# Лаговые переменные

Лаговые переменная - это переменная, значение которой мы берем не за текущий, а за отстоящий от него на определенное расстояние предыдущий момент времени.

Т.е. это переменные, взятые с запаздыванием во времени.

Лаг - величина интервала запаздывания

Варианты записи:

- lag 10
- lag = 10
- это эквивалетно **t-10**
- .shift() метод для создания лагов
- periods параметр, с помощью которого указываем порядок лага (количество периодов моментов, на которое значение должно остоять от текущего)

Пример:

```
data['Lag3'] = data['sales'].shift(periods=3)
или data['Lag3'] = data['sales'].shift(3)
```

```
In [1]: import numpy as np import pandas as pd import polars as pl import bottleneck as bn

import warnings warnings.simplefilter(action='ignore', category=FutureWarning)

import matplotlib.pyplot as plt

# настройка визуализации
% config InlineBackend.figure_format = 'retina'

# импорт необходимых классов и функций from catboost import CatBoostRegressor, Pool from sklearn.model_selection import TimeSeriesSplit from sklearn.metrics import mean_squared_error
```

Out[2]:

sales

```
2018-01-10 28002018-01-11 25002018-01-12 28902018-01-13 2610
```

Пусть у нас есть данные продаж за 12 дней. Создадим лаги с запаздыванием на

• 3, 4, 5, 6, 12 и 13 дней

```
In [3]: lags = [3, 4, 5, 6, 12, 13]
    for lag in lags:
        data[f'Lag{lag}'] = data['sales'].shift(lag)

data
```

Out[3]:

	sales	Lag3	Lag4	Lag5	Lag6	Lag12	Lag13
2018-01-09	2400	NaN	NaN	NaN	NaN	NaN	NaN
2018-01-10	2800	NaN	NaN	NaN	NaN	NaN	NaN
2018-01-11	2500	NaN	NaN	NaN	NaN	NaN	NaN
2018-01-12	2890	2400.0	NaN	NaN	NaN	NaN	NaN
2018-01-13	2610	2800.0	2400.0	NaN	NaN	NaN	NaN
2018-01-14	2500	2500.0	2800.0	2400.0	NaN	NaN	NaN
2018-01-15	2750	2890.0	2500.0	2800.0	2400.0	NaN	NaN
2018-01-16	2700	2610.0	2890.0	2500.0	2800.0	NaN	NaN
2018-01-17	2250	2500.0	2610.0	2890.0	2500.0	NaN	NaN
2018-01-18	2350	2750.0	2500.0	2610.0	2890.0	NaN	NaN
2018-01-19	2550	2700.0	2750.0	2500.0	2610.0	NaN	NaN
2018-01-20	3000	2250.0	2700.0	2750.0	2500.0	NaN	NaN

- Если порядок лага равен или превышает длину набора данных тогда получаем столбец из одних пропусков.
- Чем больше порядок лагов, тем меньше наблюдений используется при его вычислении

Теперь создадим лаги после разделения на обучающую и тестовую выборку:

- сначала удалим Lag12 и Lag13
- разобьем наш набор данных так, чтобы в тестовую выборку попали 4 последних наблюдения

Допустим, мы будем прогнозировать продажи на 4 дня вперед:

• горизонт прогнозирования - 4 дня

```
In [4]: data.drop(['Lag12', 'Lag13'], axis=1, inplace=True)
In [5]: data.head()
```

```
Out[5]:
                   sales
                         Lag3
                                Lag4 Lag5 Lag6
        2018-01-09 2400
                          NaN
                                     NaN NaN
                                NaN
        2018-01-10 2800
                          NaN
                                NaN NaN NaN
        2018-01-11 2500
                          NaN
                                NaN
                                     NaN
                                           NaN
        2018-01-12 2890 2400.0
                                NaN NaN
                                           NaN
        2018-01-13 2610 2800.0 2400.0 NaN
                                          NaN
```

```
In [6]: # задаем горизонт
HORIZON = 4

# разбиваем на обучающую и тестовую выборки
train, test = data[0:len(data)-HORIZON], data[len(data)-HORIZON:]
train
```

Out[6]:

	sales	Lag3	Lag4	Lag5	Lag6
2018-01-09	2400	NaN	NaN	NaN	NaN
2018-01-10	2800	NaN	NaN	NaN	NaN
2018-01-11	2500	NaN	NaN	NaN	NaN
2018-01-12	2890	2400.0	NaN	NaN	NaN
2018-01-13	2610	2800.0	2400.0	NaN	NaN
2018-01-14	2500	2500.0	2800.0	2400.0	NaN
2018-01-15	2750	2890.0	2500.0	2800.0	2400.0
2018-01-16	2700	2610.0	2890.0	2500.0	2800.0

In [7]: test

Out[7]:

	sales	Lag3	Lag4	Lag5	Lag6
2018-01-17	2250	2500.0	2610.0	2890.0	2500.0
2018-01-18	2350	2750.0	2500.0	2610.0	2890.0
2018-01-19	2550	2700.0	2750.0	2500.0	2610.0
2018-01-20	3000	2250.0	2700.0	2750.0	2500.0

#### Обрати внимание:

• чем выше порядок лага - тем более ранние наблюдения обучающей выборки попадают в тестовую выборку

Обрати внимание Не все лаги в тесте используют наблюдения обучающего набора:

• лаг 3 "залез" в тест, что является **подсматриванием в будущее**, которая нам не известно (обведено красным овалом)

	34103	Lago	Lag-	Lago	Lago
2018-01-09	2400	NaN	NaN	NaN	NaN
2018-01-10	2800	NaN	NaN	NaN	NaN
2018-01-11	2500	NaN	NaN	NaN	NaN
2018-01-12	2890	2400.0	NaN	NaN	NaN
2018-01-13	2610	2800.0	2400.0	NaN	NaN
2018-01-14	2500	2500.0	2800.0	2400.0	NaN
2018-01-15	2750	2890.0	2500.0	2800.0	2400.0
2018-01-16	2700	2610.0	2890.0	2500.0	2800.0
	sales	Lag3	Lag4	Lag5	Lag6
2018-01-17	2250	2500.0	2610.0	2890.0	2500.0
2018-01-18	2350	2750.0	2500.0	2610.0	2890.0
2018-01-19	2550	2700.0	2750.0	2500.0	2610.0
2018-01-20	3000	2250.0	2700.0	2750.0	2500.0

sales

Рис. 21 При вычислении лага 3 произошла «протечка»

#### важно

• Лаговые переменные необходимо создавать так, чтобы они не проникали в тестовый набор

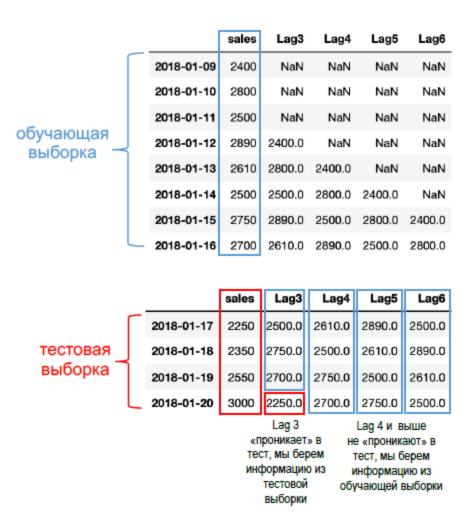


Рис. 22 Лаги, у которых порядок меньше горизонта прогнозирования, «залазят» в тест

Лаги вида  $L_{t-k}$  лучше создавать так, чтобы:

- k был равен или превышал горизонт прогнозирования
- иначе в тестовом наборе будем иметь NaN когда лаг залезет в тест

### Напишем функцию загрузки

- загружаем данные
- читаем столбец с датами как индекс
- выполняем парсинг дат

## Для избежания протечек при вычислении лагов применяют два способа:

- Значение зависимой переменной в наблюдениях исходного набора, которые будут соотстветствовать будущей тестовой выборке (новому набору данных) заменяют значениями NaN
- берем обучающую выборку и удлиняем ее на длину горизонта прогнозирования:
  - зависимая переменная в наблюдениях, соответствующая новым временным меткам (т.е. тестовой выборке/наборе новых данных) получает значения NaN

### Первый способ избежания протечек:

- значения в наблюдениях, которые будут приходится на тестовую выборку (последние 4 наблюдения исходного набора) заменяем на значения NaN
- Затем формируем лаги порядка 1, 2, 3, 4, 5

# заменяем последние 4 наблюдения на NaN data['sales'].iloc[-HORIZON:] = np.NaN

In [12]:

```
data.head()
 In [9]:
Out[9]:
                     sales
                            Lag3
                                   Lag4 Lag5 Lag6
          2018-01-09 2400
                            NaN
                                   NaN
                                        NaN
                                              NaN
          2018-01-10 2800
                            NaN
                                   NaN
                                        NaN
                                              NaN
          2018-01-11
                     2500
                            NaN
                                   NaN
                                        NaN
                                              NaN
          2018-01-12 2890 2400.0
                                   NaN
                                        NaN
                                              NaN
          2018-01-13 2610 2800.0 2400.0
                                        NaN
                                              NaN
          data.tail()
In [10]:
Out[10]:
                     sales
                            Lag3
                                   Lag4
                                         Lag5
                                                Lag6
          2018-01-16 2700 2610.0
                                 2890.0 2500.0
                                               2800.0
          2018-01-17 2250 2500.0 2610.0 2890.0
                                              2500.0
          2018-01-18 2350 2750.0 2500.0 2610.0
                                              2890.0
          2018-01-19 2550 2700.0 2750.0 2500.0
          2018-01-20 3000 2250.0 2700.0 2750.0 2500.0
         HORIZON = 4
In [11]:
          data['sales'].iloc[-HORIZON:]
         2018-01-17
                         2250
Out[11]:
         2018-01-18
                         2350
         2018-01-19
                         2550
         2018-01-20
                         3000
         Name: sales, dtype: int64
```

C:\Users\User\AppData\Local\Temp\ipykernel 19948\1987524369.py:2: SettingWithCopyWarnin

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user

A value is trying to be set on a copy of a slice from a DataFrame

```
guide/indexing.html#returning-a-view-versus-a-copy
data['sales'].iloc[-HORIZON:] = np.NaN
```

#### Создаем лаги:

- лаги меньше горизонта прогнозирования получат пропуски в наблюдениях, приходящихся на тест
- чем меньше порядок лага горизонта тем больше пропусков в тесте

Тестовая выборка

Out[13]:

	sales	Lag3	Lag4	Lag5	Lag6	Lag1	Lag2
2018-01-09	2400.0	NaN	NaN	NaN	NaN	NaN	NaN
2018-01-10	2800.0	NaN	NaN	NaN	NaN	2400.0	NaN
2018-01-11	2500.0	NaN	NaN	NaN	NaN	2800.0	2400.0
2018-01-12	2890.0	2400.0	NaN	NaN	NaN	2500.0	2800.0
2018-01-13	2610.0	2800.0	2400.0	NaN	NaN	2890.0	2500.0
2018-01-14	2500.0	2500.0	2800.0	2400.0	NaN	2610.0	2890.0
2018-01-15	2750.0	2890.0	2500.0	2800.0	2400.0	2500.0	2610.0
2018-01-16	2700.0	2610.0	2890.0	2500.0	2800.0	2750.0	2500.0
2018-01-17	NaN	2500.0	2610.0	2890.0	2500.0	2700.0	2750.0
2018-01-18	NaN	2750.0	2500.0	2610.0	2890.0	NaN	2700.0
2018-01-19	NaN	2700.0	2750.0	2500.0	2610.0	NaN	NaN
2018-01-20	NaN	NaN	2700.0	2750.0	2500.0	NaN	NaN

#### Видим что лаги:

- у который порядок меньше горизонта прогнозирования "залезают" в тестовую выборку и получают пропуски в наблюдениях, приходящихся на тестовую выборку
- чем больше горизонт прогнозирования превышает порядок лага тем больше пропусков будет в тестовой выборке

#### Так срабатывает защита от протечек

## Второй способ защиты от протечек

- удлиняем обучающую выборку на длину горизонта прогнозирования
- зависимую переменную при новых временных метках = NaN

Обучающая выборка - первые 8 наблюдений

```
In [14]: data = load_data()
    train = data.iloc[:-HORIZON]
    train
```

Out[14]:

sales

```
      2018-01-09
      2400

      2018-01-10
      2800

      2018-01-11
      2500

      2018-01-12
      2890

      2018-01-13
      2610

      2018-01-14
      2500

      2018-01-15
      2750

      2018-01-16
      2700
```

Напишем функцию, которая:

- удлиняет обучающую выборку на длину горизонта прогнозирования
- присваивает зависимой переменной при новых метках NaN

```
def calculate lags(train, target, horizon, lags range,
In [15]:
                           freq='D', aggregate=False):
             Создает лаги в обучающей и тестовой выборках.
             Параметры
             train:
                Обучающий набор
             target:
                Название завимисой переменной
             horizon:
                Горизонт прогнозирования
             lags range:
                Диапазон значений порядка лагов
             freq: str, значение по умолчанию `D`
                 Частота временного ряда
             aggregate: bool, значение по умолчанию False
                 Вычисляет аггрегированный лаг
             if min(lags range) < horizon:</pre>
                 warnings.warn(f"\nКоличество периодов для лагов нужно задавать\n"
                              f"равным или больше горизонта прогнозирования")
             if pd. version >= '1.4':
                 # создаем метки времени для горизонта
                 future dates = pd.date range(start=train.index[-1],
                                              periods=horizon + 1,
                                              freq=freq,
                                              inclusive='right')
             else:
                 # создаем метки времени для горизонта
                 future dates = pd.date range(start=train.index[-1],
                                             periods=horizon + 1,
                                             freq=freq,
                                             closed='right')
             # формирукм новый удлиненный индекс
             new index = train.index.append(future dates)
             # выполняем переиндексацию
             new df = train.reindex(new index)
```

```
# создаем лаги
for i in lags range:
    new df[f'Lag {i}'] = new df[target].shift(i)
if aggregate and min(lags range) >= horizon:
    # вычисляем агрегированный лаг
    new df['Agg Lag'] = new df[new df.filter(
        regex='Lag').columns].mean(axis=1)
train = new df.iloc[:-horizon]
test = new df.iloc[-horizon:]
return train, test
```

#### Разбор функции calculate\_lags

**2018-01-10** 2800.0

**2018-01-11** 2500.0

```
train
In [16]:
Out[16]:
                   sales
         2018-01-09 2400
         2018-01-10 2800
         2018-01-11 2500
         2018-01-12 2890
         2018-01-13 2610
         2018-01-14 2500
         2018-01-15 2750
         2018-01-16 2700
In [17]: horizon = 4
         freq = 'D'
         future dates = pd.date range(start=train.index[-1],
                                                periods=horizon + 1,
                                                freq=freq,
                                                inclusive='right')
         future dates
         DatetimeIndex(['2018-01-17', '2018-01-18', '2018-01-19', '2018-01-20'], dtype='datetime6
Out[17]:
         4[ns]', freq='D')
In [18]: new_index = train.index.append(future dates)
         new index
         DatetimeIndex(['2018-01-09', '2018-01-10', '2018-01-11', '2018-01-12',
Out[18]:
                         '2018-01-13', '2018-01-14', '2018-01-15', '2018-01-16',
                         '2018-01-17', '2018-01-18', '2018-01-19', '2018-01-20'],
                        dtype='datetime64[ns]', freq=None)
         new df = train.reindex(new index)
In [19]:
         new df
Out[19]:
                    sales
         2018-01-09 2400.0
```

```
2018-01-15 2750.0
          2018-01-16 2700.0
          2018-01-17
                       NaN
          2018-01-18
                       NaN
          2018-01-19
                       NaN
          2018-01-20
                       NaN
In [20]:
          for i in [1,2,3,4,5]:
              new df[f'Lag {i}'] = new df['sales'].shift(i)
          new df
Out[20]:
                      sales
                             Lag_1
                                    Lag_2
                                           Lag_3
                                                  Lag_4
                                                         Lag_5
          2018-01-09 2400.0
                              NaN
                                            NaN
                                                   NaN
                                                          NaN
                                     NaN
          2018-01-10 2800.0 2400.0
                                     NaN
                                            NaN
                                                   NaN
                                                          NaN
          2018-01-11 2500.0 2800.0
                                   2400.0
                                            NaN
                                                   NaN
                                                          NaN
          2018-01-12 2890.0 2500.0
                                   2800.0
                                          2400.0
                                                   NaN
                                                          NaN
          2018-01-13 2610.0
                            2890.0
                                   2500.0
                                          2800.0
                                                  2400.0
                                                          NaN
          2018-01-14 2500.0
                            2610.0
                                   2890.0
                                          2500.0
                                                  2800.0 2400.0
          2018-01-15 2750.0
                            2500.0
                                   2610.0
                                           2890.0
                                                  2500.0
                                                         2800.0
          2018-01-16 2700.0
                            2750.0 2500.0
                                          2610.0
                                                  2890.0 2500.0
          2018-01-17
                       NaN
                            2700.0
                                   2750.0
                                          2500.0
                                                 2610.0 2890.0
          2018-01-18
                       NaN
                              NaN
                                   2700.0
                                          2750.0 2500.0 2610.0
          2018-01-19
                       NaN
                              NaN
                                     NaN
                                          2700.0 2750.0 2500.0
          2018-01-20
                                     NaN
                                            NaN 2700.0 2750.0
                       NaN
                              NaN
          new df[new df.filter(regex='Lag').columns]
In [21]:
Out[21]:
                      Lag_1
                             Lag_2
                                    Lag_3
                                           Lag_4 Lag_5
          2018-01-09
                       NaN
                              NaN
                                     NaN
                                            NaN
                                                   NaN
          2018-01-10 2400.0
                              NaN
                                     NaN
                                            NaN
                                                   NaN
          2018-01-11 2800.0 2400.0
                                     NaN
                                            NaN
                                                   NaN
          2018-01-12 2500.0 2800.0 2400.0
                                                   NaN
                                            NaN
          2018-01-13 2890.0 2500.0
                                   2800.0
                                          2400.0
                                                   NaN
          2018-01-14 2610.0 2890.0 2500.0 2800.0 2400.0
          2018-01-15 2500.0 2610.0 2890.0
                                          2500.0 2800.0
          2018-01-16 2750.0 2500.0 2610.0 2890.0 2500.0
          2018-01-17 2700.0 2750.0 2500.0 2610.0 2890.0
```

**2018-01-12** 2890.0

**2018-01-13** 2610.0

2018-01-14 2500.0

```
2018-01-19
                                      2700.0 2750.0 2500.0
                         NaN
                                NaN
           2018-01-20
                         NaN
                                NaN
                                        NaN 2700.0 2750.0
           new df['Agg Lag'] = new df[new df.filter(regex='Lag').columns].mean(axis=1)
In [22]:
           new df
Out[22]:
                        sales
                                      Lag_2
                                              Lag_3
                                                     Lag_4
                                                             Lag_5
                                                                       Agg_Lag
                               Lag_1
                       2400.0
           2018-01-09
                                NaN
                                        NaN
                                               NaN
                                                       NaN
                                                              NaN
                                                                           NaN
           2018-01-10 2800.0
                                                                    2400.000000
                              2400.0
                                                       NaN
                                        NaN
                                               NaN
                                                              NaN
                       2500.0
           2018-01-11
                              2800.0
                                      2400.0
                                               NaN
                                                       NaN
                                                              NaN
                                                                    2600.000000
           2018-01-12 2890.0
                              2500.0
                                      2800.0
                                             2400.0
                                                       NaN
                                                              NaN
                                                                    2566.666667
           2018-01-13 2610.0
                              2890.0
                                      2500.0
                                             2800.0
                                                     2400.0
                                                              NaN
                                                                    2647.500000
           2018-01-14 2500.0
                              2610.0
                                      2890.0
                                             2500.0
                                                     2800.0
                                                            2400.0
                                                                    2640.000000
           2018-01-15 2750.0
                                             2890.0
                              2500.0
                                      2610.0
                                                     2500.0
                                                            2800.0
                                                                    2660.000000
           2018-01-16 2700.0
                                      2500.0
                                             2610.0
                                                     2890.0
                                                            2500.0
                              2750.0
                                                                    2650.000000
           2018-01-17
                              2700.0
                                      2750.0
                                             2500.0
                                                     2610.0
                                                            2890.0
                                                                    2690.000000
                         NaN
           2018-01-18
                                      2700.0
                                             2750.0
                                                    2500.0
                                                            2610.0
                                                                    2640.000000
                         NaN
                                NaN
           2018-01-19
                         NaN
                                             2700.0
                                                    2750.0
                                                            2500.0
                                                                    2650.000000
                                NaN
                                        NaN
           2018-01-20
                         NaN
                                NaN
                                        NaN
                                               NaN 2700.0 2750.0 2725.000000
           train = new df.iloc[:-horizon]
In [23]:
           test = new df.iloc[-horizon:]
           train
In [24]:
Out[24]:
                        sales
                               Lag_1
                                      Lag_2
                                              Lag_3
                                                     Lag_4
                                                             Lag_5
                                                                       Agg_Lag
           2018-01-09
                       2400.0
                                NaN
                                        NaN
                                               NaN
                                                       NaN
                                                              NaN
                                                                           NaN
                      2800.0
           2018-01-10
                              2400.0
                                                                    2400.000000
                                        NaN
                                               NaN
                                                       NaN
                                                              NaN
           2018-01-11
                       2500.0
                              2800.0
                                      2400.0
                                               NaN
                                                       NaN
                                                              NaN
                                                                    2600.000000
           2018-01-12 2890.0
                              2500.0
                                      2800.0 2400.0
                                                       NaN
                                                              NaN
                                                                    2566.666667
           2018-01-13 2610.0
                              2890.0
                                      2500.0
                                             2800.0
                                                     2400.0
                                                                    2647.500000
                                                              NaN
           2018-01-14 2500.0
                                             2500.0
                                                            2400.0
                              2610.0
                                      2890.0
                                                     2800.0
                                                                    2640.000000
           2018-01-15 2750.0
                              2500.0
                                             2890.0
                                                     2500.0
                                                            2800.0
                                      2610.0
                                                                    2660.000000
                                                                    2650.000000
           2018-01-16 2700.0 2750.0 2500.0 2610.0 2890.0 2500.0
In [25]:
           test
Out[25]:
                       sales
                             Lag_1
                                     Lag_2
                                            Lag_3
                                                    Lag_4
                                                           Lag_5 Agg_Lag
           2018-01-17
                                    2750.0
                                            2500.0
                       NaN
                             2700.0
                                                   2610.0
                                                           2890.0
                                                                     2690.0
           2018-01-18
                                    2700.0
                                            2750.0
                                                  2500.0
                                                           2610.0
                                                                     2640.0
                       NaN
                               NaN
```

2700.0 2750.0 2500.0

2650.0

2018-01-18

2018-01-19

NaN

NaN

NaN

NaN

2700.0 2750.0 2500.0 2610.0

**2018-01-20** NaN NaN NaN NaN 2700.0 2750.0 2725.0

#### Применим функцию и создадим лаги порядка 3, 4, 5

```
In [28]: data = load data()
         train = data.iloc[:-HORIZON]
         tr, tst = calculate lags(train, target='sales',
                                   horizon=4, lags range=range(3,6),
                                   freq='D', aggregate=False)
         C:\Users\User\AppData\Local\Temp\ipykernel 19948\1679622827.py:24: UserWarning:
         Количество периодов для лагов нужно задавать
         равным или больше горизонта прогнозирования
           warnings.warn(f"\nКоличество периодов для лагов нужно задавать\n"
         # смотрим лаги в обучающей выборке
In [29]:
         tr
Out[29]:
                     sales
                           Lag_3 Lag_4 Lag_5
         2018-01-09 2400.0
                            NaN
                                  NaN
                                         NaN
         2018-01-10 2800.0
                                         NaN
                            NaN
                                  NaN
         2018-01-11 2500.0
                            NaN
                                  NaN
                                        NaN
         2018-01-12 2890.0 2400.0
                                  NaN
                                         NaN
         2018-01-13 2610.0 2800.0 2400.0
                                         NaN
         2018-01-14 2500.0 2500.0 2800.0 2400.0
         2018-01-15 2750.0 2890.0 2500.0 2800.0
         2018-01-16 2700.0 2610.0 2890.0 2500.0
         # смотрим лаги в тесте выборке
In [30]:
         tst
Out[30]:
                    sales Lag_3 Lag_4 Lag_5
         2018-01-17 NaN 2500.0 2610.0 2890.0
         2018-01-18 NaN 2750.0 2500.0 2610.0
         2018-01-19 NaN 2700.0 2750.0 2500.0
         2018-01-20 NaN
                         NaN 2700.0 2750.0
```

Видим, что Lag\_3 получает пропуск, когда пытается использовать информацию тестовой выборки

#### Про агрегированные лаги

- можно взять простое среднее лагов:  $(L_{t-7} + L_{t-14} + L_{t-21})/3$
- можно усреднять лаги с использованием различных весов
- можно брать медиану, стандартное отклонение значений лагов

```
In [31]: # создаем лаги и агрегированный лаг
# для обучающей и тестовой выборок
data = load_data()
train = data.iloc[:-HORIZON]
```

```
In [32]: tr
```

Out[32]: sales Lag\_4 Lag\_5 Agg\_Lag

2018-01-09 2400.0 NaN NaN NaN

**2018-01-10** 2800.0 NaN NaN NaN **2018-01-11** 2500.0 NaN NaN NaN

**2018-01-12** 2890.0 NaN NaN NaN

**2018-01-13** 2610.0 2400.0 NaN 2400.0

**2018-01-14** 2500.0 2800.0 2400.0 2600.0

**2018-01-15** 2750.0 2500.0 2800.0 2650.0

**2018-01-16** 2700.0 2890.0 2500.0 2695.0

```
In [33]: tst
```

### Out[33]: sales Lag\_4 Lag\_5 Agg\_Lag

2018-01-17	NaN	2610.0	2890.0	2750.0
2018-01-18	NaN	2500.0	2610.0	2555.0
2018-01-19	NaN	2750.0	2500.0	2625.0
2018-01-20	NaN	2700.0	2750.0	2725.0

#### Функция weighted\_average\_lag:

• вычисляет агрегированный лаг на основе взвешенного среднего лага

```
In [35]: for cnt, i in enumerate(tr.columns):
    print(cnt, i)
```

0 sales

1 Lag\_4

2 Lag\_5

3 Agg\_Lag

sales Lag 4 Lag 5 Agg Lag

	2018-01-0	9 240	0.0	NaN	NaN	NaN	
	2018-01-1	0 280	0.0	NaN	NaN	NaN	
	2018-01-1	1 250	0.0	NaN	NaN	NaN	
	2018-01-1	2 289	0.0	NaN	NaN	NaN	
	2018-01-1	3 261	0.0 2	400.0	NaN	2400.0	
	2018-01-1	4 250	0.0 2	800.0	4800.0	2600.0	
	2018-01-1	5 275	0.0 2	500.0	5600.0	2650.0	
	2018-01-1	6 270	0.0 2	890.0	5000.0	2695.0	
[36]:		sales	Lag_4	Lag_5	Agg_Lag	Weighted_	Average_Lag
	2018-01-09	2400.0	NaN	NaN	NaN		NaN
	2018-01-10	2800.0	NaN	NaN	NaN		NaN
	2018-01-11	2500.0	NaN	NaN	NaN		NaN
	2018-01-12	2890.0	NaN	NaN	NaN		NaN
	2018-01-13	2610.0	2400.0	NaN	2400.0		2400.0
	2018-01-14	2500.0	2800.0	2400.0	2600.0		3800.0
	2018-01-15	2750.0	2500.0	2800.0	2650.0		4050.0
	2018-01-16	2700.0	2890.0	2500.0	2695.0		3945.0

## Вычисление лагов в Polars

# Реализация первого способа защиты от протечек

```
In [37]: # преобразуем датафрейм pandas в датафрейм Polars
data = pd.DataFrame(new_df['sales'].copy())
display(data)
polars_data = pl.DataFrame(data)
polars_data
```

```
sales
2018-01-09 2400.0
2018-01-10 2800.0
2018-01-11 2500.0
2018-01-12 2890.0
2018-01-13 2610.0
2018-01-14 2500.0
2018-01-15 2750.0
2018-01-16 2700.0
2018-01-17
             NaN
2018-01-18
             NaN
2018-01-19
             NaN
2018-01-20
             NaN
```

Out[37]: shape: (12, 1)

Out

sales

```
164
2400.0
2800.0
2500.0
2890.0
2610.0
2750.0
2700.0
null
null
null
```

#### Особенности Polars:

- нет индекса
- пропускам соответствует значения **null**

### Out[38]: shape: (12, 6)

```
sales
        Lag_1
                Lag_2 Lag_3
                                Lag_4
                                        Lag_5
   f64
           f64
                   f64
                           f64
                                   f64
                                           f64
                          null
2400.0
          null
                  null
                                   null
                                           null
2800.0
       2400.0
                  null
                          null
                                           null
                                   null
2500.0
       2800.0 2400.0
                          null
                                   null
                                           null
2890.0
       2500.0 2800.0 2400.0
                                   null
                                           null
2610.0
       2890.0
               2500.0
                        2800.0
                                2400.0
                                           null
2500.0
       2610.0 2890.0
                        2500.0
                                2800.0 2400.0
2750.0
       2500.0
                2610.0
                        2890.0
                                2500.0
                                        2800.0
                        2610.0
                                2890.0 2500.0
2700.0
       2750.0 2500.0
       2700.0
               2750.0
                        2500.0
                                2610.0 2890.0
  null
               2700.0
                        2750.0
                                2500.0 2610.0
  null
          null
                        2700.0
                                2750.0 2500.0
  null
          null
                  null
  null
          null
                  null
                          null 2700.0 2750.0
```

Видим, что:

лаги, у которых порядок меньше горизонта прогнозирования залезают в тестовую выборку и

sales	Lag_1	Lag_2	Lag_3	Lag_4	Lag_5
f64	f64	f64	f64	f64	f64
2400.0	null	null	null	null	null
2800.0	2400.0	null	null	null	null
2500.0	2800.0	2400.0	null	null	null
2890.0	2500.0	2800.0	2400.0	null	null
2610.0	2890.0	2500.0	2800.0	2400.0	null
2500.0	2610.0	2890.0	2500.0	2800.0	2400.0
2750.0	2500.0	2610.0	2890.0	2500.0	2800.0
2700.0	2750.0	2500.0	2610.0	2890.0	2500.0
null	2700.0	2750.0	2500.0	2610.0	2890.0
null	null	2700.0	2750.0	2500.0	2610.0
null	null	null	2700.0	2750.0	2500.0
null	null	null	null	2700.0	2750.0

получают пропуски в наблюдениях

Для удобства можем добавить в Polars даты.

Возьмем даты из исходного датафрейма pandas

```
# создаем копию
In [39]:
         data2 = data.copy()
        display('Создаем копию', data2)
         # на основе индекса создаем переменную с датами
         data2['date'] = data2.index
         display('на основе индекса создаем переменную с датами', data2)
         # столбец с датами ставим первым
         first column = data2.pop('date')
         data2.insert(0, 'date', first column)
         display('столбец с датами ставим первым', data2)
         # преобразовываем датафрейм pandas в датафрейм Polars
         polars data = pl.DataFrame(data2)
         display("преобразовываем датафрейм pandas в датафрейм Polars", polars data)
         # присваиваем столбцу с датами тип Date
         polars data = polars data.with columns(pl.col('date').cast(pl.Date))
         display("Присваиваем столбцу с датами тип Date", polars data)
```

'Создаем копию'

```
2018-01-09 2400.0
2018-01-10 2800.0
2018-01-11 2500.0
2018-01-12 2890.0
2018-01-13 2610.0
2018-01-14 2500.0
```

2018-01-15	2750.0
2018-01-16	2700.0
2018-01-17	NaN
2018-01-18	NaN
2018-01-19	NaN
2018-01-20	NaN

'на основе индекса создаем переменную с датами'

	sales	date
2018-01-09	2400.0	2018-01-09
2018-01-10	2800.0	2018-01-10
2018-01-11	2500.0	2018-01-11
2018-01-12	2890.0	2018-01-12
2018-01-13	2610.0	2018-01-13
2018-01-14	2500.0	2018-01-14
2018-01-15	2750.0	2018-01-15
2018-01-16	2700.0	2018-01-16
2018-01-17	NaN	2018-01-17
2018-01-18	NaN	2018-01-18
2018-01-19	NaN	2018-01-19
2018-01-20	NaN	2018-01-20

'столбец с датами ставим первым'

	date	sales
2018-01-09	2018-01-09	2400.0
2018-01-10	2018-01-10	2800.0
2018-01-11	2018-01-11	2500.0
2018-01-12	2018-01-12	2890.0
2018-01-13	2018-01-13	2610.0
2018-01-14	2018-01-14	2500.0
2018-01-15	2018-01-15	2750.0
2018-01-16	2018-01-16	2700.0
2018-01-17	2018-01-17	NaN
2018-01-18	2018-01-18	NaN
2018-01-19	2018-01-19	NaN
2018-01-20	2018-01-20	NaN

<sup>&#</sup>x27;преобразовываем датафрейм pandas в датафрейм Polars' shape: (12, 2)

date sales

```
2018-01-09 00:00:00 2400.0
           2018-01-10 00:00:00 2800.0
           2018-01-11 00:00:00 2500.0
           2018-01-12 00:00:00 2890.0
           2018-01-13 00:00:00 2610.0
           2018-01-14 00:00:00 2500.0
           2018-01-15 00:00:00 2750.0
           2018-01-16 00:00:00 2700.0
           2018-01-17 00:00:00
                                null
           2018-01-18 00:00:00
                                null
           2018-01-19 00:00:00
                                null
           2018-01-20 00:00:00
                                null
           'Присваиваем столбцу с датами тип Date'
          shape: (12, 2)
                date sales
                         f64
                date
           2018-01-09 2400.0
           2018-01-10 2800.0
           2018-01-11 2500.0
           2018-01-12 2890.0
           2018-01-13 2610.0
           2018-01-14 2500.0
           2018-01-15 2750.0
           2018-01-16 2700.0
           2018-01-17
                        null
           2018-01-18
                        null
           2018-01-19
                        null
           2018-01-20
                        null
In [40]:
           # создаем список лагов
           lags lst = list(range(1,6))
           # создаем лаги в Polars
           for i in lags lst:
               polars data = polars data.with columns([pl.col('sales').\
                                                                shift(i).alias(f'Lag {i}')])
          polars_data
Out[40]: shape: (12, 7)
                date
                       sales Lag_1 Lag_2 Lag_3 Lag_4 Lag_5
                date
                         f64
                                f64
                                        f64
                                               f64
                                                             f64
```

datetime[ns]

f64

```
2018-01-09 2400.0
                                                    null
2018-01-10 2800.0 2400.0
                             null
                                    null
                                            null
                                                    null
2018-01-11 2500.0 2800.0 2400.0
                                    null
                                            null
                                                    null
2018-01-12 2890.0 2500.0 2800.0 2400.0
                                            null
                                                    null
2018-01-13 2610.0 2890.0 2500.0 2800.0 2400.0
                                                   null
2018-01-14 2500.0 2610.0 2890.0 2500.0 2800.0 2400.0
2018-01-15 2750.0 2500.0 2610.0 2890.0 2500.0 2800.0
2018-01-16 2700.0 2750.0 2500.0 2610.0 2890.0 2500.0
2018-01-17
              null 2700.0 2750.0 2500.0 2610.0 2890.0
2018-01-18
                     null 2700.0 2750.0 2500.0 2610.0
              null
2018-01-19
              null
                     null
                             null 2700.0 2750.0 2500.0
2018-01-20
              null
                     null
                                    null 2700.0 2750.0
```

## Второй способ защиты от протечек в Polars

```
In [42]: # создаем копию
    train2 = train.copy()
# на основе индекса создаем переменную с датами
    train2['date'] = train2.index
    train2
```

```
        Out[42]:
        sales
        date

        2018-01-09
        2400
        2018-01-09

        2018-01-10
        2800
        2018-01-10

        2018-01-11
        2500
        2018-01-11

        2018-01-12
        2890
        2018-01-12

        2018-01-13
        2610
        2018-01-13

        2018-01-14
        2500
        2018-01-14

        2018-01-15
        2750
        2018-01-15

        2018-01-16
        2700
        2018-01-16
```

```
In [46]: # столбец с датами ставим первым first_column = train2.pop('date') train2.insert(0, 'date', first_column) train2
```

```
2018-01-092018-01-0924002018-01-102018-01-1028002018-01-112018-01-1125002018-01-122018-01-1228902018-01-132018-01-1326102018-01-142018-01-142500
```

date sales

Out[46]:

```
2018-01-15 2018-01-15 2750
```

```
2018-01-16 2018-01-16 2700
          # преобразовываем датафрейм пандас в датафрейм Polars
In [48]:
          polars_train_data = pd.DataFrame(train2)
          polars train data = pl.DataFrame(train2)
          polars train data
Out[48]: shape: (8, 2)
                      date sales
                             i64
                datetime[ns]
          2018-01-09 00:00:00
                            2400
          2018-01-10 00:00:00
                            2800
          2018-01-11 00:00:00
                            2500
          2018-01-12 00:00:00
                            2890
          2018-01-13 00:00:00
                            2610
          2018-01-14 00:00:00 2500
          2018-01-15 00:00:00 2750
          2018-01-16 00:00:00 2700
          # присваиваем столбцу с датами типа Date
In [50]:
          polars train data = polars train data.with columns(pl.col('date').cast(pl.Date))
          polars train data
Out[50]: shape: (8, 2)
               date sales
                      i64
                date
          2018-01-09 2400
          2018-01-10 2800
          2018-01-11 2500
          2018-01-12 2890
          2018-01-13 2610
          2018-01-14 2500
          2018-01-15 2750
          2018-01-16 2700
```

```
train:
                Обучающий набор
             target column:
                название переменной с датами
             date column:
                Название переменнйо с датами
             horizon:
                Горизонт прогнозирования
             lags range:
                 Диапазон значений порядка лагов
             aggregate: bool, значение по умолчанию False
                 Вычисляем агрегированный лаг
             .....
             # вычисляем длину горизонта
             h = len(horizon)
             if min(lags range) < h:</pre>
                 warnings.warn(f'\nКоличество периодов для лагов нужно задвать\n'
                               f'равным или больше горизонта прогнозирования')
             # удлиняем набор на длину горизонта
             dates = polars train.select(date column).to series()
             steps = pl.Series(date column, horizon).str.strptime(pl.Date)
             final dates = dates.append(steps).to frame()
             polars df = polars train.join(final dates, how='outer',
                                          on=date column)
             # создаем лаги в Polars
             for i in lags range:
                 polars df = polars df.with columns([
                     pl.col(target column).shift(i).alias(f'Lag {i}')])
             if aggregate and min(lags range) >=h:
                 # вычисляем агрегированный лаг
                 polars df = polars df.with columns([polars df.select(
                 pl.col('^Lag .*$')).mean(axis=1).alias('Agg Lag')])
             train = polars df[0:len(final dates) - h]
             test = polars df[len(final dates) - h:]
             return train, test
In [76]: # создаем лаги для обучающей и тестовой выборок
         polars train, polars test = polars calculate lags (
                                     polars train data,
                                     'sales',
                                     'date',
                                     horizon=['2018-01-17', '2018-01-18',
                                             '2018-01-19', '2018-01-20'],
                                     lags range=range(3,6))
        C:\Users\User\AppData\Local\Temp\ipykernel 19948\1248865699.py:34: UserWarning:
        Количество периодов для лагов нужно задвать
        равным или больше горизонта прогнозирования
          warnings.warn(f'\nКоличество периодов для лагов нужно задвать\n'
In [77]: # смотрим лаги в обучающей выборке
         polars train
```

Out[77]: shape: (8, 6)

Параметры

date	sales	date_right	Lag_3	Lag_4	Lag_5
date	i64	date	i64	i64	i64
2018-01-09	2400	2018-01-09	null	null	null
2018-01-10	2800	2018-01-10	null	null	null
2018-01-11	2500	2018-01-11	null	null	null
2018-01-12	2890	2018-01-12	2400	null	null
2018-01-13	2610	2018-01-13	2800	2400	null
2018-01-14	2500	2018-01-14	2500	2800	2400
2018-01-15	2750	2018-01-15	2890	2500	2800
2018-01-16	2700	2018-01-16	2610	2890	2500

```
In [78]: # смотрим лаги в тестовой выборке polars_test
```

Out[78]: shape: (4, 6)

date	sales	date_right	Lag_3	Lag_4	Lag_5
date	i64	date	i64	i64	i64
null	null	2018-01-17	2500	2610	2890
null	null	2018-01-18	2750	2500	2610
null	null	2018-01-19	2700	2750	2500
null	null	2018-01-20	null	2700	2750

C:\Users\User\AppData\Local\Temp\ipykernel\_19948\1248865699.py:52: DeprecationWarning: T
he `axis` parameter for `DataFrame.mean` is deprecated. Use `DataFrame.mean\_horizontal()
` to perform horizontal aggregation.
 pl.col('^Lag .\*\$')).mean(axis=1).alias('Agg Lag')])

```
In [80]: polars_train
```

Out[80]: shape: (8, 6)

date	sales	date_right	Lag_4	Lag_5	Agg_Lag
date	i64	date	i64	i64	f64
2018-01-09	2400	2018-01-09	null	null	null
2018-01-10	2800	2018-01-10	null	null	null
2018-01-11	2500	2018-01-11	null	null	null
2018-01-12	2890	2018-01-12	null	null	null
2018-01-13	2610	2018-01-13	2400	null	2400.0

2018-01-14	2500	2018-01-14	2800	2400	2600.0
2018-01-15	2750	2018-01-15	2500	2800	2650.0
2018-01-16	2700	2018-01-16	2890	2500	2695.0

In [81]: polars\_test

Out[81]: shape: (4, 6)

date	sales	date_right	Lag_4	Lag_5	Agg_Lag
date	i64	date	i64	i64	f64
null	null	2018-01-17	2610	2890	2750.0
null	null	2018-01-18	2500	2610	2555.0
null	null	2018-01-19	2750	2500	2625.0
null	null	2018-01-20	2700	2750	2725.0