dtype: float64
                tormoz
```

• Удалены нулевые значения по столбцам. Взяты 24 лучших параметров, которые коррелируют с target_month.

```
🥤 [36] #заменям нан на 0 во всей таблицы
                        df_train_ = df_train_.fillna(0)
√ [37] #по столбцу убираем нулевые значения
mileage_mean = df_train_[df_train_['mileage_all'] != 0]['mileage_all'].mean()
                        \label{eq:df_train_(mileage_all')} $$ df_{train_(mileage_all').apply(lambda x: mileage_mean if x == 0 else x) $$
                                                                                                                                                                                                                                                                                                                                                                                                                        Λ Ψ Θ 🗏 Φ 🖟 🖹 :
🧴 🕟 #по столбцу убираем нулевые значения
                         axl1_l_w_flange_mean = df_train_[df_train_['axl1_l_w_flange'] != 0]['axl1_l_w_flange'].mean()
                        \label{eq:df_train_(axl1_l_w_flange') = df_train_(axl1_l_w_flange').apply(lambda x: axl1_l_w_flange_mean if x == 0 else x)} df_train_(axl1_l_w_flange') = df_train_(axl1_l_w_flange').apply(lambda x: axl1_l_w_flange_mean if x == 0 else x)} df_train_(axl1_l_w_flange').apply(lambda x: axl1_l_w_flange).apply(lambda x: axl1_l_w_flan
                        axl1_r_w_flange_mean = df_train_[df_train_['axl1_r_w_flange'] != 0]['axl1_r_w_flange'].mean()
df_train_['axl1_r_w_flange'] = df_train_['axl1_r_w_flange'].apply(lambda x: axl1_r_w_flange_mean if x == 0 else x)
                        ax12 1 w flange mean = df train [df train ['ax12 1 w flange'] != 0]['ax12 1 w flange'].mean()
                        df_train_['ax12_1_w_flange'] = df_train_['ax12_1_w_flange'].apply(lambda x: ax12_1_w_flange_mean if x == 0 else x)
                         ax12\_r\_w\_flange\_mean = df\_train\_[df\_train\_['ax12\_r\_w\_flange'] != 0]['ax12\_r\_w\_flange'].mean() \\ df\_train\_['ax12\_r\_w\_flange'] = df\_train\_['ax12\_r\_w\_flange'].apply(lambda x: ax12\_r\_w\_flange_mean if x == 0 else x) \\ \\ x = 0 else x 
                         axl3_1_w_flange_mean = df_train_[df_train_['axl3_1_w_flange'] != 0]['axl3_1_w_flange'].mean()
                        df_train_['axl3_1_w_flange'] = df_train_['axl3_1_w_flange'].apply(lambda x: axl3_1_w_flange_mean if x == 0 else x)
                        axl3_r_w_flange_mean = df_train_[df_train_['axl3_r_w_flange'] != 0]['axl3_r_w_flange'].mean()
df_train_['axl3_r_w_flange'] = df_train_['axl3_r_w_flange'].apply(lambda x: axl3_r_w_flange_mean if x == 0 else x)
                         axl4\_r\_w\_flange\_mean = df\_train\_[df\_train\_['axl4\_r\_w\_flange'] != 0]['axl4\_r\_w\_flange'].mean() \\ df\_train\_['axl4\_r\_w\_flange'] = df\_train\_['axl4\_r\_w\_flange'].apply(lambda x: axl4\_r\_w\_flange_mean if x == 0 else x) \\ \\ x = 0 else x 
                         axl1\_1\_w\_rim\_mean = df\_train\_[df\_train\_['axl1\_1\_w\_rim'] != 0]['axl1\_1\_w\_rim'].mean() \\ df\_train\_['axl1\_1\_w\_rim'] = df\_train\_['axl1\_1\_w\_rim'].apply(lambda x: axl1_1_w\_rim\_mean if x == 0 else x) 
                        axl1_r_w_rim_mean = df_train_[df_train_['axl1_r_w_rim'] != 0]['axl1_r_w_rim'].mean()
                        \label{eq:df_train_(axl1_rw_rim') = df_train_(axl1_rw_rim').apply(lambda x: axl1_rw_rim_mean if x == 0 else x)} \\
                         axl2\_1\_w\_rim\_mean = df\_train\_[df\_train\_['axl2\_1\_w\_rim'] != 0]['axl2\_1\_w\_rim'].mean() \\ df\_train\_['axl1\_1\_w\_rim'] = df\_train\_['axl2\_1\_w\_rim'].apply(lambda x: axl2\_1\_w\_rim\_mean if x == 0 else x) 
                         axl2\_r\_w\_rim\_mean = df\_train\_[df\_train\_['axl2\_r\_w\_rim'] != 0]['axl2\_r\_w\_rim'].mean() \\ df\_train\_['axl2\_r\_w\_rim'] = df\_train\_['axl2\_r\_w\_rim'].apply(lambda x: axl2\_r\_w\_rim\_mean if x == 0 else x) 
                        axl3_l_w_rim_mean = df_train_[df_train_['axl3_l_w_rim'] != 0]['axl3_l_w_rim'].mean()
                        \label{eq:df_train_(axl3_l_w_rim') = df_train_(axl3_l_w_rim').apply(lambda x: axl3_l_w_rim_mean if x == 0 else x)} \\
                       axl3_r_w_rim_mean = df_train_[df_train_['axl3_r_w_rim'] != 0]['axl3_r_w_rim'].mean()
df_train_['axl3_r_w_rim'] = df_train_['axl3_r_w_rim'].apply(lambda x: axl3_r_w_rim_mean if x == 0 else x)
                         ax14\_1\_w\_rim\_mean = df\_train\_[df\_train\_['ax14\_1\_w\_rim'] != 0]['ax14\_1\_w\_rim'].mean() \\ df\_train\_['ax14\_1\_w\_rim'] = df\_train\_['ax14\_1\_w\_rim'].apply(lambda x: ax14\_1\_w\_rim\_mean if x == 0 else x) 
                        max_corr = df_train_.corr()['target_day'].apply(lambda x:abs(x)).sort_values(ascending=False,axis=0)
                          max_corr.iloc[:26]
                                                                                    1.000000

→ target_day

                          target_month
prob_end_month
                                                                                      0.582258
                          ost prob
                                                                                      0.157307
                                                                                       0.076598
                          day_run
                          mean_run
                                                                                      0.069795
                         tipvozd__3
axl1_r_w_flange
axl4_r_w_flange
axl2_1_w_flange
                                                                                       0.044693
                                                                                      0.036771
                                                                                      0.034141
                                                                                      0.033077
                          isload
                                                                                      0.032756
                           cnsi_probeg_dr
                                                                                        0.032540
                          cnsi probeg kr
                                                                                      0.030562
                          axl3_r_w_flange
axl2_r_w_flange
                                                                                      0.029426
                                                                                      0.028705
                          telega
                                                                                      0.025706
                           ax13_l_w_flange
                                                                                       0.025053
                          cnsi volumek
                                                                                      0.024210
                                                                                        0.022651
                          axl1_l_w_flange
                                                                                      0.022611
                           kuzo
                                                                                       0 021447
                          rod_id
                                                                                      0.021301
                           tormoz
                                                                                      0.019480
                           tipvozd__6
                          norma_km
axl3_l_w_rim
                                                                                      0.016718
                          Name: target_day, dtype: float64
```

Обучение модели:

Для прогнозирования выбытия вагона в плановый ремонт применяется модель К ближайших соседей (KNN). Количество соседей выбралось 3 в связи с лучшей сходимостью.

```
▼ Обучение модели
50 [52] my_random_state = 1234
√ [54] x_train, x_test, y_train, y_test = train_test_split(x, y_day, test_size=0.2, random_state = my_random_state)
KNeighborsClassifier
      KNeighborsClassifier(n_neighbors=3)
5 [56] y_pred = KNN_day.predict(x_test)
     print(accuracy_score(y_test, y_pred))
      0.9891752577319588
√ [57] x_train, x_test, y_train, y_test = train_test_split(x, y_month, test_size=0.2, random_state = my_random_state)
(58] KNN_month = KNeighborsClassifier(n_neighbors=3)
KNN_month.fit(x, y_month)
            KNeighborsClassifier
      KNeighborsClassifier(n_neighbors=3)
                                                                                                             ↑ ↓ ፡○ 📮 💠 🖟 📋 :
y_pred = KNN_month.predict(x_test)
      print(accuracy_score(y_test, y_pred))
      0.9690721649484536
```

Обработка данных для предсказания производилась так же, как и для обучения модели. Перед прогнозированием проверялось условие на наличие ремонта вагона в прошлом месяце. Если в прошлом месяце у вагона производился плановый ремонт, то в текущем месяце его не будет.

Проведение предсказания модели

0.4004992141102256

```
[ ] for i in pred['wagnum'].values:
        #если в прошлом месяце был ремонт, то в текущем его не будет
       if train[train['wagnum']== i]['target_month'].values[0] == 1:
          pred.loc[pred['wagnum']==i, 'target_month'] = 0
pred.loc[pred['wagnum']==i, 'target_day'] = 0
          x = df_pred_[df_pred_['wagnum'] == i].drop('wagnum', axis=1)
          pred_day = KNN_day.predict(x)
          if pred_day == 1:
            pred.loc[pred['wagnum']==i, 'target_month'] = int(pred_day[0])
            pred.loc[pred['wagnum']==i, 'target_day'] = int(pred_day[0])
            pred_month = KNN_month.predict(x)
            pred.loc[pred['wagnum']==i, 'target_month'] = int(pred_month[0])
pred.loc[pred['wagnum']==i, 'target_day'] = int(pred_day[0])
     pred
                           month target_month target_day
              wagnum
         0
               33361 2023-03-01
                                                0
         1
               33364 2023-03-01
                                                0
                                                             0
               33366 2023-03-01
                                                             0
               33358 2023-03-01
                                                             0
         3
         4
               33349 2023-03-01
                                                0
                                                             0
         ...
      33702
               17621 2023-03-01
                                                0
                                                             0
       33703
               25045 2023-03-01
                                                0
                                                             0
               27156 2023-03-01
      33704
                                                0
                                                             0
               21361 2023-03-01
      33705
                                                0
                                                             0
      33706
                8061 2023-03-01
     33707 rows x 4 columns
[ ] pred.to_csv('/content/y_pred_submit.csv')
[ ] print(calc_f1_score('/content/y_pred_submit.csv', '/content/drive/MyDrive/DataWagon/y_predict.csv'))
```

В результате предсказания общий счет модели равен 40%. Счет модели складывался из половины точности предсказания на месяц и половины точности предсказания на день.

Функция для оценки на приватном лидерборде (все данные)

```
🥤 🕩 # Функция для оценки на приватном лидерборде (все данные)
       def calc_f1_score(test_url: str, prediction_url: str) -> float:
           true_labels = pd.read_csv(test_url)
           pred labels = pd.read csv(prediction url)
           # Таргет для месячного прогноза
           true_labels_month = true_labels['target_month'].values
           pred_labels_month = pred_labels['target_month'].values
           # Таргет для 10 дневного прогноза
           true_labels_day = true_labels['target_day'].values
           pred_labels_day = pred_labels['target_day'].values
           # Посчитаем метрику для месяца и 10 дней
           score_month = f1_score(true_labels_month, pred_labels_month)
           score_day = f1_score(true_labels_day, pred_labels_day)
           # Посчитаем метрику с весом для двух таргетов
           score = 0.5 * score_month + 0.5 * score_day
           return score
```

Заключение:

Исходя из работы можно сделать вывод, что модель имеет удовлетворительную прогностическую способность и для дальнейшего ее использования в производстве необходима доработка, например, обучение на большем количестве данных, так как в представленном наборе в большом количестве строк отсутствовали значения.